



Università degli studi di Udine

Towards a Logic Programming Tool for Cancer Analysis

This is a pre print version of the following article:

*Original*

Towards a Logic Programming Tool for Cancer Analysis / Alice, Tarzariol; Agostino, Dovier; Alberto, Policriti. - ELETTRONICO. - 1949(2017), pp. 361-375.

*Availability:*

This version is available <http://hdl.handle.net/11390/1122838> since 2018-01-15T08:30:06Z

*Publisher:*

*Published*

DOI:

*Terms of use:*

The institutional repository of the University of Udine (<http://air.uniud.it>) is provided by ARIC services. The aim is to enable open access to all the world.

*Publisher copyright*

(Article begins on next page)

# Towards a Logic Programming tool for cancer analysis

Alice Tarzariol, Agostino Dovier, and Alberto Policriti

Dipartimento di Scienze Matematiche, Informatiche e Fisiche  
Università degli Studi di Udine, Italy

**Abstract.** The main goal of this work is to propose a pipeline capable to analyze a data collection of temporally qualified mutation profiles. A front-end and a post-processor are implemented. The front-end basically transforms the input data retrieved from medical databases into a set of facts. The post-processor allows the user to visualize the computation results as graphs. The reasoning core is based on Answer Set Programming that, allowing to deal with NP or  $\Sigma_2^P$  properties, is capable of deducing complex info from data. The system is modular: the reasoning tool can be replaced by any logic programming tool (e.g., Prolog, ILP) for different kinds of data analysis.

**Keywords:** Bioinformatics, Big Data Analysis, ASP.

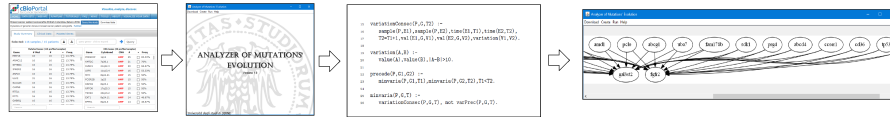
## 1 Introduction

Cancer is a game of mutation accumulation in the genome of an organism, that is the set of its genes. One of the most delicate aspects of this game is the fact that, while data is produced by today’s technology at an unprecedented rate it is, however, extremely noisy data. More precisely, even though *any* mutation might be triggering significant biological processes involved in cancer development, only some of the mutations are classified as “drivers” by experts in the field. The remaining—large—majority of mutations must be classified as “passengers”, as they are probably accumulating only as a consequence of the (devastating) side-effects of a more basic biological mechanisms already at work.

Tackling a classification problem for the collection of detected mutations in a tissue requires a multi-faceted approach. In order to build reliable filters suitable to define the *blueprint* of the disease under study it is essential to integrate different kinds of knowledge with brute force analysis.

The main goal of this work is to propose a pipeline capable to analyze a data collection of *temporally qualified* mutation profiles and synthesize a dependency graph. The information inferred by cross comparisons performed “locally” on input data, combined with information coming from external sources, allows to deduce the output graph that should represent a set of putative dependencies.

The idea is to prove the potential of the employed programming technologies in mining data that, when analyzed over realistically long sequences of time steps, rapidly generates (computationally) heavy collections of putative evolution paths.



**Fig. 1.** The proposed pipeline: from the WEB Database data is retrieved and converted into ASP input facts by a Java interface. The ASP engine elaborates the data. Finally, ASP output is converted in Graph format

ASP allows to encode compactly and elegantly problems that belongs to NP or even to  $\Sigma_2^P$ . Encodings allows great flexibility. In particular, the already implemented strategies are parametrically designed and use thresholds fixing conditions on the addition and removal of edges. Moreover, as it is always the case when dealing with biological data, many auxiliary problems—mostly related with I/O—must be faced. The data processing pipeline described here is interfaced with specialized languages and file formats for the field.

The paper is organized as follows: some related work are discussed in Sec. 2, the main ideas of our approach are presented in Sec. 3. A discussion on the data set used in the applications we have developed for extracting/converting data are presented in Sec. 4, and a brief analysis of the results found is shown in Sec. 5. Finally, some conclusions are drawn in Sec. 6.

## 2 Related work

Cancer data processing is, obviously, a very active and lively field in many scientific environments. We just report here on some work using techniques stemming from the logic programming community. We assume the reader has basic knowledge of logic and logic programming.

### 2.1 Inductive Logic Programming

Inductive Logic Programming (ILP) has been succesfully used to analyze big data and, in particular, to analyze cancer databases. ILP allows to implement a “structured” form of machine learning, by inferring a first-order theory that “explains” a set of input ground facts (defining extensionally, possibly partially, some predicates). Let us briefly recall the main ideas behind ILP (see, e.g., [14] or [5] for a more general approach to logical and relational learning).

Assume to have a set of positive observations  $O^+$  and a set of negative observations  $O^-$ . We can assume  $O^+$  and  $O^-$  as two disjoint sets of ground atoms. Assume also to have a (possibly empty) background theory (a set of clauses)  $P$ , the general inductive problem is the one of finding a set of hypotheses  $H$  (clauses) such that  $P \wedge H \models O^+$  and such that for all  $\ell \in O^-$  it holds that  $P \wedge H \not\models \ell$ .

Several strategies have been investigated and implemented by ILP systems to identify sets  $H$  made by clauses as much *general* as possible. Generality in this

context can be defined in terms of logical entailment: if clauses  $r_1$  and  $r_2$  are such that  $P \wedge r_1 \models r_2$ , then  $r_1$  is more general of  $r_2$  w.r.t.  $P$ .<sup>1</sup> Of course, this implies that if  $P \wedge r_2 \models B$  (for some  $B \subseteq O^+$ ), then  $P \wedge r_1 \models B$ . Therefore,  $r_1$  is preferable to  $r_2$  in the set of hypotheses  $H$  (unless it introduces errors, namely, there is  $\ell \in O^-$  s.t.  $P \wedge r_1 \models \ell$  and  $P \wedge r_2 \not\models \ell$ ).

In presence of large amounts of data that might be affected from errors since the beginning (e.g., medical data), one might accept a set  $H$  that might be not complete, namely it “explains” only a subset  $G$  (good) of  $O^+$ :  $P \wedge H \models G$ , and that might be not correct, namely such that  $P \wedge H \models \ell$  for  $\ell$  in a subset  $W$  (wrong) of  $O^-$ . In this case, of course, the ILP algorithm aims at maximizing the size of  $G$  and at minimizing the size of  $W$ .

An extension of ILP is PILP (probabilistic inductive logic programming) [18] where facts (observables) and rules (background theory) can be annotated by a probability. In this case a semantics based on probabilistic inference rules is adopted; the handling of uncertainty is therefore handled in a natural way.

Inductive Logic Programming has been used to analyze cancer data. For instance in [19] the authors participated to a world-wide carcinogenicity prediction competition (organized by the National Toxicology Program in the USA) using the ILP system Progol [15]. In the first round of the competition Progol produced the highest predictive accuracy of any automatic system participating the test.

Progol have been also used in [17] to infer general properties on pancreatic cancer and to allow early detection of this kind of cancer. For instance, one of the inductively computed predicates indicates that the measurements of *CEA* and *Elastase I* are very useful for the detection of pancreatic cancer.

They use data from 438 cases. The classification is based on the observation of the lymph node metastatic status and tumour differentiation status. They use a set of available lab data (e.g., CEA, CA19-9, Glucose, Elastase I, Serum Amylase, ...). After having identified the most promising features using standard feature selection criteria data is divided in three groups (low, normal, and high), as usually made by domain experts. Each patient record in the database represents an abnormality of a patient, split into positive and negative. Data is translated into a set of facts, and Progol is executed in order to rank the rules by their capability of explaining the database.

In [1] the authors propose an inductive logic programming approach to the problem of modeling evolution patterns for breast cancer. Data is obtained from a cohort of 124 patients at different stages. Data and *background knowledge* are expressed in logic programming. A set of hypotheses built on this knowledge using refinement, least general generalization, inverse resolution, and most specific corrections is computed.

---

<sup>1</sup> For definite clause programs  $\models$  is the standard logical consequence operator. For programs with negation its meaning becomes instead “if in every stable model of  $P \wedge r_1$ , the clause  $r_2$  holds.”

## 2.2 Answer Set Programming

Answer Set Programming has been used in many Bioinformatics applications. For an overview, see [4]. Among the many of them we would like to point out the approaches to the phylogenetic reconstruction (see, e.g., [7]). The origin of this work comes in fact from the desiderata, presented in the 2016 edition of the conference CILC [3], of extending the ASP tools for phylogenetic reconstruction to build a phylogenetic reconstruction of cancer evolution. If the cancer is related to a mutation of a gene present both in a human and in another specie (e.g., a rat), this may justify a deeper study of the evolution of that cancer in the other specie's body. Or, it can explain why a drug working for the other specie does not work in a human. Of course for doing that we need much more temporal data.

There are other approaches to biological data analysis based on ASP. For instance, in [12] the authors use ASP to compute an over-approximation of the set of Boolean networks which fit best with experimental data and provide the corresponding encodings. Of course, the approach can be used to explain the experimental data in a similar way as made by ILP systems, but using a purely logical approach, as we do in this paper.

## 3 The general approach

The general idea of our approach is simple and can be summarized as an attempt to infer a graph whose nodes are mutations (in genes) and whose arcs represent the causality relation between mutations, that is their most plausible (temporal) sequence. Data consists of variant allele frequencies and thresholds are used to accept/discard a temporal dependency.

The initial steps consist in a deterministic analysis followed by a blind search entirely left to the inferential engine. A pre-processing phase is necessary in order to correctly use data. A file MAF (Mutation Annotation Format) [13] is taken in input and the HUGO (HUMAN Genome Organization nomenclature) [11] gene symbols are extracted. In addition, experiment code, depth of readings—of the tumor sample in support of the variable allele—, and total depth of tumor sample readings, are uploaded. The ratio of the last two values is then calculated, thus obtaining variant allele frequencies for each sample gene. These numbers are then rounded to integers in order to make them more easily manageable by the ASP engine. The temporal information relative to the order of findings is extracted from the experiment identification, along with data that allow to classify cases—in the case of study, different organoid tissues. At this point the ASP system (groundersolver) Clingo [9] is executed on the data to compute the predicates explained in the code (according to the stable model semantics). A post-processing of Clingo's output is also implemented in order to depict the computation results as a graph.

### 3.1 ASP code over simple data

In this section we will see our ASP code and how it works over some example instances. Once filtered, input data is stored by two predicates, defined exten-

sionally by facts. Each patient can be involved in more than one experiment; the ternary predicate `sample` relates the patient ID with its experiment IDs and their times (given the experiments ordering), while the ternary predicate `val` reports each instance of mutated gene in any experiment and its value:

```

1 sample(pat-A,exp-A1,1). sample(pat-A,exp-A2,2). sample(pat-A,exp-A3,3).
2 sample(pat-A,exp-A4,4). sample(pat-B,exp-B1,1). sample(pat-B,exp-B2,2).
3 sample(pat-B,exp-B3,3). sample(pat-B,exp-B4,4). sample(pat-C,exp-C1,1).
4 sample(pat-C,exp-C2,2). sample(pat-C,exp-C3,3). sample(pat-C,exp-C4,4).
5 sample(pat-C,exp-C5,5).
6 val(exp-A1,a,3027). val(exp-A1,b,3027). val(exp-A2,a,1245).
7 val(exp-A2,c,3245). val(exp-A3,b,1245). val(exp-A3,c,1234).
8 val(exp-A4,a,2555). val(exp-A4,b,1324). val(exp-A4,c,1254).
9 val(exp-A4,d,1092). val(exp-B1,d,3027). val(exp-B1,e,3027).
10 val(exp-B2,d,1245). val(exp-B2,a,3245). val(exp-B3,a,1245).
11 val(exp-B3,b,1234). val(exp-B4,a,1334). val(exp-B4,b,1334).
12 val(exp-C1,b,3027). val(exp-C2,a,1245). val(exp-C2,b,3245).
13 val(exp-C3,a,1415). val(exp-C3,z,1415). val(exp-C4,z,9145).
14 val(exp-C4,d,9145). val(exp-C5,d,145).

```

In the first five rows we define the values of the samples of three patients `pat-A`, `pat-B` and `pat-C`, with the (temporal) order of samples. Subsequently, starting at row six, we indicate mutated genes values. For the sake of simplicity we named values `a,b,c ...`.

A number of “domain” predicates are computed from the input data (needed to limit the grounding), as follows:

```

1 gene(G) :- val(_,G,_).
2 patient(P) :- sample(P,_,_).
3 value(V) :- val(_,_,V).
4 time(0..T) :- sample(_,_,T), not sample(_,_,T+1).

```

The predicate `gene(G)` forces value `G` to appear as second parameter of `val` and, therefore, will hold for each gene uploaded in our samples. Analogous considerations for `patient(P)` and `value(V)`. `time(T)`, instead, is instanced with a list of values that goes from 0 to a maximum temporal value contained in the predicates `sample(P,E,T)`.

The following predicates allow to retrieve information from the data. Here we report some examples. Being ASP capable to deal with the whole NP class, and even the  $\Sigma_2^P$  class if disjunctive heads are used, one can capture/infer a wide family of properties. We will be back on this issue in Sec. 6.

```

1 variationConsec(P,G,T2) :-
2     sample(P,E1,T1),sample(P,E2,T2),T2=T1+1,
3     val(E1,G,V1),val(E2,G,V2),variation(V1,V2).
4 variation(A,B) :-
5     value(A),value(B),|A-B|>10.
6 precedes(P,G1,G2) :-
7     firstVariation(P,G1,T1),firstVariation(P,G2,T2),T1<T2.

```

```

8 firstVariation(P,G,T) :-
9     variationConsec(P,G,T), not varPrec(P,G,T).
10 varPrec(P,G,T) :-
11     variationConsec(P,G,T1), variationConsec(P,G,T),T1<T.
12 precedeKTimes(G1,G2,K) :-
13     gene(G1),gene(G2),K=#count{patient(P):precedes(P,G1,G2)}.
14 reachable(G1,G2) :-
15     maybePrecede(G1,G2).
16 reachable(G1,G2) :-
17     maybePrecede(G1,G3), reachable(G3,G2).
18 maybePrecede(G1,G2) :-
19     precedeKTimes(G1,G2,K1),precedeKTimes(G2,G1,K2),K1>K2,
20     not reachable(G2,G1).
21 #show maybePrecede/2.

```

variationConsec(P,G,T2) holds whenever for patient P there are two consecutive samples (the second one with time T2) with the same gene G with two values  $V_1$  e  $V_2$  that satisfy variation, namely their difference in absolute value is greater than 10. This value is the variation sensibility, in this simple example we choose a low parameter in order to take always the consecutive variation of the mutated gene. In our example we can infer that:

```

variationConsec(pat-A,a,2).  variationConsec(pat-A,c,3).
variationConsec(pat-A,b,4).  variationConsec(pat-A,c,4).
variationConsec(pat-B,d,2).  variationConsec(pat-B,a,3).
variationConsec(pat-B,a,4).  variationConsec(pat-B,b,4).
variationConsec(pat-C,b,2).  variationConsec(pat-C,a,3).
variationConsec(pat-C,z,4).  variationConsec(pat-C,d,5).

```

Before defining the main predicate of this example, namely, the predicate precedes(P,G1,G2), let us see the rule varPrec(P,G,T) at row 10: it holds if a gene G of a certain patient P has some variations in the sample at time T and also in some previous samples—and therefore it is *not* the first variation detected. This predicate is used in the body of firstVariation(P,G,T): the predicate inferring the minimum time T for which we have a variation of gene G in P.

At this point we can illustrate precedes(P,G1,G2) in row 6: it is true if, for a given patient, the first variation of G1 happens before the first variation of G2. Notice that it is not necessary to force that the two genes are different as it is implicit in the fact that the time of their first variation is different. So from the previous example we can infer:

```

precedes(pat-A,a,c).      precedes(pat-A,a,b).      precedes(pat-A,c,b).
precedes(pat-B,d,a).      precedes(pat-B,d,b).      precedes(pat-B,a,b).
precedes(pat-C,b,a).      precedes(pat-C,b,z).      precedes(pat-C,a,z).
precedes(pat-C,b,d).      precedes(pat-C,a,d).      precedes(pat-C,z,d).

```

Predicate precedeKTimes(G1,G2,K) stores in K the number of patients for which the gene G1 precedes the gene G2, while maybePrecede(G1,G2) is inferred

if G1 precedes G2 more frequently than G2 precedes G1. We report below the inferred instances of `precedeKTimes` in which the guard K is greater than zero (namely there is at least one case in which G1 precedes G2).

```
precedeKTimes(b,a,1). precedeKTimes(d,a,1). precedeKTimes(a,b,2).
precedeKTimes(c,b,1). precedeKTimes(d,b,1). precedeKTimes(a,c,1).
precedeKTimes(a,d,1). precedeKTimes(b,d,1). precedeKTimes(z,d,1).
precedeKTimes(a,z,1). precedeKTimes(b,z,1).
```

For a and b, there are 2 patients that exhibit variation of a before b, while just one exhibits the opposite variation. Hence we deduce `maybePrecede(a,b)`, while in case of a and d we don't infer any *most likely* order.

The final result in our example is:

```
maybePrecede(a,b).      maybePrecede(a,c).      maybePrecede(a,z).
maybePrecede(b,z).      maybePrecede(c,b).      maybePrecede(z,d).
```

In order to guarantee that the direct graph resulting from `maybePrecede` is acyclic, we check that the number of times the variation of the first gene precedes the second is strictly greater than the opposite order and we also impose the absence of the predicate `reachable` that states that already exists a path between two genes in the graph inferred.

The visualisation of the graph obtained is shown in Figure 2. The table containing the data (the variant allele frequency value has been omitted for simplicity) where the pairs highlighted are the variations revealed. In Figure 3 (right) it is shown the result automatically produced by Graphviz, after the Clingo output has been converted in the dot edges format Figure 3 (left).

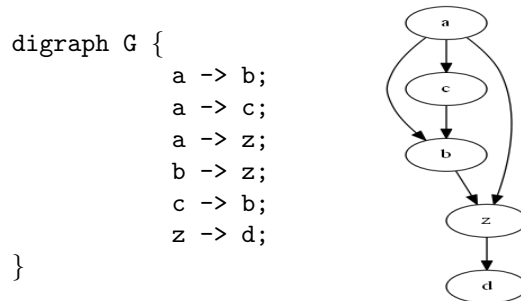
	Pat-A	Pat-B	Pat-C
1	a b	d e	b
2	a c	a d	a b
3	b c	a b	a z
4	a b c d	a b	d z
5	-	-	d

Fig. 2. Summarising table

## 4 Data sets

We tested our ASP program over the—real and publicly available—data set *Breast cancer patient xenografts* [6] available from <http://www.cbioportal.org/>.





**Fig. 3.** Graph from `maybePrecede`, in dot edges format (left), and as visualized by Graphviz (right)

This particular study was chosen as it provides different temporal analysis of the same patient’s cancer. The cBioPortal repository—as well as many other similar sites—is currently adding new case studies every year. In order to facilitate usage we, therefore, also implemented a (simple) graphical user-interface in Java providing functionalities to automate the data extraction, instances creation, as well as Clingo execution. Subsequently the ASP solver output is processed to display results by using the open source program GraphViz (Graph Visualization Software) [10].

The open source platform cBioPortal provides data from many cancer genomics data-sets (by now there are 150 of them) allowing their download, analysis, and visualisation. It was originally developed at Memorial Sloan Kettering Cancer Center (MSK) and now its software is available under an open source license via GitHub at <https://github.com/cBioPortal>. The portal is designed to bring cancer researchers near to the complexities of human genome data, allowing a rapid, intuitive, and precise visualization of expression profiles as well as clinical attributes from large-scale cancer genomics projects. A brief tutorial of the tools provided by cBioPortal is available on the web site; further documentation is available at <https://cbioportal.readthedocs.io/en/latest/README.html>.

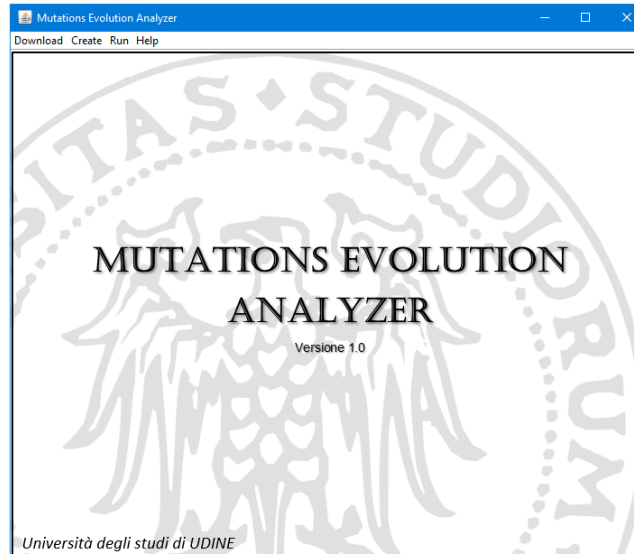
#### 4.1 Data download

From the website it is possible to query and/or download data from one or more studies satisfying given criteria (for instance, studies related to specific gene’s values) or it is possible to download the entire archive. From the home page it is also possible to query databases specified within studies, select the genetic profile, the single patient or a group of them and then specify collections of genes of interest. Users can submit even more specific queries by using the Onco Query Language (OQL), for which a brief guide is available on the website. The same information can be obtained through web API and library in R and Matlab. However, the faster way to visualize all genes from a study, consists in directly downloading the whole study database from <https://github.com/cBioPortal/datahub/tree/master/public>,



## 5 Mutations evolution analyzer

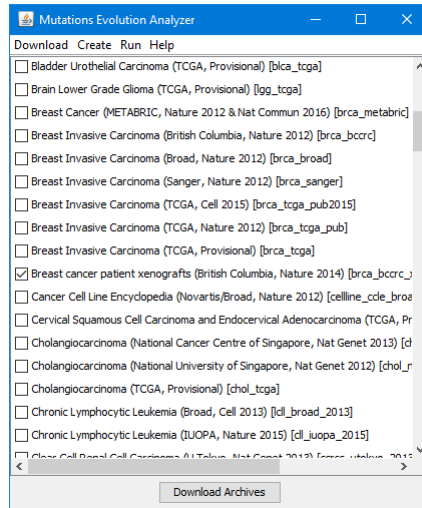
In order to get the complete archives stored in cBioPortal, automatically analyse data, and visualize the graph of variations found, we built a Java GUI called **Mutations Evolution Analyzer** (Figure 4) that contains a menu with four items: **Download**, **Create**, **Run**, and **Help**.



**Fig. 4.** Analyzer of Mutations' Evolution

The item **Download** sends a query to the web interface of the Cancer Genomic Data Server (CGDS) and, if available, returns an updated list of the studies, allowing the download over the site <https://github.com/cBioPortal/datahub/tree/master/public> of one or more of them in a user specified directory—see Figure 5.

Once the archive has been unzipped, we find *data\_mutations\_extended* file. By clicking on **Create** and then **Selected fields** the program allows to get, from a MAF file, the values of the selected fields (see Figure 6). To speed up selection, there is a checkbox called *Select standard field*, that automatically selects the fields used for our analysis: *Hugo\_Symbol*, *Tumor\_Sample\_Barcode*, *t\_alt\_count*, and *t\_depth*. In particular, for *Breast cancer patient xenografts*, the identification field is obtained by the union of the patient ID, the type of the sample (tumor or xenograft), the position in temporal order of the sample, the genotype studied, and whether it is relative to the whole genome or if it is targeted. For instance, sample *SA429X1A-targeted* is relative to patient *SA429*, it has been obtained over a xenograft (X means xenograft), is the first sample of this type (1)—as samples have a sequential number that identifies their temporal order—, applies



**Fig. 5.** Download of one ore more studies

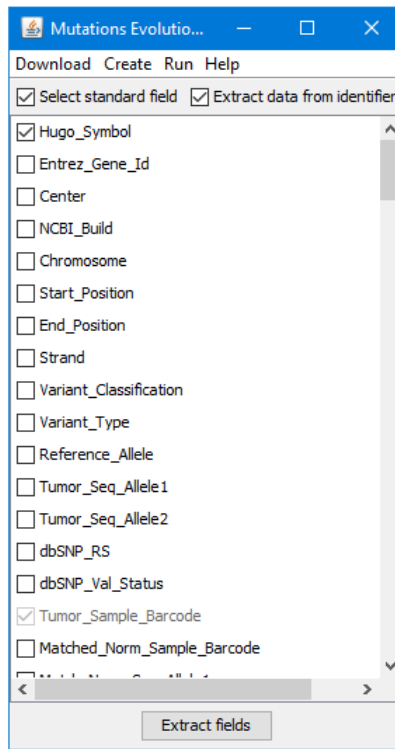
to the genotype  $A^2$ , for which it has been analyzed a particular subsequence of the genome (**targeted**). In order to get this single information directly available (instead of searching them every time from the *Tumor\_Sample\_Barcode*) it is sufficient to select the checkbox named *Extract data from identifier* that adds those new fields to the selected ones, with headers named *patientID*, *Type*, *Time*, *Genotype*, and *Wide*.

Subsequently our program creates a new file containing the selected fields, as well as a new one called *ValGene*, obtained by normalizing the variant allele frequency—namely the ratio between *t\_alt\_count* e *t\_depth*—, and then multiplies it by a reasonable value (we choose  $10^7$  in order to consider the first six digits after decimal point) rounding the result. We perform these additional operations to cast variant allele frequencies to integers, as it is easier to work with them in ASP with respect to the numbers in floating point.

Once the file containing the necessary fields has been obtained, by clicking on **Create** -> **ASP file**, a graphical interface asking the number of genes that we want analyze is shown. In Figure 7 the reader can see the default value. On click data are ordered, identifying the most frequent genes, taking the first *k* elements, with *k* as specified. An **Unknown** gene is a gene whose sequence is not known, so our program ignores it, as there is not necessarily a correlation between different unknown genes.

For the ASP syntax we use lowercase letters for all literal values, and replace characters “\_” and “.” by “\_”.

<sup>2</sup> For our usage it is sufficient to know that different letters correspond to different types. Refer to [6] for a full (and cleaner) explanation.



**Fig. 6.** Create → Selected fields

Then some Java code creates instances and adds them to the main ASP program under construction. Inessential modifications to adapt the analysis to the data under the study are performed at this stage in order, for example, to discriminate sample from same patient but with different genotypes: for each sample we add a fact for the type, wide, and genotype, in particular we consider the last one. Here there are the predicates deduced from sample **SA429X1A-targeted**:

```

1 sample(sa429,sa429x1a_targeted,1).
2 type(sa429x1a_targeted,x).
3 wide(sa429x1a_targeted,targeted).
4 genotype(sa429x1a_targeted,a).

```

As for the other rules we get:

```

1 variationConsec(P,GT,G,T2) :-
2     sample(P,E1,T1),sample(P,E2,T2),T2=T1+1,
3     val(E1,G,V1),val(E2,G,V2),variation(V1,V2),
4     genotype(E1,GT),genotype(E2,GT).
5 precedes(P,GT,G1,G2) :-

```





