





# Penalised regression with multiple sources of prior effects

Armin Rauschenberger<sup>1,\*</sup> , Zied Landoulsi<sup>1</sup> ,  
Mark A. van de Wiel<sup>2,3,†</sup> , Enrico Glaab<sup>1,†</sup> 

16 December 2022

<sup>1</sup>Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, Esch-sur-Alzette, Luxembourg. <sup>2</sup>Department of Epidemiology and Data Science (EDS), Amsterdam University Medical Centers (Amsterdam UMC), Amsterdam, The Netherlands <sup>3</sup>Medical Research Council Biostatistics Unit (MRC BSU), University of Cambridge, Cambridge, United Kingdom \*To whom correspondence should be addressed. <sup>†</sup>Mark A. van de Wiel and Enrico Glaab share senior authorship.

**In many high-dimensional prediction or classification tasks, complementary data on the features are available, e.g. prior biological knowledge on (epi)genetic markers. Here we consider tasks with numerical prior information that provide an insight into the importance (weight) and the direction (sign) of the feature effects, e.g. regression coefficients from previous studies. We propose an approach for integrating multiple sources of such prior information into penalised regression. If suitable co-data are available, this improves the predictive performance, as shown by simulation and application. The proposed method is implemented in the R package ‘transreg’ (<https://github.com/lcsb-bds/transreg>).**

**keywords:** transfer learning; co-data; prior information; ridge regression

## Background

For many biomedical prediction or classification studies, there is a previous study with a similar target and a similar high-dimensional feature space, e.g. hundreds of microRNAs (miRNAs), thousands of genes, or millions of single-nucleotide polymorphisms (SNPs). Given a trained model from a previous study, we could use it to obtain predicted values or predicted probabilities for the study of interest, but these predictions are only reliable if the two studies have the same target, the same features, and the same population. However, we expect the feature-target effects from two studies to be strongly correlated in more situations: slightly different targets (e.g. disease status vs disease stage), slightly different features (e.g. imperfectly overlapping feature space, different measurement technique), slightly different populations (e.g. hospitalised vs non-hospitalised patients), or even different modelling approaches (e.g. simple regression vs multiple regression). As it is challenging to estimate feature-target effects in high-dimensional settings, it might be advantageous to use results from previous studies as prior information for the study of interest.

Consider two prediction or classification problems, each one with a target vector and a feature matrix (samples in the rows, features in the columns). Suppose that both feature matrices cover the same features (each column in the first matrix corresponds to a column in the second matrix). In two special cases, the two problems reduce to a single problem: (i) If both problems have the same target and concern samples from the same population, they are in essence one ‘single-target’ problem (combine target vectors and feature matrices by rows, respectively), potentially with batch effects. (ii) If both problems concern the same samples, they are in essence one ‘multi-target’ problem (combine target vectors by columns, feature matrices are the same). In other cases, however, the two problems do not reduce to a single problem. Then we are in a potential transfer learning setting (Table 1).

In such settings - two or more regression problems with related targets and matched features - it might be possible to transfer information from one problem to another. If the regression problems are sufficiently

related to each other, we expect their regression coefficients to be correlated (positively or negatively). When fitting the regression model of interest, we could therefore account for the estimated regression coefficients from the other model. Transferring information on the importance and the direction of the feature effects, we could potentially increase the predictive performance.

Jiang et al. (2016) proposed the prior lasso to account for prior information in high-dimensional predictive modelling. Their method involves a preprocessing step and a weighting step. In the preprocessing step, the prior information is used to predict the target from the features. They present a solution for one set of prior effects from a closely related study (multiplying the feature matrix by the prior effects), but extensions to multiple sets of prior effects or loosely related studies may be feasible. Let  $\mathbf{y}$  represent the target and let  $\hat{\mathbf{y}}_{\text{prior}}$  represent the fitted values based on the prior information. In the weighting step, they minimise the penalised combined likelihood  $L(\mathbf{x}, \mathbf{y}; \boldsymbol{\beta}) + \eta L(\mathbf{x}, \hat{\mathbf{y}}_{\text{prior}}; \boldsymbol{\beta}) - \rho(\lambda; \boldsymbol{\beta})$  with respect to the coefficients  $\boldsymbol{\beta}$ , where  $\eta \geq 0$  (balance) and  $\lambda \geq 0$  (regularisation). If the balancing hyperparameter  $\eta$  is larger than zero, the prior predictions  $\hat{\mathbf{y}}_{\text{prior}}$  influence the estimation of the parameters  $\boldsymbol{\beta}$ .

Dhruba (2021) proposed a transfer learning method based on distribution mapping. Even if features or targets follow different distributions in two data sets, it is possible to build a predictive model using the first data set and make predictions for the second data set. Requiring matched features and targets in the source data set and unmatched features and targets in the target data set, their method transfers (i) features from the target to the source domain and (ii) predictions from the source to the target domain. By contrast, we consider transfer learning settings with matched features and targets in the target data set.

Tian and Feng (2022) proposed and implemented transfer learning for ridge and lasso regression. Their transfer learning algorithm involves two steps: (i) Estimating common coefficients for the target data set and the transferable source data sets ( $\hat{\boldsymbol{\omega}}$ ). (ii) Estimating the deviations from the common coefficients to the target coefficients ( $\hat{\boldsymbol{\delta}}$ ). Both steps together lead to the estimated target coefficients ( $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\omega}} + \hat{\boldsymbol{\delta}}$ ). Before applying their transfer learning algorithm, Tian and Feng (2022) apply a transferable source detection algorithm to exclude source data sets that are too different from the target data set. This avoids that non-transferable sources render the common coefficients misleading for the target data set ('negative transfer'). In the case of lasso regularisation in the two steps, there is sparsity in the common estimates as well as in the deviations from the common estimates to the target estimates (and thereby also in the target estimates).

The method from Tian and Feng (2022) requires not only the target data set but also the source data set(s). However, data protection regulations or restrictive data sharing policies might prevent researchers from accessing a source data set, or the available storage or processing capacity might be insufficient for analysing massive source data sets. There is therefore a need for transfer learning methods that do not require the source data but only the (anonymised) complementary data (co-data) derived from the source data. Such methods allow us to exploit summary statistics from external studies, e.g.  $p$ -values and effect sizes from a genome-wide association study (GWAS), to increase the predictive performance in the study of interest.

We propose a two-step transfer learning method, modelling with and without co-data in the first step and combining different models in the second step. Unless the source and target data sets are very similar, the coefficients from the source data set(s) will not fit well to the target data set. We therefore propose to calibrate these coefficients - preserving their signs and their order - so that they can be transferred from the source data set(s) to the target data set. Additionally, we also estimate the coefficients directly from the target data set, ignoring the co-data. The calibrated coefficients from the source data set(s) as well as the estimated coefficients from the target data set allow us to predict the outcome from the features. Finally, we combine the linear predictors from the models with and without co-data and calculate either predicted values (linear regression) or predicted probabilities (logistic regression).

In a related transfer learning setting, prior information is only available on the importance but not on the direction of the feature effects, i.e. with complementary data consisting of prior weights rather than prior effects. In the generalised linear model framework, the weighted lasso (Bergersen et al., 2011), the feature-weighted elastic net (Tay et al., 2022), and penalised regression with differential shrinkage (Zeng et al., 2021) account for prior weights in the penalty function, through feature-specific penalty factors or feature-specific regularisation parameters. Adaptive group-regularised ridge regression (van de Wiel et al., 2016) is not only applicable to categorical co-data but also to numerical co-data (prior weights), by the means

Table 1: Abstract representation of the data set of interest (without asterisk, black) and an additional data set (with asterisk, grey). Single-target learning (left): same targets, same features, different samples (from one population). Multi-target learning (centre): different targets, same features, same samples. Transfer learning (right): same or different targets, matched features, different samples (from one or two populations).

$$\begin{array}{ccc}
 \textit{Single-target learning} & \textit{Multi-target learning} & \textit{Transfer learning} \\
 \left( \begin{array}{c} y_1 \\ \vdots \\ y_n \\ y_1^* \\ \vdots \\ y_m^* \end{array} \right) \Leftarrow \left( \begin{array}{ccc} x_{11} & \cdots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{np} \\ x_{11}^* & \cdots & x_{1p}^* \\ \vdots & & \vdots \\ x_{m1}^* & \cdots & x_{mp}^* \end{array} \right) & \left( \begin{array}{cc} y_1 & y_1^* \\ \vdots & \vdots \\ y_n & y_n^* \end{array} \right) \Leftarrow \left( \begin{array}{ccc} x_{11} & \cdots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{np} \end{array} \right) & \left( \begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \right) \Leftarrow \left( \begin{array}{ccc} x_{11} & \cdots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{np} \end{array} \right) \\
 & & \Downarrow \quad \Downarrow \quad \Downarrow \\
 & & \left( \begin{array}{c} y_1^* \\ \vdots \\ y_m^* \end{array} \right) \Leftarrow \left( \begin{array}{ccc} x_{11}^* & \cdots & x_{1p}^* \\ \vdots & & \vdots \\ x_{m1}^* & \cdots & x_{mp}^* \end{array} \right)
 \end{array}$$

of creating groups of features from numerical co-data and forcing the group-penalties to be monotonically decreasing. An extension from van Nee et al. (2021) makes this approach even more suitable for numerical co-data. For single sources of co-data, it might be possible to extend these methods to prior information on the importance as well as the direction of feature effects by imposing sign constraints on the coefficients. Prior weights have not only been exploited in regression analysis, e.g. co-data moderated random forests (te Beest et al., 2017) adapt the sampling probabilities of the features to the prior weights.

## Method

### Model

Suppose one target and  $p$  features are available for  $n$  samples. We index the samples by  $i$  in  $\{1, \dots, n\}$  and the features by  $j$  in  $\{1, \dots, p\}$ . Our aim is to estimate the generalised linear model

$$\mathbb{E}[y_i] = h^{-1} \left( \beta_0 + \sum_{j=1}^p \beta_j x_{ij} \right).$$

For any sample  $i$ , the model expresses the expected value of its target ( $y_i$ ) as a function of its features ( $x_{i1}, \dots, x_{ip}$ ). The link function  $h(\cdot)$  depends on the family of distributions for the target (Gaussian: identity, binomial: logit, Poisson: log). In the linear predictor,  $\beta_0$  represents the unknown intercept, and  $\beta_j$  represents the unknown slope of feature  $j$  (i.e. the effect of the feature on the linear predictor of the target). Given the estimated intercept  $\hat{\beta}_0^*$  and the estimated slopes  $\{\hat{\beta}_1^*, \dots, \hat{\beta}_p^*\}$ , we could predict the target of previously unseen samples:

$$\hat{y}_i = h^{-1} \left( \hat{\beta}_0^* + \sum_{j=1}^p \hat{\beta}_j^* x_{ij} \right).$$

### Co-data

Suppose  $m$  sources of co-data are available, indexed by  $k$  in  $\{1, \dots, m\}$ . Let  $z_{jk}$  indicate the prior effect from source  $k$  for feature  $j$ . Our method is designed for quantitative co-data that provide an insight into the importance (absolute value) and the direction (sign) of the feature effects. Each set of prior effects ( $-1 \leq z_{ok} \leq 1$ ) is assumed to be positively correlated with the true coefficients ( $\text{cor}(z_{ok}, \beta) > 0$ ). For any source of co-data, the prior effects may be re-scaled (not re-centred), for example to the interval from  $-1$  to  $+1$ . In other words, the proposed method is invariant under multiplication of the prior effects by a positive scalar ( $z \rightarrow c \times z$  where  $c > 0$ ). We explain in the next section why this is important.

It might seem trivial to also allow for co-data that only provide an insight into the importance but not the direction of the feature effects (i.e. prior weights instead of prior effects). Each set of prior weights ( $0 \leq z_{ok} \leq 1$ ) is assumed to be positively correlated with the true absolute coefficients ( $\text{cor}(z_{ok}, |\beta|) > 0$ ). To obtain prior effects, one might want to assign the signs of the Spearman correlation coefficients between the target and the features to the prior weights. However, marginal effects and conditional effects can have opposite signs. If we wanted to extend our approach to prior weights, we would have to discover the signs inside the calibration procedure (see below), which would be related to high-dimensional regression with binary coefficients (Gamarnik and Zadik, 2017).

## Base-learners with co-data

Suppose we are in a transfer learning setting with two prediction or classification problems. For simplicity, we assume that the features do not differ in scale between the two problems. For illustration, we consider two artificial situations (where we would not use transfer learning in practice): (i) If both problems concern the same target on the same scale and the samples come from the same population, we could use the estimated regression coefficients from one problem to make predictions for the other problem. (ii) If the two problems concern the same target on different scales, we could also recycle the estimated regression coefficients, but we would have to adjust for the different scales.

When transferring estimated regression coefficients from one problem to another problem, it might not only be necessary to change their scale but it might also be beneficial to change their shape. For example, it might be that for one problem weak and strong effects matter, while for the other problem only strong effects matter. We should therefore also be able to make differences between small coefficients more or less important than those between large coefficients. We propose two calibration methods, namely exponential and isotonic calibration, to adapt the prior information to the data. For each source of co-data  $k$ , both calibration methods estimate the model

$$\mathbb{E}[y_i] = h^{-1} \left( \alpha_k + \sum_{j=1}^p \gamma_{jk} x_{ij} \right),$$

where the calibrated prior effects  $\{\hat{\gamma}_{1k}, \dots, \hat{\gamma}_{pk}\}$  depend on the initial prior effects  $\{z_{1k}, \dots, z_{pk}\}$ . The difference between exponential and isotonic scaling is how the former depend on the latter.

- exponential calibration: Let  $\gamma_{jk} = \theta_k \text{sign}(z_{jk}) |z_{jk}|^{\tau_k}$ , for  $j$  in  $\{1, \dots, p\}$ , where the factor  $\theta_k$  and the exponent  $\tau_k$  are non-negative real numbers ( $\theta_k \geq 0$  and  $\tau_k \geq 0$ ). We first fit one simple non-negative regression for different values of  $\tau_k$  (i.e. estimate  $\alpha_k$  and  $\theta_k$  given  $\tau_k$ ), and then optimise  $\tau_k$ . Once  $\theta_k$  and  $\tau_k$  have been estimated, the initial prior effects  $z_{jk}$  determine the final prior effects  $\hat{\gamma}_{jk} = \hat{\theta}_k \text{sign}(z_{jk}) |z_{jk}|^{\hat{\tau}_k}$ , for all  $j$  in  $\{1, \dots, p\}$ . The estimated factor  $\hat{\theta}_k$  and the estimated exponent  $\hat{\tau}_k$  allow the model to change the scale and the shape of the prior effects. For example,  $\hat{\theta}_k = 0$  sets them to zero,  $|\hat{\theta}_k| < 1$  makes them smaller,  $|\hat{\theta}_k| > 1$  makes them larger,  $\hat{\tau}_k = 0$  sets them to the same value,  $\hat{\tau}_k < 1$  makes (absolutely) large ones more similar, and  $\hat{\tau}_k > 1$  makes (absolutely) small ones more similar. If one or more sets of prior effects might be negatively associated with the true coefficients, we could remove the non-negativity constraints from the simple regressions (allowing  $\hat{\theta}_k < 0$  to invert the signs of the prior effects).
- isotonic calibration: We estimate  $\{\gamma_{1k}, \dots, \gamma_{pk}\}$  under the constraint that the signs of the initial prior effects  $z_{jk}$  determine the signs of the final prior effects  $\hat{\gamma}_{jk}$  (i.e.  $\hat{\gamma}_{jk} = 0 |z_{jk} = 0$ ,  $\hat{\gamma}_{jk} \geq 0 |z_{jk} > 0$ ,  $\hat{\gamma}_{jk} \leq 0 |z_{jk} < 0$ ) and under the constraint that the order of the initial prior effects determines the order of the final prior effects (i.e.  $\hat{\gamma}_{jk} \geq \hat{\gamma}_{lk} |z_{jk} \geq z_{lk}$ ,  $\hat{\gamma}_{jk} \leq \hat{\gamma}_{lk} |z_{jk} \leq z_{lk}$ ), for all  $j$  and  $l$  in  $\{1, \dots, p\}$ . If one or more sets of prior effects might be negatively associated with the true coefficients, we could fit each model with these constraints and the inverted constraints, and then select the better fit.

To make optimisation more efficient, we rewrite the sign- and order-constrained problem as a sign-constrained problem (see Table 4 in the Appendix). For each source of co-data, we order the columns of the feature matrix by increasing values of the prior effects. Suppose the first  $q$  columns correspond to negative prior effects and the last  $p - q$  columns correspond to non-negative prior effects. We take

the cumulative sum of the feature columns from left to right for the former (columns 1 to  $q$ ) and from right to left for the latter (columns  $p$  to  $q + 1$ ). We then estimate the coefficients on the left under non-positivity constraints, and those on the right under non-negativity constraints. Formally, the model equals

$$\mathbb{E}[y_i] = h^{-1} \left( \alpha_k + \sum_{j=1}^p \delta_{jk} w_{ij} \right),$$

where  $w_{ij} = \sum_{l=1}^j x_{i(l)}$  and  $\delta_{jk} \leq 0$  for  $j$  in  $\{1, \dots, q\}$ , and  $w_{ij} = \sum_{l=j}^p x_{i(l)}$  and  $\delta_{jk} \geq 0$  for  $j$  in  $\{q + 1, \dots, p\}$ , with the subscript within brackets indicating the order of the prior effects. The linear predictor of the sign-constrained model, i.e.  $\alpha_k + \sum_{j=1}^p \delta_{jk} w_{ij}$ , is equivalent to the linear predictor of the order-constrained model, i.e.  $\alpha_k + \sum_{j=1}^p \gamma_{(j)k} x_{i(j)}$ , because

$$\begin{aligned} \sum_{j=1}^q \delta_{jk} w_{ij} &= \sum_{j=1}^q \delta_{jk} \left( \sum_{l=1}^j x_{i(l)} \right) = \sum_{j=1}^q \left( \sum_{l=j}^q \delta_{lk} \right) x_{i(j)} = \sum_{j=1}^q \gamma_{(j)k} x_{i(j)}, \\ \sum_{j=q+1}^p \delta_{jk} w_{ij} &= \sum_{j=q+1}^p \delta_{jk} \left( \sum_{l=j}^p x_{i(l)} \right) = \sum_{j=q+1}^p \left( \sum_{l=q+1}^j \delta_{lk} \right) x_{i(j)} = \sum_{j=q+1}^p \gamma_{(j)k} x_{i(j)}. \end{aligned}$$

After estimating the coefficients of the sign-constrained model by maximum likelihood, we therefore estimate those of the order-constrained model by  $\hat{\gamma}_{(j)k} = \sum_{l=j}^q \hat{\delta}_{lk}$  for  $j$  in  $\{1, \dots, q\}$  and  $\hat{\gamma}_{(j)k} = \sum_{l=q+1}^j \hat{\delta}_{lk}$  for  $j$  in  $\{q + 1, \dots, p\}$ .

While exponential calibration involves three unknown parameters, namely the intercept  $\alpha_k$ , the factor  $\theta_k$  and the exponent  $\tau_k$ , isotonic calibration involves  $1 + p$  unknown parameters, namely the intercept  $\alpha_k$  and the slopes  $\gamma_k = \{\gamma_{1k}, \dots, \gamma_{pk}\}$ , for each set of co-data. Figure 1 shows the difference between exponential and isotonic calibration in several empirically assessed scenarios.

After calibration, we pre-assess the utility of each set of co-data. To do this, we calculate the residuals (depending on the family of distributions) between the fitted and the observed targets. We suggest to retain a set of co-data only if the residuals are significantly smaller than those of the intercept-only model (one-sided Wilcoxon signed-rank test) at the nominal 5% level ( $p$ -value  $\leq 0.05$ ).

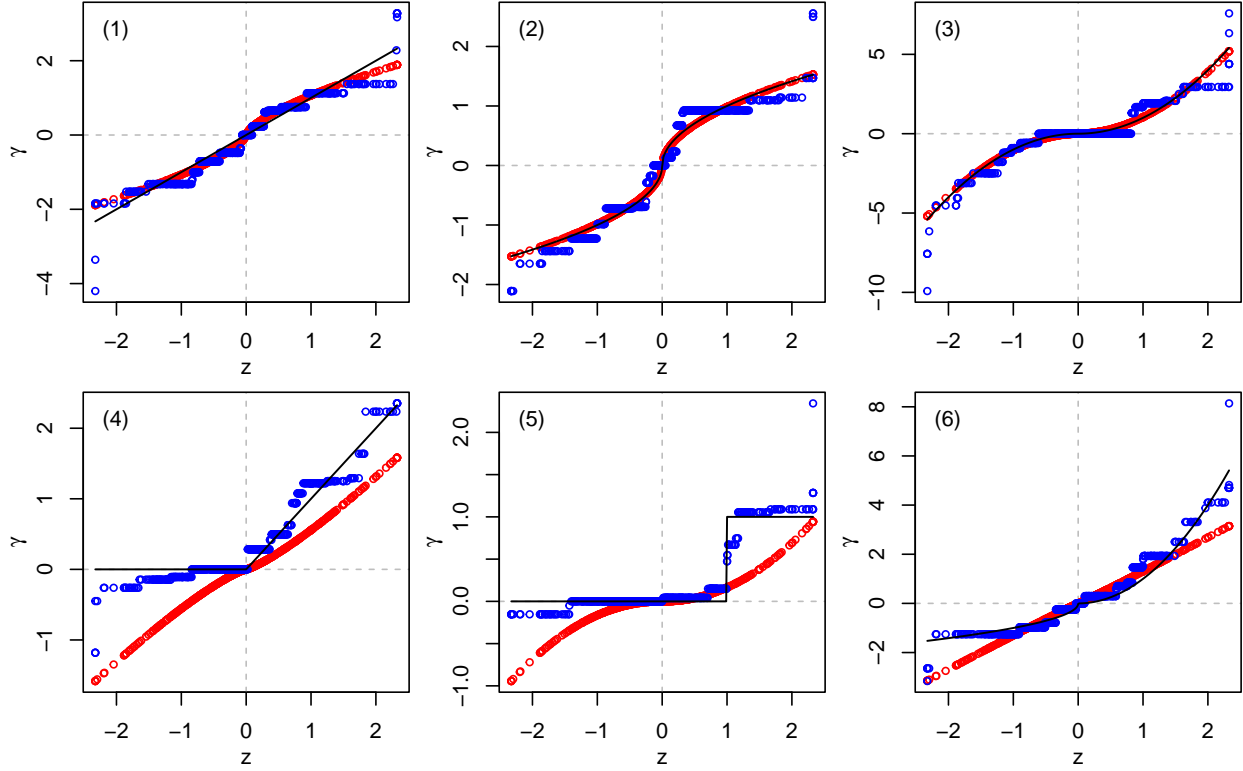


Figure 1: Final prior effects  $\gamma$  ( $y$ -axis) against initial prior effects  $z$  ( $x$ -axis), under exponential calibration (red) and isotonic calibration (blue). The black line corresponds to perfect calibration ( $\gamma = \beta$ ). We simulated the feature matrix  $\mathbf{X}$  from a standard Gaussian distribution ( $n = 200$ ,  $p = 500$ ) and the initial prior effects  $z$  from a trimmed standard Gaussian distribution (trimmed below the 1% and above the 99% quantile). We set the true coefficients to (1)  $\beta = z$ , (2)  $\beta = \text{sign}(z)\sqrt{|z|}$ , (3)  $\beta = \text{sign}(z)z^2$ , (4)  $\beta = \mathbb{I}[z > 0]z$ , (5)  $\beta = \mathbb{I}[z > 1]$ , or (6)  $\beta = -\mathbb{I}[z \leq 0]\sqrt{|z|} + \mathbb{I}[z > 0]z^2$ . And we simulated the response vector  $\mathbf{y}$  from Gaussian distributions with the means  $\boldsymbol{\eta}$  and the variance  $\text{Var}(\boldsymbol{\eta})$ , where  $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$ . While exponential calibration performs slightly better in the first three scenarios (top), isotonic calibration performs much better in the last three scenarios (bottom).

## Base-learners without co-data

We also fit the model without any co-data. We estimate the coefficients by maximising the penalised likelihood:

$$\hat{\boldsymbol{\beta}} = \operatorname{argmax}_{\boldsymbol{\beta}} \{L(\mathbf{x}; \boldsymbol{\beta}) - \rho(\lambda; \boldsymbol{\beta})\} ,$$

where  $L(\mathbf{x}, \boldsymbol{\beta})$  is the likelihood and  $\rho(\lambda; \boldsymbol{\beta})$  is the penalty. The likelihood depends on the family of distributions (Gaussian, binomial, Poisson), and the penalty can be the ridge ( $L_2$ ) or the lasso ( $L_1$ ) penalty. The penalty shrinks the squared (ridge) or absolute (lasso) slopes  $\{\beta_1, \dots, \beta_p\}$  towards zero (without penalising the intercept  $\beta_0$ ). We denote the estimated intercept by  $\hat{\beta}_0$  and the estimated slopes by  $\{\hat{\beta}_1, \dots, \hat{\beta}_p\}$ .

## Cross-validation

We split the samples into ten folds to perform 10-fold internal cross-validation. In each iteration, we fit the models to nine included folds and predict the target for the excluded fold.

Let the  $n \times m$  matrix  $\hat{\mathbf{H}}^{(0, cv)}$  represent the feature-dependent part of the cross-validated linear predictors from the models with co-data. Specifically, the entry in row  $i$  (sample) and column  $k$  (source of co-data) equals

$$\eta_{ik}^{(0, cv)} = 0 \times \hat{\alpha}_k^{-\kappa(i)} + \sum_{j=1}^p \hat{\gamma}_{jk}^{-\kappa(i)} x_{ij} ,$$

where the superscript  $-\kappa(i)$  indicates that the (ignored) intercept  $\alpha_k$  and the slopes  $\gamma_{jk}$  for  $j$  in  $\{1, \dots, p\}$  are estimated without using the fold of sample  $i$ , as in Rauschenberger and Glaab (2021).

The models without any co-data do not only have  $1 + p$  unknown parameters, namely the intercept  $\beta_0$  and the slopes  $\{\beta_1, \dots, \beta_p\}$ , but also the unknown hyperparameter  $\lambda$ . In each iteration, we fit this model for a decreasing sequence of 100 values for the regularisation parameter  $\lambda$ , indexed by  $l$  in  $\{1, \dots, 100\}$ , using the computationally efficient approach from Friedman et al. (2010, `glmnet`).

Accordingly, let the  $n \times 100$  matrix  $\hat{\mathbf{H}}^{(1, cv)}$  represent the cross-validated linear predictors from the model without co-data. Specifically, the entry in row  $i$  (sample) and column  $l$  (regularisation parameter) equals

$$\eta_{il}^{(1, cv)} = \hat{\beta}_0^{-\kappa(i), l} + \sum_{j=1}^p \hat{\beta}_j^{-\kappa(i), l} x_{ij} ,$$

where the superscripts  $-\kappa(i)$  and  $l$  indicate that the intercept  $\beta_0$  and the slopes  $\{\beta_1, \dots, \beta_p\}$  are estimated without using the fold of sample  $i$  and given the regularisation parameter  $\lambda_l$ .

To optimise the predictive performance of the co-data independent model, we would select the  $\lambda$  that minimises the cross-validated loss ( $\lambda_{\min}$ ). As we base our predictions not only on the co-data independent model but also on the co-data dependent model(s),  $\lambda_{\min}$  might be too small. The reason is that the co-data might be informative to the extent that the co-data independent model requires more penalisation. We could let the meta-learner select the optimal  $\lambda$  from the whole sequence, but this might render the inclusion and exclusion of co-data dependent models unstable. Our ad-hoc solution is to include the optimal regularisation parameter for the co-data independent model ( $\lambda_{\min}$ ) and a slightly larger one ( $\lambda_{1se}$ ). The latter is given by the one-standard-error rule, which increases  $\lambda$  until the cross-validated loss equals its minimum plus one standard error.

We concatenate  $\hat{\mathbf{H}}^{(0, cv)}$  with the columns of  $\hat{\mathbf{H}}^{(1, cv)}$  that correspond to  $\lambda_{\min}$  and  $\lambda_{1se}$  to obtain the  $n \times (m + 2)$  matrix  $\hat{\mathbf{H}}^{(cv)}$ . The first  $m$  columns correspond to the models with co-data, and the last two columns correspond to the model without co-data.

## Meta-learner

We combine the base-learners with and without co-data by stacked generalisation (Wolpert, 1992), on the level of the linear predictors (Rauschenberger et al., 2021). In the meta-layer, we regress the target on the

cross-validated linear predictors from the base-layer:

$$\mathbb{E}[y_i] = h^{-1} \left( \omega_0 + \sum_{k=1}^{m+2} \omega_k \hat{H}_{ik}^{(cv)} \right).$$

Leaving the intercept unrestricted ( $-\infty < \omega_0 < +\infty$ ) but imposing the lower bound zero on the slopes ( $\omega_1 \geq 0, \dots, \omega_{m+2} \geq 0$ ), we estimate these coefficients under lasso regularisation. Due to the feature selection property of the lasso, a source of co-data can be excluded ( $\hat{\omega}_k = 0$ ) or included ( $\hat{\omega}_k > 0$ ), where  $k$  in  $\{1, \dots, m\}$ . Similarly, the models without co-data can be excluded ( $\hat{\omega}_k = 0$ ) or included ( $\hat{\omega}_k > 0$ ), where  $k = m + 1$  for the model with  $\lambda_{\min}$  and  $k = m + 2$  for the model with  $\lambda_{\text{lse}}$ . The estimated slopes function as weights for the co-data dependent models ( $\hat{\omega}_1, \dots, \hat{\omega}_m$ ) and for the co-data independent models ( $\hat{\omega}_{m+1}, \hat{\omega}_{m+2}$ ). Thus, we do not only select sources but also weight them according to their relevance.

## Interpretation

The coefficients  $\hat{\beta}_{\min}$  and  $\hat{\beta}_{\text{lse}}$  give insight into the feature-target effects estimated without co-data, with  $\hat{\beta}_{\min,j}$  and  $\hat{\beta}_{\text{lse},j}$  representing the effect of feature  $j$ , where  $j$  in  $\{1, \dots, p\}$ . The coefficients  $\hat{\omega}$  give insight into the importance of the sources of co-data, with  $\hat{\omega}_k$  representing the importance of source  $k$ , where  $k$  in  $\{1, \dots, m\}$ . For a previously unseen sample  $i$ , the predicted value is:

$$\hat{y}_i = h^{-1} \left( \hat{\omega}_0 + \sum_{k=1}^{m+2} \hat{\omega}_k \hat{H}_{ik} \right) = h^{-1} \left( \hat{\beta}_0^* + \sum_{j=1}^p \hat{\beta}_j^* x_{ij} \right),$$

where  $\hat{\beta}_0^* = \hat{\omega}_0 + \hat{\omega}_{m+1} \hat{\beta}_{\min,0} + \hat{\omega}_{m+2} \hat{\beta}_{\text{lse},0}$

and  $\hat{\beta}_j^* = \left( \sum_{k=1}^m \hat{\omega}_k \hat{\gamma}_{jk} \right) + \hat{\omega}_{m+1} \hat{\beta}_{\min,j} + \hat{\omega}_{m+2} \hat{\beta}_{\text{lse},j}$ .

Thus, the estimated effect for a feature ( $\hat{\beta}_j^*$ ) is a weighted sum of estimated coefficients with co-data ( $\hat{\gamma}_{j1}, \dots, \hat{\gamma}_{jm}$ ) and the estimated coefficients without co-data ( $\hat{\beta}_{\min,j}, \hat{\beta}_{\text{lse},j}$ ).

Sparse models (few non-zero coefficients) are often considered to be more interpretable than dense models (many non-zero coefficients). While the original coefficients are dense ( $\sum_{j=1}^p \mathbb{I}[\hat{\beta}_j \neq 0] = p$ ) or sparse ( $\sum_{j=1}^p \mathbb{I}[\hat{\beta}_j \neq 0] \ll p$ ) depending on the choice between ridge and lasso regularisation, the weights may contain some zeros due to significance filtering or lasso regularisation ( $\sum_{k=1}^{m+2} \mathbb{I}[\hat{\omega}_k \neq 0] \leq m + 2$ ). As soon as one set of dense prior effects is selected, however, the combined coefficients also become dense ( $\sum_{j=1}^p \mathbb{I}[\hat{\beta}_j^* \neq 0] \lesssim p$ ). This means that the feature selection property of the lasso is not maintained. We should therefore choose between ridge and lasso regularisation (i) to make the model without co-data more predictive or interpretable (ii) or to make the model with co-data more predictive (iii) but not to make the model with co-data more interpretable.

## Extension

In some applications, prior information might be reliable for some features but unreliable for other features. Although the base-learners with co-data might still be predictive, the meta-learner (weighted average of the base-learners with and without co-data) might be not more predictive than the base-learner without co-data. The reason is that the meta-learner assigns the same weight to all prior effects, rather than more weight to reliable prior effects and less weight to unreliable prior effects. The same problem occurs if prior information is available for some features but missing for other features. We therefore propose an alternative approach for applications with partially informative sources of co-data.

In the following, we use the term ‘meta-features’ for the cross-validated linear predictors from the base learners with co-data. Each meta-feature - one column of the  $n \times m$  matrix  $\hat{H}^{(0,cv)}$  - corresponds to one



source of co-data. In the meta-layer, we regress the target on the meta-features and the base-features:

$$\mathbb{E}[y_i] = h^{-1} \left( \beta_0 + \sum_{k=1}^m \omega_k \hat{H}_{ik}^{(0, cv)} + \sum_{j=1}^p \beta_j x_{ij} \right),$$

with non-negativity constraints for the weights for the meta-features ( $\omega_1 \geq 0, \dots, \omega_m \geq 0$ ) but without constraints for the intercept ( $\beta_0$ ) and the slopes for the base-features ( $\beta_1, \dots, \beta_p$ ).

We estimate the weights for the meta-features and the slopes for the base-features using penalised maximum likelihood:

$$\{\hat{\omega}, \hat{\beta}\} = \underset{\{\omega, \beta\}}{\operatorname{argmax}} \{L(\mathbf{x}; \omega, \beta) - \rho(\lambda; \beta)\},$$

where  $L(\mathbf{x}; \omega, \beta)$  is the likelihood and  $\rho(\lambda; \beta)$  is the penalty. We do not penalise the weights for the meta-features ( $m \ll n$ ) but only the slopes for the base-features ( $p \gg n$ ). The more sources of co-data are available, the more it becomes necessary to penalise their weights. But then the weights  $\omega$  and the slopes  $\beta$  might need differential penalisation, for example a lasso penalty for the meta-features (selection of sources) and a ridge penalty for the base-features (many small effects). To make this computationally efficient, we would need a fast cross-validation procedure for multiple penalties (cf. van de Wiel et al., 2021) with non-negativity constraints (meta-features) and mixed lasso and ridge penalisation (meta-features vs base-features). This extension is therefore only applicable in settings with few sources of co-data.

The predicted value for a previously unseen sample  $i$  is

$$\hat{y}_i = h^{-1} \left( \hat{\beta}_0 + \sum_{k=1}^m \hat{\omega}_k \hat{H}_{ik}^{(0)} + \sum_{j=1}^p \hat{\beta}_j x_{ij} \right) = h^{-1} \left( \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j^* x_{ij} \right),$$

where  $\hat{\beta}_j^* = \left( \sum_{k=1}^m \hat{\omega}_k \hat{\gamma}_{jk} \right) + \hat{\beta}_j$ .

As the coefficients  $\beta$  are shrunk towards zero but the coefficients  $\omega$  are not penalised, the combined coefficients  $\beta^*$  are shrunk towards the calibrated prior effects. When the regularisation parameter tends to infinity ( $\lambda \rightarrow \infty$ ), the estimated deviations from the calibrated prior effects approach zero ( $\hat{\beta} \rightarrow 0$ ) and the combined estimates approach the calibrated prior effects ( $\hat{\beta}^* \rightarrow \hat{\gamma}$ ). Lasso regularisation ensures sparsity in the deviations from the calibrated prior effects ( $\sum_{j=1}^p \mathbb{I}[\hat{\beta}_j \neq 0] \ll p$ ) - in contrast to ridge regularisation - but not in the combined coefficients ( $\sum_{j=1}^p \mathbb{I}[\hat{\beta}_j^* \neq 0] \lesssim p$ ). As the combined coefficients may deviate more from unreliable than from reliable calibrated prior effects, this extension is suitable for partially informative co-data. As opposed to ‘standard stacking’, we refer to this extension as ‘simultaneous stacking’.

## Simulation

We performed two simulation studies to compare the predictive performance between our transfer learning method and the one from Tian and Feng (2022). In contrast to the method from Tian and Feng (2022), which requires the feature-target effects in the target and the source data set(s) to be *positively correlated* and on the *same scale* (i.e.  $\beta_{\text{target}} \approx \beta_{\text{source}}$ ), our method also allows for negatively correlated effects and for effects on different scales (i.e.  $\beta_{\text{target}} \approx c \times \beta_{\text{source}}$ ). Although it is possible to overcome this restriction by inverting the target (Gaussian:  $\mathbf{y}_{\text{source}} \rightarrow -\mathbf{y}_{\text{source}}$ , binomial:  $\mathbf{y}_{\text{source}} \rightarrow 1 - \mathbf{y}_{\text{source}}$ ), by re-scaling the target (Gaussian:  $\mathbf{y}_{\text{source}} \rightarrow 1/c \times \mathbf{y}_{\text{source}}$ ), or by inverting or re-scaling the features ( $\mathbf{X}_{\text{source}} \rightarrow c \times \mathbf{X}_{\text{source}}$ ), we believe it is more user-friendly to directly allow for negative correlations and different scales. To ensure a fair comparison between the two methods, we simulate positively correlated effects on the same scale. Furthermore, although the method from Tian and Feng (2022) is in theory also suitable for mixed response types, the current version of the related R package `glmtrans` requires the source data set(s) and the target data set to have the same response type (Gaussian, binomial, Poisson). We therefore always simulate the same response type in the source and target domains.

## External simulation

We use the simulation approach from Tian and Feng (2022). In each iteration, we call the function `glmtrans::models` with the arguments (1) family of distributions: `family="gaussian"` (default) or `family="binomial"`, (2) source or target data sets: `type="all"` (default), (3) difference between source and target coefficients: `h=5` (default) or `h=250`, (4) number of source data sets: `K=5` (default), (5) sample size for target data set: `n.target=100` (default), (6) sample size for each source data set: `n.source=150` (default), (7) number of non-zero coefficients: `s=15` (default) or `s=50`, (8) number of features: `p=1000` (default), number of transferable source data sets: `Ka=1`, `Ka=3` or `Ka=K=5` (default).

The simulation from Tian and Feng (2022) involves the following steps:

- **Features:** The correlation between features  $i$  and  $j$  is set to  $\Sigma_{ij} = 0.5^{|i-j|}$ , where  $i$  and  $j$  in  $\{1, \dots, p\}$ . Let  $\Sigma$  represent the correlation matrix and let  $\Sigma = \mathbf{R}^T \mathbf{R}$  represent its Cholesky decomposition, where  $\mathbf{R}$  is an upper triangular matrix. For the target data set ( $n_0 = 100$ ) and each source data set ( $n_1 = \dots = n_5 = 150$ ), the  $n_0 \times p$  matrix  $\mathbf{X}_0 = \mathbf{E}_0 \mathbf{R}$  and the  $n_k \times p$  matrices  $\mathbf{X}_k = \mathbf{E}_k \mathbf{R}$  for  $k$  in  $\{1, \dots, 5\}$  represent the features, where the  $n_0 \times p$  matrix  $\mathbf{E}_0$  and the  $n_k \times p$  matrices  $\mathbf{E}_k$  contain Gaussian noise.
- **Coefficients:** Let  $\beta_j$  represent the effect of feature  $j$ , for  $j$  in  $\{1, \dots, p\}$ , and denote the  $p$ -dimensional coefficient vectors by  $\beta_0$  for the target data set and  $\{\beta_1, \dots, \beta_5\}$  for the source data sets. For the *target* data set, the first  $s$  elements are set to  $\beta_j = 0.5$  (causal) and the last  $p - s$  elements are set to  $\beta_j = 0$  (non-causal). For *transferable source* data sets, the first  $s$  elements are set to  $\beta_j = 0.5 + (-1)^{z_j} h/p$  and the last  $p - s$  elements are set to  $\beta_j = (-1)^{z_j} h/p$ , where  $z_j$  is a realisation of  $z_j \sim \text{Bernoulli}(0.5)$ . For *non-transferable source* data sets, the first  $s$  elements are set to  $\beta_j = (-1)^{z_j} 2h/p$  (non-causal), the next  $s$  elements are set to  $\beta_j = 0.5 + (-1)^{z_j} 2h/p$  (causal), and the last  $p - 2s$  elements are a random sample of  $s$  causal and  $p - 3s$  non-causal elements generated in the same way. To obtain a non-transferable source, it would be sufficient to randomly select causal elements rather than inverting causal and non-causal elements (indices 1 to  $2s$ ).
- **Targets:** In the *Gaussian* case, the target vector is the  $n$ -dimensional vector  $\mathbf{y}_0 = \mathbf{X}_0 \beta_0 + \epsilon_0$  for the target data set, and the  $n$ -dimensional vector  $\mathbf{y}_k = 0.5 + \mathbf{X}_k \beta_k + \epsilon_k$  for the source data sets, where the  $n$ -dimensional vectors  $\{\epsilon_0, \dots, \epsilon_5\}$  contain Gaussian noise. In the *binomial* case, let  $\mathbf{p}_0 = 1/(1 + \exp(-\mathbf{X}_0 \beta_0))$  for the target data set and  $\mathbf{p}_k = 1/(1 + \exp(-0.5 - \mathbf{X}_k \beta_k))$  for the source data sets. The  $n$ -dimensional vectors  $\{\mathbf{y}_0, \dots, \mathbf{y}_5\}$  are the target vectors, with each element following a Bernoulli distribution with the probability given by  $\{\mathbf{p}_0, \dots, \mathbf{p}_5\}$ .

## Internal simulation

The simulation from Tian and Feng (2022) uses the same effect size for all causal features and a decreasing correlation structure with a fixed base. We therefore designed our own data-generating mechanism (i) to simulate different effect sizes for different causal features ( $\beta_j \in \mathbb{R}$  instead of  $\beta_j \in \{0, 0.5\}$ ) and (ii) to vary the strength of correlation between features ( $\sigma_{ij} = \rho_x^{|i-j|}$  instead of  $\sigma_{ij} = 0.5^{|i-j|}$ ).

Our simulation involves the following steps:

- **Features:** Setting the mean of feature  $i$  to  $\mu_i = 0$ , the variance of feature  $i$  to  $\sigma_{ii} = 1$ , and the covariance between features  $i$  and  $j$  to  $\sigma_{ij} = \rho_x^{|i-j|}$ , for all  $i$  and  $j$  in  $\{1, \dots, p\}$ , we simulate multiple feature matrices from the multivariate Gaussian distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\Sigma$ , namely the  $n_0 \times p$  feature matrix  $\mathbf{X}_0$  for the target data set, and the  $n_{1/2/3} \times p$  feature matrices  $\{\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3\}$  for the source data sets ( $n_0 = 100, n_1 = n_2 = n_3 = 150, p = 1000$ ).
- **Coefficients:** Setting the mean and the variance for data set  $k$  to  $\mu_k = 0$  and  $\sigma_{kk} = 1$ , for all  $k$  in  $\{0, 1, 2, 3\}$ , and the covariance between data sets  $k$  and  $l$  to  $\sigma_{kl} = 0$  if either  $k$  or  $l$  equals 1, or to  $\sigma_{kl} = \rho_\beta$  if both  $k$  and  $l$  are in  $\{0, 2, 3\}$ , we simulate two  $p \times 4$  matrices from the multivariate Gaussian distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\Sigma$ , namely  $\mathbf{B}_1$  and  $\mathbf{B}_2$ . We define the coefficients as  $\mathbf{B} = \mathbf{B}_1 \mathbb{I}[\mathbf{B}_2 > \phi^{-1}(1 - \pi)]$ , where  $\phi$  is the Gaussian cumulative distribution function and  $\pi$  equals 0.2 (dense) or 0.05 (sparse), and denote the  $p$ -dimensional coefficient vectors by  $\beta_0$  for the target data

set and by  $\{\beta_1, \beta_2, \beta_3\}$  for the source data sets. While one set of coefficients is non-transferable ( $\beta_1$ ), we transform the transferable sets of coefficients with  $\beta_2 \rightarrow \text{sign}(\beta_2)|\beta_2|^2$  and  $\beta_3 \rightarrow \text{sign}(\beta_3)\sqrt{|\beta_3|}$ .

- **Targets:** For the target data set, we compute  $\mathbf{z}_0 = \mathbf{X}_0\beta_0$  and standardise  $\mathbf{z}_0$  to obtain  $\mathbf{z}_0^*$ . For the source data sets, we proceed similarly to obtain  $\{\mathbf{z}_1^*, \mathbf{z}_2^*, \mathbf{z}_3^*\}$ . The simulated targets equal  $\mathbf{y}_k = h^{-1}(\sqrt{w}\mathbf{z}_k^* + \sqrt{1-w}\epsilon_k)$ , where  $h(\cdot)$  is a link function and  $\epsilon_k$  follows a standard Gaussian distribution, for  $k$  in  $\{0, \dots, 3\}$ . Given  $0 \leq w \leq 1$ , we have  $\text{Var}(\sqrt{w}\mathbf{z}_k^* + \sqrt{1-w}\epsilon_k) = w\text{Var}(\mathbf{z}_k^*) + (1-w)\text{Var}(\epsilon_k) = 1$ . While  $h(\cdot)$  is the identity link in the Gaussian case, it is the logit link in the binomial case, where the simulated probabilities are rounded to simulated classes.

## Simulation results

In addition to the target data set, the method from Tian and Feng (2022) requires the source data sets, while our method requires the prior effects derived from the source data sets. Since the two methods have different requirements, we first simulate the source data sets (for the competing method) and then derive the prior effects from the simulated source data sets (for the proposed method). As prior effects for the proposed method, we use the estimated regression coefficients from the source data sets. We choose the type of regularisation for both methods subject to the simulation setting, namely ridge regularisation for dense settings (**glmtrans**:  $\alpha = 0$ ; **transreg**:  $\alpha_{\text{source}} = \alpha_{\text{target}} = 0$ ) and lasso or lasso-like elastic net regularisation for sparse settings (**glmtrans**:  $\alpha = 1$ ; **transreg**:  $\alpha_{\text{source}} = 0.95, \alpha_{\text{target}} = 1$ ). The idea of the lasso-like elastic net regularisation is to render the prior information more stable.

In each simulation setting, we simulate 100 training samples and 10 000 testing samples (hold-out) for the target data set. Tables 2 and 3 show the testing loss in the external and internal simulation, under exponential and isotonic calibration. We observe that transfer learning with **glmtrans** and **transreg** leads to an improvement with respect to **glmnet**. Comparing the two different calibration approaches and the two different stacking approaches, we do not observe systematic differences.

Table 2: Testing loss in external simulation, as a percentage of the one from prediction by the mean. Settings: number of transferable source data sets ( $K_a$ ), differences between source and target coefficients ( $h$ ), dense setting with ridge regularisation ( $s = 50$ ,  $\alpha = 0$ ) or sparse setting with lasso regularisation ( $s = 15$ ,  $\alpha = 1$ ), family of distribution ('gaussian' or 'binomial'). These parameters determine (i) the average Pearson correlations among the features in the target data set ( $\bar{\rho}_x$ ) and (ii) the maximum Pearson correlation between the coefficients in the target data set and the coefficients in the source data sets ( $\max(\hat{\rho}_\beta)$ ). Methods: regularised regression (**glmnet**), competing transfer learning method (**glmtrans**), proposed transfer learning method (**transreg**) with exponential/isotonic calibration and standard/simultaneous stacking. In each setting, the colour black highlights methods that are more predictive than regularised regression without transfer learning (**glmnet**), and an underline highlights the most predictive method.

$K_a$	$h$	$\alpha$	family	$\bar{\rho}_x$	$\max(\hat{\rho}_\beta)$	glmnet	glmtrans	exp.sta	exp.sim	iso.sta	iso.sim
1	5	0	gaussian	0.01	1.00	73.6	50.3	37.3	35.7	31.6	<u>29.7</u>
3	5	0	gaussian	0.01	1.00	"	38.9	20.0	18.2	15.1	<u>13.6</u>
5	5	0	gaussian	0.01	1.00	"	23.7	13.8	12.7	10.3	<u>9.5</u>
1	250	0	gaussian	0.01	0.40	"	<u>39.5</u>	61.7	65.2	59.7	58.4
3	250	0	gaussian	0.01	0.40	"	<u>36.0</u>	48.0	46.0	42.6	42.0
5	250	0	gaussian	0.01	0.40	"	34.1	38.4	36.7	30.4	<u>27.8</u>
1	5	1	gaussian	0.01	1.00	23.3	<u>12.7</u>	14.3	16.3	13.8	14.5
3	5	1	gaussian	0.01	1.00	"	<u>9.7</u>	11.0	11.6	10.2	10.0
5	5	1	gaussian	0.01	1.00	"	<u>9.6</u>	10.8	11.4	10.0	9.8
1	250	1	gaussian	0.01	0.26	"	58.4	19.3	28.4	<u>18.8</u>	28.7
3	250	1	gaussian	0.01	0.28	"	34.5	<u>18.7</u>	28.6	18.8	32.1
5	250	1	gaussian	0.01	0.28	"	34.5	<u>18.8</u>	29.3	19.2	35.9
1	5	0	binomial	0.01	1.00	93.0	91.7	79.4	81.3	<u>76.0</u>	76.0
3	5	0	binomial	0.01	1.00	"	85.2	<u>63.7</u>	65.9	64.7	65.0
5	5	0	binomial	0.01	1.00	"	79.4	62.2	64.8	<u>61.9</u>	63.1
1	250	0	binomial	0.01	0.39	"	89.5	90.6	89.7	87.4	<u>85.4</u>
3	250	0	binomial	0.01	0.43	"	88.6	80.5	82.3	<u>75.8</u>	76.4
5	250	0	binomial	0.01	0.44	"	84.8	82.5	84.7	77.8	<u>77.5</u>
1	5	1	binomial	0.01	1.00	77.4	73.4	65.6	70.9	<u>65.1</u>	70.5
3	5	1	binomial	0.01	1.00	"	63.9	60.5	60.8	<u>58.3</u>	61.0
5	5	1	binomial	0.01	1.00	"	60.2	57.3	58.4	57.6	<u>55.8</u>
1	250	1	binomial	0.01	0.20	"	85.5	<u>77.1</u>	77.4	77.1	77.4
3	250	1	binomial	0.01	0.25	"	79.6	<u>77.1</u>	77.4	77.1	77.4
5	250	1	binomial	0.01	0.26	"	76.7	77.1	77.4	<u>76.5</u>	82.3

Table 3: Testing loss in internal simulation, as a percentage of the one from prediction by the mean. Settings: correlation parameter for features ( $\rho_x$ ), correlation parameter for coefficients ( $\rho_\beta$ ), dense setting with ridge regularisation ( $\pi = 20\%$ ,  $\alpha = 0$ ) or sparse setting with lasso regularisation ( $\pi = 5\%$ ,  $\alpha = 1$ ), family of distribution ('gaussian' or 'binomial'). These parameters determine (i) the average Pearson correlations among the features in the target data set ( $\bar{\rho}_x$ ) and (ii) the maximum Pearson correlation between the coefficients in the target data set and the coefficients in the source data sets ( $\max(\hat{\rho}_\beta)$ ). Methods: regularised regression (**glmnet**), competing transfer learning method (**glmtrans**), proposed transfer learning method (**transreg**) with exponential/isotonic calibration and standard/simultaneous stacking. In each setting, the colour black highlights methods that are more predictive than regularised regression without transfer learning (**glmnet**), and an underline highlights the most predictive method.

$\rho_x$	$\rho_\beta$	$\alpha$	family	$\bar{\rho}_x$	$\max(\hat{\rho}_\beta)$	glmnet	glmtrans	exp.sta	exp.sim	iso.sta	iso.sim
0.95	0.70	0	gaussian	0.04	0.44	<u>76.5</u>	79.7	82.5	79.5	82.5	79.5
0.99	0.70	0	gaussian	0.18	0.43	<u>63.8</u>	68.3	67.2	65.7	65.1	64.4
0.95	0.85	0	gaussian	0.04	0.62	78.0	87.9	75.6	75.6	75.0	<u>74.5</u>
0.99	0.85	0	gaussian	0.18	0.64	63.7	64.2	63.1	<u>63.0</u>	63.0	63.2
0.95	0.99	0	gaussian	0.04	0.87	70.3	74.1	65.9	<u>65.5</u>	65.6	66.2
0.99	0.99	0	gaussian	0.18	0.87	62.7	64.4	58.4	59.7	<u>58.1</u>	61.1
0.95	0.70	1	gaussian	0.04	0.28	78.6	<u>78.3</u>	78.6	78.4	78.6	78.4
0.99	0.70	1	gaussian	0.18	0.14	59.4	<u>59.3</u>	59.0	59.4	<u>58.5</u>	59.3
0.95	0.85	1	gaussian	0.04	0.39	73.4	<u>70.7</u>	71.8	72.3	72.5	72.9
0.99	0.85	1	gaussian	0.18	0.58	65.0	<u>65.9</u>	66.2	65.6	64.2	<u>62.4</u>
0.95	0.99	1	gaussian	0.04	0.87	88.2	78.0	74.8	78.4	<u>73.8</u>	76.2
0.99	0.99	1	gaussian	0.18	0.77	65.0	62.8	61.5	60.4	60.2	<u>59.8</u>
0.95	0.70	0	binomial	0.04	0.36	90.2	98.4	90.9	93.0	<u>90.0</u>	91.4
0.99	0.70	0	binomial	0.18	0.39	84.9	90.0	85.7	<u>84.2</u>	86.1	84.7
0.95	0.85	0	binomial	0.04	0.61	90.4	<u>89.3</u>	90.2	90.4	90.3	90.2
0.99	0.85	0	binomial	0.18	0.54	83.3	<u>82.7</u>	83.8	83.0	84.0	83.8
0.95	0.99	0	binomial	0.04	0.88	94.9	<u>92.7</u>	91.5	93.0	<u>90.7</u>	92.5
0.99	0.99	0	binomial	0.18	0.89	78.4	78.2	77.3	77.8	<u>77.2</u>	78.4
0.95	0.70	1	binomial	0.04	0.54	94.0	90.4	89.8	<u>89.4</u>	90.2	89.6
0.99	0.70	1	binomial	0.18	0.29	92.0	92.5	94.3	92.0	92.1	<u>90.3</u>
0.95	0.85	1	binomial	0.04	0.28	105.5	104.4	<u>96.6</u>	97.6	96.6	97.6
0.99	0.85	1	binomial	0.18	0.57	86.8	<u>85.1</u>	85.9	87.5	86.2	85.5
0.95	0.99	1	binomial	0.04	0.92	100.0	94.4	90.1	88.2	87.3	<u>86.8</u>
0.99	0.99	1	binomial	0.18	0.82	89.9	89.7	89.4	86.8	86.5	<u>86.1</u>

# Applications

## External applications

First, we consider an adapted version of the application on cervical cancer from van de Wiel et al. (2016). The aim is to transfer information from a methylation study with biopsy samples to another methylation study with self-collected samples in order to better discriminate between low-grade and high-grade precursor lesions. Specifically, we transfer the signs of the effect sizes and the  $p$ -values from the source data set to the target data set ( $n = 44$  samples,  $p = 9491$  features). We then examine whether transfer learning increases the predictive performance of ridge regression, which is more predictive than lasso regression in this application. For comparison, we consider the methods from Tay et al. (2022, `fwelnet`) and van Nee et al. (2021, `ecpc`). While the proposed method exploits information on the importance and direction of the effects (co-data:  $-\text{sign}(\text{coef}) \log_{10}(p\text{-value})$ ), the other two methods only exploit information on their importance (co-data:  $-\log_{10}(p\text{-value})$ ). After 10 repetitions of 10-fold cross-validation, we observe that transfer learning (not with exponential but with isotonic calibration) often increases the predictive performance of ridge regression (`transreg.exp.sta`: 0/10, `transreg.exp.sim`: 4/10, `transreg.iso.sta`: 7/10, `transreg.iso.sim`: 10/10, `fwelnet`: 7/10, `ecpc`: 5/10). We also observe that exploiting information on the importance as well as the direction of the effects (`transreg`) is more beneficial than exploiting information on the importance of the effects only (`fwelnet`, `ecpc`), as can be seen in the mean change in cross-validated loss (`transreg.exp.sta`: +6.56%, `transreg.exp.sim`: +0.43%, `transreg.iso.sta`: -2.50%, `transreg.iso.sim`: -9.25%, `fwelnet`: -0.12%, `ecpc`: -1.52%). Here, isotonic calibration outperforms exponential calibration, and simultaneous stacking outperforms standard stacking. A potential explanation for the large difference in performance between exponential and isotonic calibration is that positive effects might be more important than negative effects in this application, for a biological reason (methylation increases the probability of cancer) and a statistical reason (effects of overexpression are easier to detect than those of underexpression). While exponential calibration behaves symmetrically for negative and positive prior effects, isotonic calibration can shrink negative prior effects towards zero.

Second, we consider an adapted version of the application on pre-eclampsia from Tay et al. (2022). Measurements of  $p = 1125$  plasma proteins are available for  $n = 166$  patients at multiple time points ( $48 \times 2 + 8 \times 3 + 20 \times 4 + 74 \times 5 + 16 \times 6 = 666$ ). The aim is to transfer information from late time points (gestational age  $> 20$  weeks) to early time points (gestational age  $\leq 20$  weeks). We repeatedly split the patients into one source data set and one target data set. Patients with only late time points are always in the source data set, and other patients are randomly allocated to the source and the target data set. (Note that this application is somewhat artificial, as it might be better to drop transfer learning in favour of using all earliest time points in the regression of interest.) Using the source data set, we estimate two logistic regression models under ridge regularisation, once using the early time points and once using all time points. For each patient, all time points are assigned to the same cross-validation fold, and the weight is split evenly among the time points. We then use the two sets of estimated regression coefficients as co-data for the target data set. In the regression for the target data set, we only include the earliest time point of each patient. Using 10-fold cross-validation, we estimate the predictive performance of ridge regression with and without transfer learning. After repeating source-target splitting and cross-validation 10 times, we observe that transfer learning tends to decrease the cross-validated logistic deviance (`transreg.exp.sta`: 8/10, `transreg.exp.sim`: 7/10, `transreg.iso.sta`: 7/10, `transreg.iso.sim`: 9/10, `fwelnet`: 5/10, `ecpc`: 6/10). It is more beneficial to share information not only on the importance but also the direction of the effects, according to the mean change in cross-validated logistic deviance (`transreg.exp.sta`: -2.61%, `transreg.exp.sim`: -4.29%, `transreg.iso.sta`: -3.67%, `transreg.iso.sim`: -8.33%, `fwelnet`: +0.04%, `ecpc`: -6.87%). Simultaneous stacking again outperforms standard stacking, but exponential and isotonic calibration show a similar performance.

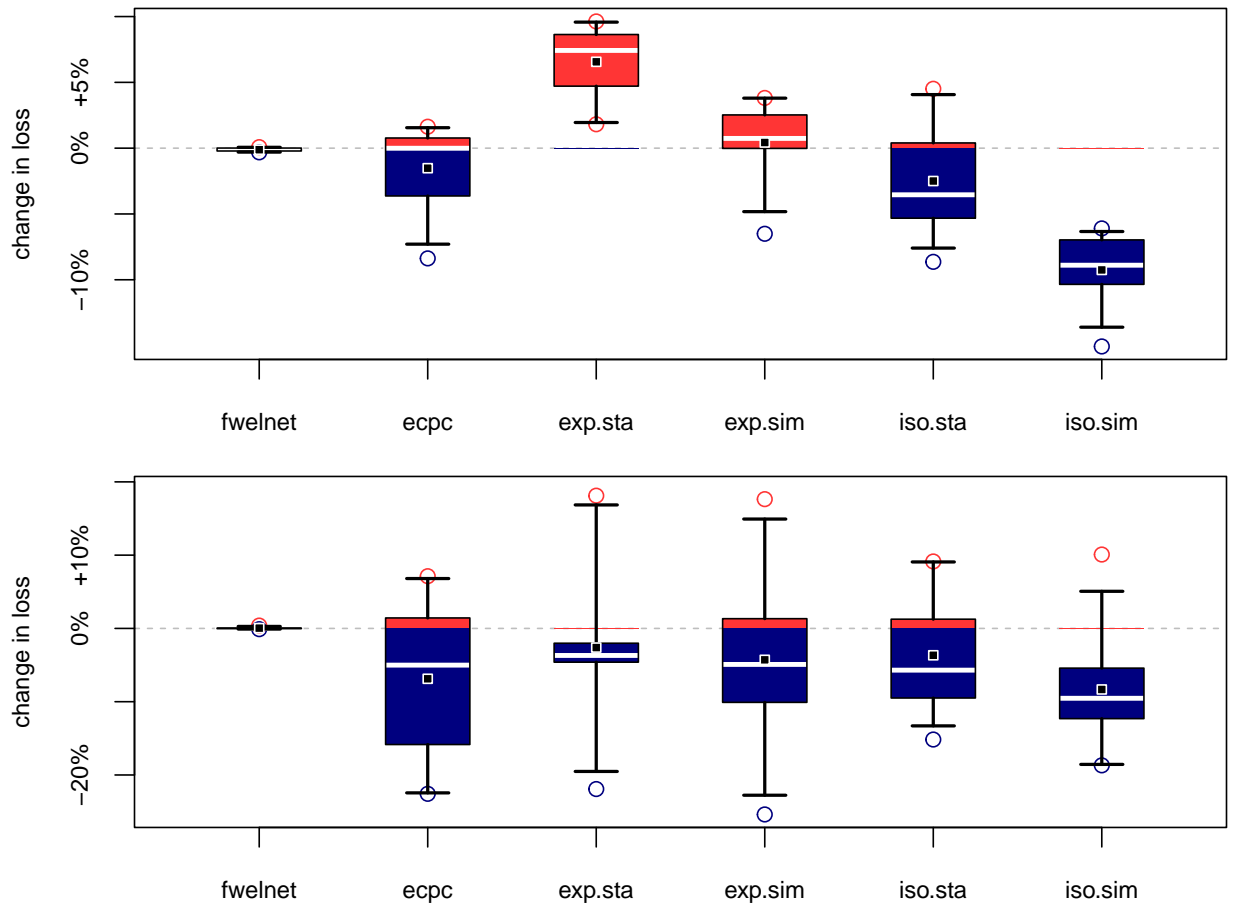


Figure 2: Percentage change in cross-validated logistic deviance from ridge regression to other methods, for 10 repetitions of 10-fold cross-validation. Top: application on cervical cancer. Bottom: application on pre-eclampsia. From left to right: co-data learning (`fwelnet`, `ecpc`), transfer learning with exponential/isotonic calibration and standard/simultaneous stacking.

## Internal application

In this application, we transfer information from a meta-analysis of genome-wide association studies on Parkinson’s disease (PD-GWAS, Nalls et al., 2019) to the Luxembourg Parkinson’s study (LUXPARK, Hipp et al., 2018). The aim is to classify samples into Parkinson’s disease (PD) patients and healthy controls based on single-nucleotide polymorphisms (SNPs).

At the time of our study, the LUXPARK data set included genotyping and clinical data of 790 PD cases and 766 healthy controls. DNA samples were genotyped using the NeuroChip array (Blauwendraat et al., 2017). Quality control steps of genotyping data were conducted according to the standard procedures reported previously (Pavelka et al., 2022). Missing genotyping data were imputed using the reference panel from the Haplotype Reference Consortium (release 1.1) on the Michigan Imputation Server (Das et al., 2016) [RRID:ID\_017579], with a filter for imputation quality ( $r^2 > 0.3$ ).

As common SNPs exhibit weak effects on PD, the sample size is likely insufficient to train a highly predictive model. However, publicly available summary statistics from the largest-to-date PD-GWAS (with around 38 000 cases and 1 400 000 controls from European ancestry) (Nalls et al., 2019) might serve as prior information on the SNP effects. For each SNP, these summary statistics are the combined results from simple logistic regression of the PD status on the SNP, namely the estimated slope (logarithmic odds ratio), its standard error, and the associated  $p$ -value. Importantly, the LUXPARK cohort was not part of the PD-GWAS, meaning that the prior information comes from independent data. As the LUXPARK cohort and the PD-GWAS cohorts have a similar ethnic background, the prior information might allow us to increase the predictive performance.

The two lists of SNPs - from the LUXPARK genotyping data (target data set) and the PD-GWAS summary statistics (source data set) - are partially overlapping. SNP data are high-dimensional and strongly correlated. From each block of SNPs in the target data set (250 kb window), we retain the most significant one and those that are in weak pairwise linkage disequilibrium with it ( $r^2 < 0.1$ ). Next, we only retain the SNPs appearing also in the source data set. These two filtering steps together reduce the dimensionality in the target data set from around 18 million SNPs to 196 018 SNPs. We code the SNP data for dominant effects, with 0 meaning no alternate allele (0/0) and 1 meaning one or two alternate alleles (0/1 or 1/1).

It seems that the results from the source data set are informative, because 5.80% of the  $p$ -values are nominally significant at the 0.05 level (11 377 out of 196 018), 77 are significant at a false discovery rate of 5% (Benjamini-Hochberg), and 35 are significant at a family-wise error rate of 5% (Holm-Bonferroni). As SNPs with a low minor allele frequency might have large effect sizes but insignificant  $p$ -values, we base the prior effects not on the estimated coefficients ( $\hat{\beta}$ ) but on the signed logarithmic  $p$ -values ( $-\text{sign}(\hat{\beta}) \log_{10}(p)$ ). For each SNP, we compared the reference and the alternate alleles between the two data sets: (i) If both data sets have the same reference allele and the same alternate allele, the signed logarithmic  $p$ -value from the source data set becomes the prior effect for the target data set. (ii) If the reference allele of each data set is the alternate allele of the other data set (swapped alleles), we invert the sign of the signed logarithmic  $p$ -value. (iii) And if the two data sets have two different sets of alleles (multiallelic SNP), we set the prior effect to zero.

Rather than using the 196 018 SNPs for predictive modelling in the target data set, we also filter them based on their significance in the source data set (which is already a type of transfer learning). For each cut-off in  $\{5 \times 10^{-2}, 5 \times 10^{-3}, \dots, 5 \times 10^{-10}\}$ , we exclude all SNPs above and include all SNPs below. This means that for the target data set, we retain a specific number of the most significant SNPs from the source data set. For each significance cut-off, we compare three modelling approaches:

- *Uninformed approach*: We use logistic regression with ridge or lasso penalisation to model the PD status based on the included SNPs. All included SNPs are treated equally, irrespective of their estimated effect in the source data set.
- *Naïve transfer learning*: After calculating for each sample the sum across the signed logarithmic  $p$ -values from the source data set multiplied by the SNPs from the target data set, we fit a simple logistic regression of the PD status on this sum.
- *Transfer learning*: The proposed transfer learning approach uses the signed logarithmic  $p$ -values from the source data set as prior effects for the target data set.



Figure 3 shows the predictive performance of modelling with estimated effects (uninformed approach), with prior effects (naïve transfer learning), or with both (transfer learning). We obtained the results with repeated nested cross-validation (10 repetitions, 10 external folds, 10 internal folds), using the same folds for all methods. If the significance cut-off is very strict, leading to a small number of significant SNPs, transfer learning does not improve the predictive performance of ridge and lasso regression. In these low-dimensional settings with many fewer SNPs than samples, prior information on the SNPs is not helpful. But otherwise transfer learning does improve the predictive performance of ridge and lasso regression. This holds for all four flavours of the proposed transfer learning method (exponential vs isotonic calibration, standard vs simultaneous stacking), but isotonic calibration works considerably better than exponential calibration and simultaneous stacking works marginally better than standard stacking.

Depending on the significance cut-off determining the number of significant SNPs, the performance of naïve transfer learning can be as high as the one of transfer learning with isotonic calibration. In these cases, the prior effects are predictive to the extent that it is not even necessary to estimate any effects. An explanation for the high performance of naïve transfer learning might be (i) the large sample size in the source data set for testing the marginal effects of the SNPs together with (ii) the linkage disequilibrium clumping leading to a selection of relatively independent SNPs.

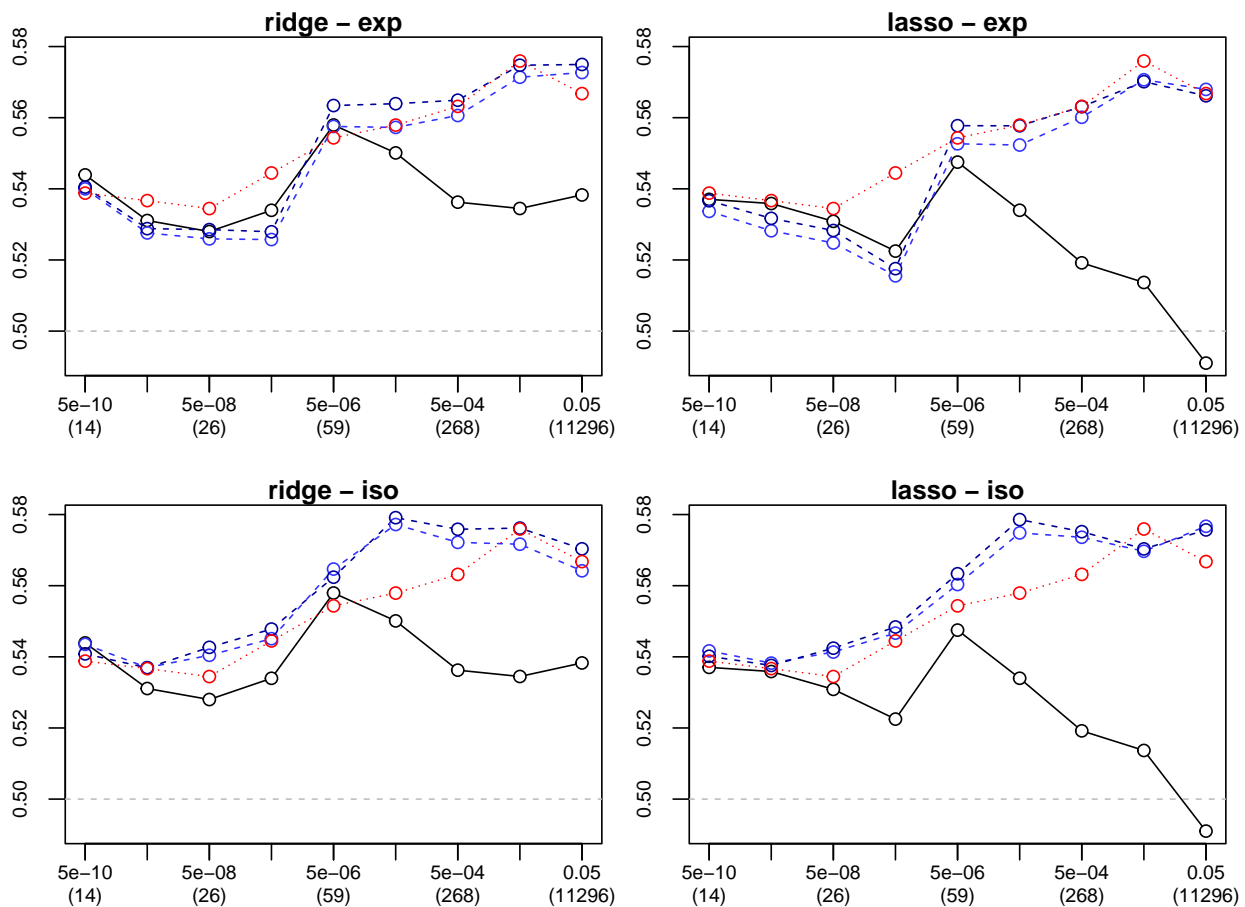


Figure 3: Mean cross-validated concordance index (C-index/AUC) from 10 times 10-fold cross-validation ( $y$ -axis) against  $p$ -value cutoff ( $x$ -axis) for regression without (solid line) and with (dashed lines) transfer learning (bright blue: standard stacking, dark blue: simultaneous stacking), under either ridge (left) or lasso (right) regularisation and either exponential (top) or isotonic (bottom) scaling. The numbers within brackets indicate the dimensionality and the dotted line is for naïve transfer learning.

## Discussion

We proposed a two-step transfer learning method for exploiting estimated coefficients from related studies to improve the predictive performance in the study of interest. First, we adapt the prior effects from the source data sets to the target data set, either with exponential or isotonic calibration. While exponential calibration is more robust to outliers (only three free parameters), isotonic calibration is more flexible (only maintains order of prior effects). We expect the former to be superior if the prior effects are close to the true effects, and the latter to be superior if there is no exponential relationship. Second, we combine the calibrated prior effects with information from the observed data, based on two variants of stacked generalisation. While the first variant (standard stacking) is more suitable if there are many sources of co-data (‘averaging calibrated prior effects and estimated effects’), the second variant (simultaneous stacking) is more suitable if there is one source of co-data with partially unreliable or partially missing prior effects (‘shrinking combined effects towards calibrated prior effects’).

The proposed transfer learning method allows for multiple sources of prior information. It does not require the source data set(s) but only the prior effects derived from the source data set(s). It therefore allows researchers to transfer predictive information from one study to another without requesting access to sensitive data. The proposed method has a competitive predictive performance with an existing but less flexible method (see simulation). In the case of closely related tasks, accounting for prior effects with transfer learning seems to be more beneficial than accounting for prior weights with co-data methods (see application). We therefore believe that the proposed method could tackle many biomedical predictions problems with one or more sets of prior effects.

In some applications, only one type of prior information derived from the source data sets is available. In other applications, multiple types of prior information are available (or the source data sets themselves). Then we can choose from multiple types of prior information. If the source and target data sets have the same feature space, estimated coefficients from penalised regression might be a reasonable choice. If the feature spaces are different, however, it is problematic that (i) lasso regression erratically selects among correlated features and (ii) ridge regression distributes weight among correlated features. This means that the presence or absence of additional correlated features in the source data sets might change the prior information on the features of interest. The same problem arises under contamination of a subset of features (van de Wiel et al., 2016). We therefore expect that signed logarithmic  $p$ -values ( $-\text{sign}(\hat{\beta}) \log_{10}(p)$ ) will often be more informative than estimated coefficients ( $\hat{\beta}$ ).

## Funding

This work was supported by the Luxembourg National Research Fund (FNR) for the ERA-Net ERACOSysMed JTC-2 project PD-Strat [INTER/11651464] and by the European Union’s Horizon 2020 research and innovation programme for the project DIGI-PD [ERAPERMED 2020-314]. The National Centre of Excellence in Research on Parkinson’s Disease (NCER-PD) is funded by the FNR [FNR/NCER13/BM/11264123].

## Acknowledgements

We are grateful to Quentin KLOPFENSTEIN for helpful discussions on differential penalisation, and to the responsible and reproducible research (R3) initiative for the pre-publication check. Data used in the preparation of this manuscript were obtained from NCER-PD. All participants provided written informed consent, and the study has been approved by the National Ethics Board (CNER Ref: 201407/13). We would like to thank all participants of the Luxembourg Parkinson’s Study for their important support of our research. Furthermore, we acknowledge the joint effort of the NCER-PD consortium members from the partner institutions Luxembourg Centre for Systems Biomedicine, Luxembourg Institute of Health, Centre Hospitalier de Luxembourg, and Laboratoire National de Santé generally contributing to the Luxembourg Parkinson’s Study as listed below: Alexander HUNDT<sup>2</sup>, Alexandre BISDORFF<sup>5</sup>, Amir SHARIFY<sup>2</sup>, Anne GRÜNEWALD<sup>1</sup>, Anne-Marie HANFF<sup>2</sup>, Armin RAUSCHENBERGER<sup>1</sup>, Beatrice NICOLAI<sup>3</sup>, Brit MOLLENHAUER<sup>12</sup>, Camille BELLORA<sup>2</sup>, Carlos VEGA MORENO<sup>1</sup>, Chouaib MEOUNI<sup>2</sup>, Christophe TREFOIS<sup>1</sup>, Claire PAULY<sup>1,3</sup>, Clare MACKAY<sup>10</sup>, Clarissa GOMES<sup>1</sup>, Daniela BERG<sup>11</sup>, Daniela ESTEVES<sup>2</sup>, Deborah MCINTYRE<sup>2</sup>, Dheeraj REDDY BOBBILI<sup>1</sup>, Eduardo ROSALES<sup>2</sup>, Ekaterina SOBOLEVA<sup>1</sup>, Elisa GÓMEZ DE LOPE<sup>1</sup>, Elodie THIRY<sup>3</sup>, Enrico GLAAB<sup>1</sup>, Estelle HENRY<sup>2</sup>, Estelle SANDT<sup>2</sup>, Evi WOLLSCHIED-LENGELING<sup>1</sup>, Françoise MEISCH<sup>1</sup>, Friedrich MÜHLSCHLEGEL<sup>4</sup>, Gaël HAMMOT<sup>2</sup>, Geeta ACHARYA<sup>2</sup>, Gelani ZELIMKHANOV<sup>3</sup>, Gessica CONTESOTTO<sup>2</sup>, Giuseppe ARENA<sup>1</sup>, Gloria AGUAYO<sup>2</sup>,

Guilherme MARQUES<sup>2</sup>, Guy BERCEM<sup>3</sup>, Guy FAGHERAZZI<sup>2</sup>, Hermann THIEN<sup>2</sup>, Ibrahim BOUSSAAD<sup>1</sup>, Inga LIEPELT<sup>11</sup>, Isabel ROSETY<sup>1</sup>, Jacek JAROSLAW LEBIODA<sup>1</sup>, Jean-Edouard SCHWEITZER<sup>1</sup>, Jean-Paul NICOLAY<sup>19</sup>, Jean-Yves FERRAND<sup>2</sup>, Jens SCHWAMBORN<sup>1</sup>, Jérôme GRAAS<sup>2</sup>, Jessica CALMES<sup>2</sup>, Jochen KLUCKEN<sup>1,2,3</sup>, Johanna TROUET<sup>2</sup>, Kate SOKOLOWSKA<sup>2</sup>, Kathrin BROCKMANN<sup>11</sup>, Katrin MARCUS<sup>13</sup>, Katy BEAUMONT<sup>2</sup>, Kirsten RUMP<sup>1</sup>, Laura LONGHINO<sup>3</sup>, Laure PAULY<sup>1</sup>, Liliana VILAS BOAS<sup>3</sup>, Linda HANSEN<sup>1,3</sup>, Lorieza CASTILLO<sup>2</sup>, Lukas PAVELKA<sup>1,3</sup>, Magali PERQUIN<sup>2</sup>, Maharshi VYAS<sup>1</sup>, Manon GANTENBEIN<sup>2</sup>, Marek OSTASZEWSKI<sup>1</sup>, Margaux SCHMITT<sup>2</sup>, Mariella GRAZIANO<sup>17</sup>, Marijus GIRAITIS<sup>2,3</sup>, Maura MINELLI<sup>2</sup>, Maxime HANSEN<sup>1,3</sup>, Mesele VALENTI<sup>2</sup>, Michael HENEKA<sup>1</sup>, Michael HEYMAN<sup>2</sup>, Michel MITTELBRONN<sup>1,4</sup>, Michel VAILLANT<sup>2</sup>, Michele BASSIS<sup>1</sup>, Michele HU<sup>8</sup>, Muhammad ALI<sup>1</sup>, Myriam ALEXANDRE<sup>2</sup>, Myriam MENSTER<sup>2</sup>, Nadine JACOBY<sup>18</sup>, Nico DIEDERICH<sup>3</sup>, Olena TSURKALENKO<sup>2</sup>, Olivier TERWINDT<sup>1,3</sup>, Patricia MARTINS CONDE<sup>1</sup>, Patrick MAY<sup>1</sup>, Paul WILMES<sup>1</sup>, Paula Cristina LUPU<sup>2</sup>, Pauline LAMBERT<sup>2</sup>, Piotr GAWRON<sup>1</sup>, Quentin KLOPFENSTEIN<sup>1</sup>, Rajesh RAWAL<sup>1</sup>, Rebecca TING JIIN LOO<sup>1</sup>, Regina BECKER<sup>1</sup>, Reinhard SCHNEIDER<sup>1</sup>, Rejko KRÜGER<sup>1,2,3</sup>, Rene DONDELINGER<sup>5</sup>, Richard WADE-MARTINS<sup>9</sup>, Robert LISZKA<sup>14</sup>, Romain NATT<sup>3</sup>, Rosalina RAMOS LIMA<sup>2</sup>, Roseline LENTZ<sup>7</sup>, Rudi BALLING<sup>1</sup>, Sabine SCHMITZ<sup>1</sup>, Sarah NICKELS<sup>1</sup>, Sascha HERZINGER<sup>1</sup>, Sinthuja PACHCHEK<sup>1</sup>, Soumyabrata GHOSH<sup>1</sup>, Stefano SAPIENZA<sup>1</sup>, Sylvia HERBRINK<sup>6</sup>, Tainá MARQUES<sup>1</sup>, Thomas GASSER<sup>11</sup>, Ulf NEHRBASS<sup>2</sup>, Valentin GROUES<sup>1</sup>, Venkata SATAGOPAM<sup>1</sup>, Victoria LORENTZ<sup>2</sup>, Walter MAETZLER<sup>15</sup>, Wei GU<sup>1</sup>, Wim AMMERLANN<sup>2</sup>, Yohan JAROZ<sup>1</sup>, Zied LANDOULSI<sup>1</sup>. <sup>1</sup>Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Esch-sur-Alzette, Luxembourg, <sup>2</sup>Luxembourg Institute of Health, Strassen, Luxembourg, <sup>3</sup>Centre Hospitalier de Luxembourg, Strassen, Luxembourg, <sup>4</sup>Laboratoire National de Santé, Dudelange, Luxembourg, <sup>5</sup>Centre Hospitalier Emile Mayrisch, Esch-sur-Alzette, Luxembourg, <sup>6</sup>Centre Hospitalier du Nord, Ettelbrück, Luxembourg, <sup>7</sup>Parkinson Luxembourg Association, Leudelange, Luxembourg, <sup>8</sup>Oxford Parkinson's Disease Centre, Nuffield Department of Clinical Neurosciences, University of Oxford, Oxford, UK, <sup>9</sup>Oxford Parkinson's Disease Centre, Department of Physiology, Anatomy and Genetics, University of Oxford, Oxford, UK, <sup>10</sup>Oxford Centre for Human Brain Activity, Wellcome Centre for Integrative Neuroimaging, Department of Psychiatry, University of Oxford, Oxford, UK, <sup>11</sup>Center of Neurology and Hertie Institute for Clinical Brain Research, Department of Neurodegenerative Diseases, University Hospital Tübingen, Tübingen, Germany, <sup>12</sup>Paracelsus-Elena-Klinik, Kassel, Germany, <sup>13</sup>Ruhr-University of Bochum, Bochum, Germany, <sup>14</sup>Westpfalz-Klinikum GmbH, Kaiserslautern, Germany, <sup>15</sup>Department of Neurology, University Medical Center Schleswig-Holstein, Kiel, Germany, <sup>16</sup>Department of Neurology Philipps, University Marburg, Marburg, Germany, <sup>17</sup>Association of Physiotherapists in Parkinson's Disease Europe, Esch-sur-Alzette, Luxembourg, <sup>18</sup>Private practice, Ettelbruck, Luxembourg, <sup>19</sup>Private practice, Luxembourg, Luxembourg.

## Reproducibility

The R package `transreg` is available on [GitHub](https://github.com/lcsb-bds/transreg) (<https://github.com/lcsb-bds/transreg>), with the code for the simulations and the applications in a vignette (<https://lcsb-bds.github.io/transreg/>). We obtained our results using R 4.2.2 [RRID:ID\_001905] on a physical machine (aarch64-apple-darwin20, macOS Monterey 12.6). Data for the application on cervical cancer are available from van de Wiel et al. (2016), in the R package `GRridge` in the data set 'dataVerlaat' (source data: Farkas et al., 2013, target data: van de Wiel et al., 2016). Data for the application on pre-eclampsia are available from Erez et al. (2017), in the supporting file 'pone.0181468.s001.csv'. For the application on Parkinson's disease, the source data are available from Nalls et al. (2019), in the online file 'nallsEtAl2019\_excluding23andMe\_allVariants.tab', and the target data are available upon request (`request.ncer-pd@uni.lu`). Information on reproducibility is also available on a frozen page (doi: 10.17881/hczj-3297).

## Author contributions

EG acquired funding. AR and MvdW developed the method. AR analysed the data and drafted the manuscript. ZL processed the SNP data and critically revised the internal application. MvdW and EG critically revised the manuscript. All authors read and approved the final manuscript.

## Appendix

Table 4: Isotonic calibration. The aim is to (i) estimate the effects of the features under sign and order constraints determined by  $q$  negative and  $p-q$  non-negative prior effects, i.e. estimate  $\gamma_{1k}, \dots, \gamma_{pk}$  for  $x_1, \dots, x_p$  under  $\hat{\gamma}_{jk} = 0|z_{jk} = 0$ ,  $\hat{\gamma}_{jk} \geq 0|z_{jk} > 0$ ,  $\hat{\gamma}_{jk} \leq 0|z_{jk} < 0$ ,  $\hat{\gamma}_{jk} \geq \hat{\gamma}_{lk}|z_{jk} \geq z_{lk}$ , and  $\hat{\gamma}_{jk} \leq \hat{\gamma}_{lk}|z_{jk} \leq z_{lk}$ . This can be solved by (ii) estimating the effects of the features ordered by the co-data under sign and order constraints, i.e. estimate  $\gamma_{(1),k}, \dots, \gamma_{(p),k}$  for  $x_{(1)}, \dots, x_{(p)}$  under  $\hat{\gamma}_{(j),k} \leq 0|j \leq p$ ,  $\hat{\gamma}_{(j),k} \geq 0|j > p$ , and  $\hat{\gamma}_{(1),k} \leq \dots \leq \hat{\gamma}_{(p),k}$ . This in turn can be solved by (iii) estimating the effects of the combined features under sign constraints, i.e. estimate  $\delta_1, \dots, \delta_p$  for  $w_1, \dots, w_p$  under  $\hat{\delta}_j \leq 0|j \leq p$  and  $\hat{\delta}_j \geq 0|j > p$ . Our algorithm receives the original features and the prior effects (i), orders the features by the prior effects (ii), combines the features (iii), estimates the effects of the combined features (iii), calculates the estimated effects of the ordered features (ii), and returns the estimated effects of the original features (i).

(i)	$x_{\circ,1}, z_{1,k}$	$x_{\circ,2}, z_{2,k}$	$\dots$	$x_{\circ,q-1}, z_{q-1,k}$	$x_{\circ,q}, z_{q,k}$
(ii)	$x_{\circ,(1)}$	$x_{\circ,(2)}$	$\dots$	$x_{\circ,(q-1)}$	$x_{\circ,(q)}$
(iii)	$w_{\circ,1} =$ $x_{\circ,(1)}$	$w_{\circ,2} =$ $x_{\circ,(1)} + x_{\circ,(2)}$	$\dots$	$w_{\circ,q-1} =$ $x_{\circ,(1)} + \dots + x_{\circ,(q-1)}$	$w_{\circ,q} =$ $x_{\circ,(1)} + \dots + x_{\circ,(q)}$
(iii)	$\hat{\delta}_{1,k}$	$\hat{\delta}_{2,k}$	$\dots$	$\hat{\delta}_{q-1,k}$	$\hat{\delta}_{q,k}$
(ii)	$\hat{\gamma}_{(1),k} =$ $\hat{\delta}_{1,k} + \dots + \hat{\delta}_{q,k}$	$\hat{\gamma}_{(2),k} =$ $\hat{\delta}_{2,k} + \dots + \hat{\delta}_{q,k}$	$\dots$	$\hat{\gamma}_{(q-1),k} =$ $\hat{\delta}_{q-1,k} + \hat{\delta}_{q,k}$	$\hat{\gamma}_{(q),k} =$ $\hat{\delta}_{q,k}$
(i)	$\hat{\gamma}_{1,k}$	$\hat{\gamma}_{2,k}$	$\dots$	$\hat{\gamma}_{q-1,k}$	$\hat{\gamma}_{q,k}$
(i)	$x_{\circ,q+1}, z_{q+1,k}$	$x_{\circ,q+2}, z_{q+2,k}$	$\dots$	$x_{\circ,p-1}, z_{p-1,k}$	$x_{\circ,p}, z_{p,k}$
(ii)	$x_{\circ,(q+1)}$	$x_{\circ,(q+2)}$	$\dots$	$x_{\circ,(p-1)}$	$x_{\circ,(p)}$
(iii)	$w_{\circ,q+1} =$ $x_{\circ,(q+1)} + \dots + x_{\circ,(p)}$	$w_{\circ,q+2} =$ $x_{\circ,(q+2)} + \dots + x_{\circ,(p)}$	$\dots$	$w_{\circ,p-1} =$ $x_{\circ,(p-1)} + x_{\circ,(p)}$	$w_{\circ,p} =$ $x_{\circ,(p)}$
(iii)	$\hat{\delta}_{q+1,k}$	$\hat{\delta}_{q+2,k}$	$\dots$	$\hat{\delta}_{p-1,k}$	$\hat{\delta}_{p,k}$
(ii)	$\hat{\gamma}_{(q+1),k} =$ $\hat{\delta}_{q+1,k}$	$\hat{\gamma}_{(q+2),k} =$ $\hat{\delta}_{q+1,k} + \hat{\delta}_{q+2,k}$	$\dots$	$\hat{\gamma}_{(p-1),k} =$ $\hat{\delta}_{q+1,k} + \dots + \hat{\delta}_{p-1,k}$	$\hat{\gamma}_{(p),k} =$ $\hat{\delta}_{q+1,k} + \dots + \hat{\delta}_{p,k}$
(i)	$\hat{\gamma}_{q+1,k}$	$\hat{\gamma}_{q+2,k}$	$\dots$	$\hat{\gamma}_{p-1,k}$	$\hat{\gamma}_{p,k}$

## References

- Bergersen, L. C., Glad, I. K., and Lyng, H. (2011). Weighted lasso with data integration. Statistical Applications in Genetics and Molecular Biology, 10(1):39. doi: 10.2202/1544-6115.1703.
- Blauwendraat, C., Faghri, F., Pihlstrom, L., Geiger, J. T., Elbaz, A., Lesage, S., Corvol, J.-C., May, P., Nicolas, A., Abramzon, Y., et al. (2017). NeuroChip, an updated version of the NeuroX genotyping platform to rapidly screen for variants associated with neurological diseases. Neurobiology of Aging, 57:247.e9–247.e13. doi: 10.1016/j.neurobiolaging.2017.05.009.
- Das, S., Forer, L., Schönherr, S., Sidore, C., Locke, A. E., Kwong, A., Vrieze, S. I., Chew, E. Y., Levy, S., McGue, M., et al. (2016). Next-generation genotype imputation service and methods. Nature Genetics, 48(10):1284–1287. doi: 10.1038/ng.3656.
- Dhruba, S. R. (2021). Application of advanced machine learning based approaches in cancer precision medicine. PhD thesis. <https://hdl.handle.net/2346/88264>. (R package `DMTL`).
- Erez, O., Romero, R., Maymon, E., Chaemsaitong, P., Done, B., Pacora, P., Panaitescu, B., Chaiworapongsa, T., Hassan, S. S., and Tarca, A. L. (2017). The prediction of late-onset preeclampsia: Results from a longitudinal proteomics study. PloS one, 12(7):e0181468. doi: 10.1371/journal.pone.0181468.
- Farkas, S. A., Milutin-Gašperov, N., Grce, M., and Nilsson, T. K. (2013). Genome-wide DNA methylation assay reveals novel candidate biomarker genes in cervical cancer. Epigenetics, 8(11):1213–1225. doi: 10.4161/epi.26346.
- Friedman, J. H., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. Journal of Statistical Software, 33(1):1–22. doi: 10.18637/jss.v033.i01. (R package `glmnet`).
- Gamarnik, D. and Zadik, I. (2017). High-dimensional regression with binary coefficients. Estimating squared error and a phase transition. In Conference on Learning Theory, volume 65 of Proceedings of Machine Learning Research, pages 948–953. doi: <http://proceedings.mlr.press/v65/david17a>.
- Hipp, G., Vaillant, M., Diederich, N. J., Roomp, K., Satagopam, V. P., Banda, P., Sandt, E., Mommaerts, K., Schmitz, S. K., Longhino, L., et al. (2018). The luxembourg parkinson’s study: a comprehensive approach for stratification and early diagnosis. Frontiers in Aging Neuroscience, 10:326. doi: 10.3389/fnagi.2018.00326.
- Jiang, Y., He, Y., and Zhang, H. (2016). Variable selection with prior information for generalized linear models via the prior LASSO method. Journal of the American Statistical Association, 111(513):355–376. doi: 10.1080/01621459.2015.1008363. (R script `pLASSO`).
- Nalls, M. A., Blauwendraat, C., Vallerga, C. L., Heilbron, K., Bandres-Ciga, S., Chang, D., Tan, M., Kia, D. A., Noyce, A. J., Xue, A., et al. (2019). Identification of novel risk loci, causal insights, and heritable risk for parkinson’s disease: a meta-analysis of genome-wide association studies. The Lancet Neurology, 18(12):1091–1102. doi: 10.1016/S1474-4422(19)30320-5.
- Pavelka, L., Rauschenberger, A., Landoulsi, Z., Pachchek, S., May, P., Glaab, E., and Krüger, R. (2022). Age at onset as stratifier in idiopathic parkinson’s disease - effect of ageing and polygenic risk score on clinical phenotypes. npj Parkinson’s Disease, 8:102. doi: 10.1038/s41531-022-00342-7.
- Rauschenberger, A. and Glaab, E. (2021). Predicting correlated outcomes from molecular data. Bioinformatics, 37(21):3889–3895. doi: 10.1093/bioinformatics/btab576. (R package `joinet`).
- Rauschenberger, A., Glaab, E., and van de Wiel, M. A. (2021). Predictive and interpretable models via the stacked elastic net. Bioinformatics, 37(14):2012–2016. doi: 10.1093/bioinformatics/btaa535. (R package `starnet`).
- Tay, J. K., Aghaeepour, N., Hastie, T., and Tibshirani, R. (2022). Feature-weighted elastic net: using “features of features” for better prediction. Statistica Sinica, in press. doi: 10.5705/ss.202020.0226. (R package `fwelnet`).

- te Beest, D. E., Mes, S. W., Wilting, S. M., Brakenhoff, R. H., and van de Wiel, M. A. (2017). Improved high-dimensional prediction with Random Forests by the use of co-data. BMC Bioinformatics, 18:584. doi: 10.1186/s12859-017-1993-1. (R package `CoRF`).
- Tian, Y. and Feng, Y. (2022). Transfer learning under high-dimensional generalized linear models. Journal of the American Statistical Association, in press. doi: 10.1080/01621459.2022.2071278. (R package `glmtrans`).
- van de Wiel, M. A., Lien, T. G., Verlaat, W., van Wieringen, W. N., and Wilting, S. M. (2016). Better prediction by use of co-data: adaptive group-regularized ridge regression. Statistics in Medicine, 35(3):368–381. doi: 10.1002/sim.6732. (R package `GRridge`).
- van de Wiel, M. A., van Nee, M. M., and Rauschenberger, A. (2021). Fast cross-validation for multi-penalty high-dimensional ridge regression. Journal of Computational and Graphical Statistics, 30(4):835–847. doi: 10.1080/10618600.2021.1904962. (R package `multiridge`).
- van Nee, M. M., Wessels, L. F., and van de Wiel, M. A. (2021). Flexible co-data learning for high-dimensional prediction. Statistics in Medicine, 40(26):5910–5925. doi: 10.1002/sim.9162. (R package `ecpc`).
- Wolpert, D. H. (1992). Stacked generalization. Neural Networks, 5(2):241–259. doi: 10.1016/S0893-6080(05)80023-1.
- Zeng, C., Thomas, D. C., and Lewinger, J. P. (2021). Incorporating prior knowledge into regularized regression. Bioinformatics, 37(4):514–521. doi: 10.1093/bioinformatics/btaa776. (R package `xtune`).