



Multiple criteria decision analysis for HTA across four EU member states: piloting the Advance Value Framework

LSE Research Online URL for this paper: <http://eprints.lse.ac.uk/102122/>

Version: Accepted Version

Article:

Angelis, Aris, Linch, Mark, Montibeller, Gilberto, Molina Lopez, Maria Teresa, Zawada, Anna, Orzel, Kinga, Arickx, Francis, Espin, Jaime and Kanavos, Panos (2019) Multiple criteria decision analysis for HTA across four EU member states: piloting the Advance Value Framework. *Social Science and Medicine*. ISSN 0037-7856

<https://doi.org/10.1016/j.socscimed.2019.112595>

Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Multiple Criteria Decision Analysis for HTA across four EU Member States: piloting the Advance Value Framework

Angelis A ¹, Linch M ², Montibeller G ³, Molina Lopez T ⁴, Zawada A ⁵, Orzel K ⁵, Arickx F ⁶, Espin J ^{7,8}, Kanavos P ¹

1. LSE Health and Department of Health Policy, London School of Economics, London, UK
2. University College London Cancer Institute and University College London Hospital, London, UK
3. School of Business and Economics, Loughborough University, Loughborough, UK
4. Andalusian Agency for Health Technology Assessment, Sevilla, Spain
5. Agency for Health Technology Assessment and Tariff System, Warsaw, Poland
6. National Institute for Health and Disability Insurance, Brussels, Belgium
7. Andalusian School of Public Health, Instituto de Investigación Biosanitaria (ibs.GRANADA), Granada, Spain
8. CIBER de Epidemiología y Salud Pública (CIBERESP), Madrid, Spain

Corresponding authors:

Aris Angelis, LSE Health and Department of Health Policy, London School of Economics, Houghton Street, London, WC2A 2AE, UK; Email: a.n.angelis@lse.ac.uk; Tel: +44 207955 6842

Panos Kanavos, LSE Health and Department of Health Policy, London School of Economics, Houghton Street, London, WC2A 2AE, UK; Email: p.g.kanavos@lse.ac.uk; Tel: +44 207955 6802

Abstract

Multiple Criteria Decision Analysis (MCDA) has emerged as a likely methodology for Health Technology Assessment (HTA). However limited empirical evidence is available on its use by decision-makers and only as part of single-setting exercises, without cross-country studies available. This pilot study applies the Advance Value Framework (AVF), an MCDA methodology for HTA based on multi-attribute value theory, through a series of case studies with decision-makers in four countries, to explore its feasibility and compare their value preferences and results.

The AVF was applied in the evaluation of three drugs for metastatic, castrate resistant, prostate cancer (abiraterone, cabazitaxel and enzalutamide in the post-chemotherapy indication). Decision conferences were organised in four European countries in collaboration with their HTA or health insurance organisations by engaging relevant assessors and experts: Sweden (TLV), Andalusia/Spain (AETSA), Poland (AOTMiT) and Belgium (INAMI-RIZIV). Participant value preferences, including performance scoring and criteria weighting, were elicited through a facilitated decision-analysis modelling approach using the MACBETH technique.

Between 6 and 11 criteria were included in the value model of each country, allocated across four criteria domains; Therapeutic Benefit criteria consistently ranked first across countries in their relative importance. Consistent drug rankings were observed in all settings, with enzalutamide generating the highest overall weighted preference value (WPV) score, followed by abiraterone and cabazitaxel; dividing drugs' overall WPV scores by their costs produced the lowest "cost-per-unit of value" for enzalutamide, followed for abiraterone and cabazitaxel. These results contrast the HTA recommendations and pricing decisions in real life.

Overall, although some differences in value preferences were observed between countries, drug rankings remained the same. The MCDA methodology employed could act as a decision support tool in HTA, due to the transparency in the construction of value preferences in a collaborative manner.

Keywords

Multiple Criteria Decision Analysis (MCDA); Health Technology Assessment (HTA);
Advance Value Framework (AVF); decision conference; value assessment; decision making;
pharmaceuticals; oncology;

Acknowledgements

This paper has been conducted in the context of the Advance-HTA project, which has received funding from the European Commission – DG Research (grant agreement number: 305983). The views reflected in the paper do not represent the views of the European Commission.

We are grateful to Carlos Bana e Costa, Larry Phillips, and Douglas Lundin for constructive feedback and suggestions on earlier versions of the paper. This research was made possible through the participation of a number of experts and stakeholders in four decisions conferences that took place in collaboration with the Swedish Dental and Pharmaceutical Benefits Agency (TLV) in Stockholm, the Andalusian Health Technology Assessment Agency (AETSA) in Seville, the Polish Agency (AOTMiT) in Warsaw, and the Belgian Health Insurance Fund (INAMI-RIZIV) in Brussels. We are indebted to all participants for their time and involvement. Thanks are also due to Olina Efthymiadou and Erica Visintin for valuable research assistance with drug costs and innovation spill-over effect calculations. Finally, we would like to thank all partners of the Advance-HTA consortium for their feedback throughout the project. All outstanding errors are our own.

Highlights

- An MCDA value framework was piloted with HTA decision-makers in four EU countries
- The value drivers of three prostate cancer drugs and their importance were analysed
- Decision-maker value preferences were elicited during four decision conferences
- Value rankings of treatment options were similar and consistent across countries
- The proposed MCDA methodology has prospects to act as a decision support tool

Introduction

In recent years, the introduction of new and costly health technologies, particularly in oncology, combined with moderate health gains, has sparked extensive debate on their value for patients and health care systems, how this value should be assessed and what should be the evaluation criteria informing coverage decisions (Cohen, 2017; Linley & Hughes, 2013). The debate has been fuelled by diverging coverage recommendations across settings for several medicines, often related to diseases associated with high morbidity and mortality (Clement et al., 2009; Faden et al., 2009; Nicod & Kanavos, 2012). Difference in opinion often arises in resource allocation decisions amongst different stakeholders, attributable, at least in part, to current evaluation methodologies not adequately capturing different notions of value (Drummond et al., 2013); this includes, for example, the Quality Adjusted Life Year (QALY), whose use in economic evaluations can at times be regarded as blunt and insufficient, among others, because it may not adequately reflect important value aspects in a variety of disease areas (Nancy Devlin & Lorgelly, 2017; Efthymiadou et al., 2019; Wouters et al., 2015). Given the limited consideration of value in traditional economic evaluations, additional parameters have been included in value assessments; however, this is often done in a non-systematic or ad-hoc manner, which may impact the transparency of decision-making processes (Angelis et al., 2018) and lead to inconsistencies in drug coverage decisions.

A growing body of literature is increasingly debating the use of highly expensive new drugs, which are perceived to bring marginal added clinical benefit on the grounds of poor value-for-money and high budget impact (Nadler et al., 2006; Shih et al., 2013; Sulmasy & Moy, 2014). Rising drug prices and the need to understand the importance of different evaluation criteria have catalysed the generation of numerous “value frameworks” aiming to inform payers, clinicians and patients on the assessment of new medicines, required for making coverage and treatment selection decisions (Anderson et al., 2014; Bach, 2015; Cherny et al., 2015; Schnipper et al., 2015). Although this is an important step towards a more inclusive value-based assessment approach (Malone et al., 2016), aspects of these frameworks

may be based on weak or ad hoc methodologies, which could potentially result in misleading recommendations or decisions (Angelis & Kanavos, 2016a).

In response to some of the concerns raised above, multiple criteria decision analysis (MCDA) has emerged as an alternative to traditional economic evaluation techniques with the prospects of addressing some of their limitations in Health Technology Assessment (HTA) (Angelis et al., 2016; NJ Devlin & Sussex, 2011; Mireille M. Goetghebeur et al., 2008; Kanavos & Angelis, 2013; Marsh et al., 2014; Radaelli et al., 2014; J Sussex et al., 2013b; Thokala, 2011), but also for eliciting stakeholder preferences and facilitating treatment selection (Danner et al., 2011; Ijzerman et al., 2008; Tervonen et al., 2015). A number of MCDA empirical studies have explored the question of value in a number of therapeutic areas, often simulating hypothetical HTA settings (Angelis et al., 2017; M. M. Goetghebeur et al., 2010; Jon Sussex et al., 2013a; Wagner et al., 2017). However, very few studies have explored the same issue by eliciting the preferences of HTA agencies and sitting decision makers and only in single-case exercises (Angelis, 2018; Jaramillo et al., 2016; Tony et al., 2011). To the best of our knowledge, no study has ever compared the value preferences of decision-makers across multiple settings using a full MCDA methodology.

By engaging HTA agencies and health insurance organisations in four EU Member States, we applied the Advance Value Framework (AVF), a recently developed multi-criteria value framework applicable to HTA (Angelis & Kanavos, 2016b; Angelis & Kanavos, 2017), to assess the value of a number of treatment options indicated for metastatic castrate resistant prostate cancer (mCRPC) following first line chemotherapy. This indication was selected because of its high disease burden and the availability of several new and expensive biologic drugs, making it a highly relevant appraisal topic for several HTA agencies.

This is to our knowledge the first cross-country, complete MCDA [pilot](#) exercise, eliciting value preferences of sitting decision-makers from different HTA agencies for the same drug treatments while considering identical sets of evidence. The two main research questions of the study relate to testing the feasibility of this MCDA methodology for HTA

decision-makers, and to observing any differences in their value perceptions as reflected through the consistency of drugs' value rankings, including value trade-offs.

Methods

Methodological Framework

An MCDA approach based on Multi-Attribute Value Theory (MAVT) was adopted (Keeney & Raiffa, 1993; von Winterfeldt & Edwards, 1986), involving the phases of problem structuring, model building, model assessment, model appraisal, and development of action plans (Angelis & Kanavos, 2016b). A series of facilitated workshops were organised taking the form of decision conferences (Phillips, 2007), adopting a facilitated decision analysis modelling approach (Franco & Montibeller, 2010b; Phillips & Phillips, 1993), in collaboration with decision-makers from four HTA agencies and health insurance bodies: the Dental and Pharmaceutical Benefits Agency (TLV, Sweden), the Andalusian Health Technology Assessment Agency (AETSA, Spain), the Agency for Health Technology Assessment and Tariff System (AOTMiT, Poland), and the National Health Insurance Agency (INAMI-RIZIV, Belgium). The agencies in these countries were selected in order to represent a set of organisations with different governance structure (arms' length HTA agency, e.g. AOTMiT, TLV and AETSA, vs integrated HTA function, e.g. INAMI-RIZIV) and responsibilities (regulatory, e.g. TLV, vs advisory AOTMiT and AETSA). This research was undertaken in the context of Advance-HTA, an EU-funded project focusing on HTA methodological advancements (London School of Economics, 2019), and all four HTA organisations were contacted to participate under the auspices of the project.

The methodological process used in terms of the design, implementation and analysis, is aligned with the ISPOR good practice guidelines on the use of MCDA for health care decisions (Marsh et al., 2016).

Problem structuring: Clinical Practice and Scope of the Exercise

Prostate cancer is the second most commonly diagnosed cancer in men globally and the most frequently diagnosed cancer among men in developed countries; it is the fifth leading cause of cancer death globally (Torre, 2015). Death rates have been decreasing in the majority of developed countries, which has mainly been attributed to improved treatment and/or early detection (Center et al., 2012).

The decision context relates to the assessment of value of second line treatments for mCRPC based on the approved European Medicines Agency (EMA) indication (EMA, 2016a, b, c), the subsequently defined scope of Technology Appraisals (TAs) by a number of HTA agencies and the ESMO guidelines (Horwich et al., 2013; NICE, 2012a, b, 2014; TLV, 2014, 2015a).

The first treatment to demonstrate a survival benefit for mCRPC patients was docetaxel chemotherapy in combination with prednisolone when compared to mitoxantrone in combination with prednisolone (Berthold et al., 2008; Tannock et al., 2004). Subsequently, new therapeutic agents have been tested in the post-chemotherapy setting with considerable success. Abiraterone, a steroid synthesis inhibitor, in combination with prednisolone showed a 3.9-month improvement in survival compared to prednisolone alone in patients pre-treated with docetaxel (14.8 vs 10.9 months, HR 0.65, $p < 0.001$) (de Bono et al., 2011). Similarly, enzalutamide, an androgen receptor antagonist, showed a 4.8-month improvement in survival (18.4 vs 13.6 months, HR 0.63, $p < 0.001$) compared to placebo alone in the same patient group (Scher et al., 2012). Cross-resistance appears to exist between abiraterone and enzalutamide meaning that patients are unlikely to derive clinical benefit by switching from one to the other agent (Bianchini et al., 2014; Lortot et al., 2013). The third agent that is widely used following progression on docetaxel is cabazitaxel, a taxane chemotherapy. Cabazitaxel led to an overall survival (OS) benefit of 2.4 months (15.1 vs 12.7 months, HR 0.70, $p < 0.0001$) compared to mitoxantrone (de Bono et al., 2010). Given this therapeutic landscape for patients with mCRPC who have progressed on first line docetaxel chemotherapy, characterised by an availability of different treatments and the apparent cross-

resistance between some of them, we adopt post-chemotherapy mCRPC as the decision context for the application of the AVF methodology.

Model Building: Advance Value Tree adaptation, treatments compared and reference levels

The model building phase comprised a number of tasks, notably the Advance Value Tree adaptation for mCRPC, the consideration of alternative drug treatments and the respective evidence, and the definition of criteria attributes and the associated ranges, all of which are discussed below. Detailed discussion on the rationale of each criterion and their value scales can be found elsewhere (Angelis & Kanavos, 2017; Angelis et al., 2017).

(a) Adaptation of the Advance Value Tree for Metastatic Prostate Cancer

At the core of AVF lies the Advance Value Tree, a hierarchical structure of evaluation criteria taking the form of a generic value tree reflecting value concerns of HTA experts and decision-makers for new medicines (Angelis & Kanavos, 2017). The Advance Value Tree consists of five criteria domains, aiming to capture the essential value attributes of new medicines in the HTA context under a prescriptive decision-aid approach. These are divided into (a) Burden of Disease (BoD); (b) Therapeutic Benefit (THE); (c) Safety Profile (SAF); (d) Innovation Level (INN); and (e) Socioeconomic Impact (SOC), summarised by the following value function:

$$Value = f(\mathbf{BoD}, \mathbf{THE}, \mathbf{SAF}, \mathbf{INN}, \mathbf{SOC}) \quad (1)$$

The Advance Value Tree was adapted into a disease-specific mCRPC value model using a bottom-up approach by comparing the characteristics of the specific drugs evaluated (Franco & Montibeller, 2010a). In consultation with a specialist medical oncologist (co-author of the paper), the generic evaluation criteria were converted into disease-specific criteria, while adhering to required criteria properties such as non-redundancy and preferential-independence (Keeney, 1992), to ensure methodological robustness and an adequate value model rooted in

decision theory. Based on the above, a preliminary mCRPC-specific value tree was produced with four criteria domains and a total of 18 criteria, each operationalised by an attribute, i.e. performance indicator, as shown in Figure 1. The BoD domain was not considered in the adaptation process on the grounds of conciseness, as all drugs were indicated for the same indication which would have identical BoD.

Criteria definitions (together with their consideration in each jurisdiction and their rankings) are provided in Table 1. The preliminary version of the mCRPC value tree was subsequently validated by decision conference participants, in line with a “socio-technical” approach, a constructive decision-aid process allowing groups of participants to interact with and learn from each other (Bana e Costa & Beinat, 2005).

<Figure 1 about here>

(b) Alternative Treatments Compared and Evidence Considered

The alternative drug options assessed in the exercise were cabazitaxel in combination with prednisolone, abiraterone in combination with prednisolone and enzalutamide monotherapy. The key evidence sources used to assess their performance included (a) the peer review publications concerning the pivotal clinical trials of the alternative treatment options that were considered for their licencing by the EMA (de Bono et al., 2011; de Bono et al., 2010; Fizazi et al., 2012; Scher et al., 2012); (b) the Product Information sections of EMA’s European Public Assessment Reports (EPAR) (Annex I and III) (EMA, 2016a, b, c); (c) the Anatomical Therapeutic Chemical (ATC) classification system indexes available through the portal of the WHO Collaborating Centre for Drug Statistics Methodology (World Health Organisation Collaborating Centre, 2016); and (d) the US National Library of Medicine clinical trials database (NIH, 2016). Additional sources of evidence included national sources (BNF, 2015; Connock et al., 2011; NICE, 2012a, b, 2014; Riensa et al., 2013) and other peer review literature (Burström et al., 2001; Collins et al., 2007; Kearns et al., 2013; Sullivan et al.,

2007), which was relevant to the study indication. Sources of evidence used relating to the performance of drugs across evaluation criteria are shown in Appendix Table A1, alongside additional information on the evidence considered.

(c) Options Performance and References Levels

By considering the performance of the alternative drug options across the value scales, “lower” (x_l) and “higher” (x_h) reference levels were defined to serve as benchmarks for the value scores of 0 and 100 respectively, acting as value anchors for constructing value functions and eliciting their relative weights (Bana e Costa & Vansnick, 1999; Keeney, 1982). The “lower” reference levels denoted a less preferred state reflecting a “satisfactory” performance level, whereas the “higher” reference levels denoted a more preferred state reflecting an “ideal” performance level.

The reference levels for the clinical attributes informing the Therapeutic and Safety criteria domains, were defined in consultation with the clinical oncologist (co-author of the paper). In principle, the rationale involved adopting the Best Supportive Care (BSC) performance as a “satisfactory” reference level, with a hypothetical 20% improvement of the best available performance acting as the “ideal” reference level (e.g. ‘overall survival’), or, alternatively, the best possible limit of the performance scale acting as an “ideal” level in cases where this was naturally restricted (e.g. ‘treatment discontinuation’). The 20% hypothetical performance improvement was selected because it was perceived to be a realistically plausible scenario for future treatment options. By considering the performance of best available option(s) among the treatments evaluated and accounting for plausible performance improvement in the near future, the value scale essentially reflected characteristics of a “global” scale to account for the performance of future options not captured in the exercise, i.e. *what is best plausible* (Belton & Stewart, 2002). Where a BSC performance was not meaningful to act as a “lower” reference level, then the lowest (i.e. worst) possible limit of the performance scale was adopted (e.g. ‘Phase 3’), or, alternatively,

20% lower than the lowest performing option was used (e.g. ‘medical costs impact’). An exception to the above was the ‘health related quality of life’ (HRQoL) attribute for which the stable disease state’s utility score was adopted as the “lower” level and the general population utility score was used as the “higher” level.

The emerging partial value function scores of the drugs for each criterion can take negative values or values higher than 100 where $v(x_{\text{lower}}) = 0$ and $v(x_{\text{higher}}) = 100$, essentially by conducting a positive linear transformation. “Lower” and “higher” reference levels for all attributes at the pre-decision conference stage and the basis of their selection are outlined in Appendix Table A2. A matrix listing the performance of drug options across the final attributes that were considered in the decision conferences, together with their reference levels, is shown in Table 2.

Model Assessment and Appraisal: Decision conferences, MCDA technique and cost calculation

The model assessment and appraisal phases comprised the tasks of conducting the decision conferences, the application of the MCDA technique for the elicitation of value preferences and cost calculation(s). These are discussed below.

(a) Decision conferences

Model assessment and model appraisal took place through a series of decision conferences (Phillips, 2007), taking the form of facilitated workshops with the participation of decision-makers, including assessors and national experts, all of whom were affiliated with the four study HTA organisations, either as members of staff or visiting external experts (their difference being in full-time employment versus part-time or visiting capacity employment). For the purposes of this study, they were both regarded as “decision-makers”, given their influence on methodological development within the agencies and on the decision outcomes of the appraisals. Across the four countries, between four (for the case of TLV) and 13 (for

the case of AOTMiT) participants were involved, typically comprising health care professionals (clinicians, pharmacists), HTA methodology experts (health economists, statisticians, HTA agency directors) and decision-makers (members of HTA appraisal committees, representatives from insurance funds and the national medicines agencies). Background material introducing the scope of the exercise in more detail was sent to the participants one week before each decision conference. Decision conferences were hosted at the head offices of the different HTA organisations between June 2015 and April 2016: Stockholm (TLV), Seville (AETSA), Warsaw (AOTMiT), and Brussels (INAMI-RIZIV).

The lead author acted as an impartial facilitator, assisted the groups' interactions and guided participants through the decision problem using the preliminary version of the mCRPC-specific value tree (Figure 1) and the relevant data. This acted as the model's starting point, based on which value judgements and preferences were elicited at the start of each decision conference while seeking group interaction and agreement (Franco & Montibeller, 2010b; Phillips, 1984; Phillips & Bana e Costa, 2007; Schein, 1999). The *Appendix* provides more information on the decision conferences.

(b) MCDA Technique

AVF adopts a value measurement MCDA methodology making use of a simple additive (i.e. linear, weighted average) value model for the aggregation of scores and weights (Angelis & Kanavos, 2017). This assumes preference independence between the different criteria, with overall value $V(\cdot)$ of an option a defined by the equation below (Keeney, 1992; von Winterfeldt & Edwards, 1986):

$$V(a) = \sum_{i=1}^m w_i v_i(a) \quad (2)$$

Where m is the number of evaluation criteria, $w_i v_i(a)$ is the weighted partial value function of evaluation criterion i for treatment a , and $V(a)$ is the overall value of a treatment a . $V(\cdot)$ is

therefore is an overall value function based on multi-attribute value theory (Keeney & Raiffa, 1993).

A value function associated with each attribute, converting the treatment performance on the attribute range to a value scale, was elicited from the participants during the decision conferences using the Measuring Attractiveness by a Categorical Based Evaluation Technique (MACBETH) questioning protocol and the M-MACBETH software (Bana e Costa & Vansnick, 1999). This protocol requires pairwise comparisons where qualitative judgements about the difference of value between different pairs of attribute levels (i.e. difference in value between x and y units on a criterion) are expressed using seven qualitative categories (i.e. no difference, very weak difference, weak difference, moderate difference, strong difference, very strong difference, or extreme difference) (Bana E Costa et al., 2012; Bana e Costa & Vansnick, 1994). MACBETH provides a constructive and user-friendly approach to generate a cardinal (interval) value scale based on the input of these qualitative pair-wise judgements, which are then converted into value scores via an optimization algorithm (Bana e Costa et al., 2016b); this approach has been widely used as a decision support tool (Bana e Costa et al., 2014; Bana e Costa et al., 2002; Bana e Costa & Oliveira, 2012; Bana e Costa & Vansnick, 1997).

Weights for a multi-attribute value function should be elicited considering the range of each attribute and the value of a “swing” between two reference levels. The weights are scaling constants that convert partial value scores into overall value scores that must reflect value trade-offs and, therefore, should not be interpreted as measurements of ‘direct importance’. An indirect (qualitative) swing weighting technique was applied to elicit relative criteria weights by first ordering the swings of each attribute and then valuing their differences using the MACBETH qualitative categories (Bana E Costa et al., 2012).

The above MACBETH-based scoring and weighting techniques were operationalised using the software M-MACBETH, (Bana e Costa & Vansnick, 1999). The software automates the additive aggregation of preference value scores and weights in order to derive overall weighted preference value (WPV) scores and also allows for sensitivity analysis on the

criteria weights. The software also enables the use of visual graphics to build a model of values, acting as a facilitation tool to inform both the design and the evaluation phases of the methodological framework (Bana e Costa et al., 2016a; Bana e Costa & Vansnick, 1999; Bana e Costa et al., 1999). More information regarding the technical details of MACBETH is available in the *Appendix*.

(c) Cost Calculation

UK list prices at ex-factory level were used as found in BNF (BNF, 2015) as a neutral benchmark in order to allow the measurement of cost(s) in a common unit across all study settings, so that overall WPV scores can then be viewed against the same cost denominator to produce comparable cost-value ratios. Access to confidential prices through risk sharing agreements was not possible. Information on the recommended dosages and treatment durations were sourced from the peer review publications of the pivotal trials and respective EPARs from EMA (de Bono et al., 2011; de Bono et al., 2010; EMA, 2016a, b, c; Scher et al., 2012). Drug administration costs for cabazitaxel were kept consistent with the respective NICE TA (NICE, 2012b), whereas for abiraterone and enzalutamide these costs were not applicable as they are orally administered.

Results

Final Value Trees, Options Performance, Criteria Weights and Value Functions

Across the four countries, decision conferences were characterised by increased interaction and extensive debate between participants, especially in cases where there was disagreement about certain values. Because the majority of participants had a shared understanding of the decision problem but also a sense of common purpose and commitment to way forward, all of which are conditions for good practice in decision conferencing, the deliberative process of each decision conference instigated a fruitful discussion and exchange of views around different criteria values and relative importance.

General consensus was reached among participants in terms of criteria consideration and model validation with no major value aspects deemed to be missing. All attributes included in each country's final mCRPC value tree, as emerged following open interaction with decision conference participants and their rankings, are shown in Table 1 (schematic illustrations of the individual value trees are shown in Appendix Figure A1). The main reason for not including a criterion attribute in the value tree was because participants considered it was non-fundamental for the evaluation, in all cases of which a zero weight was assigned. Most of the criteria attributes that were assigned a zero weight belonged in the Innovation Level domain, which comprised the highest number of criteria.

<Table 1 about here>

The performance of the drug options across the different attributes that were considered to be fundamental in the model (i.e. weight greater than zero) together with the “lower” and “higher” reference levels are shown in Table 2.

<Table 2 about here>

Between 6 (AOTMiT) and 11 (AETSA/INAMI) criteria attributes were included in the final value tree of each country, as shown in Table 3. In terms of the different criteria domains composition, the Therapeutic Benefit contained between two (TLV/AOTMiT/INAMI) and three (AETSA) criteria attributes, the Safety Profile between one (AOTMiT) and two (TLV/AETSA/INAMI), the Innovation Level between two (TLV/AOTMiT) and six (INAMI), and the Socioeconomic Impact always one.

<Table 3 about here>

During the elicitation of the ‘overall survival’ (OS) and/or ‘HRQoL’ criteria value functions, it became evident that these criteria attributes might be preference dependent. When asking participants to judge the difference in value between different increments in attribute performance (either in ‘OS’ or ‘HRQoL’), a request for clarification was raised by some of them relating to what level of performance this change was associated with on the other criterion attribute. In order to address the plausible preference-dependence observed, we combined together the two attributes in an aggregated form. The two criteria attributes were combined by multiplying the number of months in ‘OS’ and their EQ-5D utility scores in ‘HRQoL’ attributes respectively, assuming an equal (i.e. 50%) distribution of stable and progressive disease states, essentially deriving quality adjusted life months (QALMs). An example of a MACBETH value judgements matrix and its conversion into a value function for the case of the ‘OS x HRQoL’ aggregated criterion attribute in QALMs is shown in Appendix Figure A2.

There was a common set of six criteria that were considered as fundamental in all countries: (a) ‘OS x HRQoL’; (b) ‘radiographic tumour progression’ (also known as progression free survival (PFS)); (c) ‘treatment discontinuation’; (d) ‘delivery posology’; (e) ‘special instructions’; and (f) ‘medical costs impact’. This common set of criteria comprised the complete set of TLV’s value tree (n=6), whereas AOTMIT’s value tree considered ‘contraindications’ in addition (n=7). Further to these, AETSA’s value tree also considered ‘PSA response’, ‘ATCL4’, ‘Phase 3’ and ‘marketing authorisation’ (n=11), whereas INAMI’s value tree considered the same additional criteria but with ‘Phase 2’ instead of ‘PSA response’ (n=11).

Overall, the different groups of decision conferences’ participants agreed in the valuation of performance for the six common attributes that were considered across all four countries, as revealed through the elicitation of their value functions. Figure 2 plots the value scores of each drug across the six common attributes showing very similar valuations between countries.

<Figure 2 about here>

The weights of relative importance assigned to the different attributes across the four jurisdictions are shown in Figure 3. By taking into account the relative swings of the criteria attributes, i.e. the gap between the “lower” and “higher” reference levels, quantitative weights were derived for each attribute using M-MACBETH. The ‘OS x HRQoL’ aggregated criterion attribute was always assigned the highest relative weight out of 100 ([31,44] for INAMI and AETSA, respectively), followed either by ‘treatment discontinuation’ ([17,21] for AETSA and TLV, respectively) or ‘medical costs impact’ ([20,30] for INAMI and AOTMiT, respectively). Depending on the country, the third-ranked criterion was then either ‘treatment discontinuation’ (AOTMiT, INAMI), ‘medical costs impact’ (TLV), or ‘contraindications’ (AETSA) and ‘PFS’ was ranked 4th or 5th. ‘Special instructions’, although a fundamental criterion across settings, was ranked in the lowest place in 3 out of 4 settings with the ‘delivery posology’ usually at a higher position, with the exception of TLV where that order was reversed.

<Figure 3 about here>

In terms of the total weights assigned across the different criteria domains, the Therapeutic Benefit weight ranged from 40% to 54% (for AOTMiT/ INAMI and AETSA, respectively), the Safety Profile weight ranged from 20% to 33% (for AOTMiT and TLV, respectively), the Innovation Level weight ranged from 7% to 13% (for TLV and INAMI, respectively) and the Socioeconomic Impact weight ranged from 8% to 30% (for AETSA and AOTMiT, respectively) (Table 3). The above differences in relative weights reflect the different priorities of decision-makers, including the number of fundamental objectives being considered.

Overall Drug Rankings and Value-for-Money Analysis

With regards to the overall WPV scores shown in Table 4, enzalutamide consistently yielded the highest score across all four countries, always followed by abiraterone and cabazitaxel. The overall scores of abiraterone and cabazitaxel were in part influenced by a “negative” performance in the ‘treatment discontinuation’ attribute (19% and 18% respectively) which lay below the lower reference level of the scale (i.e. 10%), affecting negatively their overall value scores.

A stacked bar plot of the drugs’ overall WPV scores across all settings is shown in Figure 4. By using rounded up cost figures for enzalutamide (£24,600), abiraterone (£21,900) and cabazitaxel (£23,900, of which £22,190 related to drug cost and the remainder £1,710 to administration cost) and dividing them with overall WPV scores, their costs per MCDA value unit ranged as follows: (a) enzalutamide: £410 - £501 (for AOTMiT and AETSA, respectively); (b) abiraterone: £1,366 - £9,221 (for INAMI and TLV, respectively); and (c) cabazitaxel: £2,196 - £6,816 (for INAMI and AOTMiT, respectively) (Table 4). The overall value score of each option was driven by the fundamental objectives considered (i.e. criteria influencing the model), the criteria weights which were anchored on reference levels, and the shape of value functions which would influence the value scores.

<Table 4 about here>

<Figure 4 about here>

In terms of value-for-money, cabazitaxel was shown to be dominated by abiraterone, and was very close to being dominated by enzalutamide (i.e. a difference of £500 based on the prices used). Enzalutamide on the other hand was associated with a higher cost (a difference of £2,500 based on the prices used) and a higher overall WPV score compared to abiraterone, with a difference in score ranging between 40.4 to 52.7 value units (for AETSA and TLV, respectively). Cost benefit plots of the different options, using their overall WPV scores

versus their purchasing (plus any administration) costs across the four HTA organisations is shown in Figure 5.

<Figure 5 about here>

Similarities and differences in value perceptions across settings

By looking at Table 3 (and Figure 3) of the results, a number of similarities and differences in value preferences are observed across the four settings. The largest number of evaluation criteria were considered in Andalusia and Belgium (11 each), compared to Sweden and Poland (7 and 6, respectively), partly due to a higher number of Innovation Level criteria (5 and 6, compared to 2 each, respectively). In terms of the relative importance of criteria domains, the Therapeutic Benefit cluster consistently ranked first across all settings. The Safety Profile cluster was ranked second in three settings (except for Poland, where the Socioeconomic Impact cluster ranked higher (30% vs 20%)). The Socioeconomic Impact cluster ranked 3rd in Sweden and Belgium but 4th in Andalusia (8%). Finally, the Innovation Level cluster ranked 4th in three countries with the exception of Andalusia where it ranked 3rd (12%). The low relative importance of the Innovation Level cluster partly justifies why a hypothetical change in the final consideration of Innovation Level criteria across the different countries does not influence the ranking of the treatments, as described in the next section.

Despite the observed differences in evaluation criteria considered, the relative criteria weights assigned and the elicited value functions, the overall ranking of the treatments remained identical across countries (Table 4 and Figure 4) with enzalutamide consistently having the highest score, followed by abiraterone and cabazitaxel in all four settings.

Sensitivity and Robustness Analysis

Following each decision conference, deterministic sensitivity analysis was conducted to address parameter uncertainty on criteria weights. Specifically, changes on baseline weights

were explored to check their possible impact on treatments' overall value rankings. The results of the sensitivity analysis demonstrated that the ranking of the treatments was robust to the relative criteria weights across the different settings.

The most sensitive criterion weight, which could change enzalutamide's ranking order from first to second, was 'PFS' in the cases of INAMI and AETSA where a 10.2 and 11.1 times change (from 8.9% to 90.6% and from 8.0% to 88.5%) respectively, would be required for cabazitaxel to rank first and enzalutamide second. In other words, a higher than 10-times difference on the 'PFS' weight would be required for cabazitaxel to outperform enzalutamide, with changes of higher order required in other criteria weights for either cabazitaxel or abiraterone to rank first, in any of the study settings. Criteria weights were more sensitive with regards to the outperformance of abiraterone by cabazitaxel as the second-best treatment. Again, the most sensitive weight was for 'PFS' in the INAMI and AETSA cases, where a 2-times change (from 8.9% to 17.4% and from 8.0% to 16.7% respectively) would be needed for cabazitaxel to rank second and abiraterone third. This meant that the lowest change across criteria weights needed for an impact on treatment rankings to be observed was for the case of PFS with INAMI, where at least a 2-time difference was required for abiraterone to be outperformed. For the case of TLV and AOTMiT, the most sensitive criterion was treatment discontinuation in which a 2.6 and 3.0 times change would be needed (from 21.2% to 54.6% and from 20% to 60% respectively) for cabazitaxel to rank second-best.

The final consideration of the Innovation Level criteria cluster was explored in greater detail given that their relevance might be disputed. Removing the 'ATCL4' criterion and any spill-over effect criteria (i.e. 'Phase-2', 'Phase-3', 'MA') from the value tree of AETSA and INAMI, and any patient convenience criteria (i.e. 'delivery posology', 'special instructions') from all country value trees would not affect the treatment rankings.

Discussion and policy implications

This study is the first comparative MCDA exercise, utilising the Advance Value Framework and engaging sitting HTA decision-makers across four EU Member States to elicit and compare their preferences in the evaluation of three mCRPC treatments. In doing so, the objective was to test the usefulness-feasibility of MCDA methods for HTA decision-makers and identify differences in value perceptions.

Based on the evidence used, our results showed that the most valuable therapy for second line mCRPC was enzalutamide, followed by abiraterone and cabazitaxel. Each treatment was assessed and ranked based on their overall WPV scores, reflecting the value of their performance against a set of evaluation criteria, weighted against their relative importance. These overall scores were based on the value preferences of decision-makers that were collected via a decision conference in each setting, yielding a comprehensive and transparent, multi-dimensional benefit component. Subsequent consideration of drug costs (purchasing and administration) enabled the estimation of value-for-money in the form of “cost-per-unit of value” ratios which showed the second-ranked treatment (abiraterone) to dominate the third (cabazitaxel).

It should be noted that the constructed benefit metric excludes the cost of the treatments, i.e. the WPV score considers the impact of the technology on medical costs other than the purchasing cost of the technology. Therefore, evaluation of the treatments based solely on their overall WPV scores might not be appropriately designed to inform an HTA decision context that considers the interventions’ incremental cost per incremental benefit, but, rather, a value-based approach to reimbursement or pricing negotiation.

Attempting a comparison of the ranking achieved in this exercise with what has taken place in reality might prove challenging, partly because of how the clinical evidence was treated in the exercise, but also because it is not publicly known whether and how any of the additional value dimensions evaluated in the exercise were considered in the relevant HTA decision-making processes. In Sweden, although abiraterone’s ICER vs BSC (manufacturer

estimate of SEK820,000/QALY)(TLV, 2015a), was lower compared to enzalutamide's ICER vs BSC (TLV best estimate of SEK1,100,000/QALY)(TLV, 2014), or lower vs enzalutamide (SEK800,000/QALY)(TLV, 2015b), TLV assumed that both treatments had the same clinical effect and consequently focused on a cost-minimisation approach rather than cost-utility analysis, leading to the implementation of a confidential risk sharing agreement (RSA) as part of which discounts can be provided based on treatment duration. A similar conclusion was reached in Spain, where the Ministry of Health in its Clinical Assessment Report (Informe de Posicionamiento Terapeutico - IPT) recommended that there is no clinically relevant difference between the benefit-risk balance of enzalutamide and abiraterone, and, therefore, decisions should be guided based on drug costs (AEMPS, 2015). Pricing and reimbursement decisions are then taken by the Interministerial Committee for Pricing and Reimbursement, but the final assessment is not publicly available. At regional/hospital level, a group of hospital pharmacists conducted a full health (clinical and economic) technology assessment, where enzalutamide and abiraterone were considered to be therapeutically equivalent (GHEMA, 2016). In Poland, although AOTMiT accepted that some additional clinical effect existed for enzalutamide compared to abiraterone (mainly in secondary endpoints), it was not found to be cost-effective compared to abiraterone; however, a confidential RSA enabled a final positive recommendation by AOTMiT (AOTMiT, 2017). The final decision implemented by the Ministry of Health was to reimburse enzalutamide, similarly to the case of abiraterone (Obwieszczenie, 2017). In Belgium, following an indirect comparison no clinically relevant differences were found in the treatment outcomes of abiraterone versus enzalutamide (INAMI, 2019); eventually, a managed entry agreement (MEA) enabled reimbursement.

Consequently, and based on the evidence used to populate the MCDA model and which would inform decision-making, the hypothetical coverage decisions emerging from the ranking of the treatments based on their overall WPV scores might have been different. Given the higher overall value of enzalutamide compared to abiraterone, a cost minimisation

approach or price parity attained between the two, as inferred following the risk sharing agreements in place, might not have been justified.

One reason why our value models make slightly different predictions is because it has captured benefits that go beyond the current formal remit of HTA agencies, therefore the results should be viewed as ‘proof-of-concept’, for the purposes of testing the performance of the methodology. Furthermore, the decision context addressed in the exercise was a one-off evaluation problem within the indication of mCRPC which might contradict the operational scope of some HTA agencies and health insurance bodies relating to repeated decisions around the reimbursement of drugs across different disease areas.

The extent to which HTA decision-makers can be relied upon, or not, to reflect societal preferences when constructing their value preferences is a very important topic for discussion but not aimed to be addressed in this study. Here, we simply elicited decision-makers’ own preferences without considering whether these might be representative for society or not. In reality, evidence in Belgium suggests that health care coverage related preferences of decision-makers differ to those of the public (Cleemput et al., 2018), and therefore more research would be needed to reveal such discrepancies.

Overall, the HTA decision-makers that participated in the decision conferences provided positive feedback about the potential usefulness of the value framework and the MCDA approach in general, raising the prospects of the framework acting as a decision support tool in the evaluation of new medicines. According to participants, key advantages of the framework included the feasibility to transparently assess the performance of the options across a number of explicit evaluation criteria, while allowing the elicitation of value trade-offs (i.e. their relative importance), and its overall facilitative nature in the construction and analysis of group value preferences. Our results are in line with past evidence on a different oncology indication (Angelis et al., 2017).

Challenges of MCDA applications in HTA

The assessment across 4 settings has offered a number of important insights relating to the application of MCDA in HTA and the challenges this represents. In order for any MCDA methodology to become a useful tool for HTA decision-makers and serve their needs, certain requirements must be met: first, sound methods should be used to ensure technical requirements are fulfilled (Keeney & Raiffa, 1993); second, social aspects of the process should be treated carefully to ensure various socio-technical requirements are fulfilled (Baltussen et al., 2017); and, third, tools and guidelines should be available and tailored for the appropriate audience ensuring that best practice requirements are fulfilled (Phillips, 2017).

Among the first group of technical requirements, one key challenge of MCDA studies in HTA relates to the theoretical properties that are required for the evaluation criteria. Due to the popularity of using a simple additive (i.e. weighted average) value model, the violation of preference-independence is of particular relevance as it might undermine the validity of such models and the insights offered by the results (Marsh et al., 2018; Morton, 2017). Evidence suggests that preference dependencies might exist between health gain and disease severity (Nord et al., 2009), or between OS and HRQoL (Angelis & Kanavos, 2017). The latter also featured strongly in this study, where such a preference dependence between OS and HRQoL was detected during the decision conferences and, as a result, the two criteria attributes were combined into a common aggregated attribute. Beyond combining the two criteria into a common aggregated attribute, other more technically complex solutions exist for addressing preference dependencies, such as using other functional forms of aggregation for combining scores and weights together, such as multiplicative models (Chongtrakul et al., 2005). Furthermore, tests for identifying preference dependencies have existed for many years (Currim & Sarin, 1984; Keeney, 1992; Rodrigues et al., 2017).

Other technical challenges relate to the need for evaluation criteria to be non-overlapping so that there can be no double counting, and that criteria weights are connected to the attribute ranges. If either one of these conditions is not satisfied, criteria weights could misrepresent decision makers' true value preferences. Furthermore, a number of cognitive

biases may affect value judgments and thus appropriate elicitation protocols and de-biasing tools must be employed (Montibeller & Winterfeldt, 2015).

In order to avoid double-counting, a clear justification of their inclusion is needed, which should be on the grounds of addressing the fundamental objectives of the analysis, rather than be informed based on the existence of available evidence and data (Keeney, 1992; Keeney & Gregory, 2005). This process could be supported by the use of problem structuring tools aiming to distinguish between ‘fundamental objectives’ and ‘means objectives’ (Franco & Montibeller, 2010a), as we adopted in this exercise.

In terms of weighting, asking direct questions for the general importance of criteria are known to be one of the most common mistakes when eliciting value trade-offs (Keeney, 1992; Keeney, 2002). Instead, sound weighting procedures for the assignment of relative weights should take place in accordance with the use of explicit lower and higher reference levels (Belton & Stewart, 2002; Keeney, 2002), ideally through user-friendly indirect technique protocols that can reduce bias, similar to what we aimed for in this exercise through the explicit definition of reference levels and the implementation of the qualitative (MACBETH) swing weighting technique.

A further challenge relates to the linking of MCDA results with coverage and resource allocation decisions, possibly through the use of specific value thresholds, that can reflect the efficiency and opportunity cost of funding decisions (Sculpher et al., 2017). In economic evaluation, incremental cost effectiveness ratio (ICER) thresholds are supposed to reflect the opportunity cost of the benefit foregone elsewhere in the health care system that would have resulted from the coverage of alternative technologies (Claxton et al., 2015). Assuming that a QALY-based ICER threshold is accurate, it could be used as a benchmark to create an MCDA value threshold by extrapolating the ICER threshold in proportion to how much of the MCDA model’s weight is accounted for by non-QALY value components (Phelps & Madhavan, 2018). Alternatively, following the generation of a multi-dimensional benefit component, purchasing costs could be used to derive treatments’ cost-value ratios to

inform the resource allocation decisions assuming a fixed budget (Peacock et al., 2007), similar to our approach in this exercise with the calculation of the “cost per unit of value”.

Study limitations

The study has a number of limitations, both related to the clinical evidence used and the MCDA process followed, so results should be interpreted with caution. First, in terms of the clinical data used, there was a lack of relative treatment effects; in order to counteract that, absolute treatment effects from different clinical trials were used based on the assumption that they are directly comparable which might not be accurate even for similar patient populations in the studies. As a result, differences in the performance of the options that have been valued might in reality not be statistically significant, e.g. in OS. Ideally, one would need indirect comparisons or a network meta-analysis (NMA) through a mixed treatment comparison (Jansen et al., 2011), therefore, an evidence synthesis step would be required as part of the model-building phase; as, for example, in the case of assessing the comparative benefit-risk of statins in primary prevention (Tervonen et al., 2015) or second-generation antidepressants (van Valkenhoef et al., 2012).

Second, another clinical evidence related limitation could be that only the treatments’ impact on HRQoL of the stable disease state was assessed, because no treatment was assumed to have any effect during progression (NICE, 2014). This might not be true for other disease indications in which case the relevant HRQoL attribute would have to capture both the stable and progressive disease states.

Third, there are also a number of limitations in terms of the MCDA process adopted: one of them relates to the relatively small number of participants in some decision conferences, which could reflect a limited representation of perspectives for the purpose of informing policy-making. A group size of between seven and 15 participants is known to be ideal as they are large enough to represent all major perspectives but small enough to work towards agreement, effectively allowing for efficient group processes to emerge while preserving individuality, (Phillips & Phillips, 1993). However, capturing an all-round set of

preferences was not among the primary aims of the exercise. The value scale of the treatment discontinuation attribute and, more specifically, the “lower” reference level of “10%” could be perceived as a limitation because it influenced the negative partial value scores of two treatments whose performance was worse. This was the outcome of consultation with an oncologist, based on evidence from one of the clinical trials’ placebo-controlled arms, because it was believed to better resemble BSC used in practice; although others might have chosen a different performance level to define the “lower” reference level, the overall ranking of the treatments did not change when altering the lowest reference level to a much less preferred hypothetical performance (20% lower than the worst performing option), while keeping the weights constant.

One major advantage in MCDA, is that it can be tailor-made to reflect decision-makers’ needs, by taking into account different fundamental objectives through the consideration of a variety of criteria, reflecting their priorities (by eliciting relative weights) and representing their preferences (by eliciting value functions). However, it should be recognised that the emerging differences that have been described above, prevent the direct comparison of overall value scores for alternative options; these would require identical value trees (i.e. the same set of criteria, weights and value functions across settings), in addition to the same evidence on options performance. The ranking comparisons that we have made in this study using ordinal scales reflect these limitations.

Conclusions and implications

In this study, we tested the application of AVF, a multi-criteria value framework, in collaboration with HTA decision-makers in order to deduce its feasibility and compare results across settings, in an effort to investigate its potential usefulness and limitations for the purposes of HTA. We found that the AVF methodology can act as a valuable decision support tool because of the transparent construction of value preferences in a collaborative manner, which facilitates the evaluation processes of groups, including the elicitation of value

preferences and trade-offs. Although we observed setting-specific differences in value perceptions, the rankings of drugs remained consistent across all countries. Based on the evidence used in the exercise, a coverage decision using this method would have pointed towards a different recommendation denoting differences in value between the first two treatments, in contrast with the cost minimisation approach adopted or the price parity attained between the two in real life.

Despite a number of limitations relating to data and process issues and the existence of broader challenges with the use of MCDA in HTA due to specific methodological requirements which would need to be satisfied, the present study has demonstrated that an MCDA framework can, in fact, provide meaningful valuations of novel health technologies which, in turn, can inform coverage decisions.

The MCDA methodology adopted enabled participants in the study countries to reflect on certain value dimensions and incorporate these more explicitly in the deliberation process, supporting its use as a transparent value communication tool. Future research efforts could involve similar cross-country case studies, the advancement of MCDA methods and their alignment with HTA policy needs, or repeating the study with different participants to understand whether similarities and differences identified in this study can be replicated.

References

- AEMPS. (2015). INFORME DE POSICIONAMIENTO TERAPÉUTICO PT-ENZALUTAMIDA/V1/30072015. *Agencia Española de Medicamentos y Productos Sanitarios*.
- Anderson, J.L., Heidenreich, P.A., Barnett, P.G., Creager, M.A., Fonarow, G.C., Gibbons, R.J., et al. (2014). ACC/AHA statement on cost/value methodology in clinical practice guidelines and performance measures: a report of the American College of Cardiology/American Heart Association Task Force on Performance Measures and Task Force on Practice Guidelines. *Journal Of The American College Of Cardiology*, 63, 2304-2322.
- Angelis, A. (2018). Evaluating the benefits of new drugs in health technology assessment using multiple criteria decision analysis: a case study on metastatic prostate cancer with the dental and pharmaceuticals benefits agency (TLV) in Sweden. *MDM Policy & Practice*, 3.
- Angelis, A., & Kanavos, P. (2016a). Critique of the American Society of Clinical Oncology Value Assessment Framework for Cancer Treatments: Putting Methodologic Robustness First. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*.
- Angelis, A., & Kanavos, P. (2016b). Value-Based Assessment of New Medical Technologies: Towards a Robust Methodological Framework for the Application of Multiple Criteria Decision Analysis in the Context of Health Technology Assessment. *Pharmacoeconomics*, 34, 435-446.
- Angelis, A., & Kanavos, P. (2017). Multiple Criteria Decision Analysis (MCDA) for evaluating new medicines in Health Technology Assessment and beyond: the Advance Value Framework. *Social Science & Medicine*, 188, 137-156.
- Angelis, A., Kanavos, P., & Montibeller, G. (2016). Resource allocation and priority setting in health care: a multi-criteria decision analysis problem of value? *Global Policy*.
- Angelis, A., Lange, A., & Kanavos, P. (2018). Using health technology assessment to assess the value of new medicines: results of a systematic review and expert consultation across eight European countries. *European Journal of Health Economics*.
- Angelis, A., Montibeller, G., Hochhauser, D., & Kanavos, P. (2017). Multiple criteria decision analysis in the context of health technology assessment: a simulation exercise on metastatic colorectal cancer with multiple stakeholders in the English setting. *BMC Medical Informatics and Decision Making*, 17.
- AOTMiT. (2017). Rekomendacja nr 19/2017 z dnia 27 marca 2017 r. Prezesa Agencji Oceny Technologii Medycznych i Taryfikacji w sprawie objęcia refundacją produktu leczniczego Xtandi, enzalutamidum, kapsułki miękkie, 40 mg, 112 kaps., w ramach programu lekowego „Leczenie opornego na kastrację raka gruczołu krokowego z przerzutami (ICD-10 C-61). http://bipold.aotm.gov.pl/assets/files/zlecenia_mz/2017/008/REK/RP_19_2017_Xtandi.pdf.
- Bach, P. (2015). DrugAbacus App. Memorial Sloan Kettering Cancer Center.
- Baltussen, R., Jansen, M.P.M., Bijlmakers, L., Grutters, J., Kluytmans, A., Reuzel, R.P., et al. (2017). Value Assessment Frameworks for HTA Agencies: The

- Organization of Evidence-Informed Deliberative Processes. *Value in Health*, 20, 256-260.
- Bana e Costa, C., & Beinat, E. (2005). Model-structuring in public decision-aiding. Operational Research working papers. London: London School of Economics and Political Science.
- Bana e Costa, C., De Corte, J., & Vansnick, J. (2016a). M-MACBETH website.
- Bana e Costa, C., De Corte, J., & Vansnick, J. (2016b). On the Mathematical Foundations of MACBETH. In S. Greco, M. Ehrgott, & J. Figueira (Eds.), *Multiple Criteria Decision Analysis: State of the Art Surveys*: Springer New York.
- Bana e Costa, C., Lourenço, J., Oliveira, M., & Bana e Costa, J. (2014). A Socio-technical Approach for Group Decision Support in Public Strategic Planning: The Pernambuco PPA Case. *Group Decision & Negotiation*, 23, 5-29.
- Bana e Costa, C., & Vansnick, J. (1999). The MACBETH Approach: Basic Ideas, Software, and an Application. In N. Meskens, & M. Roubens (Eds.), *Advances in Decision Analysis* pp. 131-157): Springer Netherlands.
- Bana e Costa, C.A., Corrêa, É.C., De Corte, J.-M., & Vansnick, J.-C. (2002). Facilitating bid evaluation in Public call for tenders: a socio-technical approach. *Omega*, 30, 227.
- Bana E Costa, C.A., De Corte, J.-M., & Vansnick, J.-C. (2012). MACBETH. *International Journal of Information Technology & Decision Making*, 11, 359-387.
- Bana e Costa, C.A., Ensslin, L., Corrêa, É.C., & Vansnick, J.-C. (1999). Decision Support Systems in action: Integrated application in a multicriteria decision aid process. *European Journal of Operational Research*, 113, 315-335.
- Bana e Costa, C.A., & Oliveira, M.D. (2012). A multicriteria decision analysis model for faculty evaluation. *Omega*, 40, 424-436.
- Bana e Costa, C.A., & Vansnick, J.-C. (1994). MACBETH — An Interactive Path Towards the Construction of Cardinal Value Functions. *International Transactions in Operational Research*, 1, 489.
- Bana e Costa, C.A., & Vansnick, J.-C. (1997). Applications of the MACBETH Approach in the Framework of an Additive Aggregation Model. *Journal of Multi-Criteria Decision Analysis*, 6, 107-114.
- Belton, V., & Stewart, T. (2002). *Multiple criteria decision analysis: an integrated approach*. Dordrecht: Kluwer Academic Publishers.
- Berthold, D.R., Pond, G.R., Soban, F., de Wit, R., Eisenberger, M., & Tannock, I.F. (2008). Docetaxel plus prednisone or mitoxantrone plus prednisone for advanced prostate cancer: updated survival in the TAX 327 study. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*, 26, 242-245.
- Bianchini, D., Lorente, D., Rodriguez-Vida, A., Omlin, A., Pezaro, C., Ferraldeschi, R., et al. (2014). Antitumour activity of enzalutamide (MDV3100) in patients with metastatic castration-resistant prostate cancer (CRPC) pre-treated with docetaxel and abiraterone. *European Journal Of Cancer (Oxford, England: 1990)*, 50, 78-84.
- BNF. (2015). British National Formulary 69. <https://www.bnf.org/>.
- Burström, K., Johannesson, M., & Diderichsen, F. (2001). Swedish population health-related quality of life results using the EQ-5D. *Quality Of Life Research: An International Journal Of Quality Of Life Aspects Of Treatment, Care And Rehabilitation*, 10, 621-635.

- Center, M.M., Jemal, A., Lortet-Tieulent, J., Ward, E., Ferlay, J., Brawley, O., et al. (2012). International Variation in Prostate Cancer Incidence and Mortality Rates. *European Urology*, 61, 1079-1092.
- Cherny, N.I., Sullivan, R., Dafni, U., Kerst, J.M., Sobrero, A., Zielinski, C., et al. (2015). A standardised, generic, validated approach to stratify the magnitude of clinical benefit that can be anticipated from anti-cancer therapies: the European Society for Medical Oncology Magnitude of Clinical Benefit Scale (ESMO-MCBS). *Annals Of Oncology: Official Journal Of The European Society For Medical Oncology / ESMO*, 26, 1547-1573.
- Chongtrakul, P., Sumpradit, N., & Yoongthong, W. (2005). *ISafe and the evidence-based approach for essential medicines selection in Thailand*. pp. 18-19): Essential Drugs Monitor.
- Claxton, K., Martin, S., Soares, M., Rice, N., Spackman, E., Hinde, S., et al. (2015). Methods for the estimation of the National Institute for Health and care excellence cost- effectiveness threshold. *Health Technology Assessment*, 19, 1-503.
- Cleemput, I., Devriese, S., Kohn, L., Devos, C., van Til, J., Groothuis-Oudshoorn, C.G.M., et al. (2018). What Does the Public Want? Structural Consideration of Citizen Preferences in Health Care Coverage Decisions. *MDM Policy & Practice*, 3.
- Clement, F.M., Harris, A., Li, J.J., Yong, K., Lee, K.M., & Manns, B.J. (2009). Using effectiveness and cost-effectiveness to make drug coverage decisions: a comparison of Britain, Australia, and Canada. *JAMA*, 302, 1437-1443.
- Cohen, D. (2017). Cancer drugs: high price, uncertain value. *BMJ*, 359.
- Collins, R., Fenwick, E., Trowman, R., Perard, R., Norman, G., Light, K., et al. (2007). A systematic review and economic model of the clinical effectiveness and cost-effectiveness of docetaxel in combination with prednisone or prednisolone for the treatment of hormone-refractory metastatic prostate cancer. *Health Technology Assessment (Winchester, England)*, 11, iii.
- Connock, M., Cummins, E., Shyangdan, D., Hall, B., Grove, A., & Clarke, A. (2011). Abiraterone acetate for the treatment of metastatic, castrate-resistant prostate cancer following previous cytotoxic chemotherapy: A Single Technology Appraisal. Warwick Evidence.
- Currim, I.S., & Sarin, R.K. (1984). A Comparative Evaluation of Multiattribute Consumer Preference Models. *Management Science*, 30, 543-561.
- Danner, M., Hummel, J.M., Volz, F., van Manen, J.G., Wiegard, B., Dintsios, C.-M., et al. (2011). Integrating patients' views into health technology assessment: Analytic hierarchy process (AHP) as a method to elicit patient preferences. *International Journal Of Technology Assessment In Health Care*, 27, 369-375.
- de Bono, J.S., Logothetis, C.J., Molina, A., Fizazi, K., North, S., Chu, L., et al. (2011). Abiraterone and increased survival in metastatic prostate cancer. *The New England Journal Of Medicine*, 364, 1995-2005.
- de Bono, J.S., Oudard, S., Ozguroglu, M., Hansen, S., Machiels, J.-P., Kocak, I., et al. (2010). Prednisone plus cabazitaxel or mitoxantrone for metastatic castration-resistant prostate cancer progressing after docetaxel treatment: a randomised open-label trial. *Lancet*, 376, 1147-1154.
- Devlin, N., & Lorgelly, P. (2017). QALYs as a measure of value in cancer. *Journal of Cancer Policy*, 11, 19-25.
- Devlin, N., & Sussex, J. (2011). Incorporating multiple criteria in HTA: methods and processes. London: Office of Health Economics.

- Drummond, M., Tarricone, R., & Torbica, A. (2013). Assessing the added value of health technologies: reconciling different perspectives. *Value In Health: The Journal Of The International Society For Pharmacoeconomics And Outcomes Research*, 16, S7-S13.
- Efthymiadou, O., Mossman, J., & Kanavos, P. (2019). Health related quality of life aspects not captured by EQ-5D-5L: results from an international survey of patients. *Health Policy*, 123, 159-165.
- EMA. (2016a). Jevtana (cabazitaxel) EPAR - Product Information. European Medicines Agency.
- EMA. (2016b). Xtandi (enzalutamide) EPAR - Product Information. European Medicines Agency.
- EMA. (2016c). Zytiga (abiraterone) EPAR - Product Information. European Medicines Agency.
- Faden, R.R., Chalkidou, K., Appleby, J., Waters, H.R., & Leider, J.P. (2009). Expensive Cancer Drugs: A Comparison between the United States and the United Kingdom. *Milbank Quarterly*, 87, 789-819.
- Fasolo, B., & Bana e Costa, C.A. (2014). Tailoring value elicitation to decision makers' numeracy and fluency: Expressing value judgments in numbers or words. *Omega*, 44, 83-90.
- Fizazi, K., Scher, H.I., Molina, A., Logothetis, C.J., Chi, K.N., Jones, R.J., et al. (2012). Abiraterone acetate for treatment of metastatic castration-resistant prostate cancer: final overall survival analysis of the COU-AA-301 randomised, double-blind, placebo-controlled phase 3 study. *The Lancet. Oncology*, 13, 983-992.
- Franco, L., & Montibeller, G. (2010a). Problem Structuring for Multicriteria Decision Analysis Interventions analysis interventions. *Wiley Encyclopedia of Operations Research and Management Science*.
- Franco, L.A., & Montibeller, G. (2010b). Facilitated modelling in operational research. *European Journal of Operational Research*, 205, 489-500.
- GHEMA. (2016). ENZALUTAMIDA en cáncer próstata metastásico resistente a la castración (previo a quimioterapia). Grupo de Evaluación de Novedades, EStandarización e Investigación en Selección de Medicamentos.
- Goetghebeur, M.M., Wagner, M., Khoury, H., Levitt, R.J., Erickson, L.J., & Rindress, D. (2008). Evidence and Value: Impact on DEcisionMaking--the EVIDEM framework and potential applications. *BMC Health Services Research*, 8, 270.
- Goetghebeur, M.M., Wagner, M., Khoury, H., Rindress, D., Grégoire, J., & Deal, C. (2010). Combining multicriteria decision analysis, ethics and health technology assessment: Applying the EVIDEM decisionmaking framework to growth hormone for Turner syndrome patients. *Cost Effectiveness and Resource Allocation*, 8, <xocs:firstpage xmlns:xocs=""/>.
- Horwich, A., Parker, C., de Reijke, T., & Kataja, V. (2013). Prostate cancer: ESMO Clinical Practice Guidelines for diagnosis, treatment and follow-up. *Annals Of Oncology: Official Journal Of The European Society For Medical Oncology / ESMO*, 24 Suppl 6, vi106-vi114.
- Ijzerman, M.J., van Til, J.A., & Snoek, G.J. (2008). Comparison of two multi-criteria decision techniques for eliciting treatment preferences in people with neurological disorders. *The Patient*, 1, 265-272.
- INAMI. (2019). Remboursement de médicaments : Décisions ministérielles et rapports d'évaluation de la CRM. <https://www.inami.fgov.be/fr/programmes-web/Pages/applications-rapports-crm.aspx>.

- Jansen, J.P., Fleurence, R., Devine, B., Itzler, R., Barrett, A., Hawkins, N., et al. (2011). Interpreting indirect treatment comparisons and network meta-analysis for health-care decision making: report of the ISPOR Task Force on Indirect Treatment Comparisons Good Research Practices: part 1. *Value In Health: The Journal Of The International Society For Pharmacoeconomics And Outcomes Research*, 14, 417-428.
- Jaramillo, H.E.C., Goetghebeur, M., & Moreno-Mattar, O. (2016). TESTING MULTI-CRITERIA DECISION ANALYSIS FOR MORE TRANSPARENT RESOURCE-ALLOCATION DECISION MAKING IN COLOMBIA. 32, 307-314.
- Kanavos, P., & Angelis, A. (2013). Multiple criteria decision analysis for value based assessment of new medical technologies: a conceptual framework. The LSE Health Working Paper Series in Health Policy and Economics: London School of Economics and Political Science.
- Kearns, B., Lloyd Jones, M., Stevenson, M., & Littlewood, C. (2013). Cabazitaxel for the Second-Line Treatment of Metastatic Hormone-Refractory Prostate Cancer: A NICE Single Technology Appraisal. *PharmacoEconomics*, 31, 479-488.
- Keeney, R. (1992). *Value focused thinking: a path to creative decision making*. Cambridge: Harvard University Press.
- Keeney, R., & Raiffa, H. (1993). *Decisions with multiple objectives: preferences and value trade-offs*. Cambridge: Cambridge University Press.
- Keeney, R.L. (1982). Decision Analysis: An Overview. *Operations Research*, 30, 803.
- Keeney, R.L. (2002). Common Mistakes In Making Value Trade-Offs. *Operations Research*, 50, 935.
- Keeney, R.L., & Gregory, R.S. (2005). Selecting Attributes to Measure the Achievement of Objectives. *Operations Research*, 53, 1-11.
- Linley, W.G., & Hughes, D.A. (2013). Societal views on NICE, cancer drugs fund and value-based pricing criteria for prioritising medicines: a cross-sectional survey of 4118 adults in Great Britain. *Health Economics*, 22, 948-964.
- London School of Economics. (2019). ADVANCE_HTA Project. <http://www.lse.ac.uk/lse-health/research/projects/advance-hta>.
- Loriot, Y., Bianchini, D., Ileana, E., Sandhu, S., Patrikidou, A., Pezaro, C., et al. (2013). Antitumour activity of abiraterone acetate against metastatic castration-resistant prostate cancer progressing after docetaxel and enzalutamide (MDV3100). *Annals Of Oncology: Official Journal Of The European Society For Medical Oncology / ESMO*, 24, 1807-1812.
- Malone, D.C., Berg, N.S., Claxton, K., Garrison, L.P., Jr., Ijzerman, M., Marsh, K., et al. (2016). International Society for Pharmacoeconomics and Outcomes Research Comments on the American Society of Clinical Oncology Value Framework. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*, 34, 2936-2937.
- Marsh, K., Ijzerman, M., Thokala, P., Baltussen, R., Boysen, M., Kaló, Z., et al. (2016). Multiple Criteria Decision Analysis for Health Care Decision Making-Emerging Good Practices: Report 2 of the ISPOR MCDA Emerging Good Practices Task Force. *Value In Health: The Journal Of The International Society For Pharmacoeconomics And Outcomes Research*, 19, 125-137.

- Marsh, K., Lanitis, T., Neasham, D., Orfanos, P., & Caro, J. (2014). Assessing the Value of Healthcare Interventions Using Multi-Criteria Decision Analysis: A Review of the Literature. *Pharmacoeconomics*, 32, 345-365.
- Marsh, K., Sculpher, M., Caro, J.J., & Tervonen, T. (2018). The Use of MCDA in HTA: Great Potential, but More Effort Needed. *Value in Health*.
- Montibeller, G., & Winterfeldt, D. (2015). Cognitive and Motivational Biases in Decision and Risk Analysis. *Risk Analysis: An International Journal*, 35, 1230-1251.
- Morton, A. (2017). Treacle and Smallpox: Two Tests for Multicriteria Decision Analysis Models in Health Technology Assessment. *Value in Health*, 20, 512-515.
- Nadler, E., Eckert, B., & Neumann, P.J. (2006). Do Oncologists Believe New Cancer Drugs Offer Good Value? *Oncologist*, 11, 90-95.
- NICE. (2012a). Abiraterone for castration-resistant metastatic prostate cancer previously treated with a docetaxel-containing regimen. Technology Appraisal Guidance 259: National Institute for Health and Care Excellence.
- NICE. (2012b). Cabazitaxel for hormone-refractory metastatic prostate cancer previously treated with a docetaxel-containing regimen. Technology Appraisal Guidance 255: National Institute for Health and Care Excellence.
- NICE. (2014). Enzalutamide for metastatic hormone- relapsed prostate cancer previously treated with a docetaxel - containing regimen. Technology Appraisal Guidance 316: National Institute for Health and Care Excellence.
- Nicod, E., & Kanavos, P. (2012). Commonalities and differences in HTA outcomes: A comparative analysis of five countries and implications for coverage decisions. *Health Policy*, 108, 167-177.
- NIH. (2016). ClinicalTrials.gov. US National Institutes of Health.
- Nord, E., Daniels, N., & Kamlet, M. (2009). QALYs: Some Challenges. *Value in Health*, 12, S10-S15.
- Obwieszczenie, M. (2017). Obwieszczenie Ministra Zdrowia z dnia 25 października 2017 r. w sprawie wykazu refundowanych leków, środków spożywczych specjalnego przeznaczenia żywieniowego oraz wyrobów medycznych (DZ. URZ. Min. Zdr. 2017.105). <http://www.bip.mz.gov.pl/legislacja/akty-prawne-1/obwieszczenie-ministra-zdrowia-z-dnia-25-pazdziernika-2017-r-w-sprawie-wykazu-refundowanych-lekow-srodkow-spozywczych-specjalnego-przeznaczenia-zywieniowego-oraz-wyrobow-medycznych-na-1-listopada-20/>.
- Peacock, S.J., Richardson, J.R.J., Carter, R., & Edwards, D. (2007). Priority setting in health care using multi- attribute utility theory and programme budgeting and marginal analysis (PBMA). *Social Science and Medicine*, 64, 897-910.
- Phelps, C., & Madhavan, G. (2018). Resource allocation in decision support frameworks. *Cost Effectiveness and Resource Allocation*, 16.
- Phillips, L. (1984). A theory of requisite decision models. *Acta Psychologica*, 56, 29-48.
- Phillips, L. (2007). Decision Conferencing. In W. Edwards, R. Miles, & D. von Winterfeldt (Eds.), *Advances in Decision Analysis: From Foundations to Applications*. Cambridge: Cambridge University Press.
- Phillips, L., & Phillips, M. (1993). Facilitated Work Groups: Theory and Practice. *The Journal of the Operational Research Society*, 44.
- Phillips, L.D. (2017). Best Practice for

- MCDA in Healthcare. In K. Marsh, Goetghebeur, M., P. Thokala, & R. Baltussen (Eds.), *Multi-Criteria Decision Analysis to Support Healthcare Decisions*: Springer.
- Phillips, L.D., & Bana e Costa, C.A. (2007). Transparent prioritisation, budgeting and resource allocation with multi-criteria decision analysis and decision conferencing. *Annals of Operations Research*, 154, 51-68.
- Radaelli, G., Lettieri, E., Masella, C., Merlino, L., Strada, A., & Tringali, M. (2014). Implementation of EUnetHTA core Model® in Lombardia: the VTS framework. *International Journal Of Technology Assessment In Health Care*, 30, 105-112.
- Riemsma, R., Ramaekers, B., Tomini, F., Wolff, R., van Asselt, A., Joore, M., et al. (2013). Abiraterone for the treatment of chemotherapy naïve metastatic castration-resistant prostate cancer: a Single Technology Appraisal. Kleijnen Systematic Reviews Ltd.
- Rodrigues, T.C., Montibeller, G., Oliveira, M.D., & Bana E Costa, C.A. (2017). Modelling multicriteria value interactions with Reasoning Maps. *European Journal of Operational Research*, 258, 1054-1071.
- Schein, E. (1999). *Process consultation revisited: building the helping relationship*. Reading: Addison–Wesley.
- Scher, H.I., Fizazi, K., Saad, F., Taplin, M.-E., Sternberg, C.N., Miller, K., et al. (2012). Increased survival with enzalutamide in prostate cancer after chemotherapy. *The New England Journal Of Medicine*, 367, 1187-1197.
- Schnipper, L.E., Davidson, N.E., Wollins, D.S., Tyne, C., Blayney, D.W., Blum, D., et al. (2015). American Society of Clinical Oncology Statement: A Conceptual Framework to Assess the Value of Cancer Treatment Options. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*, 33, 2563-2577.
- Sculpher, M., Claxton, K., & Pearson, S.D. (2017). Developing a Value Framework: The Need to Reflect the Opportunity Costs of Funding Decisions. *Value in Health*, 20, 234-239.
- Shih, Y.C.T., Ganz, P.A., Aberle, D., Abernethy, A., Bekelman, J., Brawley, O., et al. (2013). Delivering high-quality and affordable care throughout the cancer care continuum. *Journal of Clinical Oncology*, 31, 4151-4157.
- Sullivan, P., Mulani, P., Fishman, M., & Sleep, D. (2007). Quality of life findings from a multicenter, multinational, observational study of patients with metastatic hormone-refractory prostate cancer. *An International Journal of Quality of Life Aspects of Treatment, Care and Rehabilitation - Official Journal of the International Society of Quality of Life Research*, 16, 571-575.
- Sulmasy, D., & Moy, B. (2014). Debating the oncologist's role in defining the value of cancer care: Our duty is to our patients. *Journal of Clinical Oncology*, 32, 4039-4041.
- Sussex, J., Rollet, P., Garau, M., Schmitt, C., Kent, A., & Hutchings, A. (2013a). A Pilot Study of Multicriteria Decision Analysis for Valuing Orphan Medicines. *Value in Health*, 16, 1163-1169.
- Sussex, J., Rollet, P., Garau, M., Schmitt, C., Kent, A., & Hutchings, A. (2013b). Multi-criteria decision analysis to value orphan medicines. London: Office of Health Economics.
- Tannock, I.F., de Wit, R., Berry, W.R., Horti, J., Pluzanska, A., Chi, K.N., et al. (2004). Docetaxel plus prednisone or mitoxantrone plus prednisone for

- advanced prostate cancer. *The New England Journal Of Medicine*, 351, 1502-1512.
- Tervonen, T., Naci, H., van Valkenhoef, G., Ades, A.E., Angelis, A., Hillege, H.L., et al. (2015). Applying Multiple Criteria Decision Analysis to Comparative Benefit-Risk Assessment: Choosing among Statins in Primary Prevention. *Medical Decision Making: An International Journal Of The Society For Medical Decision Making*, 35, 859-871.
- Thokala, P. (2011). Multi criteria decision analysis for health technology assessment: report by the decision support unit. Sheffield: University of Sheffield.
- TLV. (2014). Designation 2775/2013. Tandvårds- och läkemedelsförmånsverket.
- TLV. (2015a). Designation 4774/2014. Tandvårds- och läkemedelsförmånsverket.
- TLV. (2015b). Designation 4852/2014. Tandvårds- och läkemedelsförmånsverket.
- Tony, M., Wagner, M., Khoury, H., Rindress, D., Papastavros, T., Oh, P., et al. (2011). Bridging health technology assessment (HTA) with multicriteria decision analyses (MCDA): field testing of the EVIDEM framework for coverage decisions by a public payer in Canada. *BMC Health Services Research*, 11, 329.
- Torre, L.A. (2015). Global cancer statistics, 2012. *CA: A Cancer Journal For Clinicians*, 65, 87.
- van Valkenhoef, G., Tervonen, T., Zhao, J., de Brock, B., Hillege, H.L., & Postmus, D. (2012). Multicriteria benefit-risk assessment using network meta-analysis. *Journal Of Clinical Epidemiology*, 65, 394-403.
- von Winterfeldt, D., & Edwards, W. (1986). *Decision analysis and behavioral research*. Cambridge: Cambridge University Press.
- Wagner, M., Khoury, H., Bennetts, L., Berto, P., Ehreth, J., Badia, X., et al. (2017). Appraising the holistic value of Lenvatinib for radio-iodine refractory differentiated thyroid cancer: A multi-country study applying pragmatic MCDA. *BMC Cancer*, 17, <xocs:firstpage xmlns:xocs=""/>.
- World Health Organisation Collaborating Centre. (2016). ATC/DDD Index 2016. World Health Organisation Collaborating Centre.
- Wouters, O.J., Naci, H., & Samani, N.J. (2015). QALYs in cost-effectiveness analysis: an overview for cardiologists. *Heart (British Cardiac Society)*, 101, 1868-1873.

Appendix

Model Building: Alternative Treatments Compared and Evidence Considered

The source of evidence used for identifying the performance of options across the evaluation criteria is shown in Table A1.

Model Building: Setting Attribute Ranges and Reference Levels

For the case of clinical therapeutic attributes, “lower” reference levels were based on best standard of care (BSC) performance, coming from the median of the respective placebo arm of the *AFFIRM* trial, with the exception of the HRQoL attribute (EQ-5D utility score) that was based on the utility of stable disease with no treatment coming from past NICE TAs (NICE, 2012a, b). The “higher” reference levels were derived by adding a 20% absolute improvement to the performance level of the best performing option, besides for the case of the HRQoL attribute (EQ-5D utility score) that was based on the general Swedish population (Burström et al., 2001). The rationale was to design a value scale incorporating a “global” reference level (Belton & Stewart, 2002), reflecting an “ideal” performance (as proxied by the 20% improvement in best available performance), corresponding to the 100 anchor level of the value scale. This could also offer a flexibility margin to be able to incorporate the performance of future improved options within the same elicited value scale. Consequently, two reference levels within the attribute range were defined in most cases: i) the “lower” reference level (x_l) (i.e. BSC-based satisfactory performance), acting on the same time also as the minimum limit of the attribute range (x_*); and ii) the “higher” reference level (x_h) (i.e. 20% better than the best performing option), acting on the same time as the maximum limit of the attribute range (x^*) to give $x_* = x_l \leq x_h = x^*$.

A similar, but reverse, logic was used for setting the reference levels in the “treatment discontinuation” attribute of the safety cluster; the “lower” reference level was defined to be equal to the BSC (i.e. placebo) arm of the *AFFIRM* trial. However, contrary to the logic

adopted so far for the therapeutic benefit criteria, the “higher” reference level was not set equal to 20% worse than the best performing option (because the lower the performance, the higher the value), but rather equal to the minimum, i.e. worst possible, natural limit of the attribute scale (i.e. 0%) which was regarded as an “ideal” level. In turn, the minimum limit of the scale was derived by worsening the performance of the worst performing treatment option by 20%. A similar approach was used for setting the reference levels of the qualitative “contraindications” attribute, defining the “higher” reference level equal to the maximum (i.e. most attractive) limit of the attribute scale (i.e. none known contraindications) and the “lower” reference level equal to the minimum (i.e. least attractive) limit of the attribute scale.

For the innovation attributes, the “higher” reference level was set either equal to 20% better than the best performing option for the case of natural quantitative attributes (e.g. number of new indications for which the technology is investigated in a given clinical development stage), or equal to the maximum, i.e. best possible, limit of the scale for the case of constructed qualitative attributes (e.g. the existence of any special instructions, the technology's relative market entrance in regards to its ATC Level), reflecting a “global” versus “local” scaling approach respectively. Given that the BSC performance was irrelevant to be used as satisfactory level in the innovation attributes, and any efforts to derive a “satisfactory” level would be subjective in nature, the minimum limit of the scale for each attribute was used as a “lower” reference level. Therefore the “lower” reference level was based on the worst performance plausible as inferred from the lowest possible limit of the scales, both for the case of natural quantitative attributes (e.g. 0 number of new indications for which the technology is investigated in a given clinical development stage), and the case of constructed qualitative attributes (e.g. worst possible combination of special instructions, 5th entrance at an ATC level).

For the socioeconomics attribute (impact on direct costs), the “higher” reference level was based on the BSC's impact on cost (i.e. £0 impact on costs), given that by definition impact on costs for all treatment options are incremental to BSC, and the “lower” reference

level was derived by adding a 20% absolute increment to the worst performing option (i.e. to the one with the biggest impact on costs).

“Lower” and “higher” reference levels for all attributes at the pre-workshop stage and the basis of their selection are outlined in Table A2 (assuming no impact of luteinizing hormone-releasing hormone analogue).

Model Assessment and Appraisal: Decision Conference

On the day of each decision conference the preliminary model was validated with the participants by revising it cluster by cluster through an open discussion, seeking group consensus and adopting an iterative and interactive-model-building process where debate was encouraged and differences of opinion were actively sought.

In terms of the decision-aiding methodology used, the lead author acted as an impartial facilitator with the aim of enhancing content and process interaction, while refraining from contributing to the content of the group’s discussions, essentially guiding the group in how to think about the issues but not what to think (Phillips & Bana e Costa, 2007; Schein, 1999).

In terms of facilities, the rooms of the decision conferences had a Π-shaped meeting table for all the participants to have direct eye to eye contact, with an overhead projector screen and a second portable projector or large TV screen. The M-MACBETH software (more information provided in the MCDA Technique section of the main text and below) was operated using a laptop, the screen of which was connected to the projector, and the second screen was used to show the list of the evaluation criteria together with their “lower” and “higher” reference levels.

The decision conferences took place over a full working day or two half working days; in the former case, there was one lunch break and two coffee breaks throughout the day, whereas in the latter case only a coffee break took place around the middle of each session. In each decision conference, the day started with an overview of the MCDA methodology

adopted and the description of the preliminary version of the value tree which was then analysed cluster by cluster. At the beginning of each cluster the value tree was validated; the various criteria were explained, followed by a group discussion relating to their relevance and completeness. As a result of this iterative process, some of the criteria were not included because they were perceived as irrelevant or non-fundamental. Schematic illustrations of the final versions of the value trees are shown in Figure A1. Then, value functions were elicited for the different criteria and relative weights were assigned within the clusters. Finally, relative weights were assigned across clusters, enabling the calculation of the options' overall WPV scores.

Model Assessment and Appraisal: MCDA Technique

MACBETH uses seven semantic categories ranging between “no difference” to “extreme difference”, in order to distinguish between the value of different attribute levels. Based on these qualitative judgements of difference and, by analysing judgmental inconsistencies, it facilitates the move from ordinal preference modeling, a cognitively less demanding elicitation of preferences, to a quantitative value function. The approach has evolved through the course of theoretical research and real world practical applications, making it an interactive decision support system that facilitates decision-makers' communication. An example of the type of questioning being asked would be “What do you judge to be the difference of value between x' and x'' ?” where x' and x'' are two different attribute levels of attribute x , across the plausible range (i.e. $x^* \leq x', x'' \leq x^*$). The value judgements matrix for the Overall Survival attribute and their conversion into its value function is provided as an example in Figure A2.

Following the elicitation of value functions, criteria baseline weights can be elicited. Questions of direct importance for a criterion such as “*How important is a given criterion?*” are known to be as one of the most common mistakes when making value trade-offs because they are assessing them independent of the respective attribute ranges (Keeney, 2002). In contrast, indirect weighting technique that assess value trade-offs in tandem with the

respective ranges of attributes should be employed. For example, the quantitative swing weighting technique asks for judgments of relative value between ‘swings’ (i.e. changes from standard lower level x_* to higher reference level x^* on each x^{th} attribute) taking the form “*How would you rank the relative importance of the criteria, considering their attributes ranges relative to 100 for the highest-ranked criterion considering its range?*”. Each swing, i.e. a relative change from a lower attribute level to a higher attribute level, is valued between 0 and 100, with the most valuable swing anchored as 100 (von Winterfeldt & Edwards, 1986). Normalised weights are then calculated, as a proportion of each swing weight, so the normalised weights summed to 100%. Instead, relative attribute weights were calculated using an alternative qualitative swing weighting protocol, by using the MACBETH procedure to elicit the differences in attractiveness between the lower and higher reference levels of the different attributes, initially at individual level and then at criteria cluster level (i.e. by considering multiple attribute swings on the same time) (Bana e Costa et al., 2016b; Bana E Costa et al., 2012).

Finally criteria preference value scores and the respective weights can be combined together through an additive aggregation approach as described in equation 2 (if the adequate conditions of complete and transitive preferences are met as well as multi-attribute preferential independence conditions (von Winterfeldt & Edwards, 1986)).

The M-MACBETH software automatically performs consistency checking between the qualitative judgements expressed, and in addition a second consistency check was manually performed by the author to validate the cardinality, i.e. interval nature, of the emerging value scale. This was done by comparing the sizes of the intervals between the proposed scores and inviting participants to adjust them if necessary (Fasolo & Bana e Costa, 2014), a requirement which is essential for the application of simple additive value models.

Figure Captions

Tables and Figures

Table 1: Criteria definitions, their consideration in each jurisdiction and their ranking

Criteria Sub-Domain	Evaluation criteria	Definition	Country (competent HTA organisation)			
			Belgium (INAMI/RIZIV)	Poland (AOTMiT)	Andalusia (AETSA)	Sweden (TLV)
Criteria Domain 1: Therapeutic Benefit						
Direct endpoints	Overall survival x Health related quality of life*	The median time from treatment randomisation to death adjusted for the mean health related quality of life using the EQ-5D utility score	✓ (1 st)	✓ (1 st)	✓ (1 st)	✓ (1 st)
Indirect endpoints	Radiographic tumour progression	The median survival time on which patients have not experienced disease progression (using RECIST criteria)	✓ (5 th)	✓ (5 th)	✓ (4 th)	✓ (5 th)
	PSA response	The proportion of patients having a ≥50% reduction in PSA			✓ (8 th)	
Criteria Domain 2: Safety Profile						
Tolerability	Treatment discontinuation	The proportion of patients discontinuing treatment due to adverse events	✓ (3 rd)	✓ (3 rd)	✓ (2 nd)	✓ (2 nd)
Contra-indications & warnings	Contra-indications	The existence of any type of contra-indication accompanying the treatment	✓ (4 th)		✓ (3 rd)	✓ (4 th)
Criteria Domain 3: Innovation Level						
Type and timing of innovation	ATC Level 1	The technology's relative market entrance in regards to its ATC Level 1 (Anatomical)				
	ATC Level 2	The technology's relative market entrance in regards to its ATC Level 2 (Therapeutic)				
	ATC Level 3	The technology's relative market entrance in regards to its ATC Level 3 (Pharmacological)				
	ATC Level 4	The technology's relative market entrance in regards to its ATC Level 4 (Chemical)	✓ (6 th)		✓ (10 th)	

	ATC Level 5	The technology's relative market entrance in regards to its ATC Level 5 (Molecular)				
Spill-over effect	Phase 1	The number of new indications for which the technology is investigated in Phase 1 clinical trials				
	Phase 2	The number of new indications for which the technology is investigated in Phase 2 clinical trials	✓ (8 th)			
	Phase 3	The number of new indications for which the technology is investigated in Phase 2 clinical trials	✓ (9 th)		✓ (9 th)	
	Marketing authorisation	The number of new indications that the technology has gained an approval for at the stage of marketing authorisation	✓ (10 th)		✓ (7 th)	
Patient convenience	Delivery posology	The combination of the delivery system (RoA and dosage form) with the posology (frequency of dosing and duration of administration) of the treatment	✓ (7 th)	✓ (4 th)	✓ (6 th)	✓ (7 th)
	Special instructions	The existence of any special instructions accompanying the administration of the treatment	✓ (11 th)	✓ (6 th)	✓ (11 th)	✓ (6 th)
Criteria Domain 4: Socio-Economic Impact						
Direct costs	Medical costs impact	The impact of the technology on direct medical costs excluding the purchasing costs of the technology	✓ (2 nd)	✓ (2 nd)	✓ (5 th)	✓ (3 rd)

Notes: *: Aggregation between OS and HRQoL criteria took place due to preference-dependence leading to a combined criterion; PSA= prostate-specific antigen; ATC=Anatomical Therapeutic Chemical classification system; RoA=Route of Administration.

Source: The authors, based on DCs in Andalusia/Spain, Belgium, Poland and Sweden.

Table 2: Performance matrix and reference levels considered across the final criteria attributes

Criterion name	Attribute metric	Lower level	Abiraterone	Cabazitaxel	Enzalutamide	Higher level
Overall survival (OS)*	Months	13.6	15.8	15.1	18.4	22.1
Health Related Quality of Life (HRQoL), stable disease*	Utility (EQ-5D)	0.72	0.76	0.76***	0.76	0.82
Health Related Quality of Life (HRQoL), progressive disease*	Utility (EQ-5D)	0.64	0.64	0.64	0.64	0.82
OS X HRQoL**	Quality adjusted life months (QALMs)	9.2	11	10.5	12.8	18.1
Radiographic tumour progression, i.e. progression free survival (PFS)	Months	2.9	5.6	8.8	8.3	10.6
PSA response	% of patients	1.5	29.5	39.2	54	64.8
Treatment discontinuation	% of patients	10	19	18	8	0
Contra-indication(s)	Type of contra-indication	hyp + hep imp + low neut	hyp + hep imp	hyp + hep imp + low neut	hyp	None
ATC Level 4, i.e. chemical mechanism of action	Relative market entrance	5 th	2 nd	2 nd	1 st	1 st
Phase 2	Number of new indications	0	1	13	4	16
Phase 3	Number of new indications	0	1	2	0	2
Marketing authorisation	Number of new indications	0	0	0	0	1
Delivery posology	Type of delivery system & posology combinations	Oral, daily - one off + IV, every 3 weeks - 1 hr	Oral, daily - one off	Oral, daily - one off + IV, every 3 weeks - 1 hr	Oral, daily - one off	Oral, daily - one off

Special instructions	Type(s) of special instructions	Concomitant and/or pre- med + no food	Concomitant and/or pre- med + no food	Concomitant and/or pre- med	None	None
Medical costs impact	GBP	10,000	5,750	7,992	567	0

Notes: * Used for the calculation of the quality adjusted life months (QALMs) attribute of the aggregated OS x HRQoL criterion; ** Calculated assuming an equal 50% split in time duration between the stable disease and progressive disease states in HRQoL; *** Used the same score of the other two options as data not available; hyp = hypersensitivity; hep imp = hepatic impairment; low neut = low neutrophil count.

Source: The authors from the literature.

Table 3: Number of criteria attributes per cluster, relative weights per criteria cluster and their ranking across the four HTA settings.

<i>HTA Agency/ Criteria Clusters</i>	Sweden (TLV)			Andalusia (AETSA)			Poland (AOTMiT)			Belgium (INAMI-RIZIV)		
	Criteria numbers	Criteria weights	Criteria ranking	Criteria numbers	Criteria weights	Criteria ranking	Criteria numbers	Criteria weights	Criteria ranking	Criteria numbers	Criteria weights	Criteria ranking
Therapeutic Benefit	2	44.5	1 st	3	54.3	1 st	2	40.0	1 st	2	40.0	1 st
Safety Profile	2	33.3	2 nd	2	26.0	2 nd	1	20.0	3 rd	2	26.7	2 nd
Innovation Level	2	7.4	4 th	5	11.8	3 rd	2	10.0	4 th	6	13.3	4 th
Socioeconomic Impact	1	14.8	3 rd	1	7.9	4 th	1	30.0	2 nd	1	20.0	3 rd
Total	7	100		11	100		6	100		11	100	

Source: The authors based on input from decision conferences.

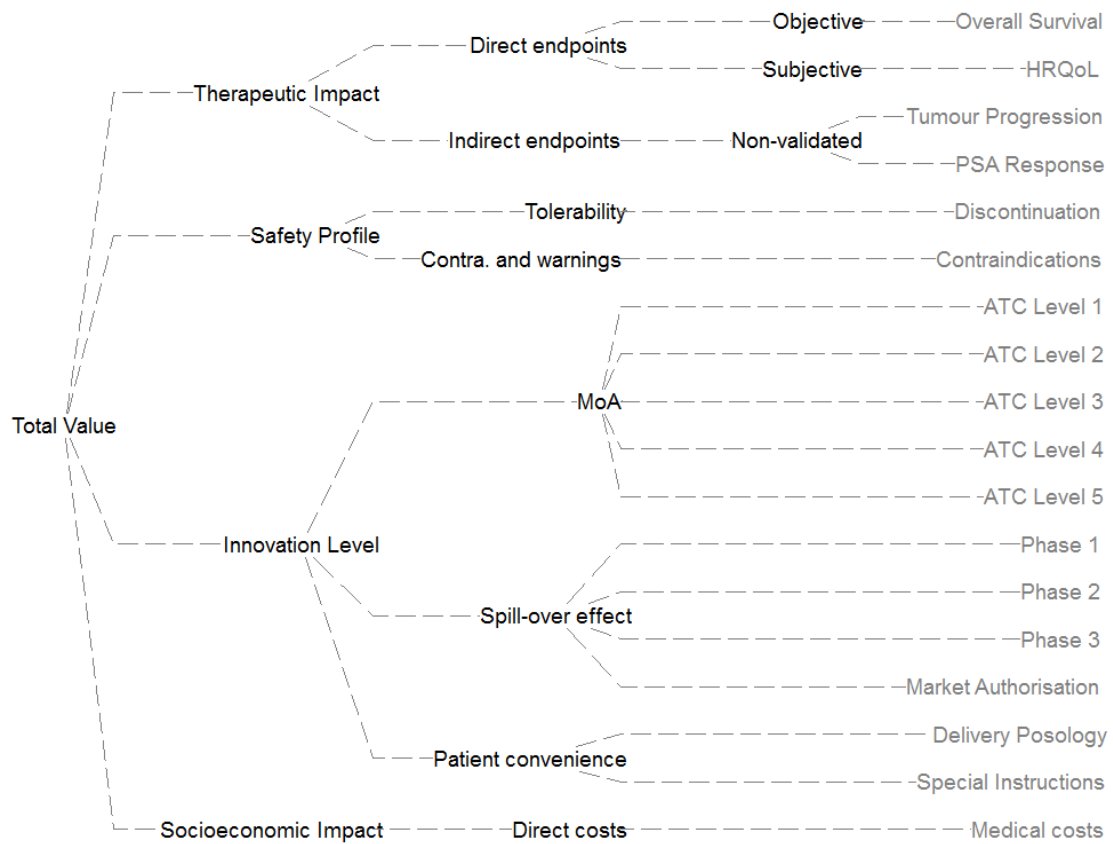
Table 4: Overall weighted preference value (WPV) scores, costs and costs per unit of value across the four HTA settings.

<i>Treatments/ HTA agency</i>	Enzalutamide		Abiraterone		Cabazitaxel	
	<i>Overall WPV score</i>	<i>Ranking per country</i>	<i>Overall WPV score</i>	<i>Ranking per country</i>	<i>Overall WPV score</i>	<i>Ranking per country</i>
Sweden (TLV)	55.1	1 st	2.4	2 nd	-3.4	3 rd
Andalusia (AETSA)	49.1	1 st	8.8	2 nd	4.4	3 rd
Poland (AOTMiT)	59.9	1 st	12.1	2 nd	3.5	3 rd
Belgium (INAMI-RIZIV)	58.6	1 st	16.0	2 nd	10.9	3 rd
<i>Costs (£)</i>	24,600		21,900		23,900	
	<i>Cost per unit of value</i>	<i>Ranking per country</i>	<i>Cost per unit of value</i>	<i>Ranking per country</i>	<i>Cost per unit of value</i>	<i>Ranking per country</i>
Sweden (TLV)	447	1 st	9,221	2 nd	N/A	3 rd
Andalusia (AETSA)	501	1 st	2,496	2 nd	5,481	3 rd
Poland (AOTMiT)	410	1 st	1,805	2 nd	6,816	3 rd
Belgium (INAMI-RIZIV)	420	1 st	1,366	2 nd	2,196	3 rd

Note: No cost-per-unit of value was calculated because of the negative overall WPV score (i.e. having a worst overall performance compared to the performance of the lower reference level), which would produce a negative cost-per-unit of value ($£23,900/(-3.4) = -7,072$) and would therefore faultily “improve” the median figure of the treatment.

Source: The authors.

Figure 1: Preliminary value tree for metastatic prostate cancer (pre-workshop).



Notes: Contra. = Contraindications; MoA = Mechanism of action; HRQoL = Health related quality of life; PSA = Prostate-specific Antigen; ATC = Anatomical therapeutic chemical; Image produced using the Hiview3 software version 3.2.0.4.

Source: The authors.

Figure 2: Criteria valuation drug profiles.

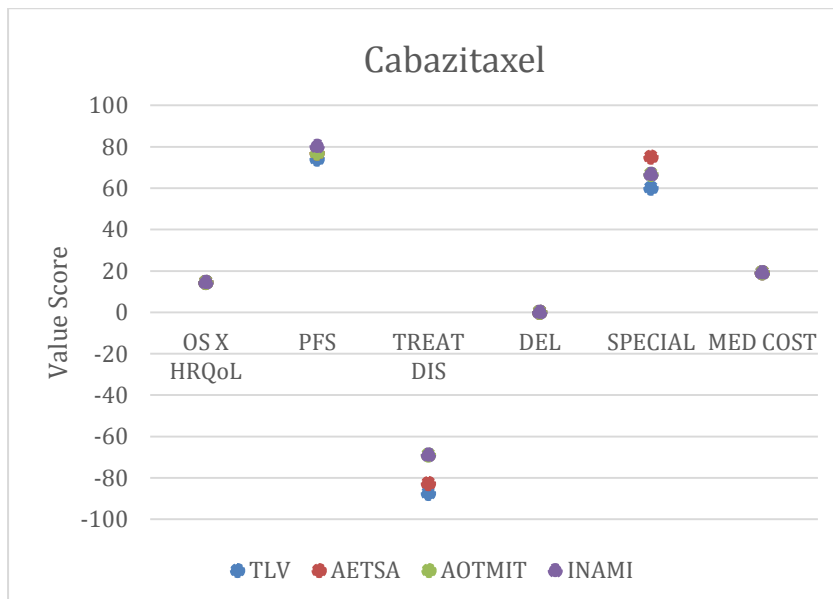
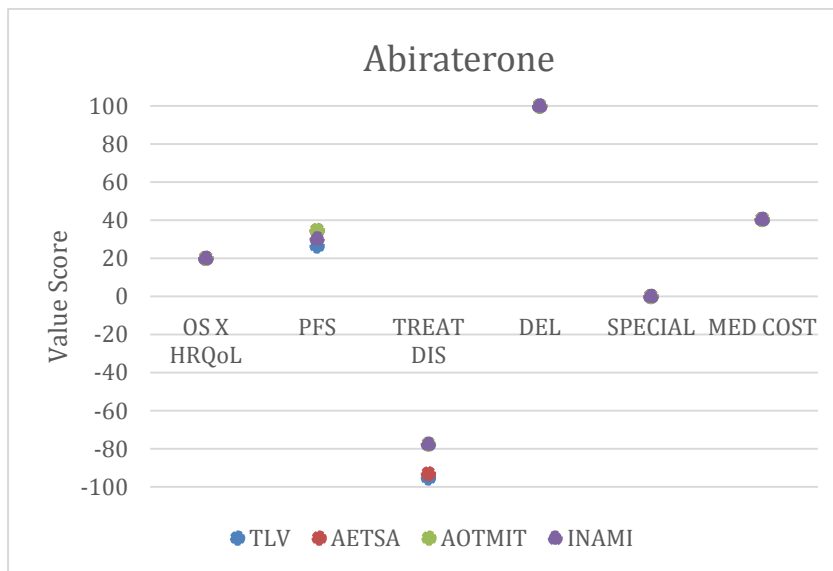
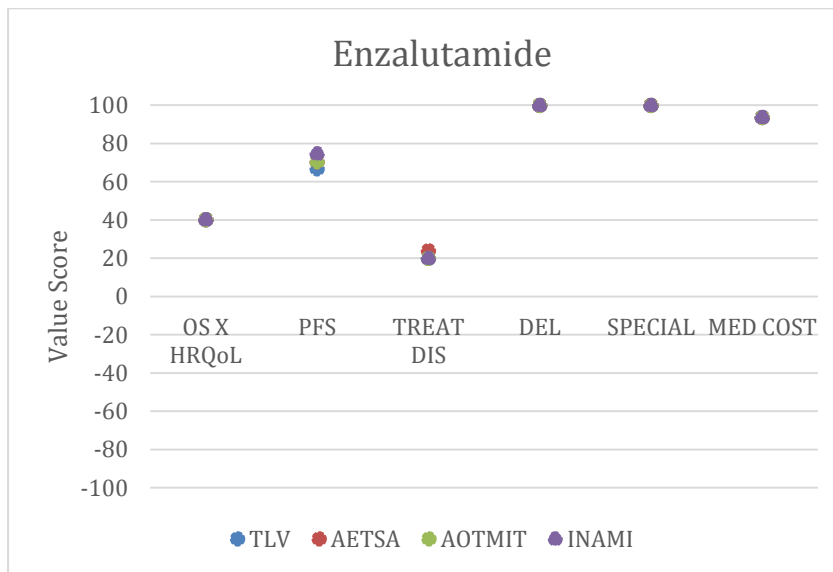


Figure 3: Relative criteria weights stacked bars across the four HTA settings.

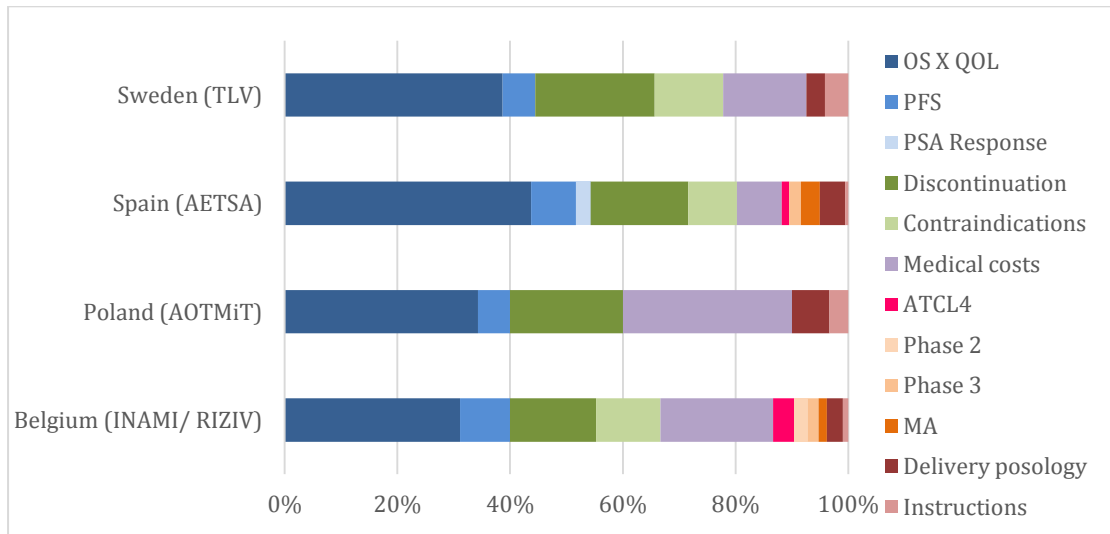
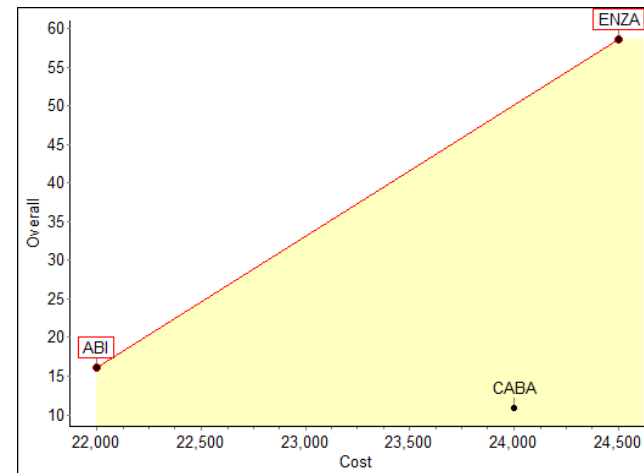
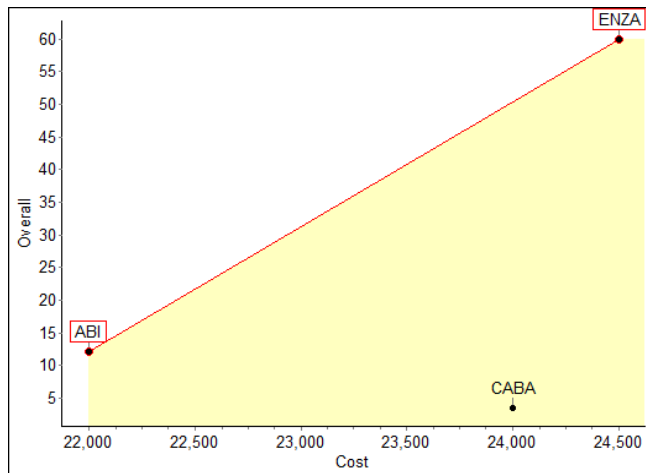
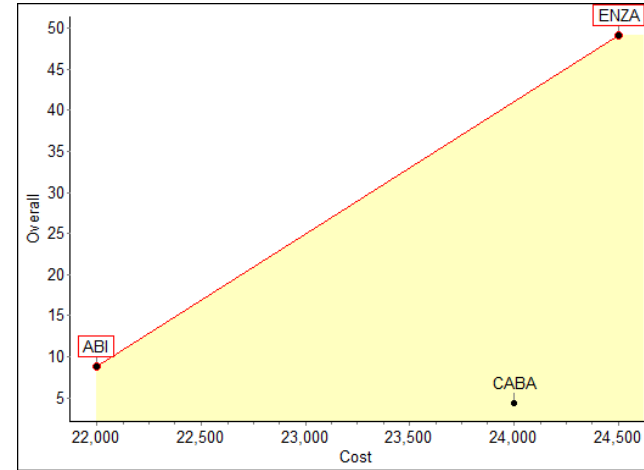
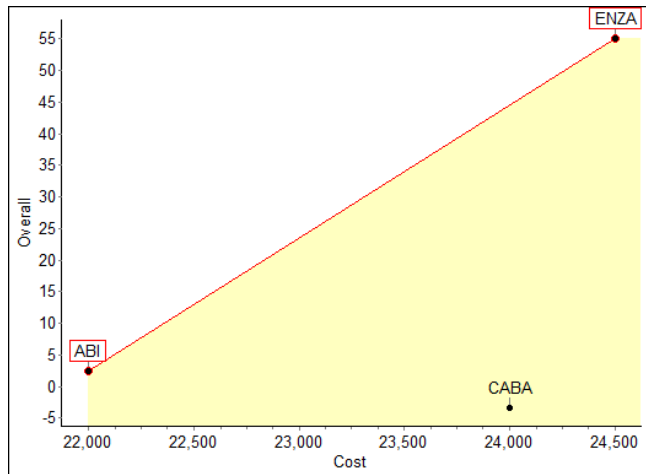


Figure 5: Cost benefit plots of treatments overall weighted preference value scores versus their purchasing costs across the four HTA settings (TLV top left, AETSA top right, AOTMiT, bottom left, INAMI bottom right).



Introduction

In recent years, the introduction of new and costly health technologies, particularly in oncology, combined with moderate health gains, has sparked extensive debate on their value for patients and health care systems, how this value should be assessed and what should be the evaluation criteria informing coverage decisions (Cohen, 2017; Linley & Hughes, 2013). The debate has been fuelled by diverging coverage recommendations across settings for several medicines, often related to diseases associated with high morbidity and mortality (Clement et al., 2009; Faden et al., 2009; Nicod & Kanavos, 2012). Difference in opinion often arises in resource allocation decisions amongst different stakeholders, attributable, at least in part, to current evaluation methodologies not adequately capturing different notions of value (Drummond et al., 2013); this includes, for example, the Quality Adjusted Life Year (QALY), whose use in economic evaluations can at times be regarded as blunt and insufficient, among others, because it may not adequately reflect important value aspects in a variety of disease areas (Nancy Devlin & Lorgelly, 2017; Efthymiadou et al., 2019; Wouters et al., 2015). Given the limited consideration of value in traditional economic evaluations, additional parameters have been included in value assessments; however, this is often done in a non-systematic or ad-hoc manner, which may impact the transparency of decision-making processes (Angelis et al., 2018) and lead to inconsistencies in drug coverage decisions.

A growing body of literature is increasingly debating the use of highly expensive new drugs, which are perceived to bring marginal added clinical benefit on the grounds of poor value-for-money and high budget impact (Nadler et al., 2006; Shih et al., 2013; Sulmasy & Moy, 2014). Rising drug prices and the need to understand the importance of different evaluation criteria have catalysed the generation of numerous “value frameworks” aiming to inform payers, clinicians and patients on the assessment of new medicines, required for making coverage and treatment selection decisions (Anderson et al., 2014; Bach, 2015; Cherny et al., 2015; Schnipper et al., 2015). Although this is an important step towards a more inclusive value-based assessment approach (Malone et al., 2016), aspects of these frameworks

may be based on weak or ad hoc methodologies, which could potentially result in misleading recommendations or decisions (Angelis & Kanavos, 2016a).

In response to some of the concerns raised above, multiple criteria decision analysis (MCDA) has emerged as an alternative to traditional economic evaluation techniques with the prospects of addressing some of their limitations in Health Technology Assessment (HTA) (Angelis et al., 2016; NJ Devlin & Sussex, 2011; Mireille M. Goetghebeur et al., 2008; Kanavos & Angelis, 2013; Marsh et al., 2014; Radaelli et al., 2014; J Sussex et al., 2013b; Thokala, 2011), but also for eliciting stakeholder preferences and facilitating treatment selection (Danner et al., 2011; Ijzerman et al., 2008; Tervonen et al., 2015). A number of MCDA empirical studies have explored the question of value in a number of therapeutic areas, often simulating hypothetical HTA settings (Angelis et al., 2017; M. M. Goetghebeur et al., 2010; Jon Sussex et al., 2013a; Wagner et al., 2017). However, very few studies have explored the same issue by eliciting the preferences of HTA agencies and sitting decision makers and only in single-case exercises (Angelis, 2018; Jaramillo et al., 2016; Tony et al., 2011). To the best of our knowledge, no study has ever compared the value preferences of decision-makers across multiple settings using a full MCDA methodology.

By engaging HTA agencies and health insurance organisations in four EU Member States, we applied the Advance Value Framework (AVF), a recently developed multi-criteria value framework applicable to HTA (Angelis & Kanavos, 2016b; Angelis & Kanavos, 2017), to assess the value of a number of treatment options indicated for metastatic castrate resistant prostate cancer (mCRPC) following first line chemotherapy. This indication was selected because of its high disease burden and the availability of several new and expensive biologic drugs, making it a highly relevant appraisal topic for several HTA agencies.

This is to our knowledge the first cross-country, complete MCDA pilot exercise, eliciting value preferences of sitting decision-makers from different HTA agencies for the same drug treatments while considering identical sets of evidence. The two main research questions of the study relate to testing the feasibility of this MCDA methodology for HTA

decision-makers, and to observing any differences in their value perceptions as reflected through the consistency of drugs' value rankings, including value trade-offs.

Methods

Methodological Framework

An MCDA approach based on Multi-Attribute Value Theory (MAVT) was adopted (Keeney & Raiffa, 1993; von Winterfeldt & Edwards, 1986), involving the phases of problem structuring, model building, model assessment, model appraisal, and development of action plans (Angelis & Kanavos, 2016b). A series of facilitated workshops were organised taking the form of decision conferences (Phillips, 2007), adopting a facilitated decision analysis modelling approach (Franco & Montibeller, 2010b; Phillips & Phillips, 1993), in collaboration with decision-makers from four HTA agencies and health insurance bodies: the Dental and Pharmaceutical Benefits Agency (TLV, Sweden), the Andalusian Health Technology Assessment Agency (AETSA, Spain), the Agency for Health Technology Assessment and Tariff System (AOTMiT, Poland), and the National Health Insurance Agency (INAMI-RIZIV, Belgium). The agencies in these countries were selected in order to represent a set of organisations with different governance structure (arms' length HTA agency, e.g. AOTMiT, TLV and AETSA, vs integrated HTA function, e.g. INAMI-RIZIV) and responsibilities (regulatory, e.g. TLV, vs advisory AOTMiT and AETSA). This research was undertaken in the context of Advance-HTA, an EU-funded project focusing on HTA methodological advancements (London School of Economics, 2019), and all four HTA organisations were contacted to participate under the auspices of the project.

The methodological process used in terms of the design, implementation and analysis, is aligned with the ISPOR good practice guidelines on the use of MCDA for health care decisions (Marsh et al., 2016).

Problem structuring: Clinical Practice and Scope of the Exercise

Prostate cancer is the second most commonly diagnosed cancer in men globally and the most frequently diagnosed cancer among men in developed countries; it is the fifth leading cause of cancer death globally (Torre, 2015). Death rates have been decreasing in the majority of developed countries, which has mainly been attributed to improved treatment and/or early detection (Center et al., 2012).

The decision context relates to the assessment of value of second line treatments for mCRPC based on the approved European Medicines Agency (EMA) indication (EMA, 2016a, b, c), the subsequently defined scope of Technology Appraisals (TAs) by a number of HTA agencies and the ESMO guidelines (Horwich et al., 2013; NICE, 2012a, b, 2014; TLV, 2014, 2015a).

The first treatment to demonstrate a survival benefit for mCRPC patients was docetaxel chemotherapy in combination with prednisolone when compared to mitoxantrone in combination with prednisolone (Berthold et al., 2008; Tannock et al., 2004). Subsequently, new therapeutic agents have been tested in the post-chemotherapy setting with considerable success. Abiraterone, a steroid synthesis inhibitor, in combination with prednisolone showed a 3.9-month improvement in survival compared to prednisolone alone in patients pre-treated with docetaxel (14.8 vs 10.9 months, HR 0.65, $p < 0.001$) (de Bono et al., 2011). Similarly, enzalutamide, an androgen receptor antagonist, showed a 4.8-month improvement in survival (18.4 vs 13.6 months, HR 0.63, $p < 0.001$) compared to placebo alone in the same patient group (Scher et al., 2012). Cross-resistance appears to exist between abiraterone and enzalutamide meaning that patients are unlikely to derive clinical benefit by switching from one to the other agent (Bianchini et al., 2014; Lortol et al., 2013). The third agent that is widely used following progression on docetaxel is cabazitaxel, a taxane chemotherapy. Cabazitaxel led to an overall survival (OS) benefit of 2.4 months (15.1 vs 12.7 months, HR 0.70, $p < 0.0001$) compared to mitoxantrone (de Bono et al., 2010). Given this therapeutic landscape for patients with mCRPC who have progressed on first line docetaxel chemotherapy, characterised by an availability of different treatments and the apparent cross-

resistance between some of them, we adopt post-chemotherapy mCRPC as the decision context for the application of the AVF methodology.

Model Building: Advance Value Tree adaptation, treatments compared and reference levels

The model building phase comprised a number of tasks, notably the Advance Value Tree adaptation for mCRPC, the consideration of alternative drug treatments and the respective evidence, and the definition of criteria attributes and the associated ranges, all of which are discussed below. Detailed discussion on the rationale of each criterion and their value scales can be found elsewhere (Angelis & Kanavos, 2017; Angelis et al., 2017).

(a) Adaptation of the Advance Value Tree for Metastatic Prostate Cancer

At the core of AVF lies the Advance Value Tree, a hierarchical structure of evaluation criteria taking the form of a generic value tree reflecting value concerns of HTA experts and decision-makers for new medicines (Angelis & Kanavos, 2017). The Advance Value Tree consists of five criteria domains, aiming to capture the essential value attributes of new medicines in the HTA context under a prescriptive decision-aid approach. These are divided into (a) Burden of Disease (BoD); (b) Therapeutic Benefit (THE); (c) Safety Profile (SAF); (d) Innovation Level (INN); and (e) Socioeconomic Impact (SOC), summarised by the following value function:

$$Value = f(\mathbf{BoD}, \mathbf{THE}, \mathbf{SAF}, \mathbf{INN}, \mathbf{SOC}) \quad (1)$$

The Advance Value Tree was adapted into a disease-specific mCRPC value model using a bottom-up approach by comparing the characteristics of the specific drugs evaluated (Franco & Montibeller, 2010a). In consultation with a specialist medical oncologist (co-author of the paper), the generic evaluation criteria were converted into disease-specific criteria, while adhering to required criteria properties such as non-redundancy and preferential-independence (Keeney, 1992), to ensure methodological robustness and an adequate value model rooted in

decision theory. Based on the above, a preliminary mCRPC-specific value tree was produced with four criteria domains and a total of 18 criteria, each operationalised by an attribute, i.e. performance indicator, as shown in Figure 1. The BoD domain was not considered in the adaptation process on the grounds of conciseness, as all drugs were indicated for the same indication which would have identical BoD.

Criteria definitions (together with their consideration in each jurisdiction and their rankings) are provided in Table 1. The preliminary version of the mCRPC value tree was subsequently validated by decision conference participants, in line with a “socio-technical” approach, a constructive decision-aid process allowing groups of participants to interact with and learn from each other (Bana e Costa & Beinat, 2005).

<Figure 1 about here>

(b) Alternative Treatments Compared and Evidence Considered

The alternative drug options assessed in the exercise were cabazitaxel in combination with prednisolone, abiraterone in combination with prednisolone and enzalutamide monotherapy. The key evidence sources used to assess their performance included (a) the peer review publications concerning the pivotal clinical trials of the alternative treatment options that were considered for their licencing by the EMA (de Bono et al., 2011; de Bono et al., 2010; Fizazi et al., 2012; Scher et al., 2012); (b) the Product Information sections of EMA’s European Public Assessment Reports (EPAR) (Annex I and III) (EMA, 2016a, b, c); (c) the Anatomical Therapeutic Chemical (ATC) classification system indexes available through the portal of the WHO Collaborating Centre for Drug Statistics Methodology (World Health Organisation Collaborating Centre, 2016); and (d) the US National Library of Medicine clinical trials database (NIH, 2016). Additional sources of evidence included national sources (BNF, 2015; Connock et al., 2011; NICE, 2012a, b, 2014; Riensa et al., 2013) and other peer review literature (Burström et al., 2001; Collins et al., 2007; Kearns et al., 2013; Sullivan et al.,

2007), which was relevant to the study indication. Sources of evidence used relating to the performance of drugs across evaluation criteria are shown in Appendix Table A1, alongside additional information on the evidence considered.

(c) Options Performance and References Levels

By considering the performance of the alternative drug options across the value scales, “lower” (x_l) and “higher” (x_h) reference levels were defined to serve as benchmarks for the value scores of 0 and 100 respectively, acting as value anchors for constructing value functions and eliciting their relative weights (Bana e Costa & Vansnick, 1999; Keeney, 1982). The “lower” reference levels denoted a less preferred state reflecting a “satisfactory” performance level, whereas the “higher” reference levels denoted a more preferred state reflecting an “ideal” performance level.

The reference levels for the clinical attributes informing the Therapeutic and Safety criteria domains, were defined in consultation with the clinical oncologist (co-author of the paper). In principle, the rationale involved adopting the Best Supportive Care (BSC) performance as a “satisfactory” reference level, with a hypothetical 20% improvement of the best available performance acting as the “ideal” reference level (e.g. ‘overall survival’), or, alternatively, the best possible limit of the performance scale acting as an “ideal” level in cases where this was naturally restricted (e.g. ‘treatment discontinuation’). The 20% hypothetical performance improvement was selected because it was perceived to be a realistically plausible scenario for future treatment options. By considering the performance of best available option(s) among the treatments evaluated and accounting for plausible performance improvement in the near future, the value scale essentially reflected characteristics of a “global” scale to account for the performance of future options not captured in the exercise, i.e. *what is best plausible* (Belton & Stewart, 2002). Where a BSC performance was not meaningful to act as a “lower” reference level, then the lowest (i.e. worst) possible limit of the performance scale was adopted (e.g. ‘Phase 3’), or, alternatively,

20% lower than the lowest performing option was used (e.g. ‘medical costs impact’). An exception to the above was the ‘health related quality of life’ (HRQoL) attribute for which the stable disease state’s utility score was adopted as the “lower” level and the general population utility score was used as the “higher” level.

The emerging partial value function scores of the drugs for each criterion can take negative values or values higher than 100 where $v(x_{\text{lower}}) = 0$ and $v(x_{\text{higher}}) = 100$, essentially by conducting a positive linear transformation. “Lower” and “higher” reference levels for all attributes at the pre-decision conference stage and the basis of their selection are outlined in Appendix Table A2. A matrix listing the performance of drug options across the final attributes that were considered in the decision conferences, together with their reference levels, is shown in Table 2.

Model Assessment and Appraisal: Decision conferences, MCDA technique and cost calculation

The model assessment and appraisal phases comprised the tasks of conducting the decision conferences, the application of the MCDA technique for the elicitation of value preferences and cost calculation(s). These are discussed below.

(a) Decision conferences

Model assessment and model appraisal took place through a series of decision conferences (Phillips, 2007), taking the form of facilitated workshops with the participation of decision-makers, including assessors and national experts, all of whom were affiliated with the four study HTA organisations, either as members of staff or visiting external experts (their difference being in full-time employment versus part-time or visiting capacity employment). For the purposes of this study, they were both regarded as “decision-makers”, given their influence on methodological development within the agencies and on the decision outcomes of the appraisals. Across the four countries, between four (for the case of TLV) and 13 (for

the case of AOTMiT) participants were involved, typically comprising health care professionals (clinicians, pharmacists), HTA methodology experts (health economists, statisticians, HTA agency directors) and decision-makers (members of HTA appraisal committees, representatives from insurance funds and the national medicines agencies). Background material introducing the scope of the exercise in more detail was sent to the participants one week before each decision conference. Decision conferences were hosted at the head offices of the different HTA organisations between June 2015 and April 2016: Stockholm (TLV), Seville (AETSA), Warsaw (AOTMiT), and Brussels (INAMI-RIZIV).

The lead author acted as an impartial facilitator, assisted the groups' interactions and guided participants through the decision problem using the preliminary version of the mCRPC-specific value tree (Figure 1) and the relevant data. This acted as the model's starting point, based on which value judgements and preferences were elicited at the start of each decision conference while seeking group interaction and agreement (Franco & Montibeller, 2010b; Phillips, 1984; Phillips & Bana e Costa, 2007; Schein, 1999). The *Appendix* provides more information on the decision conferences.

(b) MCDA Technique

AVF adopts a value measurement MCDA methodology making use of a simple additive (i.e. linear, weighted average) value model for the aggregation of scores and weights (Angelis & Kanavos, 2017). This assumes preference independence between the different criteria, with overall value $V(\cdot)$ of an option a defined by the equation below (Keeney, 1992; von Winterfeldt & Edwards, 1986):

$$V(a) = \sum_{i=1}^m w_i v_i(a) \quad (2)$$

Where m is the number of evaluation criteria, $w_i v_i(a)$ is the weighted partial value function of evaluation criterion i for treatment a , and $V(a)$ is the overall value of a treatment a . $V(\cdot)$ is

therefore is an overall value function based on multi-attribute value theory (Keeney & Raiffa, 1993).

A value function associated with each attribute, converting the treatment performance on the attribute range to a value scale, was elicited from the participants during the decision conferences using the Measuring Attractiveness by a Categorical Based Evaluation Technique (MACBETH) questioning protocol and the M-MACBETH software (Bana e Costa & Vansnick, 1999). This protocol requires pairwise comparisons where qualitative judgements about the difference of value between different pairs of attribute levels (i.e. difference in value between x and y units on a criterion) are expressed using seven qualitative categories (i.e. no difference, very weak difference, weak difference, moderate difference, strong difference, very strong difference, or extreme difference) (Bana E Costa et al., 2012; Bana e Costa & Vansnick, 1994). MACBETH provides a constructive and user-friendly approach to generate a cardinal (interval) value scale based on the input of these qualitative pair-wise judgements, which are then converted into value scores via an optimization algorithm (Bana e Costa et al., 2016b); this approach has been widely used as a decision support tool (Bana e Costa et al., 2014; Bana e Costa et al., 2002; Bana e Costa & Oliveira, 2012; Bana e Costa & Vansnick, 1997).

Weights for a multi-attribute value function should be elicited considering the range of each attribute and the value of a “swing” between two reference levels. The weights are scaling constants that convert partial value scores into overall value scores that must reflect value trade-offs and, therefore, should not be interpreted as measurements of ‘direct importance’. An indirect (qualitative) swing weighting technique was applied to elicit relative criteria weights by first ordering the swings of each attribute and then valuing their differences using the MACBETH qualitative categories (Bana E Costa et al., 2012).

The above MACBETH-based scoring and weighting techniques were operationalised using the software M-MACBETH, (Bana e Costa & Vansnick, 1999). The software automates the additive aggregation of preference value scores and weights in order to derive overall weighted preference value (WPV) scores and also allows for sensitivity analysis on the

criteria weights. The software also enables the use of visual graphics to build a model of values, acting as a facilitation tool to inform both the design and the evaluation phases of the methodological framework (Bana e Costa et al., 2016a; Bana e Costa & Vansnick, 1999; Bana e Costa et al., 1999). More information regarding the technical details of MACBETH is available in the *Appendix*.

(c) Cost Calculation

UK list prices at ex-factory level were used as found in BNF (BNF, 2015) as a neutral benchmark in order to allow the measurement of cost(s) in a common unit across all study settings, so that overall WPV scores can then be viewed against the same cost denominator to produce comparable cost-value ratios. Access to confidential prices through risk sharing agreements was not possible. Information on the recommended dosages and treatment durations were sourced from the peer review publications of the pivotal trials and respective EPARs from EMA (de Bono et al., 2011; de Bono et al., 2010; EMA, 2016a, b, c; Scher et al., 2012). Drug administration costs for cabazitaxel were kept consistent with the respective NICE TA (NICE, 2012b), whereas for abiraterone and enzalutamide these costs were not applicable as they are orally administered.

Results

Final Value Trees, Options Performance, Criteria Weights and Value Functions

Across the four countries, decision conferences were characterised by increased interaction and extensive debate between participants, especially in cases where there was disagreement about certain values. Because the majority of participants had a shared understanding of the decision problem but also a sense of common purpose and commitment to way forward, all of which are conditions for good practice in decision conferencing, the deliberative process of each decision conference instigated a fruitful discussion and exchange of views around different criteria values and relative importance.

General consensus was reached among participants in terms of criteria consideration and model validation with no major value aspects deemed to be missing. All attributes included in each country's final mCRPC value tree, as emerged following open interaction with decision conference participants and their rankings, are shown in Table 1 (schematic illustrations of the individual value trees are shown in Appendix Figure A1). The main reason for not including a criterion attribute in the value tree was because participants considered it was non-fundamental for the evaluation, in all cases of which a zero weight was assigned. Most of the criteria attributes that were assigned a zero weight belonged in the Innovation Level domain, which comprised the highest number of criteria.

<Table 1 about here>

The performance of the drug options across the different attributes that were considered to be fundamental in the model (i.e. weight greater than zero) together with the “lower” and “higher” reference levels are shown in Table 2.

<Table 2 about here>

Between 6 (AOTMiT) and 11 (AETSA/INAMI) criteria attributes were included in the final value tree of each country, as shown in Table 3. In terms of the different criteria domains composition, the Therapeutic Benefit contained between two (TLV/AOTMiT/INAMI) and three (AETSA) criteria attributes, the Safety Profile between one (AOTMiT) and two (TLV/AETSA/INAMI), the Innovation Level between two (TLV/AOTMiT) and six (INAMI), and the Socioeconomic Impact always one.

<Table 3 about here>

During the elicitation of the ‘overall survival’ (OS) and/or ‘HRQoL’ criteria value functions, it became evident that these criteria attributes might be preference dependent. When asking participants to judge the difference in value between different increments in attribute performance (either in ‘OS’ or ‘HRQoL’), a request for clarification was raised by some of them relating to what level of performance this change was associated with on the other criterion attribute. In order to address the plausible preference-dependence observed, we combined together the two attributes in an aggregated form. The two criteria attributes were combined by multiplying the number of months in ‘OS’ and their EQ-5D utility scores in ‘HRQoL’ attributes respectively, assuming an equal (i.e. 50%) distribution of stable and progressive disease states, essentially deriving quality adjusted life months (QALMs). An example of a MACBETH value judgements matrix and its conversion into a value function for the case of the ‘OS x HRQoL’ aggregated criterion attribute in QALMs is shown in Appendix Figure A2.

There was a common set of six criteria that were considered as fundamental in all countries: (a) ‘OS x HRQoL’; (b) ‘radiographic tumour progression’ (also known as progression free survival (PFS)); (c) ‘treatment discontinuation’; (d) ‘delivery posology’; (e) ‘special instructions’; and (f) ‘medical costs impact’. This common set of criteria comprised the complete set of TLV’s value tree (n=6), whereas AOTMIT’s value tree considered ‘contraindications’ in addition (n=7). Further to these, AETSA’s value tree also considered ‘PSA response’, ‘ATCL4’, ‘Phase 3’ and ‘marketing authorisation’ (n=11), whereas INAMI’s value tree considered the same additional criteria but with ‘Phase 2’ instead of ‘PSA response’ (n=11).

Overall, the different groups of decision conferences’ participants agreed in the valuation of performance for the six common attributes that were considered across all four countries, as revealed through the elicitation of their value functions. Figure 2 plots the value scores of each drug across the six common attributes showing very similar valuations between countries.

<Figure 2 about here>

The weights of relative importance assigned to the different attributes across the four jurisdictions are shown in Figure 3. By taking into account the relative swings of the criteria attributes, i.e. the gap between the “lower” and “higher” reference levels, quantitative weights were derived for each attribute using M-MACBETH. The ‘OS x HRQoL’ aggregated criterion attribute was always assigned the highest relative weight out of 100 ([31,44] for INAMI and AETSA, respectively), followed either by ‘treatment discontinuation’ ([17,21] for AETSA and TLV, respectively) or ‘medical costs impact’ ([20,30] for INAMI and AOTMiT, respectively). Depending on the country, the third-ranked criterion was then either ‘treatment discontinuation’ (AOTMiT, INAMI), ‘medical costs impact’ (TLV), or ‘contraindications’ (AETSA) and ‘PFS’ was ranked 4th or 5th. ‘Special instructions’, although a fundamental criterion across settings, was ranked in the lowest place in 3 out of 4 settings with the ‘delivery posology’ usually at a higher position, with the exception of TLV where that order was reversed.

<Figure 3 about here>

In terms of the total weights assigned across the different criteria domains, the Therapeutic Benefit weight ranged from 40% to 54% (for AOTMiT/ INAMI and AETSA, respectively), the Safety Profile weight ranged from 20% to 33% (for AOTMiT and TLV, respectively), the Innovation Level weight ranged from 7% to 13% (for TLV and INAMI, respectively) and the Socioeconomic Impact weight ranged from 8% to 30% (for AETSA and AOTMiT, respectively) (Table 3). The above differences in relative weights reflect the different priorities of decision-makers, including the number of fundamental objectives being considered.

Overall Drug Rankings and Value-for-Money Analysis

With regards to the overall WPV scores shown in Table 4, enzalutamide consistently yielded the highest score across all four countries, always followed by abiraterone and cabazitaxel. The overall scores of abiraterone and cabazitaxel were in part influenced by a “negative” performance in the ‘treatment discontinuation’ attribute (19% and 18% respectively) which lay below the lower reference level of the scale (i.e. 10%), affecting negatively their overall value scores.

A stacked bar plot of the drugs’ overall WPV scores across all settings is shown in Figure 4. By using rounded up cost figures for enzalutamide (£24,600), abiraterone (£21,900) and cabazitaxel (£23,900, of which £22,190 related to drug cost and the remainder £1,710 to administration cost) and dividing them with overall WPV scores, their costs per MCDA value unit ranged as follows: (a) enzalutamide: £410 - £501 (for AOTMiT and AETSA, respectively); (b) abiraterone: £1,366 - £9,221 (for INAMI and TLV, respectively); and (c) cabazitaxel: £2,196 - £6,816 (for INAMI and AOTMiT, respectively) (Table 4). The overall value score of each option was driven by the fundamental objectives considered (i.e. criteria influencing the model), the criteria weights which were anchored on reference levels, and the shape of value functions which would influence the value scores.

<Table 4 about here>

<Figure 4 about here>

In terms of value-for-money, cabazitaxel was shown to be dominated by abiraterone, and was very close to being dominated by enzalutamide (i.e. a difference of £500 based on the prices used). Enzalutamide on the other hand was associated with a higher cost (a difference of £2,500 based on the prices used) and a higher overall WPV score compared to abiraterone, with a difference in score ranging between 40.4 to 52.7 value units (for AETSA and TLV, respectively). Cost benefit plots of the different options, using their overall WPV scores

versus their purchasing (plus any administration) costs across the four HTA organisations is shown in Figure 5.

<Figure 5 about here>

Similarities and differences in value perceptions across settings

By looking at Table 3 (and Figure 3) of the results, a number of similarities and differences in value preferences are observed across the four settings. The largest number of evaluation criteria were considered in Andalusia and Belgium (11 each), compared to Sweden and Poland (7 and 6, respectively), partly due to a higher number of Innovation Level criteria (5 and 6, compared to 2 each, respectively). In terms of the relative importance of criteria domains, the Therapeutic Benefit cluster consistently ranked first across all settings. The Safety Profile cluster was ranked second in three settings (except for Poland, where the Socioeconomic Impact cluster ranked higher (30% vs 20%)). The Socioeconomic Impact cluster ranked 3rd in Sweden and Belgium but 4th in Andalusia (8%). Finally, the Innovation Level cluster ranked 4th in three countries with the exception of Andalusia where it ranked 3rd (12%). The low relative importance of the Innovation Level cluster partly justifies why a hypothetical change in the final consideration of Innovation Level criteria across the different countries does not influence the ranking of the treatments, as described in the next section.

Despite the observed differences in evaluation criteria considered, the relative criteria weights assigned and the elicited value functions, the overall ranking of the treatments remained identical across countries (Table 4 and Figure 4) with enzalutamide consistently having the highest score, followed by abiraterone and cabazitaxel in all four settings.

Sensitivity and Robustness Analysis

Following each decision conference, deterministic sensitivity analysis was conducted to address parameter uncertainty on criteria weights. Specifically, changes on baseline weights

were explored to check their possible impact on treatments' overall value rankings. The results of the sensitivity analysis demonstrated that the ranking of the treatments was robust to the relative criteria weights across the different settings.

The most sensitive criterion weight, which could change enzalutamide's ranking order from first to second, was 'PFS' in the cases of INAMI and AETSA where a 10.2 and 11.1 times change (from 8.9% to 90.6% and from 8.0% to 88.5%) respectively, would be required for cabazitaxel to rank first and enzalutamide second. In other words, a higher than 10-times difference on the 'PFS' weight would be required for cabazitaxel to outperform enzalutamide, with changes of higher order required in other criteria weights for either cabazitaxel or abiraterone to rank first, in any of the study settings. Criteria weights were more sensitive with regards to the outperformance of abiraterone by cabazitaxel as the second-best treatment. Again, the most sensitive weight was for 'PFS' in the INAMI and AETSA cases, where a 2-times change (from 8.9% to 17.4% and from 8.0% to 16.7% respectively) would be needed for cabazitaxel to rank second and abiraterone third. This meant that the lowest change across criteria weights needed for an impact on treatment rankings to be observed was for the case of PFS with INAMI, where at least a 2-time difference was required for abiraterone to be outperformed. For the case of TLV and AOTMiT, the most sensitive criterion was treatment discontinuation in which a 2.6 and 3.0 times change would be needed (from 21.2% to 54.6% and from 20% to 60% respectively) for cabazitaxel to rank second-best.

The final consideration of the Innovation Level criteria cluster was explored in greater detail given that their relevance might be disputed. Removing the 'ATCL4' criterion and any spill-over effect criteria (i.e. 'Phase-2', 'Phase-3', 'MA') from the value tree of AETSA and INAMI, and any patient convenience criteria (i.e. 'delivery posology', 'special instructions') from all country value trees would not affect the treatment rankings.

Discussion and policy implications

This study is the first comparative MCDA exercise, utilising the Advance Value Framework and engaging sitting HTA decision-makers across four EU Member States to elicit and compare their preferences in the evaluation of three mCRPC treatments. In doing so, the objective was to test the feasibility of MCDA methods for HTA decision-makers and identify differences in value perceptions.

Based on the evidence used, our results showed that the most valuable therapy for second line mCRPC was enzalutamide, followed by abiraterone and cabazitaxel. Each treatment was assessed and ranked based on their overall WPV scores, reflecting the value of their performance against a set of evaluation criteria, weighted against their relative importance. These overall scores were based on the value preferences of decision-makers that were collected via a decision conference in each setting, yielding a comprehensive and transparent, multi-dimensional benefit component. Subsequent consideration of drug costs (purchasing and administration) enabled the estimation of value-for-money in the form of “cost-per-unit of value” ratios which showed the second-ranked treatment (abiraterone) to dominate the third (cabazitaxel).

It should be noted that the constructed benefit metric excludes the cost of the treatments, i.e. the WPV score considers the impact of the technology on medical costs other than the purchasing cost of the technology. Therefore, evaluation of the treatments based solely on their overall WPV scores might not be appropriately designed to inform an HTA decision context that considers the interventions’ incremental cost per incremental benefit, but, rather, a value-based approach to reimbursement or pricing negotiation.

Attempting a comparison of the ranking achieved in this exercise with what has taken place in reality might prove challenging, partly because of how the clinical evidence was treated in the exercise, but also because it is not publicly known whether and how any of the additional value dimensions evaluated in the exercise were considered in the relevant HTA decision-making processes. In Sweden, although abiraterone’s ICER vs BSC (manufacturer

estimate of SEK820,000/QALY)(TLV, 2015a), was lower compared to enzalutamide's ICER vs BSC (TLV best estimate of SEK1,100,000/QALY)(TLV, 2014), or lower vs enzalutamide (SEK800,000/QALY)(TLV, 2015b), TLV assumed that both treatments had the same clinical effect and consequently focused on a cost-minimisation approach rather than cost-utility analysis, leading to the implementation of a confidential risk sharing agreement (RSA) as part of which discounts can be provided based on treatment duration. A similar conclusion was reached in Spain, where the Ministry of Health in its Clinical Assessment Report (Informe de Posicionamiento Terapeutico - IPT) recommended that there is no clinically relevant difference between the benefit-risk balance of enzalutamide and abiraterone, and, therefore, decisions should be guided based on drug costs (AEMPS, 2015). Pricing and reimbursement decisions are then taken by the Interministerial Committee for Pricing and Reimbursement, but the final assessment is not publicly available. At regional/hospital level, a group of hospital pharmacists conducted a full health (clinical and economic) technology assessment, where enzalutamide and abiraterone were considered to be therapeutically equivalent (GHEMA, 2016). In Poland, although AOTMiT accepted that some additional clinical effect existed for enzalutamide compared to abiraterone (mainly in secondary endpoints), it was not found to be cost-effective compared to abiraterone; however, a confidential RSA enabled a final positive recommendation by AOTMiT (AOTMiT, 2017). The final decision implemented by the Ministry of Health was to reimburse enzalutamide, similarly to the case of abiraterone (Obwieszczenie, 2017). In Belgium, following an indirect comparison no clinically relevant differences were found in the treatment outcomes of abiraterone versus enzalutamide (INAMI, 2019); eventually, a managed entry agreement (MEA) enabled reimbursement.

Consequently, and based on the evidence used to populate the MCDA model and which would inform decision-making, the hypothetical coverage decisions emerging from the ranking of the treatments based on their overall WPV scores might have been different. Given the higher overall value of enzalutamide compared to abiraterone, a cost minimisation

approach or price parity attained between the two, as inferred following the risk sharing agreements in place, might not have been justified.

One reason why our value models make slightly different predictions is because it has captured benefits that go beyond the current formal remit of HTA agencies, therefore the results should be viewed as ‘proof-of-concept’, for the purposes of testing the performance of the methodology. Furthermore, the decision context addressed in the exercise was a one-off evaluation problem within the indication of mCRPC which might contradict the operational scope of some HTA agencies and health insurance bodies relating to repeated decisions around the reimbursement of drugs across different disease areas.

The extent to which HTA decision-makers can be relied upon, or not, to reflect societal preferences when constructing their value preferences is a very important topic for discussion but not aimed to be addressed in this study. Here, we simply elicited decision-makers’ own preferences without considering whether these might be representative for society or not. In reality, evidence in Belgium suggests that health care coverage related preferences of decision-makers differ to those of the public (Cleemput et al., 2018), and therefore more research would be needed to reveal such discrepancies.

Overall, the HTA decision-makers that participated in the decision conferences provided positive feedback about the potential usefulness of the value framework and the MCDA approach in general, raising the prospects of the framework acting as a decision support tool in the evaluation of new medicines. According to participants, key advantages of the framework included the feasibility to transparently assess the performance of the options across a number of explicit evaluation criteria, while allowing the elicitation of value trade-offs (i.e. their relative importance), and its overall facilitative nature in the construction and analysis of group value preferences. Our results are in line with past evidence on a different oncology indication (Angelis et al., 2017).

Challenges of MCDA applications in HTA

The assessment across 4 settings has offered a number of important insights relating to the application of MCDA in HTA and the challenges this represents. In order for any MCDA methodology to become a useful tool for HTA decision-makers and serve their needs, certain requirements must be met: first, sound methods should be used to ensure technical requirements are fulfilled (Keeney & Raiffa, 1993); second, social aspects of the process should be treated carefully to ensure various socio-technical requirements are fulfilled (Baltussen et al., 2017); and, third, tools and guidelines should be available and tailored for the appropriate audience ensuring that best practice requirements are fulfilled (Phillips, 2017).

Among the first group of technical requirements, one key challenge of MCDA studies in HTA relates to the theoretical properties that are required for the evaluation criteria. Due to the popularity of using a simple additive (i.e. weighted average) value model, the violation of preference-independence is of particular relevance as it might undermine the validity of such models and the insights offered by the results (Marsh et al., 2018; Morton, 2017). Evidence suggests that preference dependencies might exist between health gain and disease severity (Nord et al., 2009), or between OS and HRQoL (Angelis & Kanavos, 2017). The latter also featured strongly in this study, where such a preference dependence between OS and HRQoL was detected during the decision conferences and, as a result, the two criteria attributes were combined into a common aggregated attribute. Beyond combining the two criteria into a common aggregated attribute, other more technically complex solutions exist for addressing preference dependencies, such as using other functional forms of aggregation for combining scores and weights together, such as multiplicative models (Chongtrakul et al., 2005). Furthermore, tests for identifying preference dependencies have existed for many years (Currim & Sarin, 1984; Keeney, 1992; Rodrigues et al., 2017).

Other technical challenges relate to the need for evaluation criteria to be non-overlapping so that there can be no double counting, and that criteria weights are connected to the attribute ranges. If either one of these conditions is not satisfied, criteria weights could misrepresent decision makers' true value preferences. Furthermore, a number of cognitive

biases may affect value judgments and thus appropriate elicitation protocols and de-biasing tools must be employed (Montibeller & Winterfeldt, 2015).

In order to avoid double-counting, a clear justification of their inclusion is needed, which should be on the grounds of addressing the fundamental objectives of the analysis, rather than be informed based on the existence of available evidence and data (Keeney, 1992; Keeney & Gregory, 2005). This process could be supported by the use of problem structuring tools aiming to distinguish between ‘fundamental objectives’ and ‘means objectives’ (Franco & Montibeller, 2010a), as we adopted in this exercise.

In terms of weighting, asking direct questions for the general importance of criteria are known to be one of the most common mistakes when eliciting value trade-offs (Keeney, 1992; Keeney, 2002). Instead, sound weighting procedures for the assignment of relative weights should take place in accordance with the use of explicit lower and higher reference levels (Belton & Stewart, 2002; Keeney, 2002), ideally through user-friendly indirect technique protocols that can reduce bias, similar to what we aimed for in this exercise through the explicit definition of reference levels and the implementation of the qualitative (MACBETH) swing weighting technique.

A further challenge relates to the linking of MCDA results with coverage and resource allocation decisions, possibly through the use of specific value thresholds, that can reflect the efficiency and opportunity cost of funding decisions (Sculpher et al., 2017). In economic evaluation, incremental cost effectiveness ratio (ICER) thresholds are supposed to reflect the opportunity cost of the benefit foregone elsewhere in the health care system that would have resulted from the coverage of alternative technologies (Claxton et al., 2015). Assuming that a QALY-based ICER threshold is accurate, it could be used as a benchmark to create an MCDA value threshold by extrapolating the ICER threshold in proportion to how much of the MCDA model’s weight is accounted for by non-QALY value components (Phelps & Madhavan, 2018). Alternatively, following the generation of a multi-dimensional benefit component, purchasing costs could be used to derive treatments’ cost-value ratios to

inform the resource allocation decisions assuming a fixed budget (Peacock et al., 2007), similar to our approach in this exercise with the calculation of the “cost per unit of value”.

Study limitations

The study has a number of limitations, both related to the clinical evidence used and the MCDA process followed, so results should be interpreted with caution. First, in terms of the clinical data used, there was a lack of relative treatment effects; in order to counteract that, absolute treatment effects from different clinical trials were used based on the assumption that they are directly comparable which might not be accurate even for similar patient populations in the studies. As a result, differences in the performance of the options that have been valued might in reality not be statistically significant, e.g. in OS. Ideally, one would need indirect comparisons or a network meta-analysis (NMA) through a mixed treatment comparison (Jansen et al., 2011), therefore, an evidence synthesis step would be required as part of the model-building phase; as, for example, in the case of assessing the comparative benefit-risk of statins in primary prevention (Tervonen et al., 2015) or second-generation antidepressants (van Valkenhoef et al., 2012).

Second, another clinical evidence related limitation could be that only the treatments’ impact on HRQoL of the stable disease state was assessed, because no treatment was assumed to have any effect during progression (NICE, 2014). This might not be true for other disease indications in which case the relevant HRQoL attribute would have to capture both the stable and progressive disease states.

Third, there are also a number of limitations in terms of the MCDA process adopted: one of them relates to the relatively small number of participants in some decision conferences, which could reflect a limited representation of perspectives for the purpose of informing policy-making. A group size of between seven and 15 participants is known to be ideal as they are large enough to represent all major perspectives but small enough to work towards agreement, effectively allowing for efficient group processes to emerge while preserving individuality, (Phillips & Phillips, 1993). However, capturing an all-round set of

preferences was not among the primary aims of the exercise. The value scale of the treatment discontinuation attribute and, more specifically, the “lower” reference level of “10%” could be perceived as a limitation because it influenced the negative partial value scores of two treatments whose performance was worse. This was the outcome of consultation with an oncologist, based on evidence from one of the clinical trials’ placebo-controlled arms, because it was believed to better resemble BSC used in practice; although others might have chosen a different performance level to define the “lower” reference level, the overall ranking of the treatments did not change when altering the lowest reference level to a much less preferred hypothetical performance (20% lower than the worst performing option), while keeping the weights constant.

One major advantage in MCDA, is that it can be tailor-made to reflect decision-makers’ needs, by taking into account different fundamental objectives through the consideration of a variety of criteria, reflecting their priorities (by eliciting relative weights) and representing their preferences (by eliciting value functions). However, it should be recognised that the emerging differences that have been described above, prevent the direct comparison of overall value scores for alternative options; these would require identical value trees (i.e. the same set of criteria, weights and value functions across settings), in addition to the same evidence on options performance. The ranking comparisons that we have made in this study using ordinal scales reflect these limitations.

Conclusions and implications

In this study, we tested the application of AVF, a multi-criteria value framework, in collaboration with HTA decision-makers in order to deduce its feasibility and compare results across settings, in an effort to investigate its potential usefulness and limitations for the purposes of HTA. We found that the AVF methodology can act as a valuable decision support tool because of the transparent construction of value preferences in a collaborative manner, which facilitates the evaluation processes of groups, including the elicitation of value

preferences and trade-offs. Although we observed setting-specific differences in value perceptions, the rankings of drugs remained consistent across all countries. Based on the evidence used in the exercise, a coverage decision using this method would have pointed towards a different recommendation denoting differences in value between the first two treatments, in contrast with the cost minimisation approach adopted or the price parity attained between the two in real life.

Despite a number of limitations relating to data and process issues and the existence of broader challenges with the use of MCDA in HTA due to specific methodological requirements which would need to be satisfied, the present study has demonstrated that an MCDA framework can, in fact, provide meaningful valuations of novel health technologies which, in turn, can inform coverage decisions.

The MCDA methodology adopted enabled participants in the study countries to reflect on certain value dimensions and incorporate these more explicitly in the deliberation process, supporting its use as a transparent value communication tool. Future research efforts could involve similar cross-country case studies, the advancement of MCDA methods and their alignment with HTA policy needs, or repeating the study with different participants to understand whether similarities and differences identified in this study can be replicated.

References

- AEMPS. (2015). INFORME DE POSICIONAMIENTO TERAPÉUTICO PT-ENZALUTAMIDA/V1/30072015. *Agencia Española de Medicamentos y Productos Sanitarios*.
- Anderson, J.L., Heidenreich, P.A., Barnett, P.G., Creager, M.A., Fonarow, G.C., Gibbons, R.J., et al. (2014). ACC/AHA statement on cost/value methodology in clinical practice guidelines and performance measures: a report of the American College of Cardiology/American Heart Association Task Force on Performance Measures and Task Force on Practice Guidelines. *Journal Of The American College Of Cardiology*, 63, 2304-2322.
- Angelis, A. (2018). Evaluating the benefits of new drugs in health technology assessment using multiple criteria decision analysis: a case study on metastatic prostate cancer with the dental and pharmaceuticals benefits agency (TLV) in Sweden. *MDM Policy & Practice*, 3.
- Angelis, A., & Kanavos, P. (2016a). Critique of the American Society of Clinical Oncology Value Assessment Framework for Cancer Treatments: Putting Methodologic Robustness First. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*.
- Angelis, A., & Kanavos, P. (2016b). Value-Based Assessment of New Medical Technologies: Towards a Robust Methodological Framework for the Application of Multiple Criteria Decision Analysis in the Context of Health Technology Assessment. *Pharmacoeconomics*, 34, 435-446.
- Angelis, A., & Kanavos, P. (2017). Multiple Criteria Decision Analysis (MCDA) for evaluating new medicines in Health Technology Assessment and beyond: the Advance Value Framework. *Social Science & Medicine*, 188, 137-156.
- Angelis, A., Kanavos, P., & Montibeller, G. (2016). Resource allocation and priority setting in health care: a multi-criteria decision analysis problem of value? *Global Policy*.
- Angelis, A., Lange, A., & Kanavos, P. (2018). Using health technology assessment to assess the value of new medicines: results of a systematic review and expert consultation across eight European countries. *European Journal of Health Economics*.
- Angelis, A., Montibeller, G., Hochhauser, D., & Kanavos, P. (2017). Multiple criteria decision analysis in the context of health technology assessment: a simulation exercise on metastatic colorectal cancer with multiple stakeholders in the English setting. *BMC Medical Informatics and Decision Making*, 17.
- AOTMiT. (2017). Rekomendacja nr 19/2017 z dnia 27 marca 2017 r. Prezesa Agencji Oceny Technologii Medycznych i Taryfikacji w sprawie objęcia refundacją produktu leczniczego Xtandi, enzalutamidum, kapsułki miękkie, 40 mg, 112 kaps., w ramach programu lekowego „Leczenie opornego na kastrację raka gruczołu krokowego z przerzutami (ICD-10 C-61).
http://bipold.aotm.gov.pl/assets/files/zlecenia_mz/2017/008/REK/RP_19_2017_Xtandi.pdf.
- Bach, P. (2015). DrugAbacus App. Memorial Sloan Kettering Cancer Center.
- Baltussen, R., Jansen, M.P.M., Bijlmakers, L., Grutters, J., Kluytmans, A., Reuzel, R.P., et al. (2017). Value Assessment Frameworks for HTA Agencies: The

- Organization of Evidence-Informed Deliberative Processes. *Value in Health*, 20, 256-260.
- Bana e Costa, C., & Beinat, E. (2005). Model-structuring in public decision-aiding. Operational Research working papers. London: London School of Economics and Political Science.
- Bana e Costa, C., De Corte, J., & Vansnick, J. (2016a). M-MACBETH website.
- Bana e Costa, C., De Corte, J., & Vansnick, J. (2016b). On the Mathematical Foundations of MACBETH. In S. Greco, M. Ehrgott, & J. Figueira (Eds.), *Multiple Criteria Decision Analysis: State of the Art Surveys*: Springer New York.
- Bana e Costa, C., Lourenço, J., Oliveira, M., & Bana e Costa, J. (2014). A Socio-technical Approach for Group Decision Support in Public Strategic Planning: The Pernambuco PPA Case. *Group Decision & Negotiation*, 23, 5-29.
- Bana e Costa, C., & Vansnick, J. (1999). The MACBETH Approach: Basic Ideas, Software, and an Application. In N. Meskens, & M. Roubens (Eds.), *Advances in Decision Analysis* pp. 131-157): Springer Netherlands.
- Bana e Costa, C.A., Corrêa, É.C., De Corte, J.-M., & Vansnick, J.-C. (2002). Facilitating bid evaluation in Public call for tenders: a socio-technical approach. *Omega*, 30, 227.
- Bana E Costa, C.A., De Corte, J.-M., & Vansnick, J.-C. (2012). MACBETH. *International Journal of Information Technology & Decision Making*, 11, 359-387.
- Bana e Costa, C.A., Ensslin, L., Corrêa, É.C., & Vansnick, J.-C. (1999). Decision Support Systems in action: Integrated application in a multicriteria decision aid process. *European Journal of Operational Research*, 113, 315-335.
- Bana e Costa, C.A., & Oliveira, M.D. (2012). A multicriteria decision analysis model for faculty evaluation. *Omega*, 40, 424-436.
- Bana e Costa, C.A., & Vansnick, J.-C. (1994). MACBETH — An Interactive Path Towards the Construction of Cardinal Value Functions. *International Transactions in Operational Research*, 1, 489.
- Bana e Costa, C.A., & Vansnick, J.-C. (1997). Applications of the MACBETH Approach in the Framework of an Additive Aggregation Model. *Journal of Multi-Criteria Decision Analysis*, 6, 107-114.
- Belton, V., & Stewart, T. (2002). *Multiple criteria decision analysis: an integrated approach*. Dordrecht: Kluwer Academic Publishers.
- Berthold, D.R., Pond, G.R., Soban, F., de Wit, R., Eisenberger, M., & Tannock, I.F. (2008). Docetaxel plus prednisone or mitoxantrone plus prednisone for advanced prostate cancer: updated survival in the TAX 327 study. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*, 26, 242-245.
- Bianchini, D., Lorente, D., Rodriguez-Vida, A., Omlin, A., Pezaro, C., Ferraldeschi, R., et al. (2014). Antitumour activity of enzalutamide (MDV3100) in patients with metastatic castration-resistant prostate cancer (CRPC) pre-treated with docetaxel and abiraterone. *European Journal Of Cancer (Oxford, England: 1990)*, 50, 78-84.
- BNF. (2015). British National Formulary 69. <https://www.bnf.org/>.
- Burström, K., Johannesson, M., & Diderichsen, F. (2001). Swedish population health-related quality of life results using the EQ-5D. *Quality Of Life Research: An International Journal Of Quality Of Life Aspects Of Treatment, Care And Rehabilitation*, 10, 621-635.

- Center, M.M., Jemal, A., Lortet-Tieulent, J., Ward, E., Ferlay, J., Brawley, O., et al. (2012). International Variation in Prostate Cancer Incidence and Mortality Rates. *European Urology*, 61, 1079-1092.
- Cherny, N.I., Sullivan, R., Dafni, U., Kerst, J.M., Sobrero, A., Zielinski, C., et al. (2015). A standardised, generic, validated approach to stratify the magnitude of clinical benefit that can be anticipated from anti-cancer therapies: the European Society for Medical Oncology Magnitude of Clinical Benefit Scale (ESMO-MCBS). *Annals Of Oncology: Official Journal Of The European Society For Medical Oncology / ESMO*, 26, 1547-1573.
- Chongtrakul, P., Sumpradit, N., & Yoongthong, W. (2005). *ISafe and the evidence-based approach for essential medicines selection in Thailand*. pp. 18-19): Essential Drugs Monitor.
- Claxton, K., Martin, S., Soares, M., Rice, N., Spackman, E., Hinde, S., et al. (2015). Methods for the estimation of the National Institute for Health and care excellence cost- effectiveness threshold. *Health Technology Assessment*, 19, 1-503.
- Cleemput, I., Devriese, S., Kohn, L., Devos, C., van Til, J., Groothuis-Oudshoorn, C.G.M., et al. (2018). What Does the Public Want? Structural Consideration of Citizen Preferences in Health Care Coverage Decisions. *MDM Policy & Practice*, 3.
- Clement, F.M., Harris, A., Li, J.J., Yong, K., Lee, K.M., & Manns, B.J. (2009). Using effectiveness and cost-effectiveness to make drug coverage decisions: a comparison of Britain, Australia, and Canada. *JAMA*, 302, 1437-1443.
- Cohen, D. (2017). Cancer drugs: high price, uncertain value. *BMJ*, 359.
- Collins, R., Fenwick, E., Trowman, R., Perard, R., Norman, G., Light, K., et al. (2007). A systematic review and economic model of the clinical effectiveness and cost-effectiveness of docetaxel in combination with prednisone or prednisolone for the treatment of hormone-refractory metastatic prostate cancer. *Health Technology Assessment (Winchester, England)*, 11, iii.
- Connock, M., Cummins, E., Shyangdan, D., Hall, B., Grove, A., & Clarke, A. (2011). Abiraterone acetate for the treatment of metastatic, castrate-resistant prostate cancer following previous cytotoxic chemotherapy: A Single Technology Appraisal. Warwick Evidence.
- Currim, I.S., & Sarin, R.K. (1984). A Comparative Evaluation of Multiattribute Consumer Preference Models. *Management Science*, 30, 543-561.
- Danner, M., Hummel, J.M., Volz, F., van Manen, J.G., Wiegard, B., Dintsios, C.-M., et al. (2011). Integrating patients' views into health technology assessment: Analytic hierarchy process (AHP) as a method to elicit patient preferences. *International Journal Of Technology Assessment In Health Care*, 27, 369-375.
- de Bono, J.S., Logothetis, C.J., Molina, A., Fizazi, K., North, S., Chu, L., et al. (2011). Abiraterone and increased survival in metastatic prostate cancer. *The New England Journal Of Medicine*, 364, 1995-2005.
- de Bono, J.S., Oudard, S., Ozguroglu, M., Hansen, S., Machiels, J.-P., Kocak, I., et al. (2010). Prednisone plus cabazitaxel or mitoxantrone for metastatic castration-resistant prostate cancer progressing after docetaxel treatment: a randomised open-label trial. *Lancet*, 376, 1147-1154.
- Devlin, N., & Lorgelly, P. (2017). QALYs as a measure of value in cancer. *Journal of Cancer Policy*, 11, 19-25.
- Devlin, N., & Sussex, J. (2011). Incorporating multiple criteria in HTA: methods and processes. London: Office of Health Economics.

- Drummond, M., Tarricone, R., & Torbica, A. (2013). Assessing the added value of health technologies: reconciling different perspectives. *Value In Health: The Journal Of The International Society For Pharmacoeconomics And Outcomes Research*, 16, S7-S13.
- Efthymiadou, O., Mossman, J., & Kanavos, P. (2019). Health related quality of life aspects not captured by EQ-5D-5L: results from an international survey of patients. *Health Policy*, 123, 159-165.
- EMA. (2016a). Jevtana (cabazitaxel) EPAR - Product Information. European Medicines Agency.
- EMA. (2016b). Xtandi (enzalutamide) EPAR - Product Information. European Medicines Agency.
- EMA. (2016c). Zytiga (abiraterone) EPAR - Product Information. European Medicines Agency.
- Faden, R.R., Chalkidou, K., Appleby, J., Waters, H.R., & Leider, J.P. (2009). Expensive Cancer Drugs: A Comparison between the United States and the United Kingdom. *Milbank Quarterly*, 87, 789-819.
- Fasolo, B., & Bana e Costa, C.A. (2014). Tailoring value elicitation to decision makers' numeracy and fluency: Expressing value judgments in numbers or words. *Omega*, 44, 83-90.
- Fizazi, K., Scher, H.I., Molina, A., Logothetis, C.J., Chi, K.N., Jones, R.J., et al. (2012). Abiraterone acetate for treatment of metastatic castration-resistant prostate cancer: final overall survival analysis of the COU-AA-301 randomised, double-blind, placebo-controlled phase 3 study. *The Lancet. Oncology*, 13, 983-992.
- Franco, L., & Montibeller, G. (2010a). Problem Structuring for Multicriteria Decision Analysis Interventions analysis interventions. *Wiley Encyclopedia of Operations Research and Management Science*.
- Franco, L.A., & Montibeller, G. (2010b). Facilitated modelling in operational research. *European Journal of Operational Research*, 205, 489-500.
- GHEMA. (2016). ENZALUTAMIDA en cáncer próstata metastásico resistente a la castración (previo a quimioterapia). Grupo de Evaluación de Novedades, EStandarización e Investigación en Selección de Medicamentos.
- Goetghebeur, M.M., Wagner, M., Khoury, H., Levitt, R.J., Erickson, L.J., & Rindress, D. (2008). Evidence and Value: Impact on DEcisionMaking--the EVIDEM framework and potential applications. *BMC Health Services Research*, 8, 270.
- Goetghebeur, M.M., Wagner, M., Khoury, H., Rindress, D., Grégoire, J., & Deal, C. (2010). Combining multicriteria decision analysis, ethics and health technology assessment: Applying the EVIDEM decisionmaking framework to growth hormone for Turner syndrome patients. *Cost Effectiveness and Resource Allocation*, 8, <xocs:firstpage xmlns:xocs=""/>.
- Horwich, A., Parker, C., de Reijke, T., & Kataja, V. (2013). Prostate cancer: ESMO Clinical Practice Guidelines for diagnosis, treatment and follow-up. *Annals Of Oncology: Official Journal Of The European Society For Medical Oncology / ESMO*, 24 Suppl 6, vi106-vi114.
- Ijzerman, M.J., van Til, J.A., & Snoek, G.J. (2008). Comparison of two multi-criteria decision techniques for eliciting treatment preferences in people with neurological disorders. *The Patient*, 1, 265-272.
- INAMI. (2019). Remboursement de médicaments : Décisions ministérielles et rapports d'évaluation de la CRM. <https://www.inami.fgov.be/fr/programmes-web/Pages/applications-rapports-crm.aspx>.

- Jansen, J.P., Fleurence, R., Devine, B., Itzler, R., Barrett, A., Hawkins, N., et al. (2011). Interpreting indirect treatment comparisons and network meta-analysis for health-care decision making: report of the ISPOR Task Force on Indirect Treatment Comparisons Good Research Practices: part 1. *Value In Health: The Journal Of The International Society For Pharmacoeconomics And Outcomes Research*, 14, 417-428.
- Jaramillo, H.E.C., Goetghebeur, M., & Moreno-Mattar, O. (2016). TESTING MULTI-CRITERIA DECISION ANALYSIS FOR MORE TRANSPARENT RESOURCE-ALLOCATION DECISION MAKING IN COLOMBIA. 32, 307-314.
- Kanavos, P., & Angelis, A. (2013). Multiple criteria decision analysis for value based assessment of new medical technologies: a conceptual framework. The LSE Health Working Paper Series in Health Policy and Economics: London School of Economics and Political Science.
- Kearns, B., Lloyd Jones, M., Stevenson, M., & Littlewood, C. (2013). Cabazitaxel for the Second-Line Treatment of Metastatic Hormone-Refractory Prostate Cancer: A NICE Single Technology Appraisal. *PharmacoEconomics*, 31, 479-488.
- Keeney, R. (1992). *Value focused thinking: a path to creative decision making*. Cambridge: Harvard University Press.
- Keeney, R., & Raiffa, H. (1993). *Decisions with multiple objectives: preferences and value trade-offs*. Cambridge: Cambridge University Press.
- Keeney, R.L. (1982). Decision Analysis: An Overview. *Operations Research*, 30, 803.
- Keeney, R.L. (2002). Common Mistakes In Making Value Trade-Offs. *Operations Research*, 50, 935.
- Keeney, R.L., & Gregory, R.S. (2005). Selecting Attributes to Measure the Achievement of Objectives. *Operations Research*, 53, 1-11.
- Linley, W.G., & Hughes, D.A. (2013). Societal views on NICE, cancer drugs fund and value-based pricing criteria for prioritising medicines: a cross-sectional survey of 4118 adults in Great Britain. *Health Economics*, 22, 948-964.
- London School of Economics. (2019). ADVANCE_HTA Project. <http://www.lse.ac.uk/lse-health/research/projects/advance-hta>.
- Loriot, Y., Bianchini, D., Ileana, E., Sandhu, S., Patrikidou, A., Pezaro, C., et al. (2013). Antitumour activity of abiraterone acetate against metastatic castration-resistant prostate cancer progressing after docetaxel and enzalutamide (MDV3100). *Annals Of Oncology: Official Journal Of The European Society For Medical Oncology / ESMO*, 24, 1807-1812.
- Malone, D.C., Berg, N.S., Claxton, K., Garrison, L.P., Jr., Ijzerman, M., Marsh, K., et al. (2016). International Society for Pharmacoeconomics and Outcomes Research Comments on the American Society of Clinical Oncology Value Framework. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*, 34, 2936-2937.
- Marsh, K., Ijzerman, M., Thokala, P., Baltussen, R., Boysen, M., Kaló, Z., et al. (2016). Multiple Criteria Decision Analysis for Health Care Decision Making-Emerging Good Practices: Report 2 of the ISPOR MCDA Emerging Good Practices Task Force. *Value In Health: The Journal Of The International Society For Pharmacoeconomics And Outcomes Research*, 19, 125-137.

- Marsh, K., Lanitis, T., Neasham, D., Orfanos, P., & Caro, J. (2014). Assessing the Value of Healthcare Interventions Using Multi-Criteria Decision Analysis: A Review of the Literature. *Pharmacoeconomics*, 32, 345-365.
- Marsh, K., Sculpher, M., Caro, J.J., & Tervonen, T. (2018). The Use of MCDA in HTA: Great Potential, but More Effort Needed. *Value in Health*.
- Montibeller, G., & Winterfeldt, D. (2015). Cognitive and Motivational Biases in Decision and Risk Analysis. *Risk Analysis: An International Journal*, 35, 1230-1251.
- Morton, A. (2017). Treacle and Smallpox: Two Tests for Multicriteria Decision Analysis Models in Health Technology Assessment. *Value in Health*, 20, 512-515.
- Nadler, E., Eckert, B., & Neumann, P.J. (2006). Do Oncologists Believe New Cancer Drugs Offer Good Value? *Oncologist*, 11, 90-95.
- NICE. (2012a). Abiraterone for castration-resistant metastatic prostate cancer previously treated with a docetaxel-containing regimen. Technology Appraisal Guidance 259: National Institute for Health and Care Excellence.
- NICE. (2012b). Cabazitaxel for hormone-refractory metastatic prostate cancer previously treated with a docetaxel-containing regimen. Technology Appraisal Guidance 255: National Institute for Health and Care Excellence.
- NICE. (2014). Enzalutamide for metastatic hormone- relapsed prostate cancer previously treated with a docetaxel - containing regimen. Technology Appraisal Guidance 316: National Institute for Health and Care Excellence.
- Nicod, E., & Kanavos, P. (2012). Commonalities and differences in HTA outcomes: A comparative analysis of five countries and implications for coverage decisions. *Health Policy*, 108, 167-177.
- NIH. (2016). ClinicalTrials.gov. US National Institutes of Health.
- Nord, E., Daniels, N., & Kamlet, M. (2009). QALYs: Some Challenges. *Value in Health*, 12, S10-S15.
- Obwieszczenie, M. (2017). Obwieszczenie Ministra Zdrowia z dnia 25 października 2017 r. w sprawie wykazu refundowanych leków, środków spożywczych specjalnego przeznaczenia żywieniowego oraz wyrobów medycznych (DZ. URZ. Min. Zdr. 2017.105). <http://www.bip.mz.gov.pl/legislacja/akty-prawne-1/obwieszczenie-ministra-zdrowia-z-dnia-25-pazdziernika-2017-r-w-sprawie-wykazu-refundowanych-lekow-srodkow-spozywczych-specjalnego-przeznaczenia-zywieniowego-oraz-wyrobow-medycznych-na-1-listopada-20/>.
- Peacock, S.J., Richardson, J.R.J., Carter, R., & Edwards, D. (2007). Priority setting in health care using multi- attribute utility theory and programme budgeting and marginal analysis (PBMA). *Social Science and Medicine*, 64, 897-910.
- Phelps, C., & Madhavan, G. (2018). Resource allocation in decision support frameworks. *Cost Effectiveness and Resource Allocation*, 16.
- Phillips, L. (1984). A theory of requisite decision models. *Acta Psychologica*, 56, 29-48.
- Phillips, L. (2007). Decision Conferencing. In W. Edwards, R. Miles, & D. von Winterfeldt (Eds.), *Advances in Decision Analysis: From Foundations to Applications*. Cambridge: Cambridge University Press.
- Phillips, L., & Phillips, M. (1993). Facilitated Work Groups: Theory and Practice. *The Journal of the Operational Research Society*, 44.
- Phillips, L.D. (2017). Best Practice for

- MCDA in Healthcare. In K. Marsh, Goetghebeur, M., P. Thokala, & R. Baltussen (Eds.), *Multi-Criteria Decision Analysis to Support Healthcare Decisions*: Springer.
- Phillips, L.D., & Bana e Costa, C.A. (2007). Transparent prioritisation, budgeting and resource allocation with multi-criteria decision analysis and decision conferencing. *Annals of Operations Research*, 154, 51-68.
- Radaelli, G., Lettieri, E., Masella, C., Merlino, L., Strada, A., & Tringali, M. (2014). Implementation of EUnetHTA core Model® in Lombardia: the VTS framework. *International Journal Of Technology Assessment In Health Care*, 30, 105-112.
- Riemsma, R., Ramaekers, B., Tomini, F., Wolff, R., van Asselt, A., Joore, M., et al. (2013). Abiraterone for the treatment of chemotherapy naïve metastatic castration-resistant prostate cancer: a Single Technology Appraisal. Kleijnen Systematic Reviews Ltd.
- Rodrigues, T.C., Montibeller, G., Oliveira, M.D., & Bana E Costa, C.A. (2017). Modelling multicriteria value interactions with Reasoning Maps. *European Journal of Operational Research*, 258, 1054-1071.
- Schein, E. (1999). *Process consultation revisited: building the helping relationship*. Reading: Addison–Wesley.
- Scher, H.I., Fizazi, K., Saad, F., Taplin, M.-E., Sternberg, C.N., Miller, K., et al. (2012). Increased survival with enzalutamide in prostate cancer after chemotherapy. *The New England Journal Of Medicine*, 367, 1187-1197.
- Schnipper, L.E., Davidson, N.E., Wollins, D.S., Tyne, C., Blayney, D.W., Blum, D., et al. (2015). American Society of Clinical Oncology Statement: A Conceptual Framework to Assess the Value of Cancer Treatment Options. *Journal Of Clinical Oncology: Official Journal Of The American Society Of Clinical Oncology*, 33, 2563-2577.
- Sculpher, M., Claxton, K., & Pearson, S.D. (2017). Developing a Value Framework: The Need to Reflect the Opportunity Costs of Funding Decisions. *Value in Health*, 20, 234-239.
- Shih, Y.C.T., Ganz, P.A., Aberle, D., Abernethy, A., Bekelman, J., Brawley, O., et al. (2013). Delivering high-quality and affordable care throughout the cancer care continuum. *Journal of Clinical Oncology*, 31, 4151-4157.
- Sullivan, P., Mulani, P., Fishman, M., & Sleep, D. (2007). Quality of life findings from a multicenter, multinational, observational study of patients with metastatic hormone-refractory prostate cancer. *An International Journal of Quality of Life Aspects of Treatment, Care and Rehabilitation - Official Journal of the International Society of Quality of Life Research*, 16, 571-575.
- Sulmasy, D., & Moy, B. (2014). Debating the oncologist's role in defining the value of cancer care: Our duty is to our patients. *Journal of Clinical Oncology*, 32, 4039-4041.
- Sussex, J., Rollet, P., Garau, M., Schmitt, C., Kent, A., & Hutchings, A. (2013a). A Pilot Study of Multicriteria Decision Analysis for Valuing Orphan Medicines. *Value in Health*, 16, 1163-1169.
- Sussex, J., Rollet, P., Garau, M., Schmitt, C., Kent, A., & Hutchings, A. (2013b). Multi-criteria decision analysis to value orphan medicines. London: Office of Health Economics.
- Tannock, I.F., de Wit, R., Berry, W.R., Horti, J., Pluzanska, A., Chi, K.N., et al. (2004). Docetaxel plus prednisone or mitoxantrone plus prednisone for

- advanced prostate cancer. *The New England Journal Of Medicine*, 351, 1502-1512.
- Tervonen, T., Naci, H., van Valkenhoef, G., Ades, A.E., Angelis, A., Hillege, H.L., et al. (2015). Applying Multiple Criteria Decision Analysis to Comparative Benefit-Risk Assessment: Choosing among Statins in Primary Prevention. *Medical Decision Making: An International Journal Of The Society For Medical Decision Making*, 35, 859-871.
- Thokala, P. (2011). Multi criteria decision analysis for health technology assessment: report by the decision support unit. Sheffield: University of Sheffield.
- TLV. (2014). Designation 2775/2013. Tandvårds- och läkemedelsförmånsverket.
- TLV. (2015a). Designation 4774/2014. Tandvårds- och läkemedelsförmånsverket.
- TLV. (2015b). Designation 4852/2014. Tandvårds- och läkemedelsförmånsverket.
- Tony, M., Wagner, M., Khoury, H., Rindress, D., Papastavros, T., Oh, P., et al. (2011). Bridging health technology assessment (HTA) with multicriteria decision analyses (MCDA): field testing of the EVIDEM framework for coverage decisions by a public payer in Canada. *BMC Health Services Research*, 11, 329.
- Torre, L.A. (2015). Global cancer statistics, 2012. *CA: A Cancer Journal For Clinicians*, 65, 87.
- van Valkenhoef, G., Tervonen, T., Zhao, J., de Brock, B., Hillege, H.L., & Postmus, D. (2012). Multicriteria benefit-risk assessment using network meta-analysis. *Journal Of Clinical Epidemiology*, 65, 394-403.
- von Winterfeldt, D., & Edwards, W. (1986). *Decision analysis and behavioral research*. Cambridge: Cambridge University Press.
- Wagner, M., Khoury, H., Bennetts, L., Berto, P., Ehreth, J., Badia, X., et al. (2017). Appraising the holistic value of Lenvatinib for radio-iodine refractory differentiated thyroid cancer: A multi-country study applying pragmatic MCDA. *BMC Cancer*, 17, <xocs:firstpage xmlns:xocs=""/>.
- World Health Organisation Collaborating Centre. (2016). ATC/DDD Index 2016. World Health Organisation Collaborating Centre.
- Wouters, O.J., Naci, H., & Samani, N.J. (2015). QALYs in cost-effectiveness analysis: an overview for cardiologists. *Heart (British Cardiac Society)*, 101, 1868-1873.

Appendix

Model Building: Alternative Treatments Compared and Evidence Considered

The source of evidence used for identifying the performance of options across the evaluation criteria is shown in Table A1.

Model Building: Setting Attribute Ranges and Reference Levels

For the case of clinical therapeutic attributes, “lower” reference levels were based on best standard of care (BSC) performance, coming from the median of the respective placebo arm of the *AFFIRM* trial, with the exception of the HRQoL attribute (EQ-5D utility score) that was based on the utility of stable disease with no treatment coming from past NICE TAs (NICE, 2012a, b). The “higher” reference levels were derived by adding a 20% absolute improvement to the performance level of the best performing option, besides for the case of the HRQoL attribute (EQ-5D utility score) that was based on the general Swedish population (Burström et al., 2001). The rationale was to design a value scale incorporating a “global” reference level (Belton & Stewart, 2002), reflecting an “ideal” performance (as proxied by the 20% improvement in best available performance), corresponding to the 100 anchor level of the value scale. This could also offer a flexibility margin to be able to incorporate the performance of future improved options within the same elicited value scale. Consequently, two reference levels within the attribute range were defined in most cases: i) the “lower” reference level (x_l) (i.e. BSC-based satisfactory performance), acting on the same time also as the minimum limit of the attribute range (x_*); and ii) the “higher” reference level (x_h) (i.e. 20% better than the best performing option), acting on the same time as the maximum limit of the attribute range (x^*) to give $x_* = x_l \leq x_h = x^*$.

A similar, but reverse, logic was used for setting the reference levels in the “treatment discontinuation” attribute of the safety cluster; the “lower” reference level was defined to be equal to the BSC (i.e. placebo) arm of the *AFFIRM* trial. However, contrary to the logic

adopted so far for the therapeutic benefit criteria, the “higher” reference level was not set equal to 20% worse than the best performing option (because the lower the performance, the higher the value), but rather equal to the minimum, i.e. worst possible, natural limit of the attribute scale (i.e. 0%) which was regarded as an “ideal” level. In turn, the minimum limit of the scale was derived by worsening the performance of the worst performing treatment option by 20%. A similar approach was used for setting the reference levels of the qualitative “contraindications” attribute, defining the “higher” reference level equal to the maximum (i.e. most attractive) limit of the attribute scale (i.e. none known contraindications) and the “lower” reference level equal to the minimum (i.e. least attractive) limit of the attribute scale.

For the innovation attributes, the “higher” reference level was set either equal to 20% better than the best performing option for the case of natural quantitative attributes (e.g. number of new indications for which the technology is investigated in a given clinical development stage), or equal to the maximum, i.e. best possible, limit of the scale for the case of constructed qualitative attributes (e.g. the existence of any special instructions, the technology's relative market entrance in regards to its ATC Level), reflecting a “global” versus “local” scaling approach respectively. Given that the BSC performance was irrelevant to be used as satisfactory level in the innovation attributes, and any efforts to derive a “satisfactory” level would be subjective in nature, the minimum limit of the scale for each attribute was used as a “lower” reference level. Therefore the “lower” reference level was based on the worst performance plausible as inferred from the lowest possible limit of the scales, both for the case of natural quantitative attributes (e.g. 0 number of new indications for which the technology is investigated in a given clinical development stage), and the case of constructed qualitative attributes (e.g. worst possible combination of special instructions, 5th entrance at an ATC level).

For the socioeconomics attribute (impact on direct costs), the “higher” reference level was based on the BSC's impact on cost (i.e. £0 impact on costs), given that by definition impact on costs for all treatment options are incremental to BSC, and the “lower” reference

level was derived by adding a 20% absolute increment to the worst performing option (i.e. to the one with the biggest impact on costs).

“Lower” and “higher” reference levels for all attributes at the pre-workshop stage and the basis of their selection are outlined in Table A2 (assuming no impact of luteinizing hormone-releasing hormone analogue).

Model Assessment and Appraisal: Decision Conference

On the day of each decision conference the preliminary model was validated with the participants by revising it cluster by cluster through an open discussion, seeking group consensus and adopting an iterative and interactive-model-building process where debate was encouraged and differences of opinion were actively sought.

In terms of the decision-aiding methodology used, the lead author acted as an impartial facilitator with the aim of enhancing content and process interaction, while refraining from contributing to the content of the group’s discussions, essentially guiding the group in how to think about the issues but not what to think (Phillips & Bana e Costa, 2007; Schein, 1999).

In terms of facilities, the rooms of the decision conferences had a Π-shaped meeting table for all the participants to have direct eye to eye contact, with an overhead projector screen and a second portable projector or large TV screen. The M-MACBETH software (more information provided in the MCDA Technique section of the main text and below) was operated using a laptop, the screen of which was connected to the projector, and the second screen was used to show the list of the evaluation criteria together with their “lower” and “higher” reference levels.

The decision conferences took place over a full working day or two half working days; in the former case, there was one lunch break and two coffee breaks throughout the day, whereas in the latter case only a coffee break took place around the middle of each session. In each decision conference, the day started with an overview of the MCDA methodology

adopted and the description of the preliminary version of the value tree which was then analysed cluster by cluster. At the beginning of each cluster the value tree was validated; the various criteria were explained, followed by a group discussion relating to their relevance and completeness. As a result of this iterative process, some of the criteria were not included because they were perceived as irrelevant or non-fundamental. Schematic illustrations of the final versions of the value trees are shown in Figure A1. Then, value functions were elicited for the different criteria and relative weights were assigned within the clusters. Finally, relative weights were assigned across clusters, enabling the calculation of the options' overall WPV scores.

Model Assessment and Appraisal: MCDA Technique

MACBETH uses seven semantic categories ranging between “no difference” to “extreme difference”, in order to distinguish between the value of different attribute levels. Based on these qualitative judgements of difference and, by analysing judgmental inconsistencies, it facilitates the move from ordinal preference modeling, a cognitively less demanding elicitation of preferences, to a quantitative value function. The approach has evolved through the course of theoretical research and real world practical applications, making it an interactive decision support system that facilitates decision-makers' communication. An example of the type of questioning being asked would be “What do you judge to be the difference of value between x' and x'' ?” where x' and x'' are two different attribute levels of attribute x , across the plausible range (i.e. $x^* \leq x', x'' \leq x^*$). The value judgements matrix for the Overall Survival attribute and their conversion into its value function is provided as an example in Figure A2.

Following the elicitation of value functions, criteria baseline weights can be elicited. Questions of direct importance for a criterion such as “*How important is a given criterion?*” are known to be as one of the most common mistakes when making value trade-offs because they are assessing them independent of the respective attribute ranges (Keeney, 2002). In contrast, indirect weighting technique that assess value trade-offs in tandem with the

respective ranges of attributes should be employed. For example, the quantitative swing weighting technique asks for judgments of relative value between ‘swings’ (i.e. changes from standard lower level x_* to higher reference level x^* on each x^{th} attribute) taking the form “*How would you rank the relative importance of the criteria, considering their attributes ranges relative to 100 for the highest-ranked criterion considering its range?*”. Each swing, i.e. a relative change from a lower attribute level to a higher attribute level, is valued between 0 and 100, with the most valuable swing anchored as 100 (von Winterfeldt & Edwards, 1986). Normalised weights are then calculated, as a proportion of each swing weight, so the normalised weights summed to 100%. Instead, relative attribute weights were calculated using an alternative qualitative swing weighting protocol, by using the MACBETH procedure to elicit the differences in attractiveness between the lower and higher reference levels of the different attributes, initially at individual level and then at criteria cluster level (i.e. by considering multiple attribute swings on the same time) (Bana e Costa et al., 2016b; Bana E Costa et al., 2012).

Finally criteria preference value scores and the respective weights can be combined together through an additive aggregation approach as described in equation 2 (if the adequate conditions of complete and transitive preferences are met as well as multi-attribute preferential independence conditions (von Winterfeldt & Edwards, 1986)).

The M-MACBETH software automatically performs consistency checking between the qualitative judgements expressed, and in addition a second consistency check was manually performed by the author to validate the cardinality, i.e. interval nature, of the emerging value scale. This was done by comparing the sizes of the intervals between the proposed scores and inviting participants to adjust them if necessary (Fasolo & Bana e Costa, 2014), a requirement which is essential for the application of simple additive value models.

Figure Captions

Tables and Figures

Table 1: Criteria definitions, their consideration in each jurisdiction and their ranking

Criteria Sub-Domain	Evaluation criteria	Definition	Country (competent HTA organisation)			
			Belgium (INAMI/RIZIV)	Poland (AOTMiT)	Andalusia (AETSA)	Sweden (TLV)
Criteria Domain 1: Therapeutic Benefit						
Direct endpoints	Overall survival x Health related quality of life*	The median time from treatment randomisation to death adjusted for the mean health related quality of life using the EQ-5D utility score	✓ (1 st)	✓ (1 st)	✓ (1 st)	✓ (1 st)
Indirect endpoints	Radiographic tumour progression	The median survival time on which patients have not experienced disease progression (using RECIST criteria)	✓ (5 th)	✓ (5 th)	✓ (4 th)	✓ (5 th)
	PSA response	The proportion of patients having a ≥50% reduction in PSA			✓ (8 th)	
Criteria Domain 2: Safety Profile						
Tolerability	Treatment discontinuation	The proportion of patients discontinuing treatment due to adverse events	✓ (3 rd)	✓ (3 rd)	✓ (2 nd)	✓ (2 nd)
Contra-indications & warnings	Contra-indications	The existence of any type of contra-indication accompanying the treatment	✓ (4 th)		✓ (3 rd)	✓ (4 th)
Criteria Domain 3: Innovation Level						
Type and timing of innovation	ATC Level 1	The technology's relative market entrance in regards to its ATC Level 1 (Anatomical)				
	ATC Level 2	The technology's relative market entrance in regards to its ATC Level 2 (Therapeutic)				
	ATC Level 3	The technology's relative market entrance in regards to its ATC Level 3 (Pharmacological)				
	ATC Level 4	The technology's relative market entrance in regards to its ATC Level 4 (Chemical)	✓ (6 th)		✓ (10 th)	

	ATC Level 5	The technology's relative market entrance in regards to its ATC Level 5 (Molecular)				
Spill-over effect	Phase 1	The number of new indications for which the technology is investigated in Phase 1 clinical trials				
	Phase 2	The number of new indications for which the technology is investigated in Phase 2 clinical trials	✓ (8 th)			
	Phase 3	The number of new indications for which the technology is investigated in Phase 2 clinical trials	✓ (9 th)		✓ (9 th)	
	Marketing authorisation	The number of new indications that the technology has gained an approval for at the stage of marketing authorisation	✓ (10 th)		✓ (7 th)	
Patient convenience	Delivery posology	The combination of the delivery system (RoA and dosage form) with the posology (frequency of dosing and duration of administration) of the treatment	✓ (7 th)	✓ (4 th)	✓ (6 th)	✓ (7 th)
	Special instructions	The existence of any special instructions accompanying the administration of the treatment	✓ (11 th)	✓ (6 th)	✓ (11 th)	✓ (6 th)
Criteria Domain 4: Socio-Economic Impact						
Direct costs	Medical costs impact	The impact of the technology on direct medical costs excluding the purchasing costs of the technology	✓ (2 nd)	✓ (2 nd)	✓ (5 th)	✓ (3 rd)

Notes: *: Aggregation between OS and HRQoL criteria took place due to preference-dependence leading to a combined criterion; PSA= prostate-specific antigen; ATC=Anatomical Therapeutic Chemical classification system; RoA=Route of Administration.

Source: The authors, based on DCs in Andalusia/Spain, Belgium, Poland and Sweden.

Table 2: Performance matrix and reference levels considered across the final criteria attributes

Criterion name	Attribute metric	Lower level	Abiraterone	Cabazitaxel	Enzalutamide	Higher level
Overall survival (OS)*	Months	13.6	15.8	15.1	18.4	22.1
Health Related Quality of Life (HRQoL), stable disease*	Utility (EQ-5D)	0.72	0.76	0.76***	0.76	0.82
Health Related Quality of Life (HRQoL), progressive disease*	Utility (EQ-5D)	0.64	0.64	0.64	0.64	0.82
OS X HRQoL**	Quality adjusted life months (QALMs)	9.2	11	10.5	12.8	18.1
Radiographic tumour progression, i.e. progression free survival (PFS)	Months	2.9	5.6	8.8	8.3	10.6
PSA response	% of patients	1.5	29.5	39.2	54	64.8
Treatment discontinuation	% of patients	10	19	18	8	0
Contra-indication(s)	Type of contra-indication	hyp + hep imp + low neut	hyp + hep imp	hyp + hep imp + low neut	hyp	None
ATC Level 4, i.e. chemical mechanism of action	Relative market entrance	5 th	2 nd	2 nd	1 st	1 st
Phase 2	Number of new indications	0	1	13	4	16
Phase 3	Number of new indications	0	1	2	0	2
Marketing authorisation	Number of new indications	0	0	0	0	1
Delivery posology	Type of delivery system & posology combinations	Oral, daily - one off + IV, every 3 weeks - 1 hr	Oral, daily - one off	Oral, daily - one off + IV, every 3 weeks - 1 hr	Oral, daily - one off	Oral, daily - one off

Special instructions	Type(s) of special instructions	Concomitant and/or pre- med + no food	Concomitant and/or pre- med + no food	Concomitant and/or pre- med	None	None
Medical costs impact	GBP	10,000	5,750	7,992	567	0

Notes: * Used for the calculation of the quality adjusted life months (QALMs) attribute of the aggregated OS x HRQoL criterion; ** Calculated assuming an equal 50% split in time duration between the stable disease and progressive disease states in HRQoL; *** Used the same score of the other two options as data not available; hyp = hypersensitivity; hep imp = hepatic impairment; low neut = low neutrophil count.

Source: The authors from the literature.

Table 3: Number of criteria attributes per cluster, relative weights per criteria cluster and their ranking across the four HTA settings.

<i>HTA Agency/ Criteria Clusters</i>	Sweden (TLV)			Andalusia (AETSA)			Poland (AOTMiT)			Belgium (INAMI-RIZIV)		
	Criteria numbers	Criteria weights	Criteria ranking	Criteria numbers	Criteria weights	Criteria ranking	Criteria numbers	Criteria weights	Criteria ranking	Criteria numbers	Criteria weights	Criteria ranking
Therapeutic Benefit	2	44.5	1 st	3	54.3	1 st	2	40.0	1 st	2	40.0	1 st
Safety Profile	2	33.3	2 nd	2	26.0	2 nd	1	20.0	3 rd	2	26.7	2 nd
Innovation Level	2	7.4	4 th	5	11.8	3 rd	2	10.0	4 th	6	13.3	4 th
Socioeconomic Impact	1	14.8	3 rd	1	7.9	4 th	1	30.0	2 nd	1	20.0	3 rd
Total	7	100		11	100		6	100		11	100	

Source: The authors based on input from decision conferences.

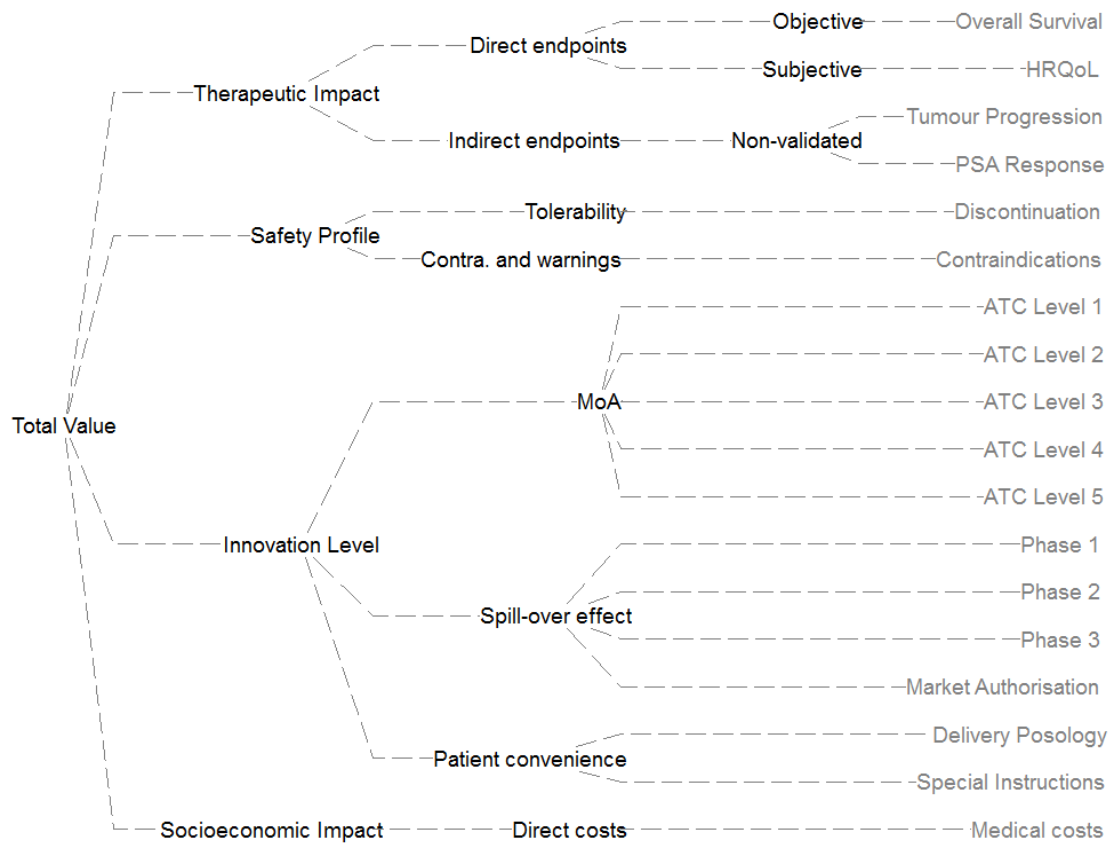
Table 4: Overall weighted preference value (WPV) scores, costs and costs per unit of value across the four HTA settings.

<i>Treatments/ HTA agency</i>	Enzalutamide		Abiraterone		Cabazitaxel	
	<i>Overall WPV score</i>	<i>Ranking per country</i>	<i>Overall WPV score</i>	<i>Ranking per country</i>	<i>Overall WPV score</i>	<i>Ranking per country</i>
Sweden (TLV)	55.1	1 st	2.4	2 nd	-3.4	3 rd
Andalusia (AETSA)	49.1	1 st	8.8	2 nd	4.4	3 rd
Poland (AOTMiT)	59.9	1 st	12.1	2 nd	3.5	3 rd
Belgium (INAMI-RIZIV)	58.6	1 st	16.0	2 nd	10.9	3 rd
<i>Costs (£)</i>	24,600		21,900		23,900	
	<i>Cost per unit of value</i>	<i>Ranking per country</i>	<i>Cost per unit of value</i>	<i>Ranking per country</i>	<i>Cost per unit of value</i>	<i>Ranking per country</i>
Sweden (TLV)	447	1 st	9,221	2 nd	N/A	3 rd
Andalusia (AETSA)	501	1 st	2,496	2 nd	5,481	3 rd
Poland (AOTMiT)	410	1 st	1,805	2 nd	6,816	3 rd
Belgium (INAMI-RIZIV)	420	1 st	1,366	2 nd	2,196	3 rd

Note: No cost-per-unit of value was calculated because of the negative overall WPV score (i.e. having a worst overall performance compared to the performance of the lower reference level), which would produce a negative cost-per-unit of value ($£23,900/(-3.4) = -7,072$) and would therefore faultily “improve” the median figure of the treatment.

Source: The authors.

Figure 1: Preliminary value tree for metastatic prostate cancer (pre-workshop).



Notes: Contra. = Contraindications; MoA = Mechanism of action; HRQoL = Health related quality of life; PSA = Prostate-specific Antigen; ATC = Anatomical therapeutic chemical; Image produced using the Hiview3 software version 3.2.0.4.

Source: The authors.

Figure 2: Criteria valuation drug profiles.

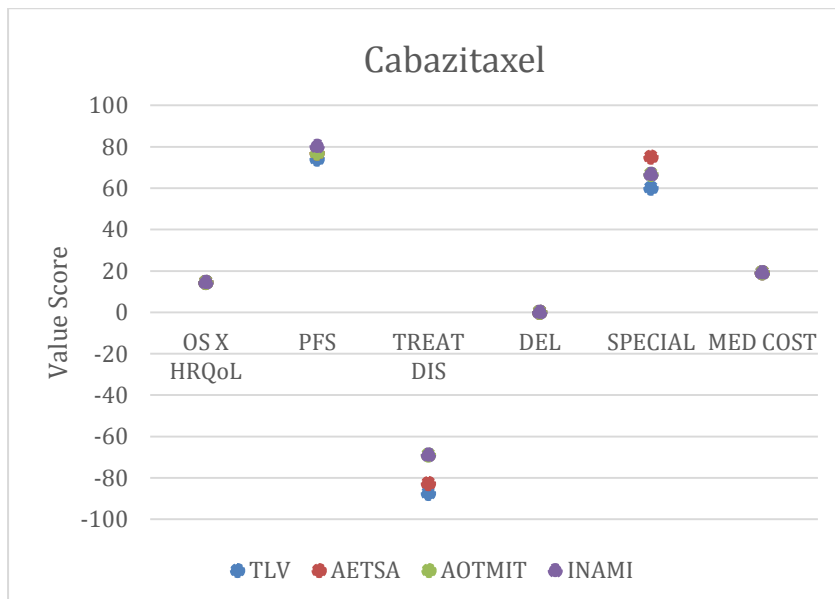
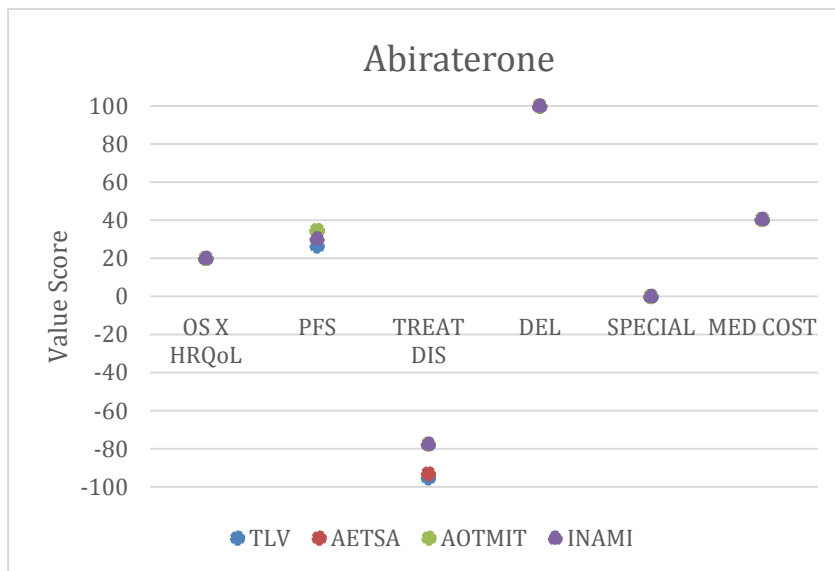
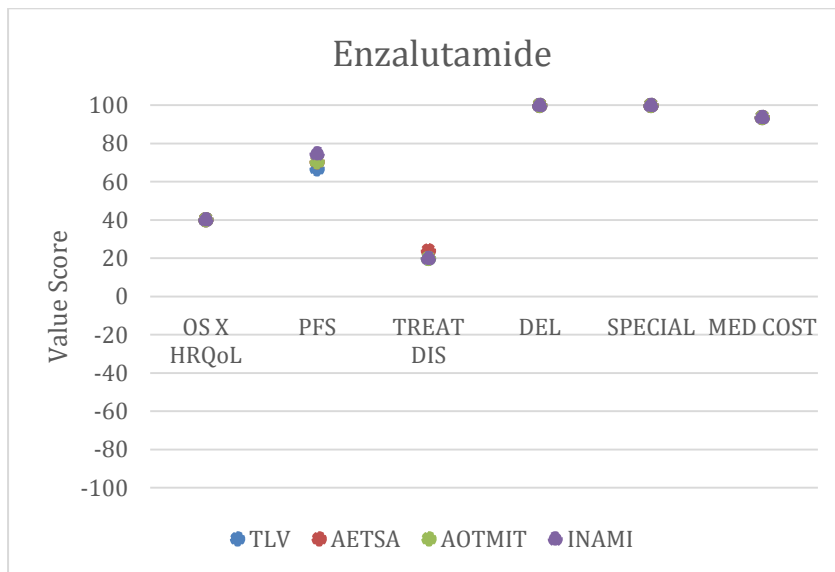


Figure 3: Relative criteria weights stacked bars across the four HTA settings.

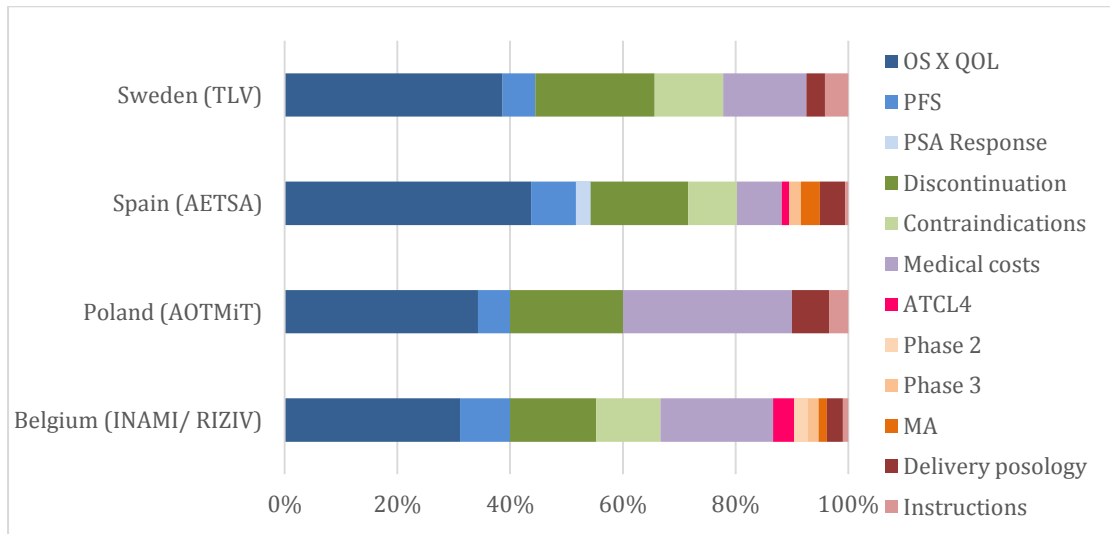


Figure 4: Stacked bar plot of treatments' overall weighted preference value scores across the four HTA settings.

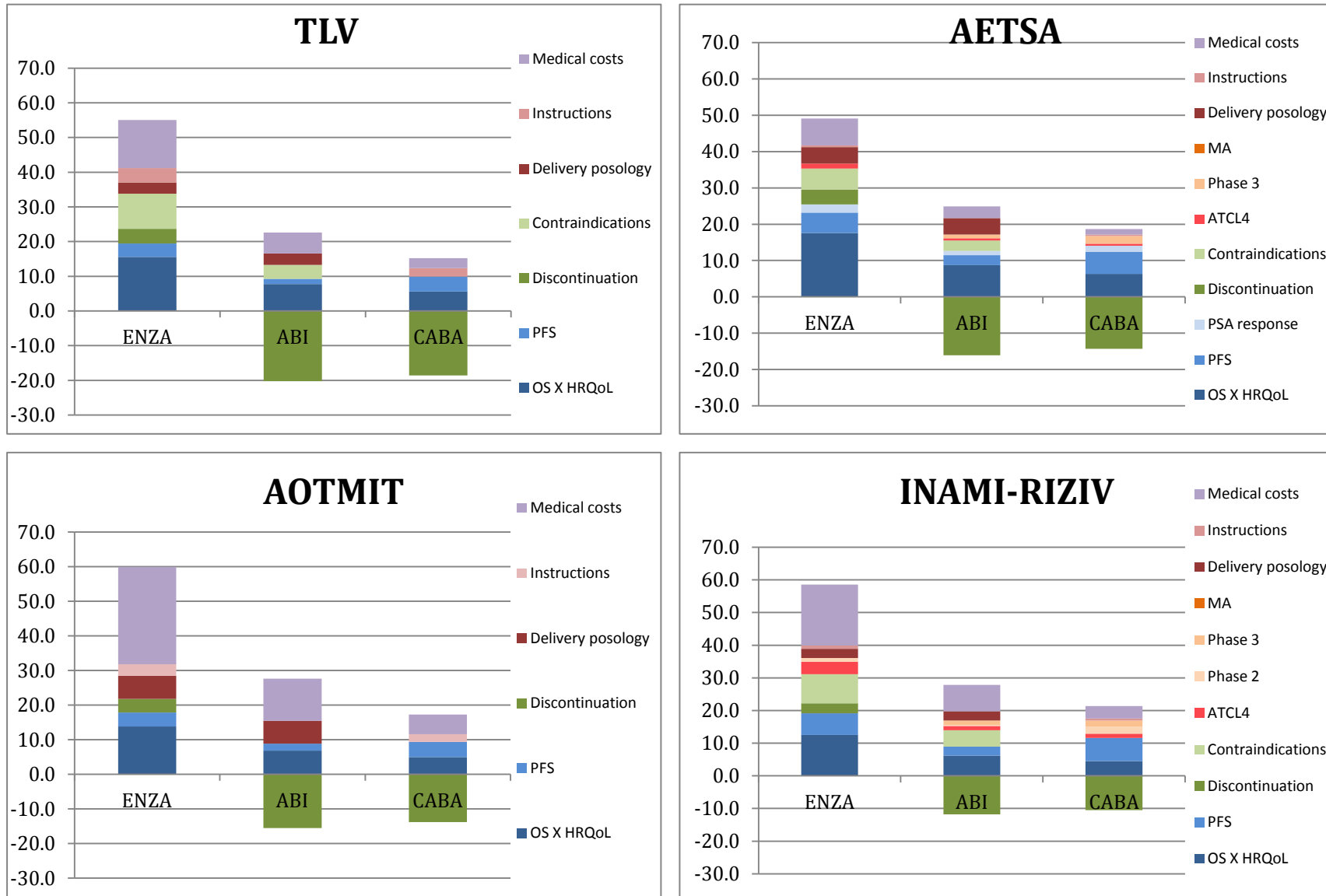
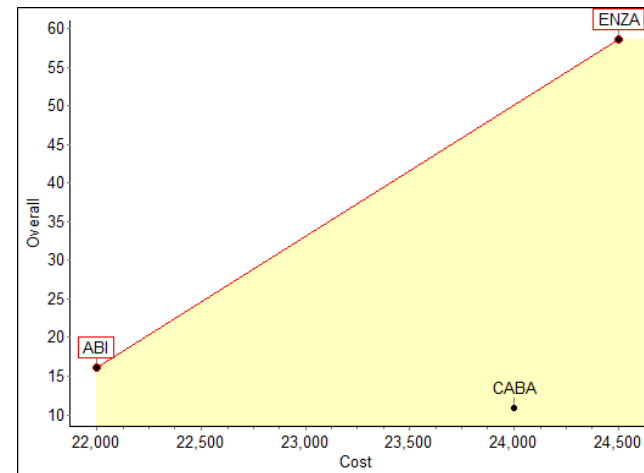
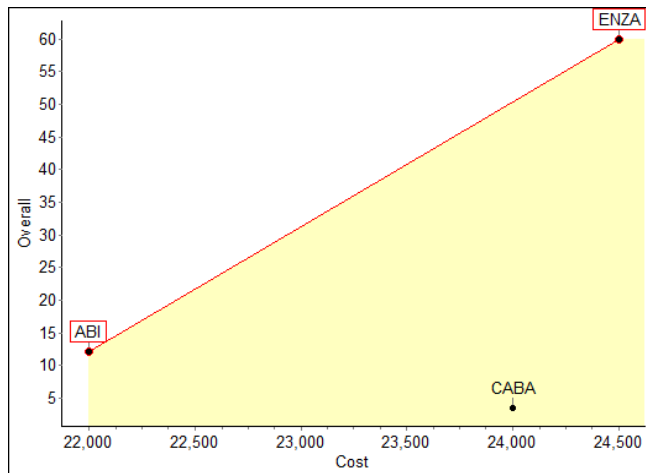
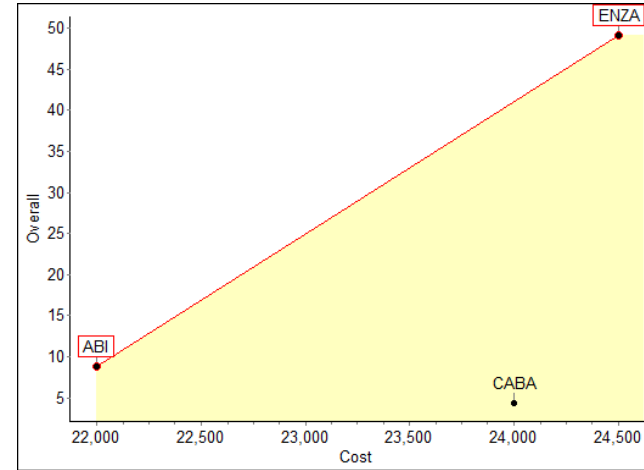
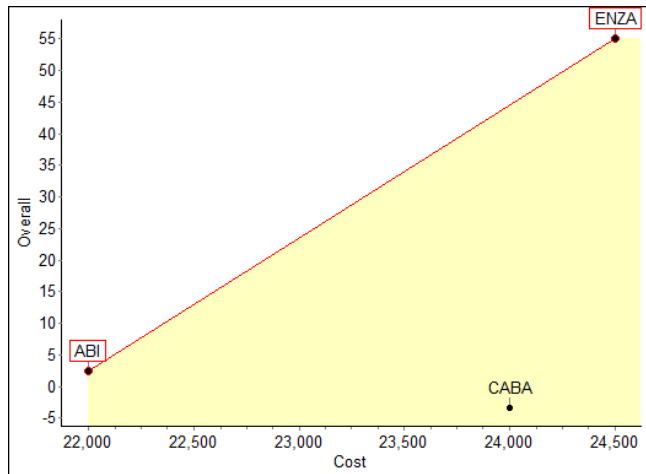


Figure 5: Cost benefit plots of treatments overall weighted preference value scores versus their purchasing costs across the four HTA settings (TLV top left, AETSA top right, AOTMiT, bottom left, INAMI bottom right).



Appendix

Tables and Figures

Table A1: Attributes definition and sources of evidence

Cluster	Attribute	Definition	Evidence source		
			Abiraterone	Cabazitaxel	Enzalutamide
THERAPEUTIC BENEFIT	Overall survival	The median time from treatment randomisation to death	de Bono et al 2011	de Bono et al 2010	Scher et al 2012
	Health related quality of life	Health related quality of life using the EQ-5D score	Sullivan et al 2007; TA 255; TA259; TA316	N/A – assumed, based on Sullivan et al 2007; TA 255; TA259; TA316	Sullivan et al 2007; TA 255; TA259; TA316
	Radiographic tumour progression	The median survival time on which patients have not experienced disease progression (using RECIST criteria)	de Bono et al 2011	de Bono et al 2010	Scher et al 2012
	PSA response	The proportion of patients having a $\geq 50\%$ reduction in PSA	Fizazi et al 2012	de Bono et al 2010	Scher et al 2013
SAFETY PROFILE	Treatment discontinuation	The proportion of patients discontinuing treatment due to AEs	de Bono et al 2011	de Bono et al 2010	Scher et al 2012
	Contra-indications	The existence of any type of contra-indication accompanying the treatment	EPAR, Prescribing info	EPAR, Prescribing info	EPAR, Prescribing info
INNOVATION LEVEL	ATC Level 1	The technology's relative market entrance in regards to its ATC Level 1 (Anatomical)	WHO ATC index	WHO ATC index	WHO ATC index
	ATC Level 2	The technology's relative market entrance in regards to its ATC Level 2 (Therapeutic)	WHO ATC index	WHO ATC index	WHO ATC index
	ATC Level 3	The technology's relative market entrance in regards to its ATC Level 3 (Pharmacological)	WHO ATC index	WHO ATC index	WHO ATC index
	ATC Level 4	The technology's relative market entrance in regards to its ATC Level 4 (Chemical)	WHO ATC index	WHO ATC index	WHO ATC index
	ATC Level 5	The technology's relative market entrance in regards to its ATC Level 5 (Molecular)	WHO ATC index	WHO ATC index	WHO ATC index

	Phase 1	The number of new indications for which the technology is investigated in Phase 1 clinical trials	ClinicalTrials.gov	ClinicalTrials.gov	ClinicalTrials.gov
	Phase 2	The number of new indications for which the technology is investigated in Phase 2 clinical trials	ClinicalTrials.gov	ClinicalTrials.gov	ClinicalTrials.gov
	Phase 3	The number of new indications for which the technology is investigated in Phase 2 clinical trials	ClinicalTrials.gov	ClinicalTrials.gov	ClinicalTrials.gov
	Marketing authorisation	The number of new indications that the technology has gained an approval for at the stage of marketing authorisation	ClinicalTrials.gov	ClinicalTrials.gov	ClinicalTrials.gov
	Delivery posology	The combination of the delivery system (RoA and dosage form) with the posology (frequency of dosing and duration of administration) of the treatment	EPAR, Prescribing info	EPAR, Prescribing info	EPAR, Prescribing info
	Special instructions	The existence of any special instructions accompanying the administration of the treatment	EPAR, Prescribing info	EPAR, Prescribing info	EPAR, Prescribing info
SOCIO-ECONOMIC IMPACT	Medical costs impact	The impact of the technology on direct medical costs excluding the purchasing costs of the technology*	BNF 69, Prescribing info, Connock et al 2011, Riemsa et al 2013, TA259	BNF 69, Prescribing info, de Bono et al 2010, TA255	BNF 69, TA316

Notes: * These costs include i) concomitant medications, ii) outpatient visits, diagnostic/laboratory tests, hospitalisations and other monitoring costs (including management AEs), and iii) terminal care.

Source: The authors.

Table A2: Pre-decision conference attribute reference levels and basis of selection

Cluster	Attribute name	Attribute metric	Lower level	Basis	Higher level	Basis
THERAPEUTIC BENEFIT	Overall survival	months	13.6	Best supportive care (BSC)	22.1	20% higher than the best performing option
	Health related quality of life	utility (EQ-5D)	0.72	Utility used for stable disease	0.82	Utility scores of general population
	Radiographic tumour progression	months	2.9	BSC	10.6	20% higher than the best performing option
	PSA response	% patients	1.5	BSC	64.8	20% higher than the best performing option
SAFETY PROFILE	Treatment discontinuation (% of patients)	% patients	10	BSC	0	Highest possible limit of the scale
	Contra-indications	types of contra-indications	Hypersensitivity + hepatic impairment + low neutrophil counts	Lowest possible limit of the scale	None known contraindications	Highest possible limit of the scale
INNOVATION LEVEL	ATC Level 1	relative market entrance	5 th	Lowest possible limit of the scale	1st	Highest possible limit of the scale
	ATC Level 2	relative market entrance	5 th	Lowest possible limit of the scale	1st	Highest possible limit of the scale
	ATC Level 3	relative market entrance	5 th	Lowest possible limit of the scale	1st	Highest possible limit of the scale
	ATC Level 4	relative market entrance	5 th	Lowest possible limit of the scale	1st	Highest possible limit of the scale
	ATC Level 5	relative market entrance	5 th	Lowest possible limit of the scale	1st	Highest possible limit of the scale

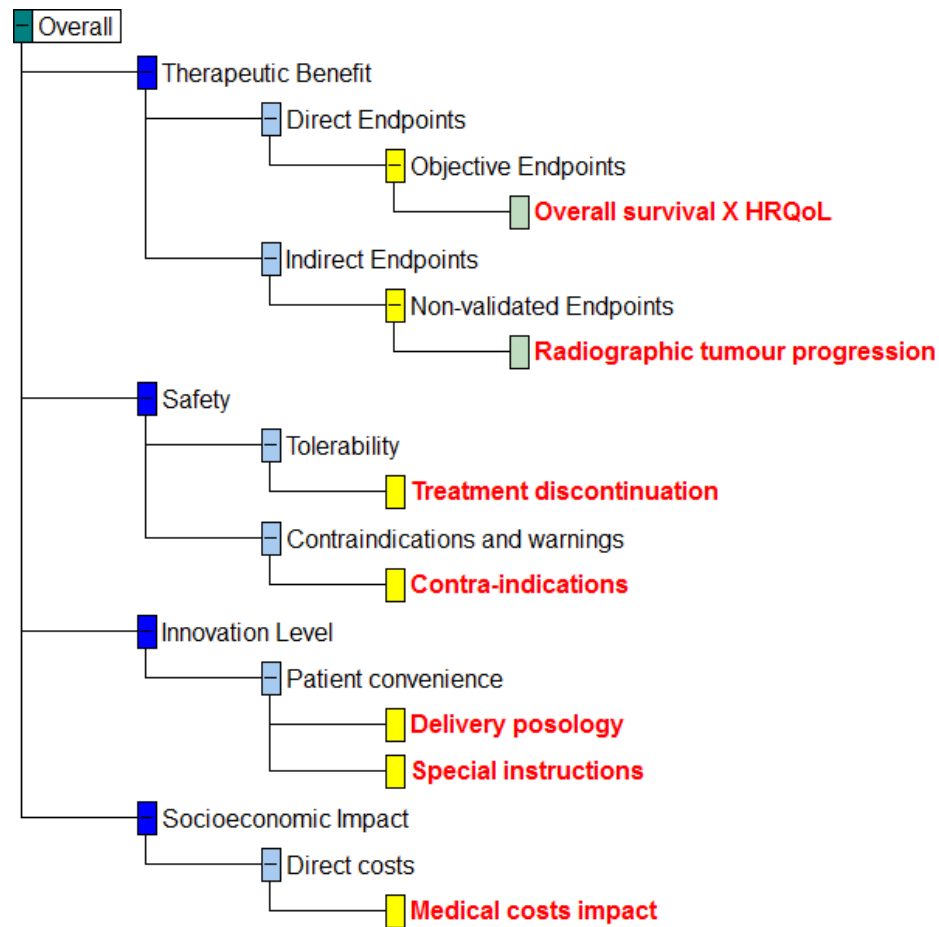
	Phase 1	number of new indications	0	Lowest possible limit of the scale	10	20% higher than the best performing option
	Phase 2	number of new indications	0	Lowest possible limit of the scale	16	20% higher than the best performing option
	Phase 3	number of new indications	0	Lowest possible limit of the scale	2	20% higher than the best performing option
	Marketing authorisation	number of new indications	0	Lowest possible limit of the scale	1	20% higher than the best performing option
	Delivery Posology	types of delivery system & posology combinations	Oral, every day - one off + IV, every 3 weeks - 1 hour*	Lowest possible limit of the scale	Oral, every day - one off*	Highest possible limit of the scale
	Special instructions	types of special instructions	No food + concomitant and/or pre-medication*	Lowest possible limit of the scale	None*	Highest possible limit of the scale
SOCIO-ECONOMIC IMPACT	Medical costs impact	GBP	10,000	20% higher than the worst performing option (rounded up)	0	BSC

Note: * Assuming no impact on Luteinizing hormone-releasing hormone (LHRH) analogue.

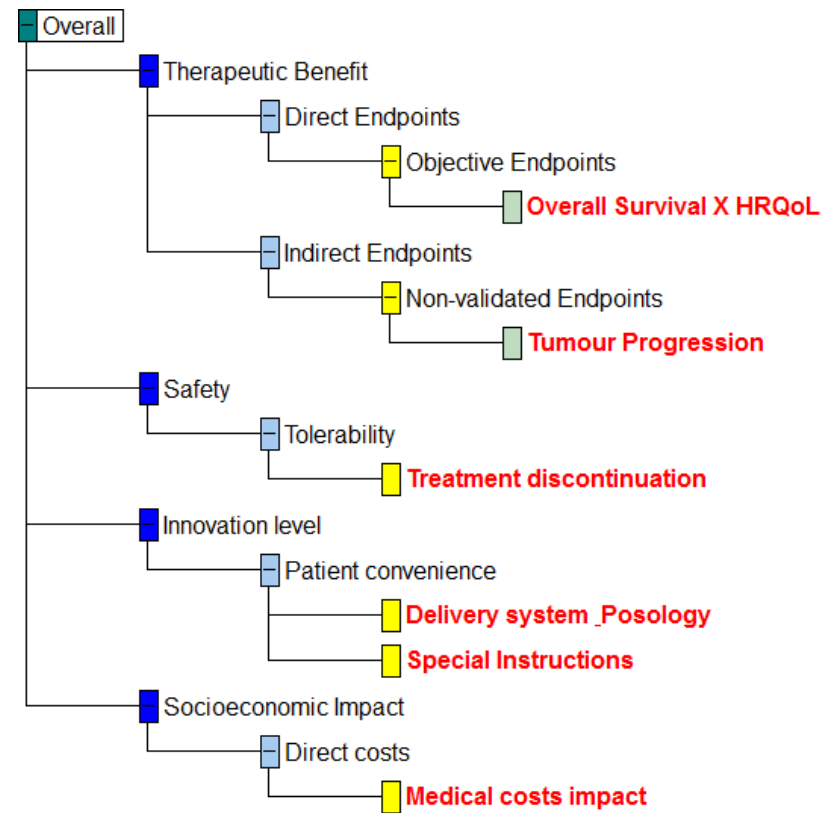
Source: The authors based on the literature.

Figure A1: Final value trees for metastatic prostate cancer across the four HTA agencies*

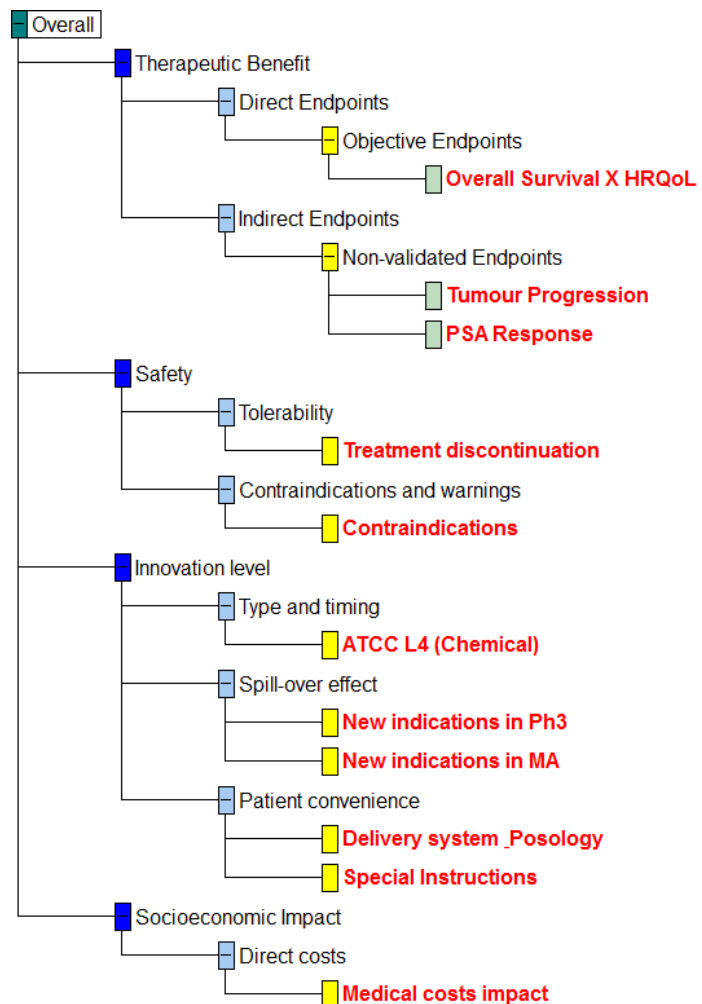
TLV



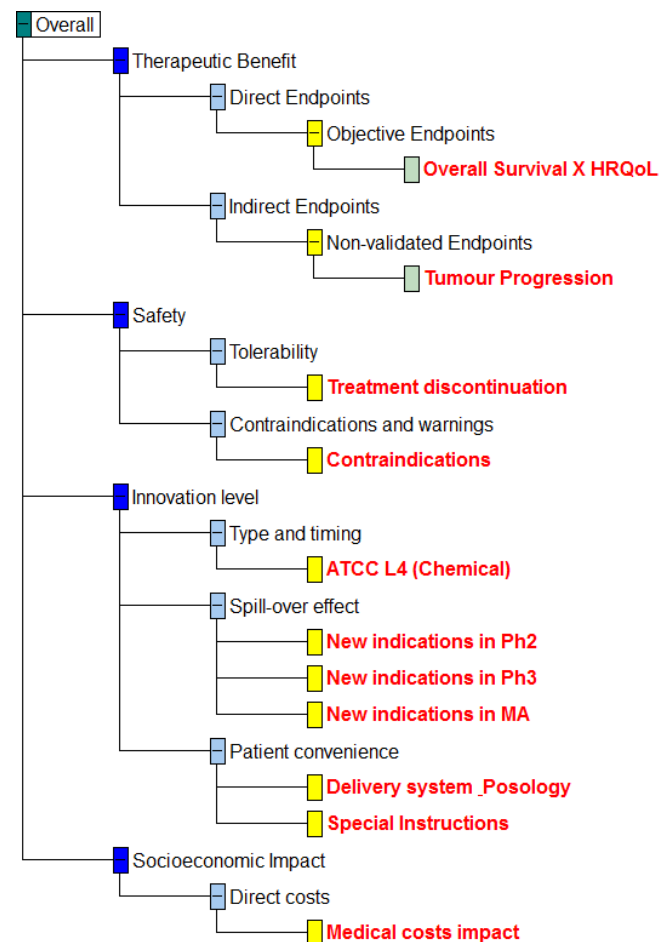
AOTMIT



AETSA

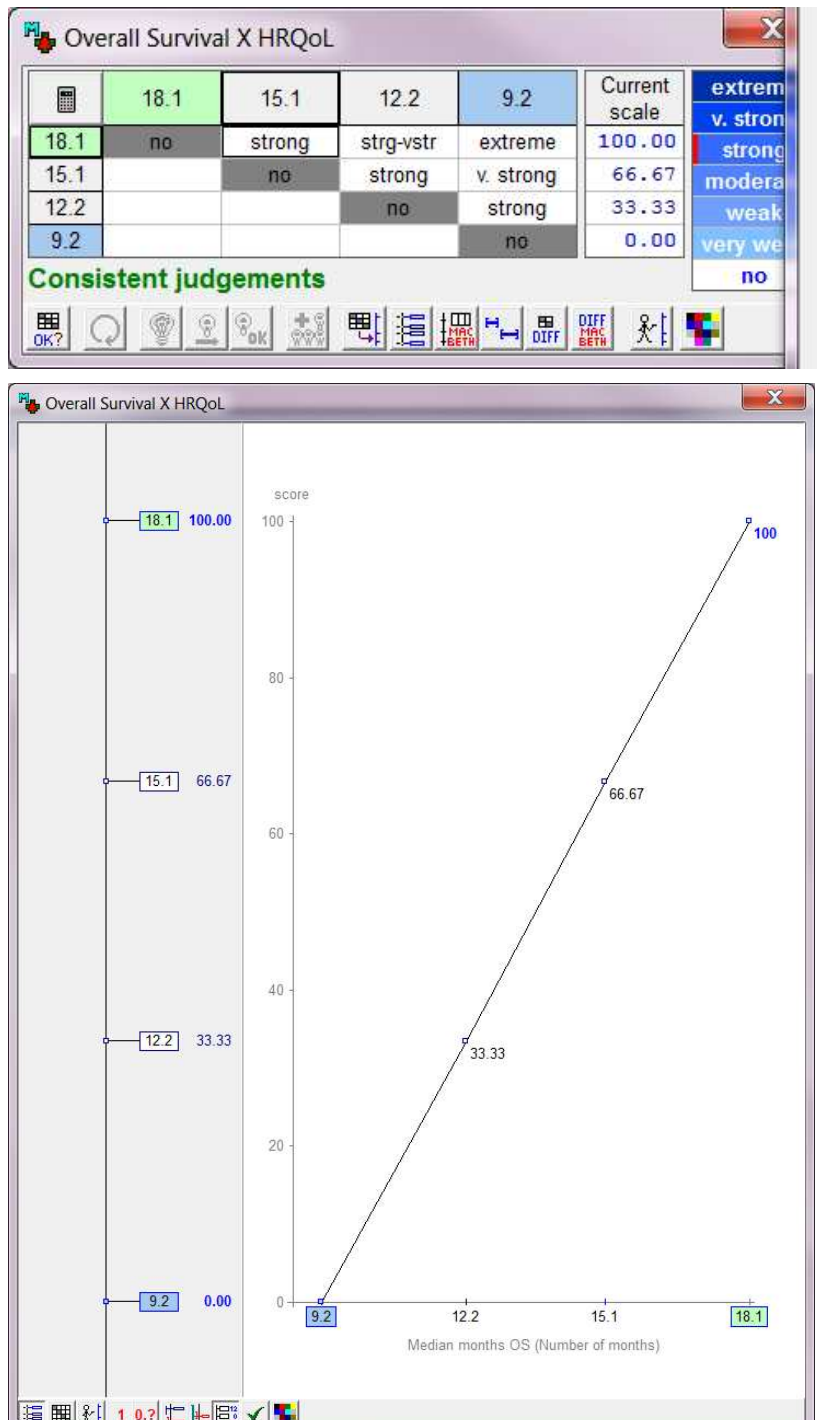


INAMI-RIZIV



* Images produced using the M-MACBETH (beta) software version 3.0.0

Figure A2: Example of value judgements matrix for the “Overall Survival x HRQoL” attribute measured in quality adjusted life months (*QALMs*) and its conversion into value functions (from the AOTMiT decision conference).



*Image produced using the M-MACBETH (beta) software version 3.0.0

Caption: In the Overall Survival x HRQoL attribute example, measured in quality adjusted life months (*QALMs*), the question asked was the following: “*What do you judge to be the difference of value between 9.2 and 18.1 QALMs? No difference, very weak, weak, moderate, strong, very strong, or extreme?*” Once a decision was reached (by consensus or majority voting), the next question came along: “*What do you judge to be the difference of value between 12.2 and 18.1 months QALMs? No difference, very weak, weak, moderate, strong, very strong, or extreme?*” The same process was followed until value judgments for all the different combinations of attribute levels were elicited, filling in the different rows from the right-hand side (i.e. lower range) to the left-hand side (i.e. higher range).

Ethics approval:

Ethics approval is not required for this paper as no personal or sensitive data from human subjects were collected.