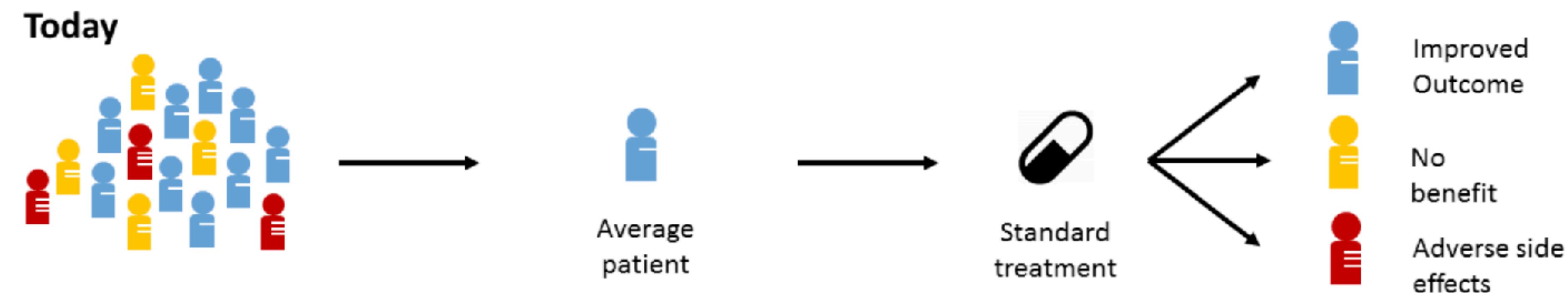


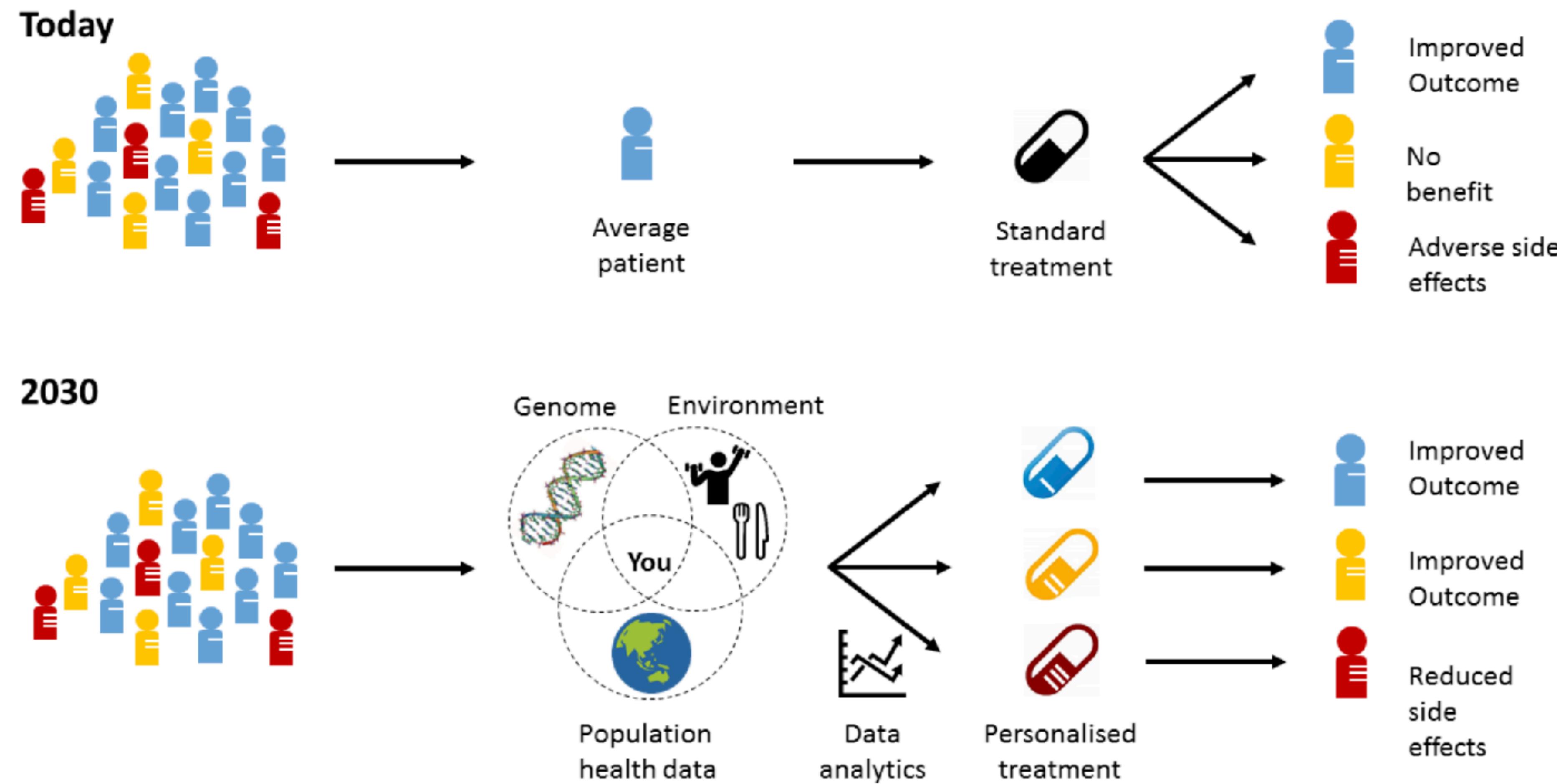
Kevin Wang

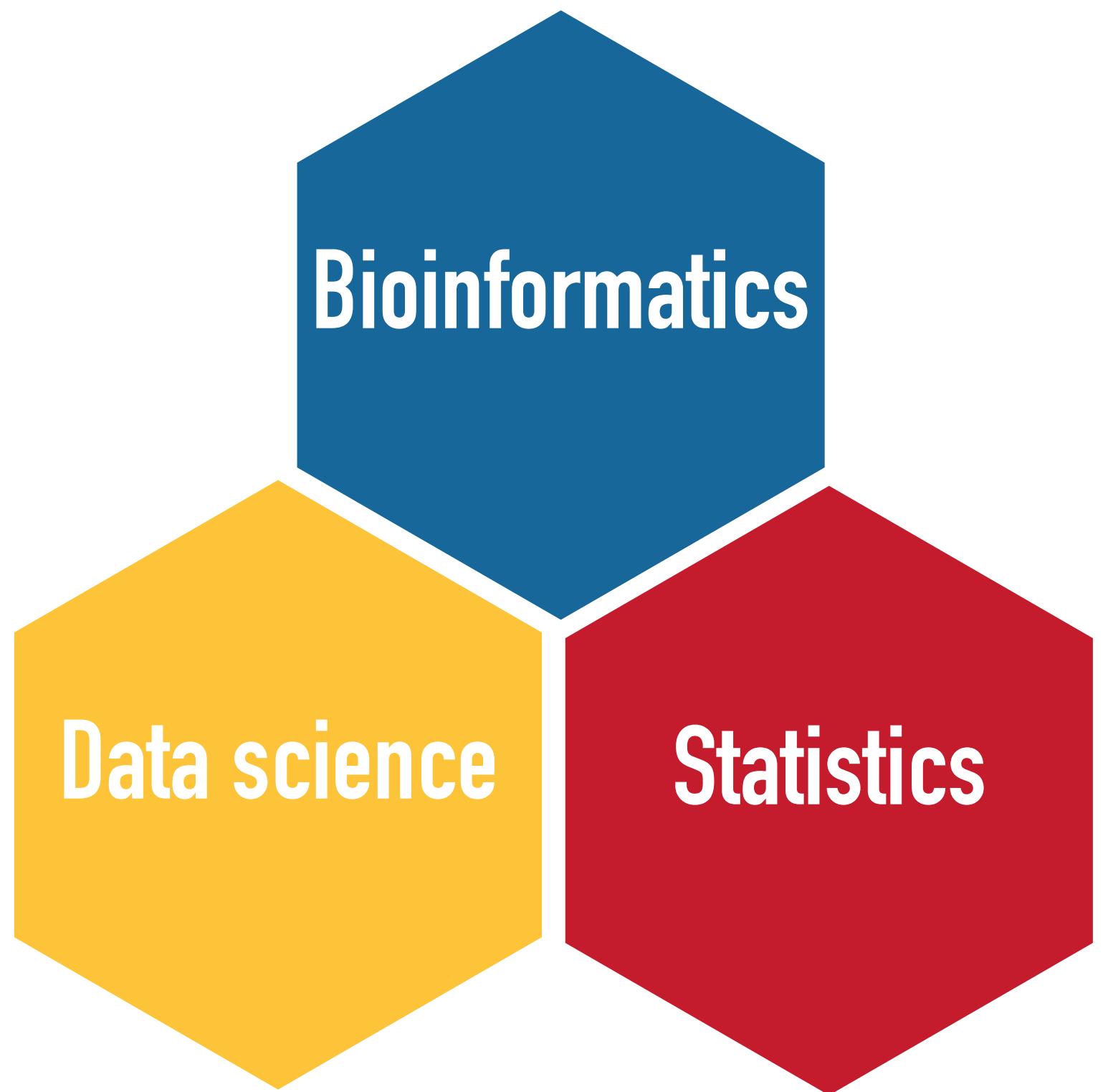
Cross-Platform Omics Prediction: a step towards precision medicine

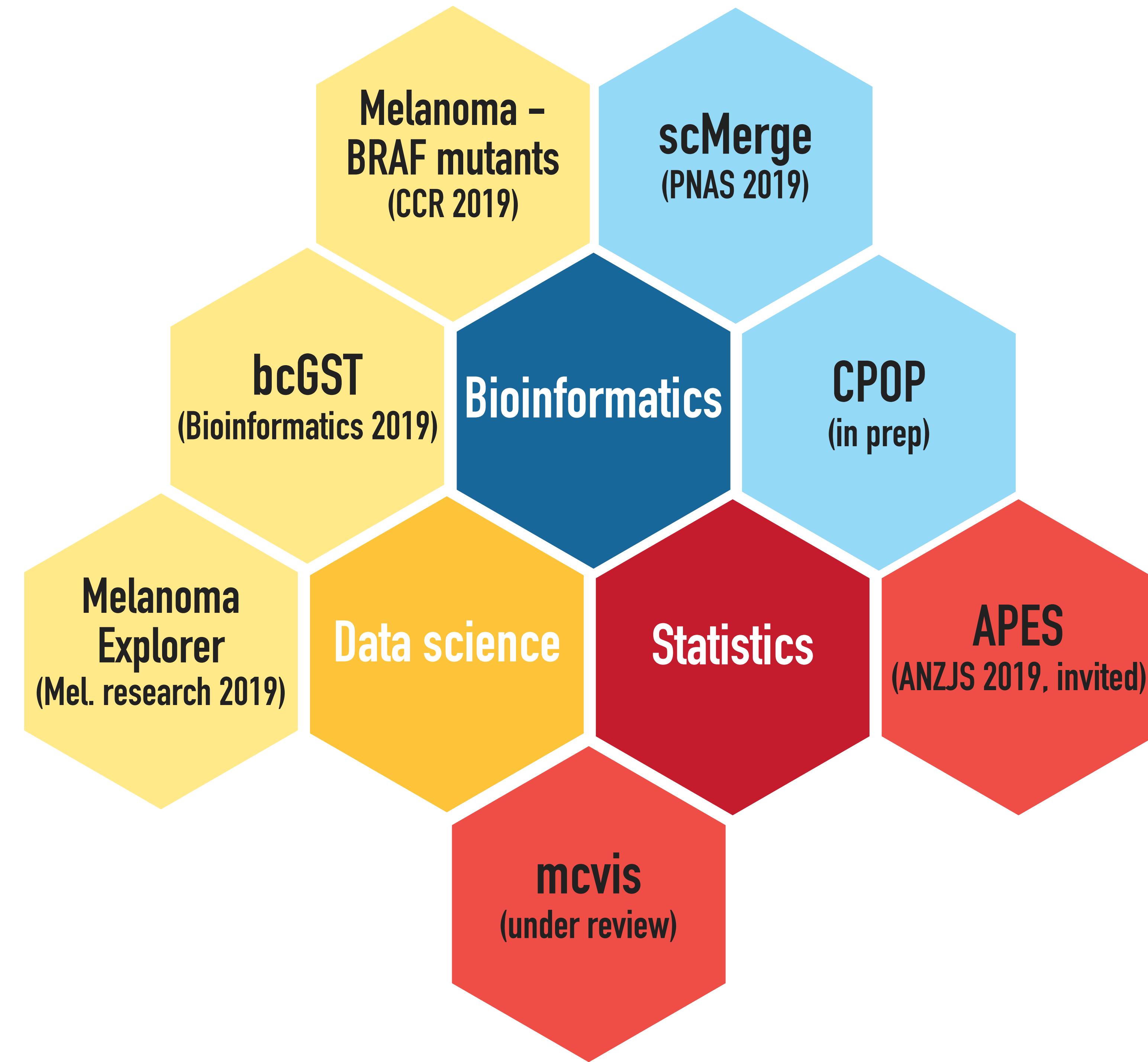
Precision medicine: predicting best cause of action using omics data



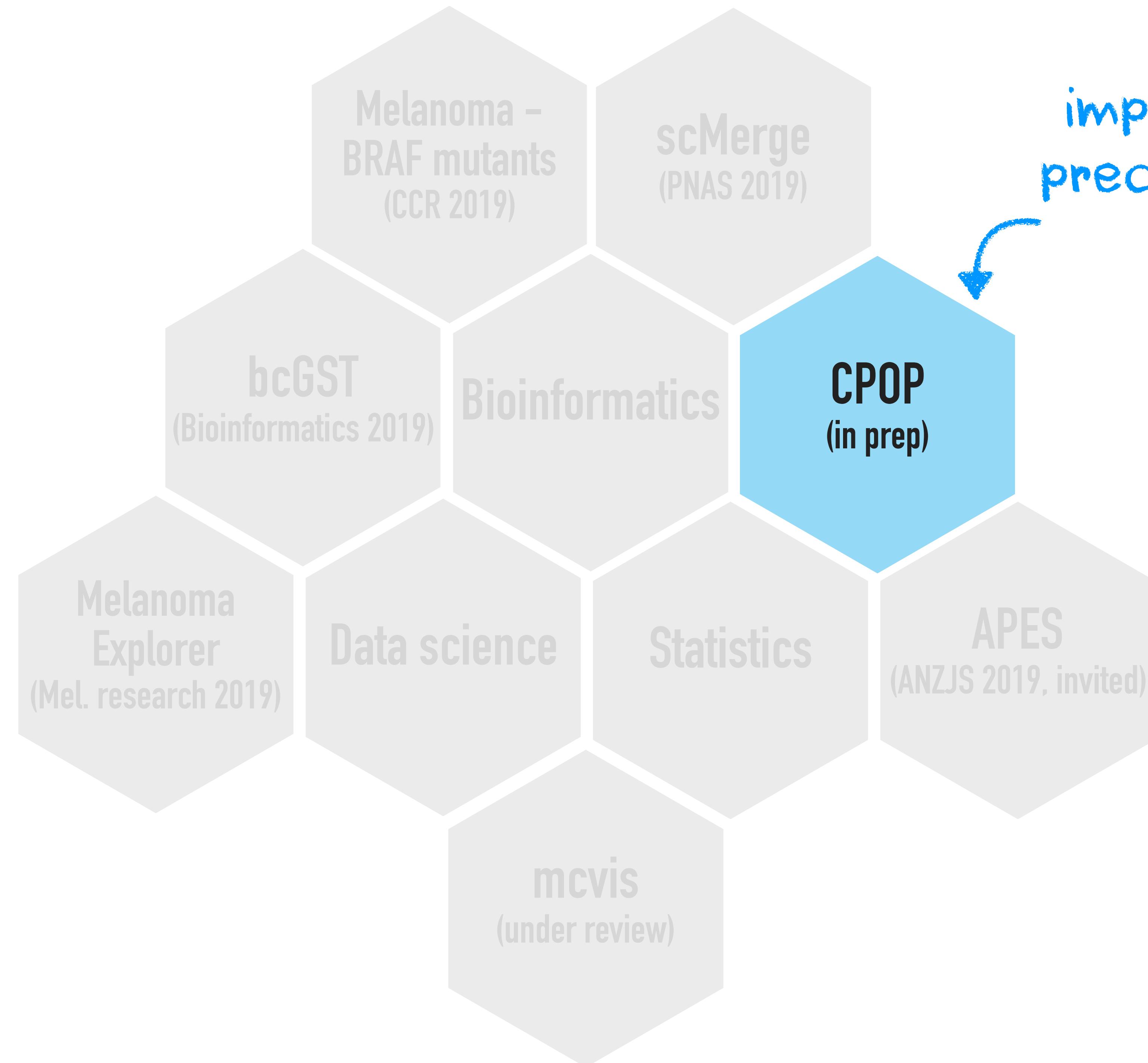
Precision medicine: predicting best cause of action using omics data





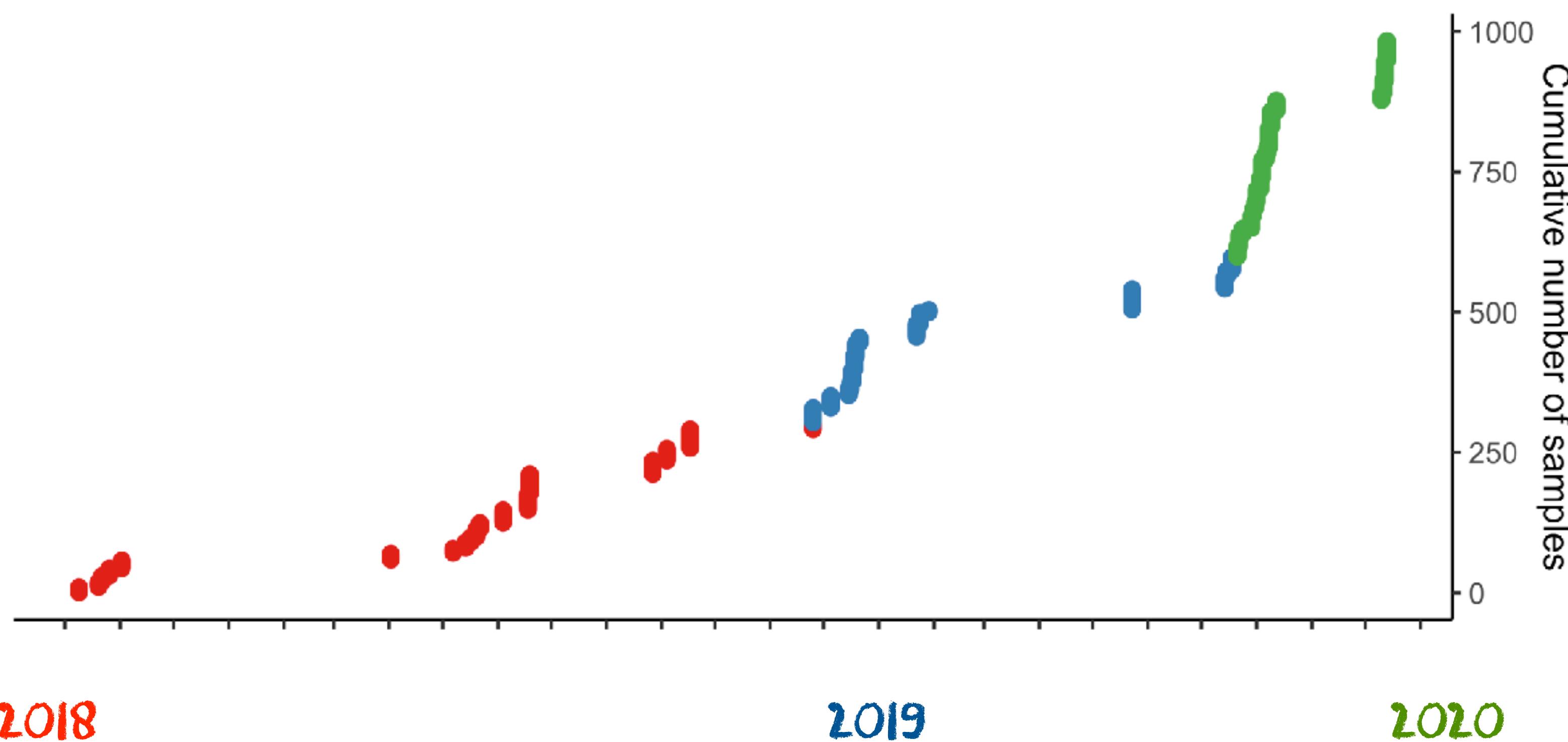


Towards
implementation
precision medicine

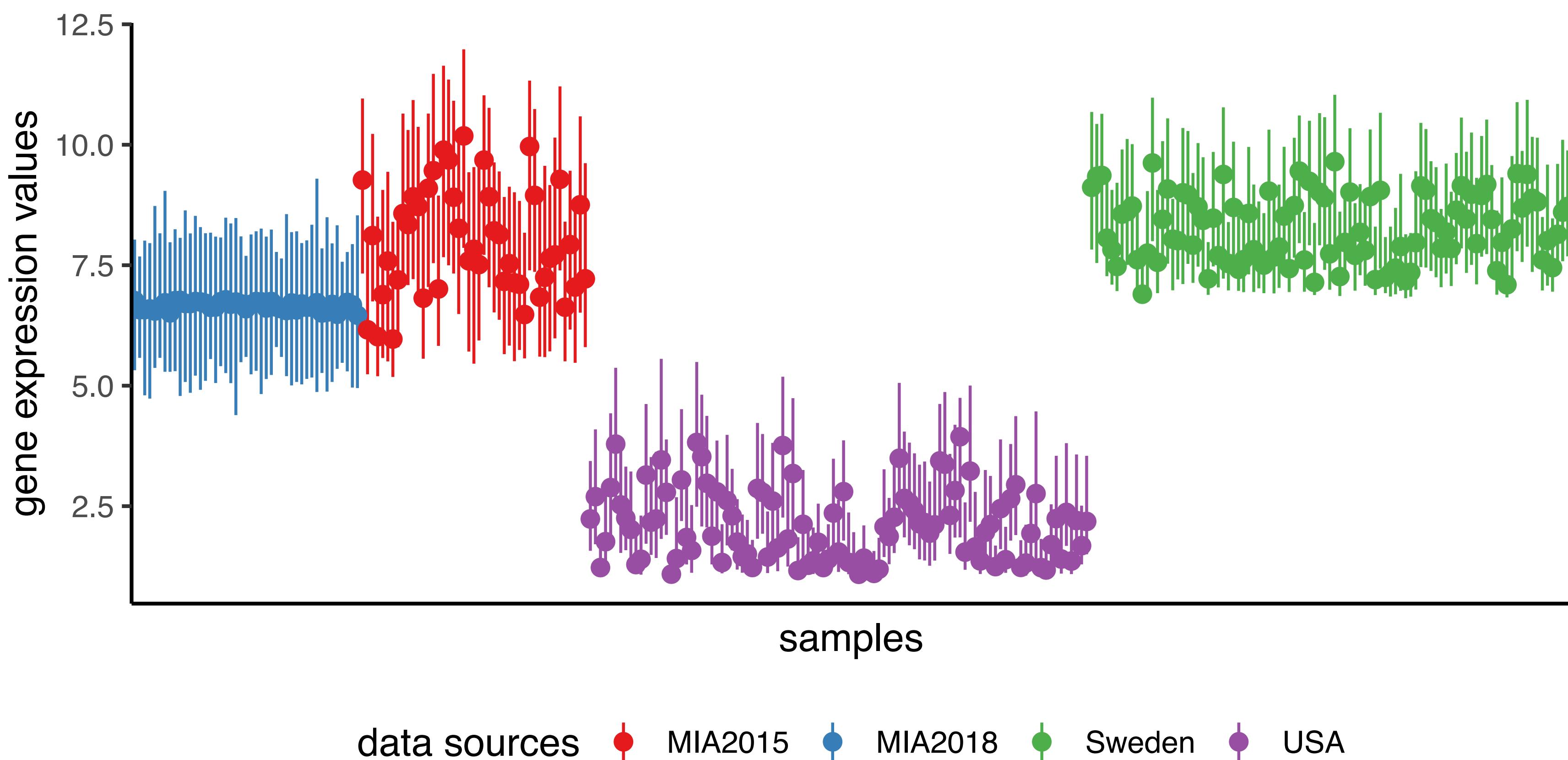


Melanoma Institute Australia

- ▶ Predict patient outcomes using gene expression:
- ▶ prospectively



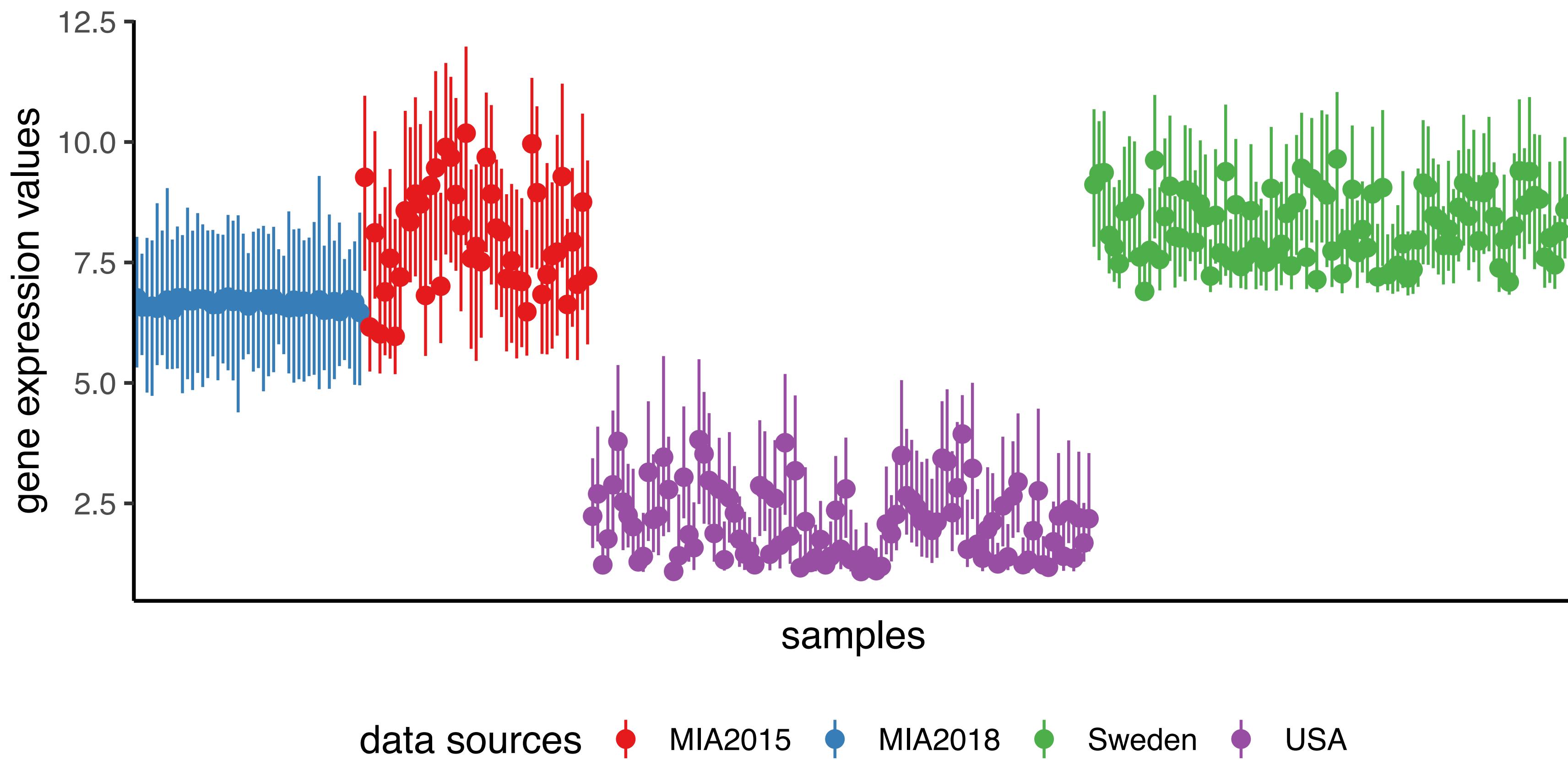
Melanoma Institute Australia



- ▶ Predict patient outcomes using gene expression:
 - ▶ prospectively
 - ▶ multi-centres
- ▶ CRE grant for implementation



Vision of a risk score



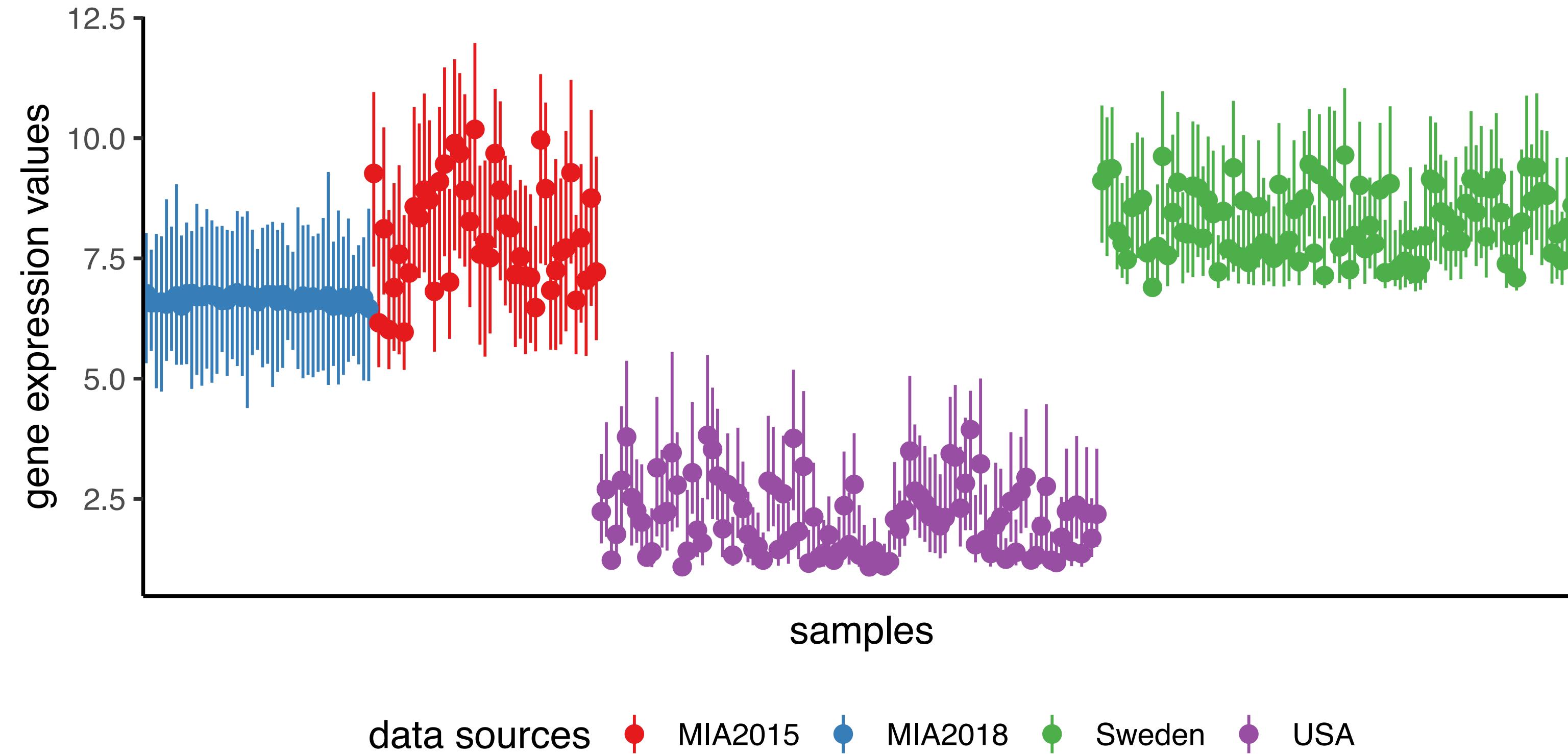
$$\hat{y} = X \hat{\beta}$$

regression-based
risk score



**Omics-based
risk score**

Omics-based clinical risk score: what is so difficult?

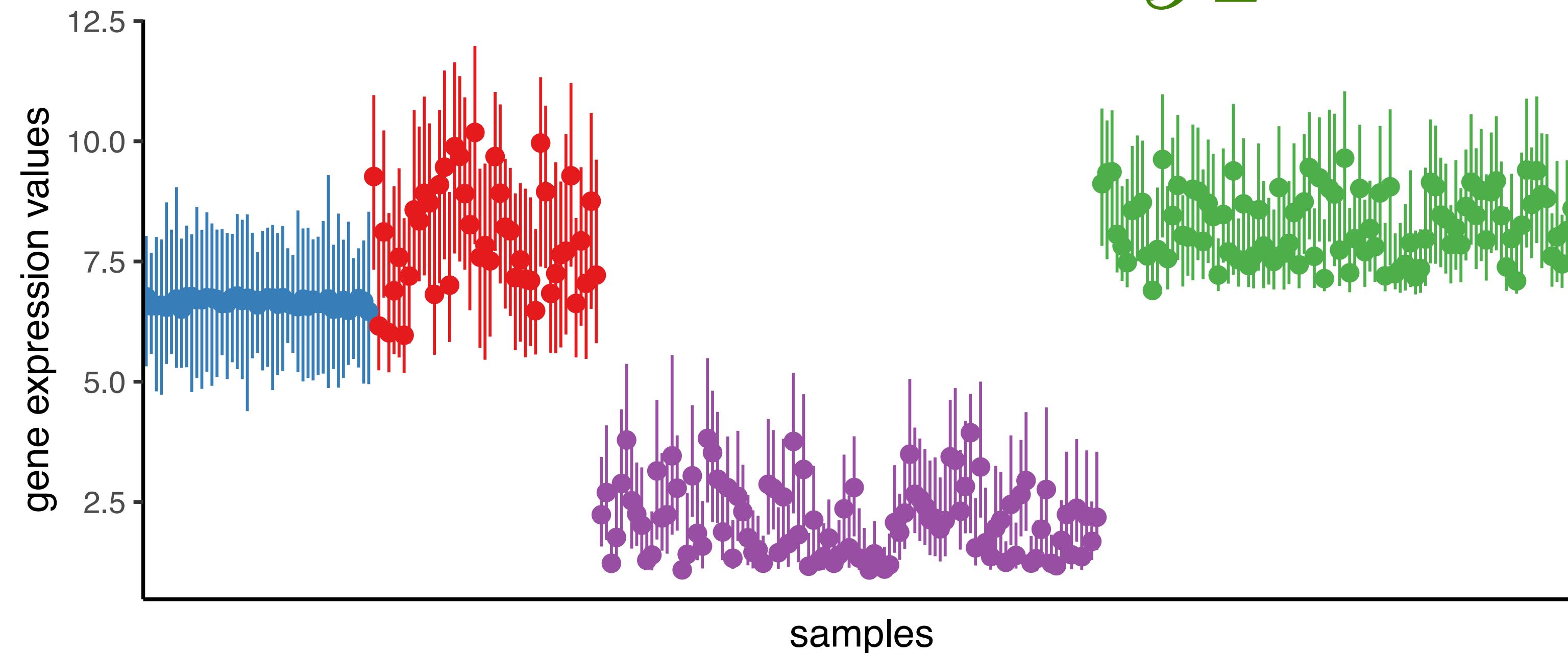


Transferability
For the same samples,
the prediction from one gene
expression platform
should be equivalent to
another platform

Omics-based clinical risk score: what is so difficult?

$$\hat{y}_1 = X_1 \hat{\beta}_1$$

$$\hat{y}_2 = X_2 \hat{\beta}_1$$



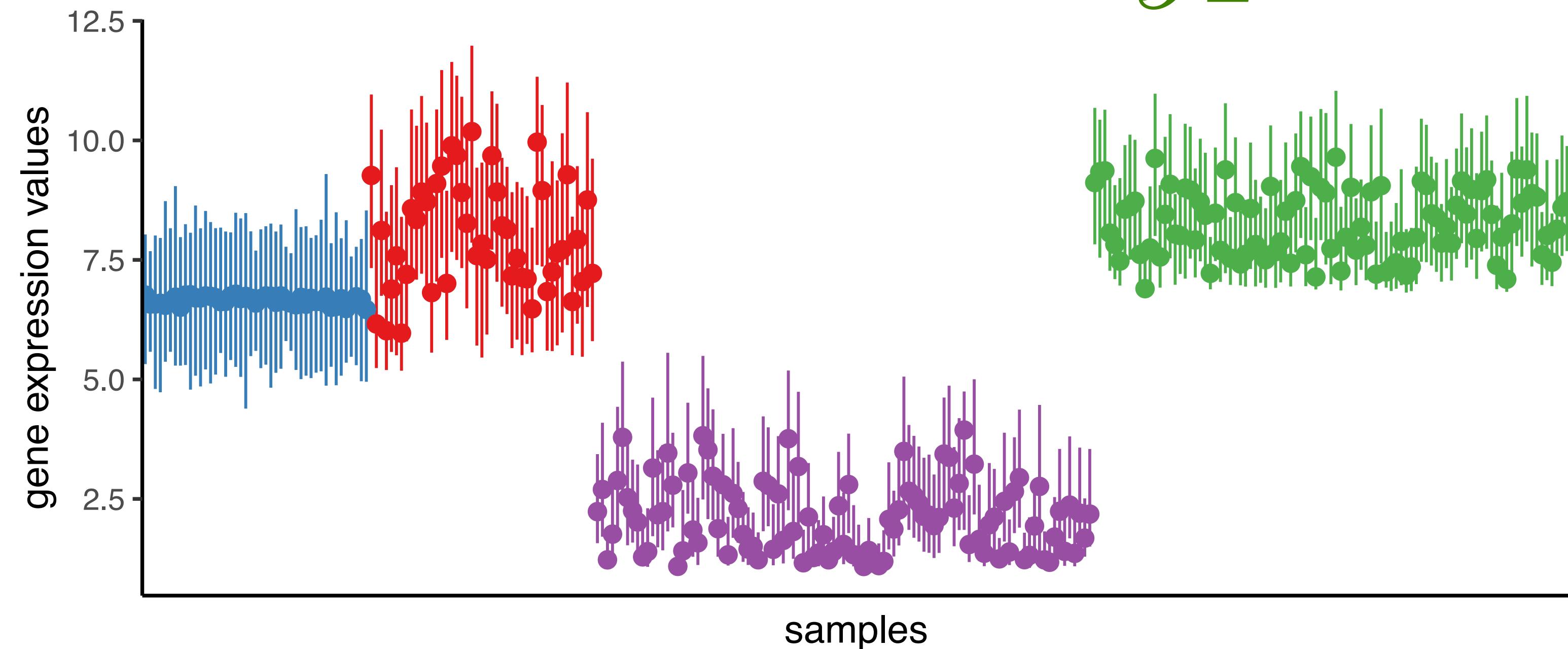
Omics features are on a
relative and unit-less scale,
a typical value in one data
can be an impossible value
on another

Omics-based clinical risk score: what is so difficult?

Assuming a
noiseless shift

$$\hat{y}_1 = X_1 \hat{\beta}_1$$

$$\hat{y}_2 = X_2 \hat{\beta}_1 = (X_1 + 1) \hat{\beta}_1$$



Omics features are on a
relative and unit-less scale,
a typical value in one data
can be an impossible value
on another

Components of a risk score

Data

$$(X_1, y_1)$$

$$(X_2, y_2)$$

Model

$$\hat{\beta}_1$$

Prediction

$$\hat{y}_1 = X_1 \hat{\beta}_1$$

$$\hat{y}_2 = X_2 \hat{\beta}_1$$

Components of a risk score

Data

$$(X_1, y_1)$$

$$(X_2, y_2)$$

No renormalisation

Model

$$\hat{\beta}_1$$

No model retraining

Prediction

$$\hat{y}_1 = X_1 \hat{\beta}_1$$

$$\hat{y}_2 = X_2 \hat{\beta}_1$$

Scale-equivalent prediction

Existing approach vs CPOP

Existing approach: data-harmonisation

Data

$$(X_1, y_1)$$

$$(X_2, y_2)$$

Harmonisation

Model

$$\hat{\beta}_1$$

Classical estimation
methods

Prediction

$$\hat{y}_1 = X_1 \hat{\beta}_1$$

$$\hat{y}_2 = X_2 \hat{\beta}_1$$

Classical model prediction

1. Data-harmonisation: normalisation or standardisation

Data

(X_1, y_1)

(X_2, y_2)

Harmonisation

Introduces strong dependence between training samples and models

1. Prospective: re-normalisation upon new single-samples
2. Multi-centres: re-training of model upon new populations

2. Model-based approach: CPOP

Data

$$(X_1, y_1) \rightarrow (\textcolor{red}{Z}_1, y_1)$$

$$(X_2, y_2) \rightarrow (\textcolor{red}{Z}_2, y_2)$$

Model

$$\hat{\beta}_1 \approx \hat{\beta}_2$$

Prediction

$$\textcolor{red}{Z}_1 \hat{\beta}_1 \approx Z_1 \hat{\beta}_2$$

$$Z_2 \hat{\beta}_1 \approx \textcolor{red}{Z}_2 \hat{\beta}_2$$

Feature transform

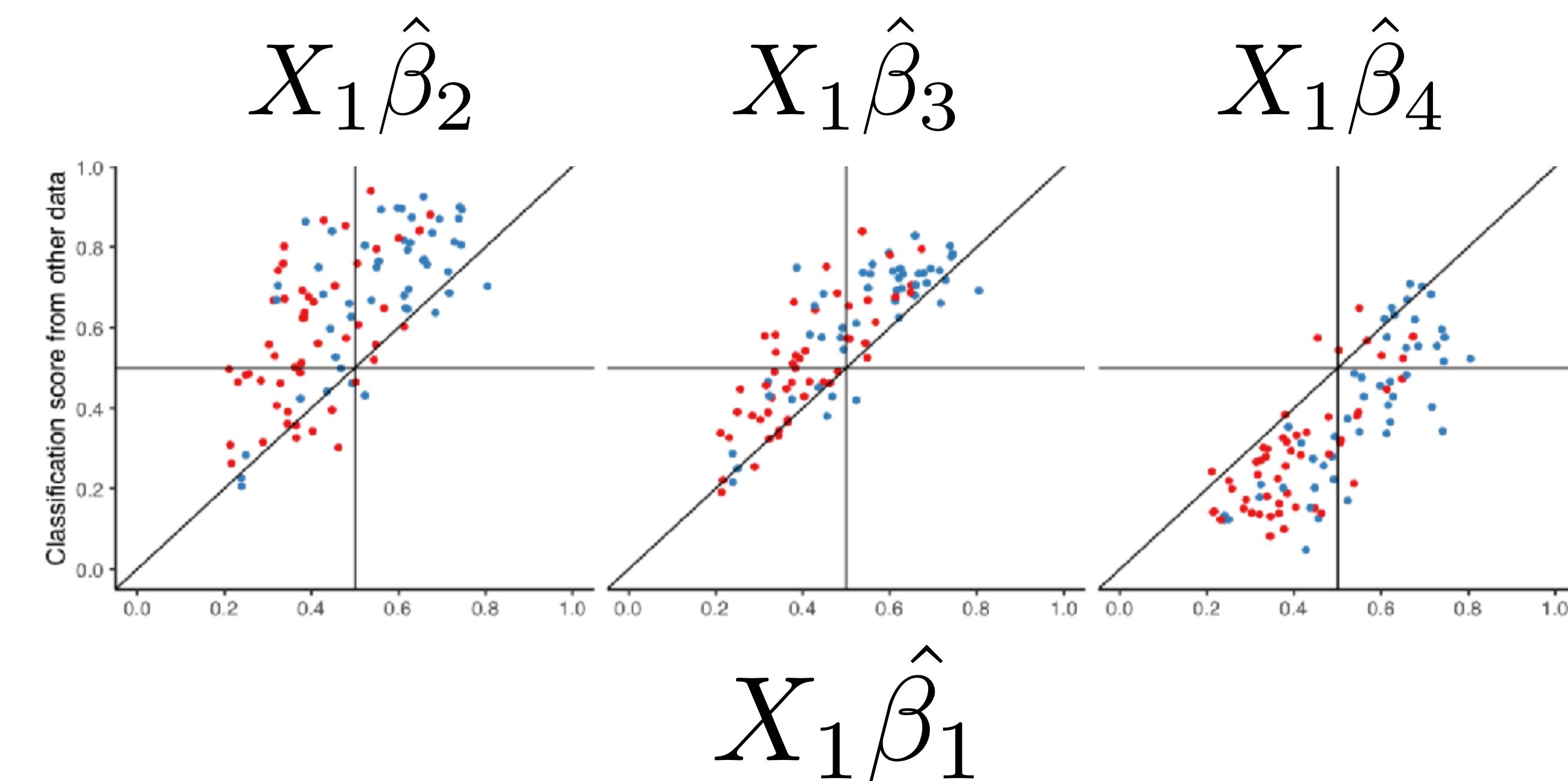
Stable estimation

Stable prediction

Statistical challenges

1. Concordance in features scaling across datasets
2. Concordance in feature selection and estimates
3. Single-patient prediction

Transferability implies that patients should be evenly scattered around the identity line



First component of CPOP: feature transform



就让我来次透彻心扉的痛
都拿走 让我再次两手空空
只有奄奄一息过
那个真正的我
他才能够诞生

Log-ratio transformation

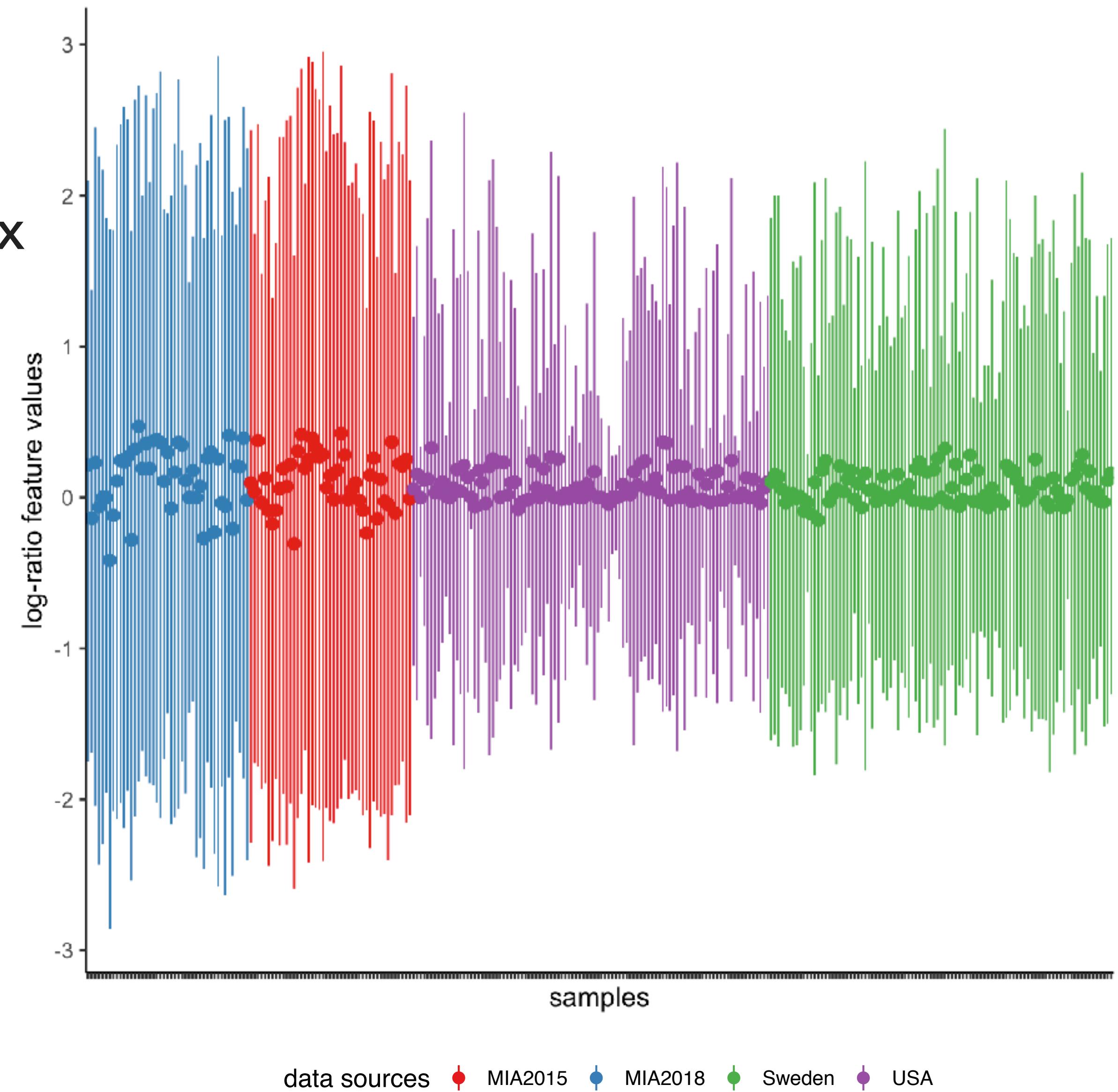
For each column in the gene expression matrix

$$X = \{x_1, \dots, x_p\} \in \mathbb{R}^{n \times p}$$

Construct $Z \in \mathbb{R}^{n \times \binom{p}{2}}$ column-wise:

$$z_j = \log\left(\frac{x_l}{x_m}\right)$$

$$\text{for } 1 \leq l < m \leq p$$

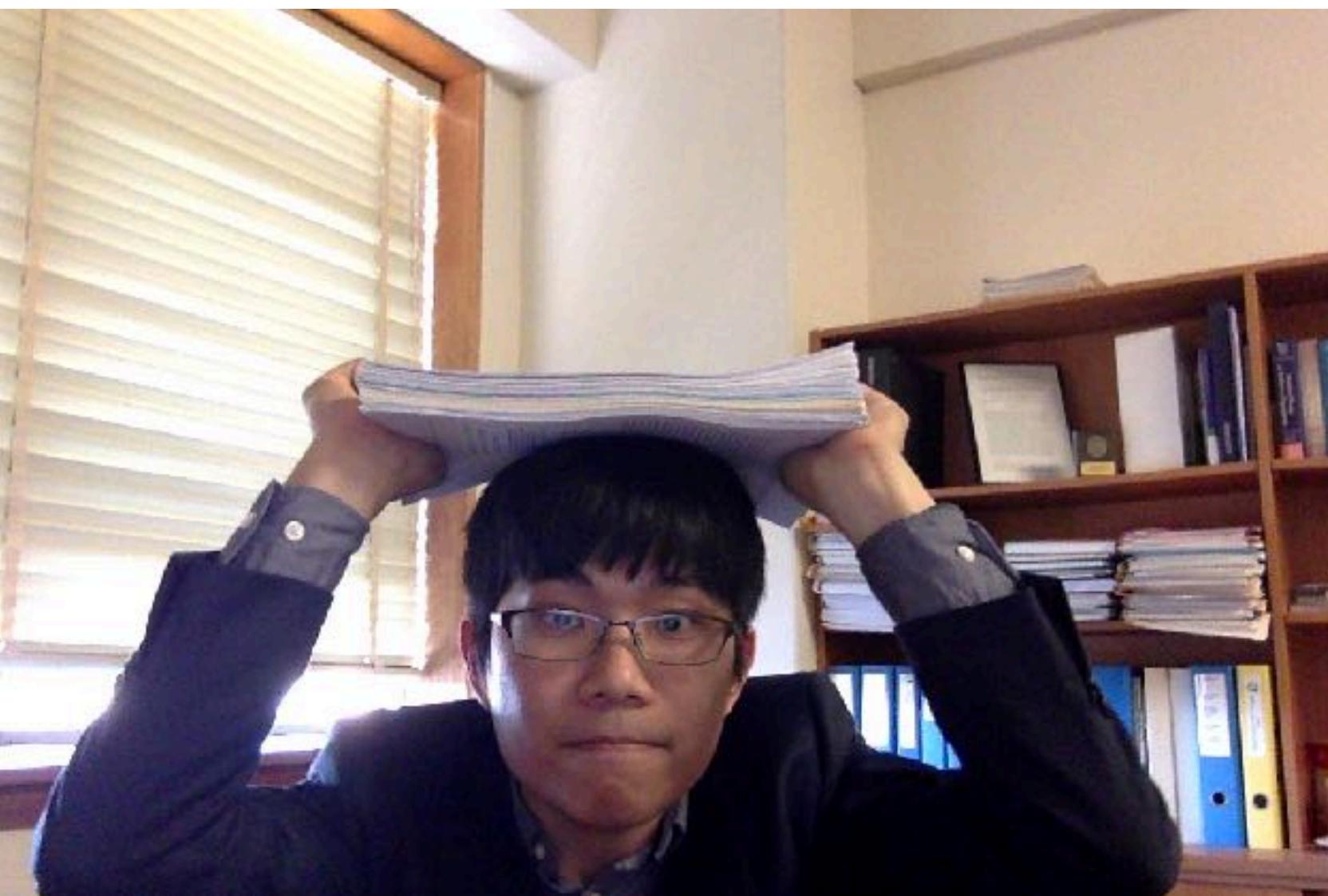


Why log-ratios?

- ▶ Within-sample standardisation avoids re-normalisation and model re-training
- ▶ Under-used in patient outcome prediction
- ▶ Potentials for further method developments

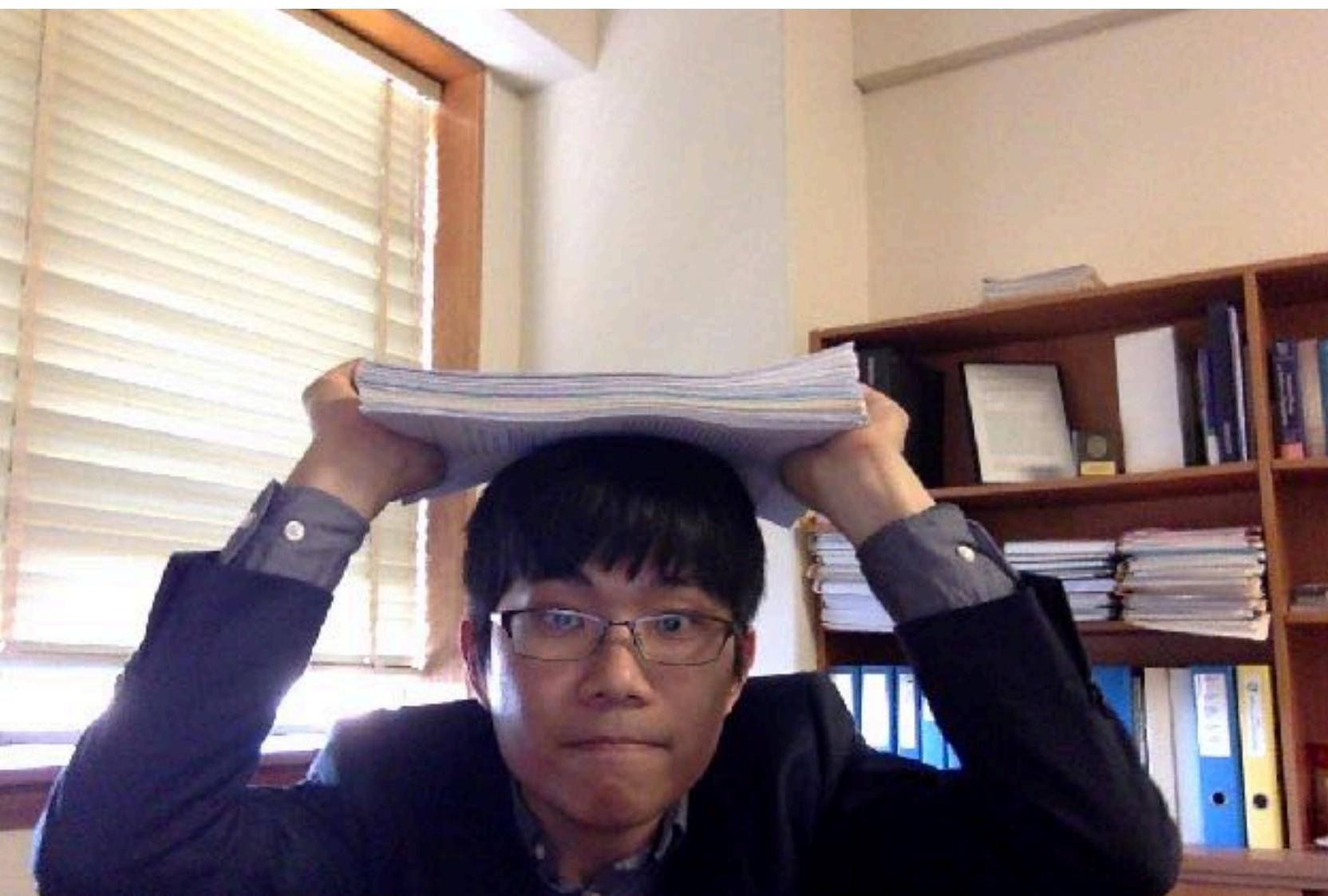
Why log-ratios?

- ▶ Within-sample standardisation avoids re-normalisation and model re-training
- ▶ Under-used in patient outcome prediction
- ▶ Potentials for further method developments

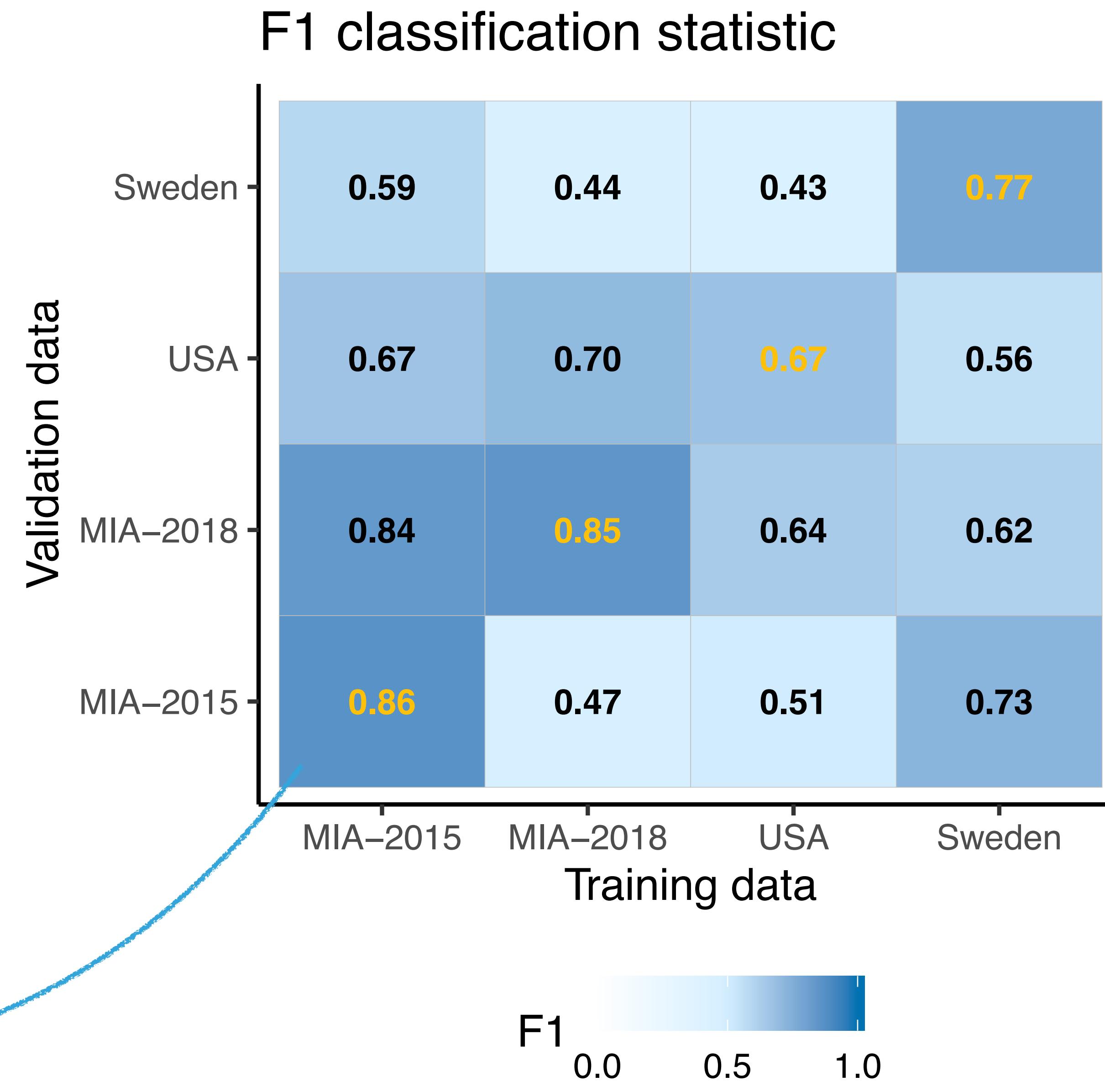
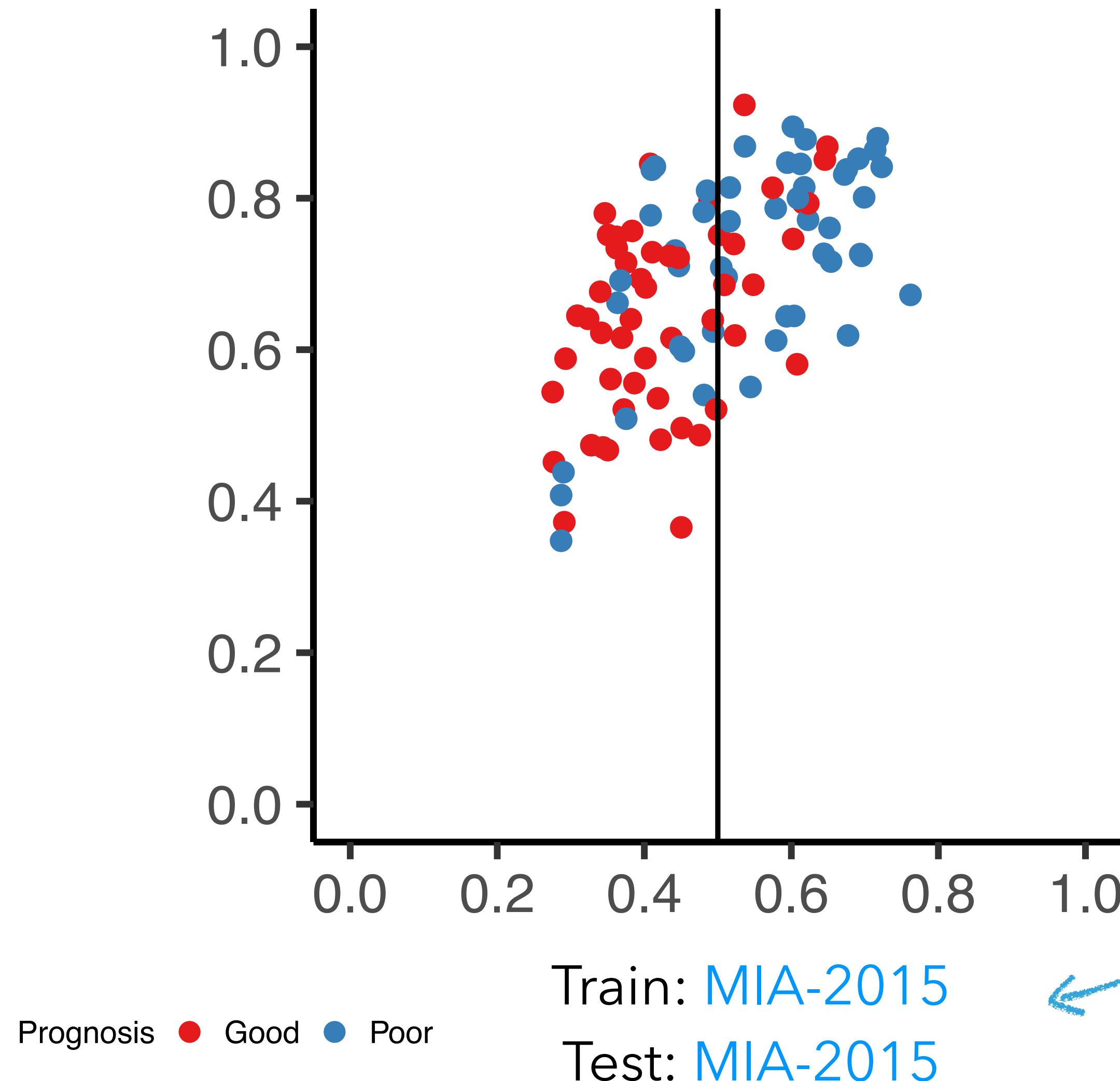


Why log-ratios?

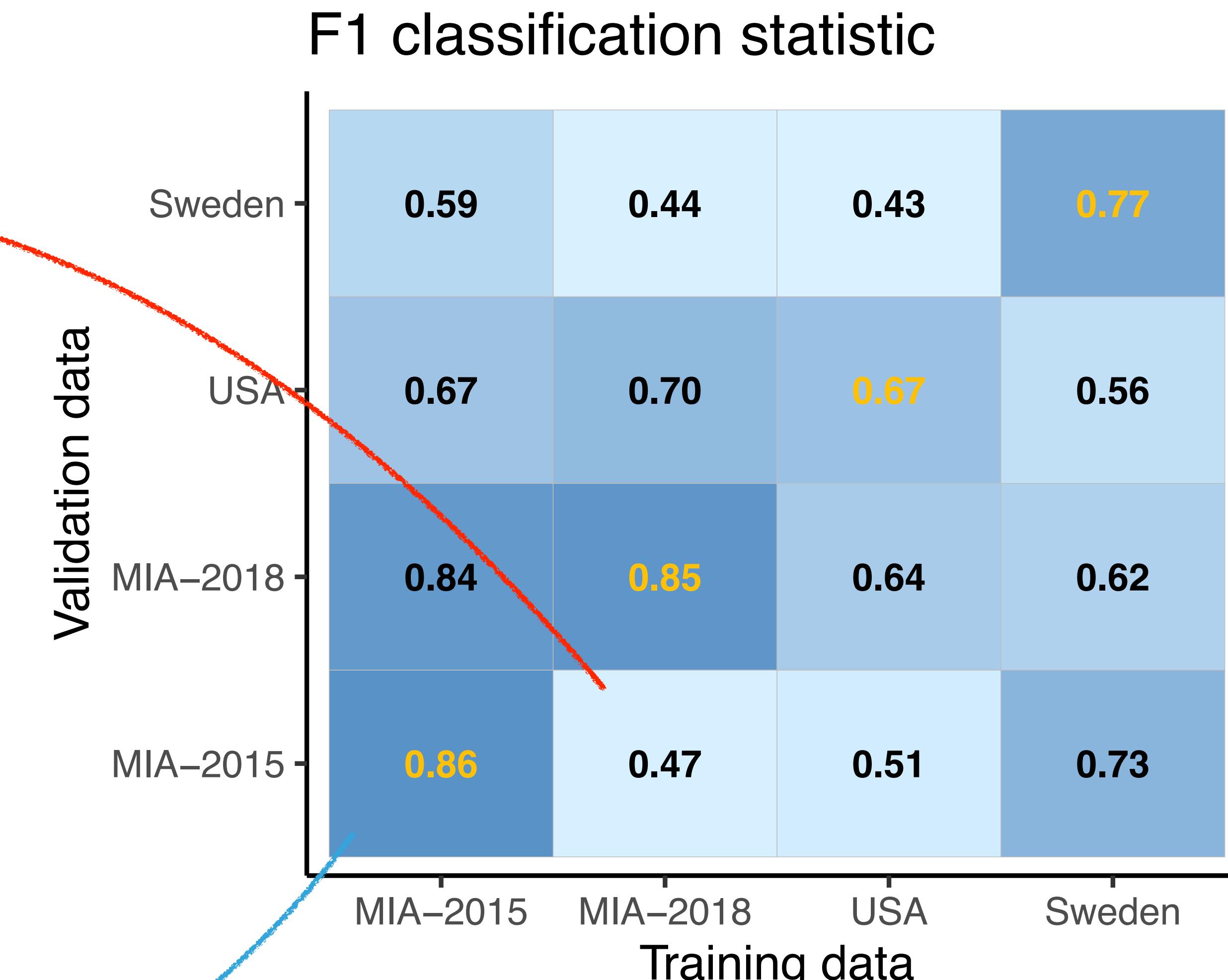
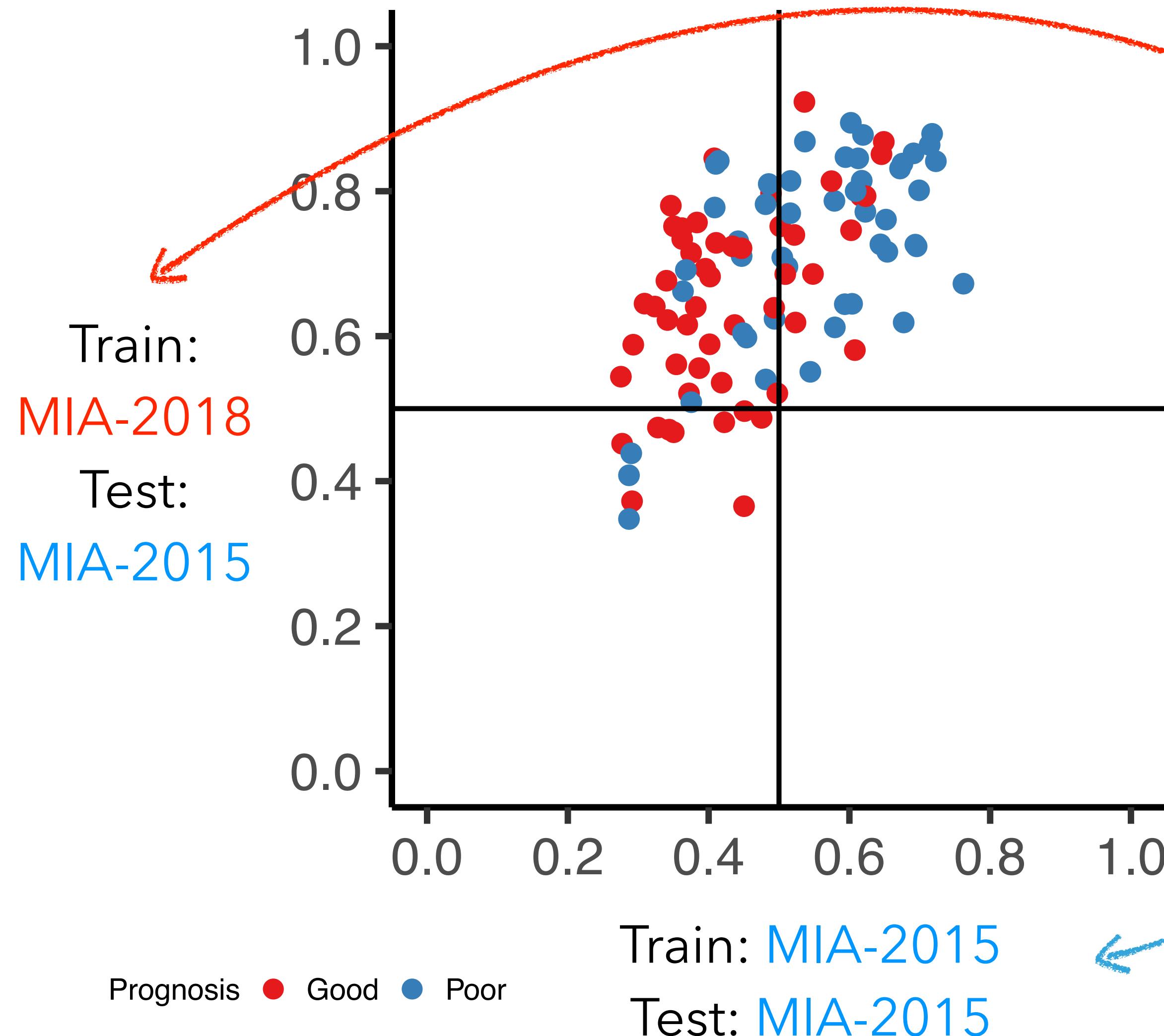
- ▶ Within-sample standardisation avoids re-normalisation and model re-training
- ▶ Under-used in patient outcome prediction
- ▶ Potentials for further method developments



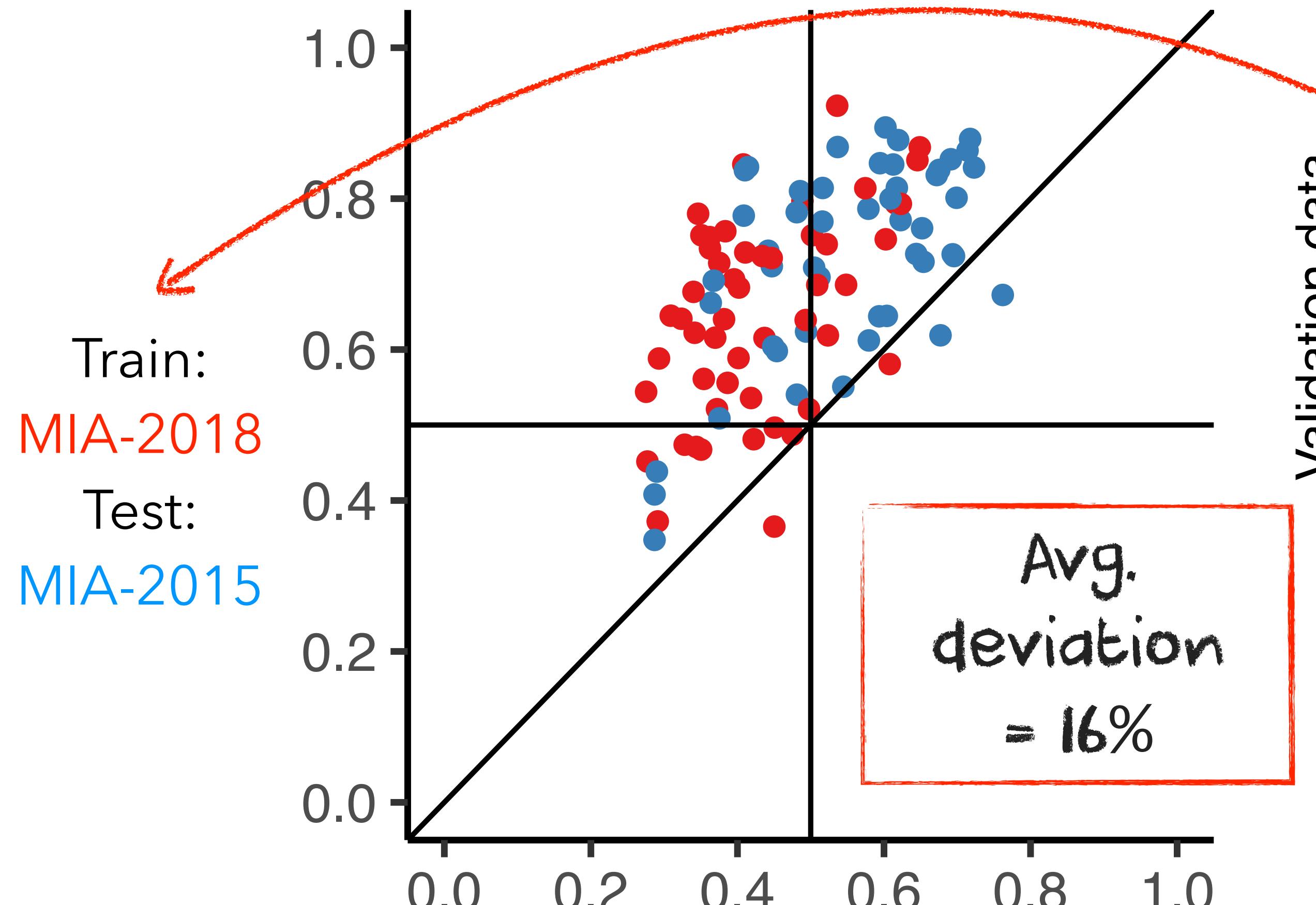
Is log-ratio enough?



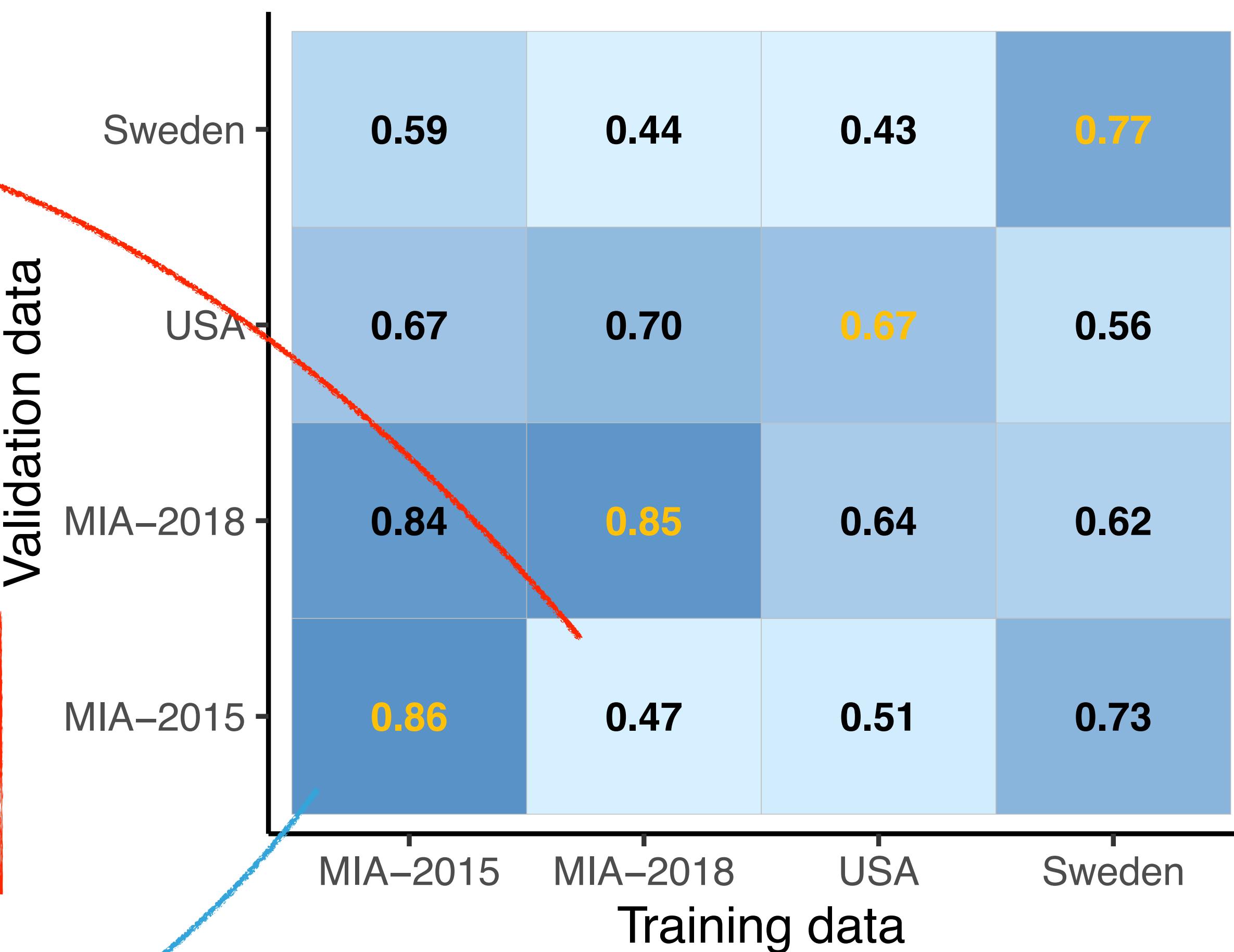
Is log-ratio enough?



Is log-ratio enough?

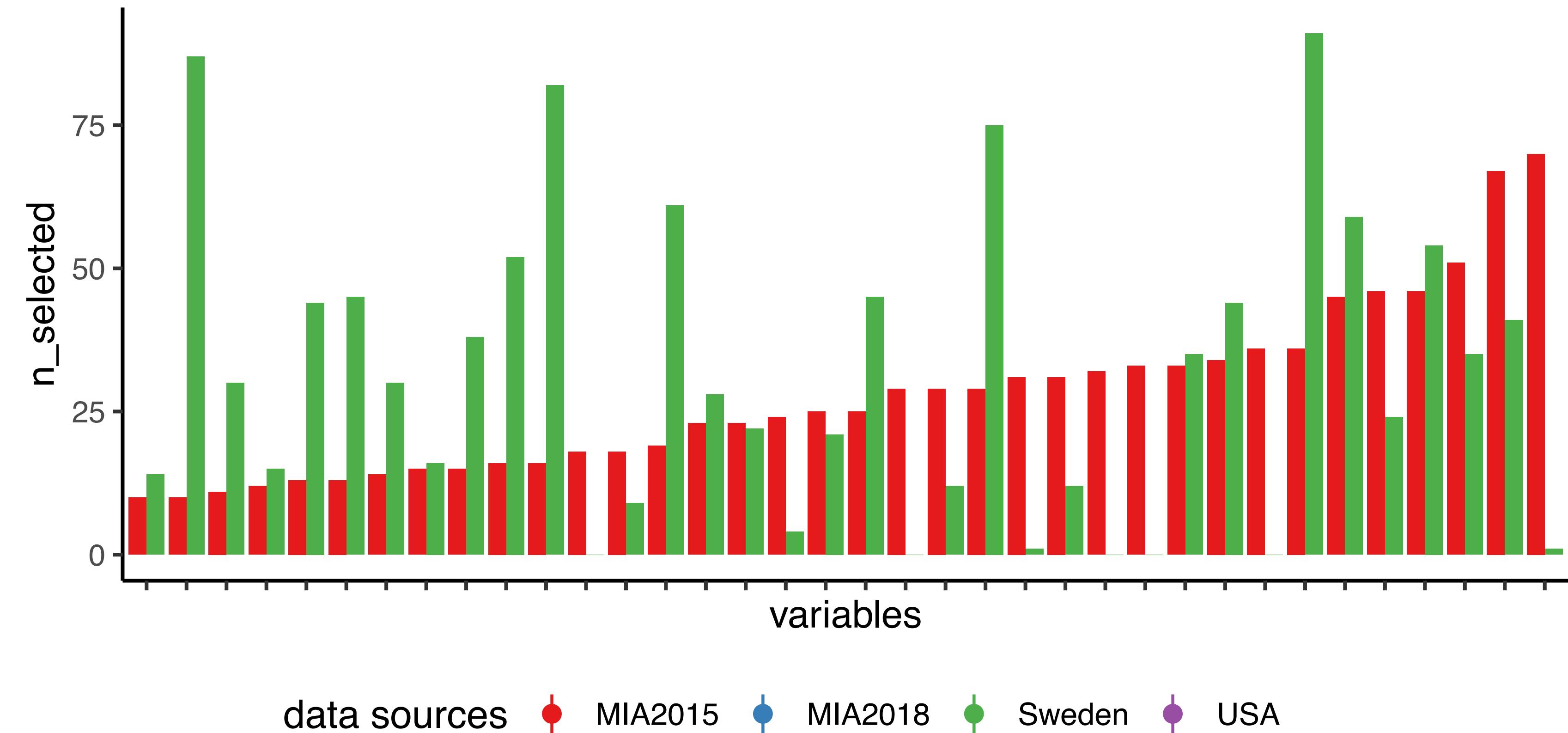


F1 classification statistic



Lasso variable selection is not stable!

- ▶ Bach 2008
 - ▶ Meinshausen & Bühlmann 2010
 - ▶ Lim & Yu 2016
- all pointed out that Lasso is unstable under cross-validation



Second component of CPOP: stable feature selection and estimation



我曾经毁了我的一切
只想永远地离开
我曾经堕入无边黑暗
想挣扎无法自拔
我曾经像你像他像那野草野花
绝望着也渴望着
也哭也笑也平凡着

CPOP flowchart

Data

$$(X_1, y_1) \rightarrow (Z_1, y_1)$$

$$(X_2, y_2) \rightarrow (Z_2, y_2)$$

Model

$$\hat{\beta}_1 \approx \hat{\beta}_2$$

Prediction

$$Z_1 \hat{\beta}_1 \approx Z_1 \hat{\beta}_2$$

$$Z_2 \hat{\beta}_1 \approx Z_2 \hat{\beta}_2$$

Feature transform

Stable estimation

Stable prediction

Weighted Elastic Net

$$\hat{\beta}(y, Z) = \underset{\beta \in \mathbb{R}^{\binom{p}{2}}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - z_i^\top \beta)^2 + \lambda \sum_{j=1}^q w_j \left[\frac{(1-\alpha)}{2} \beta_j^2 + \alpha |\beta_j| \right]$$

Modelling on log-ratio

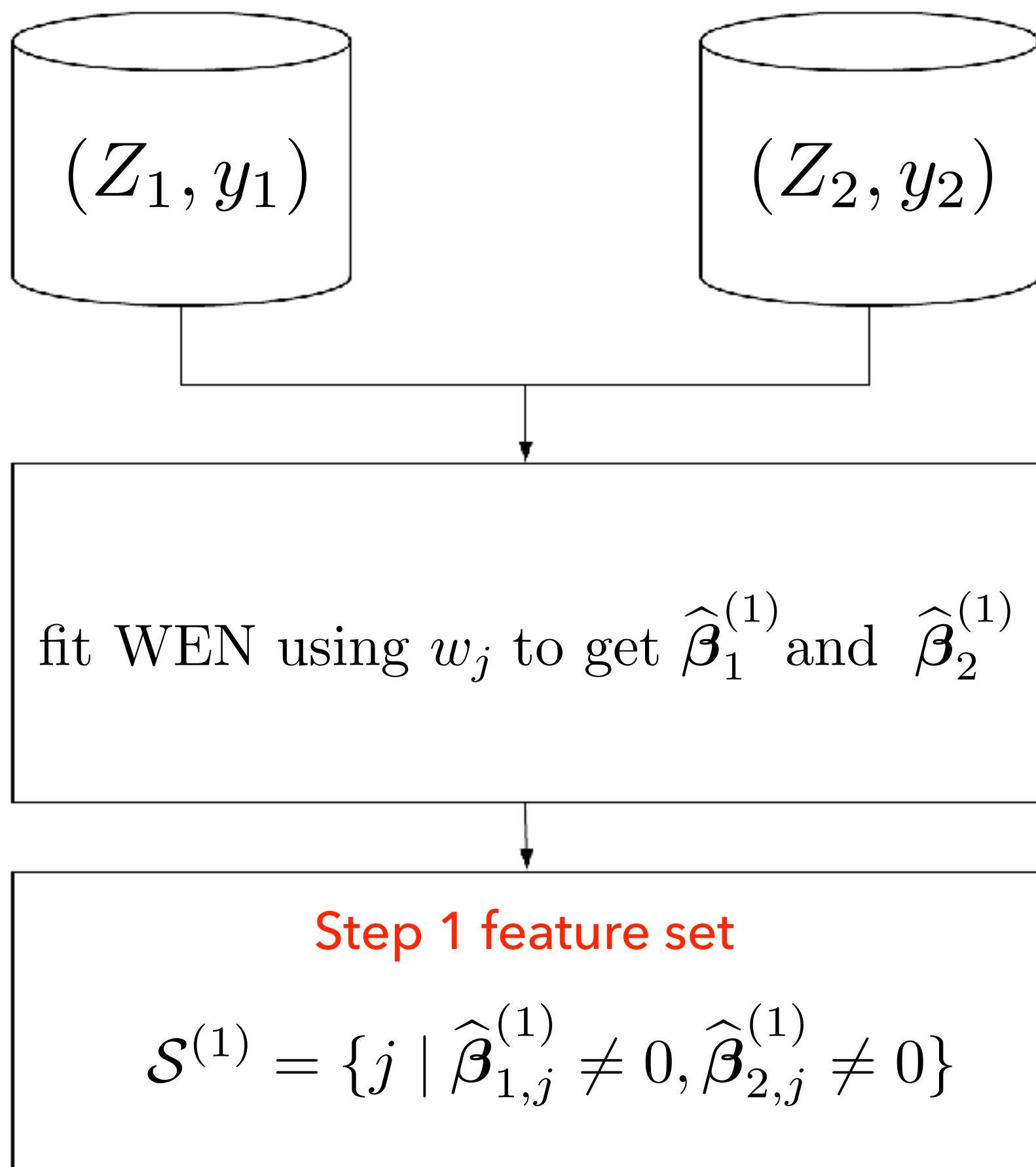
loss function,
could be
logistic or

Cox

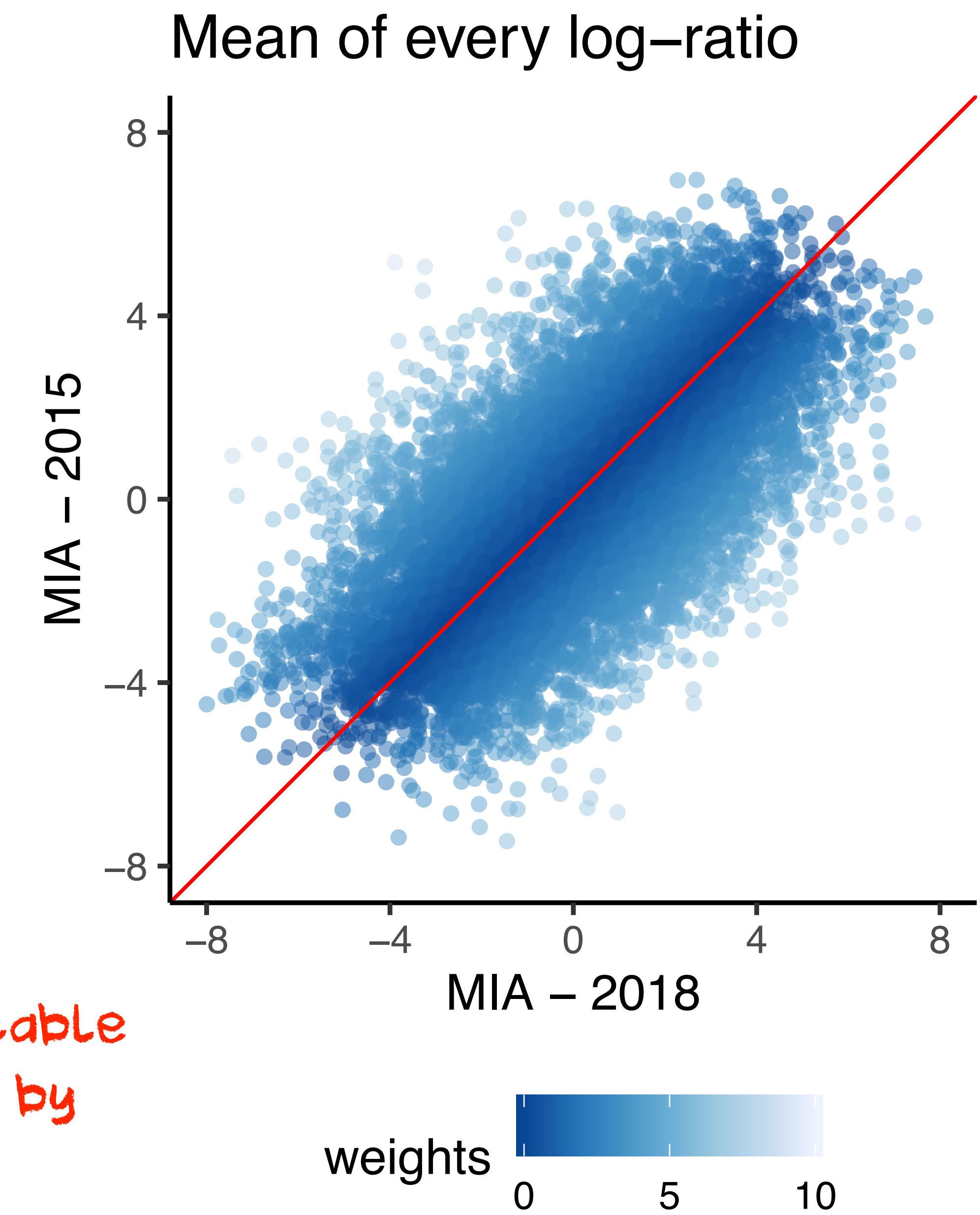
A mix of L1
and L2
penalties

Weights on each feature,
proportional to the
stability of features

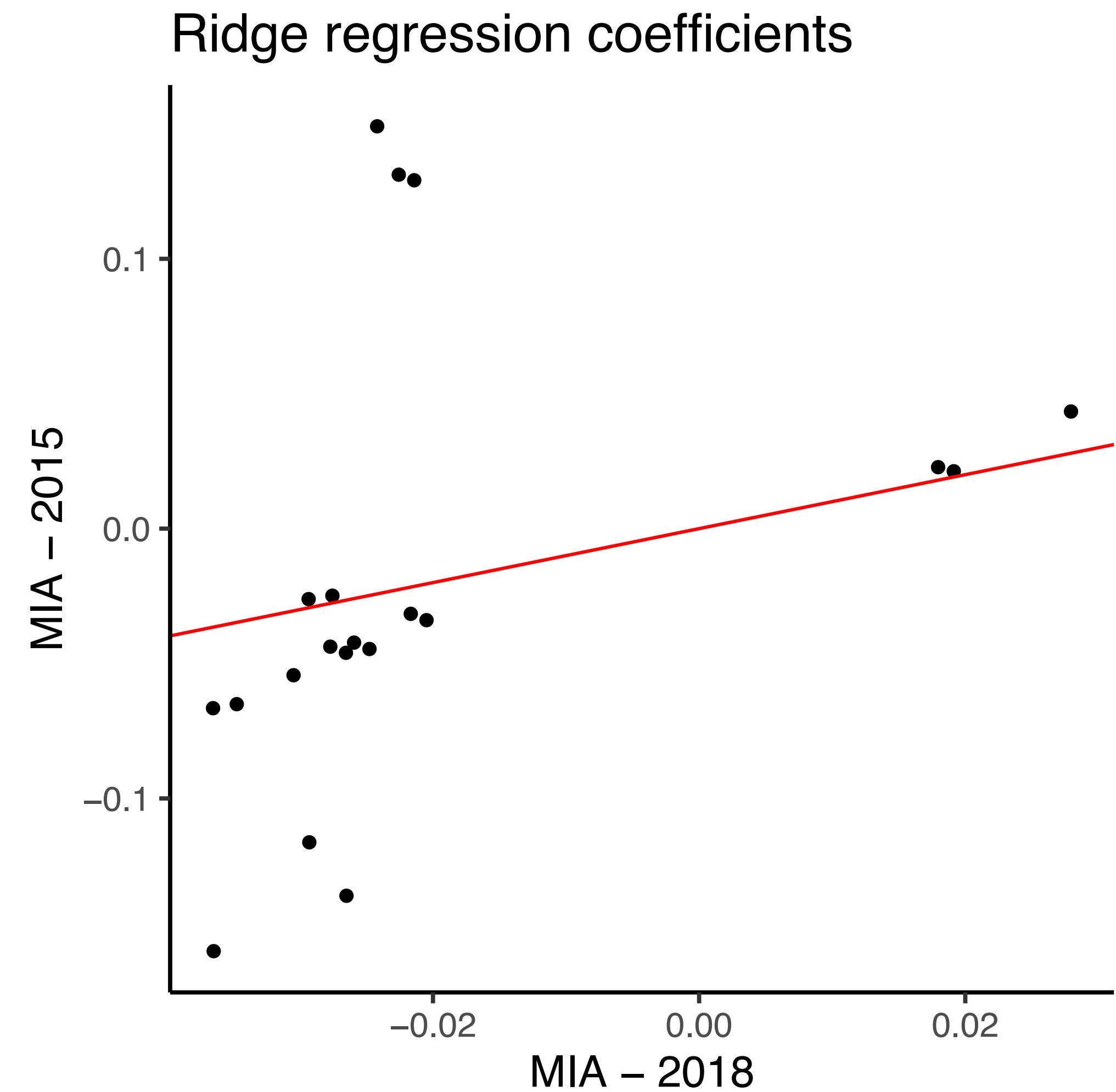
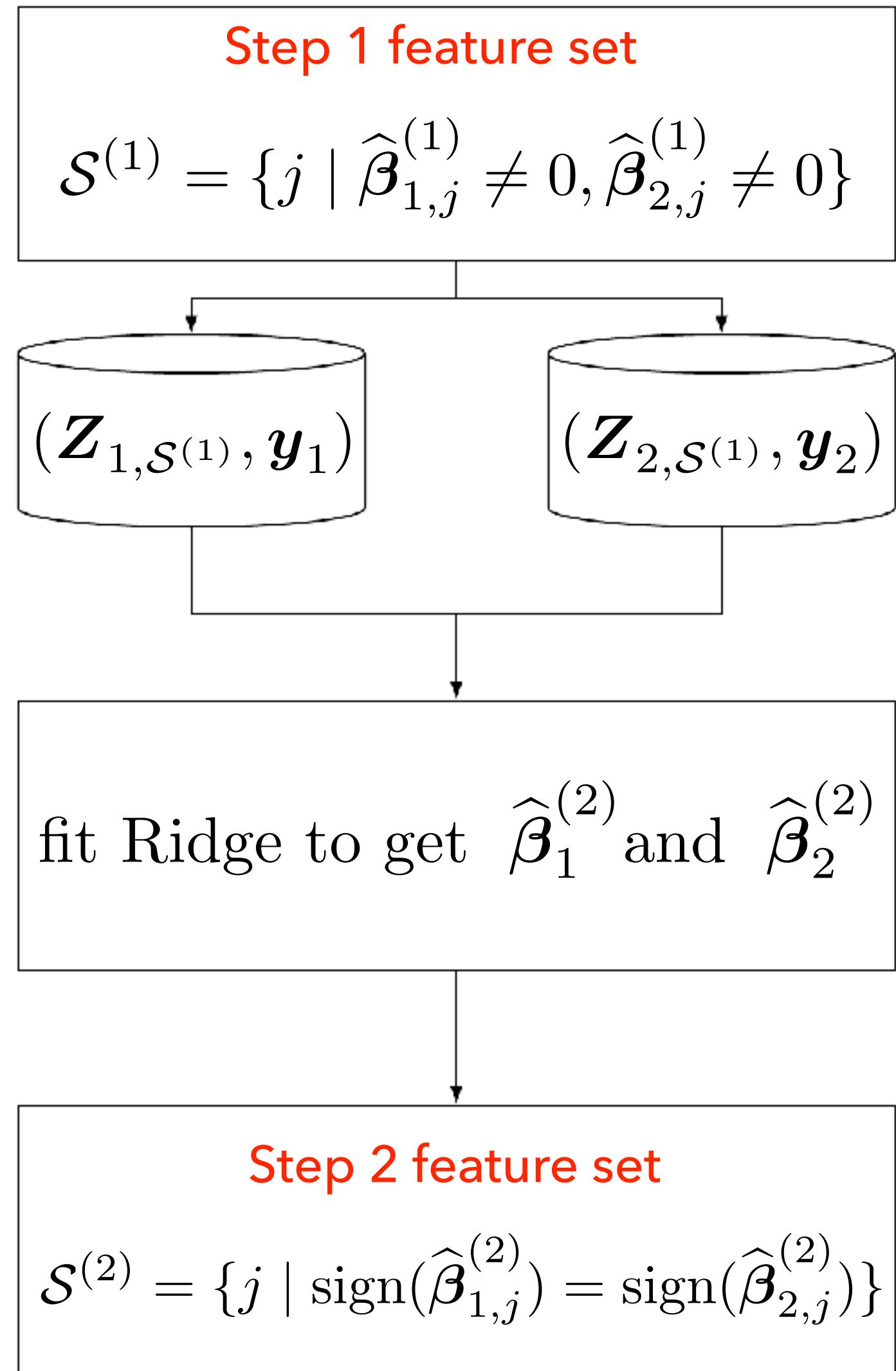
Step 1: feature selection stability



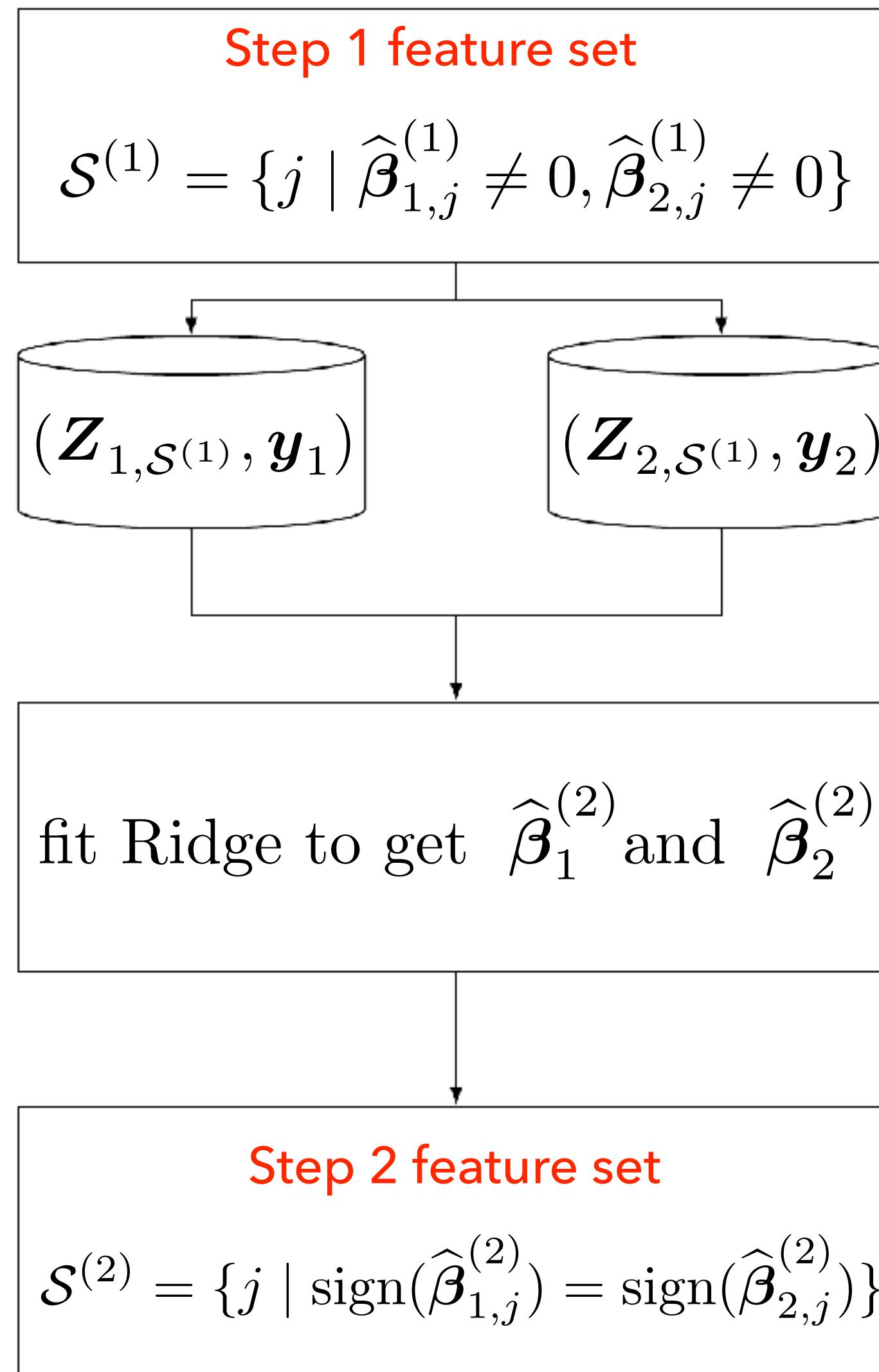
informative and stable
features agreed by
both data



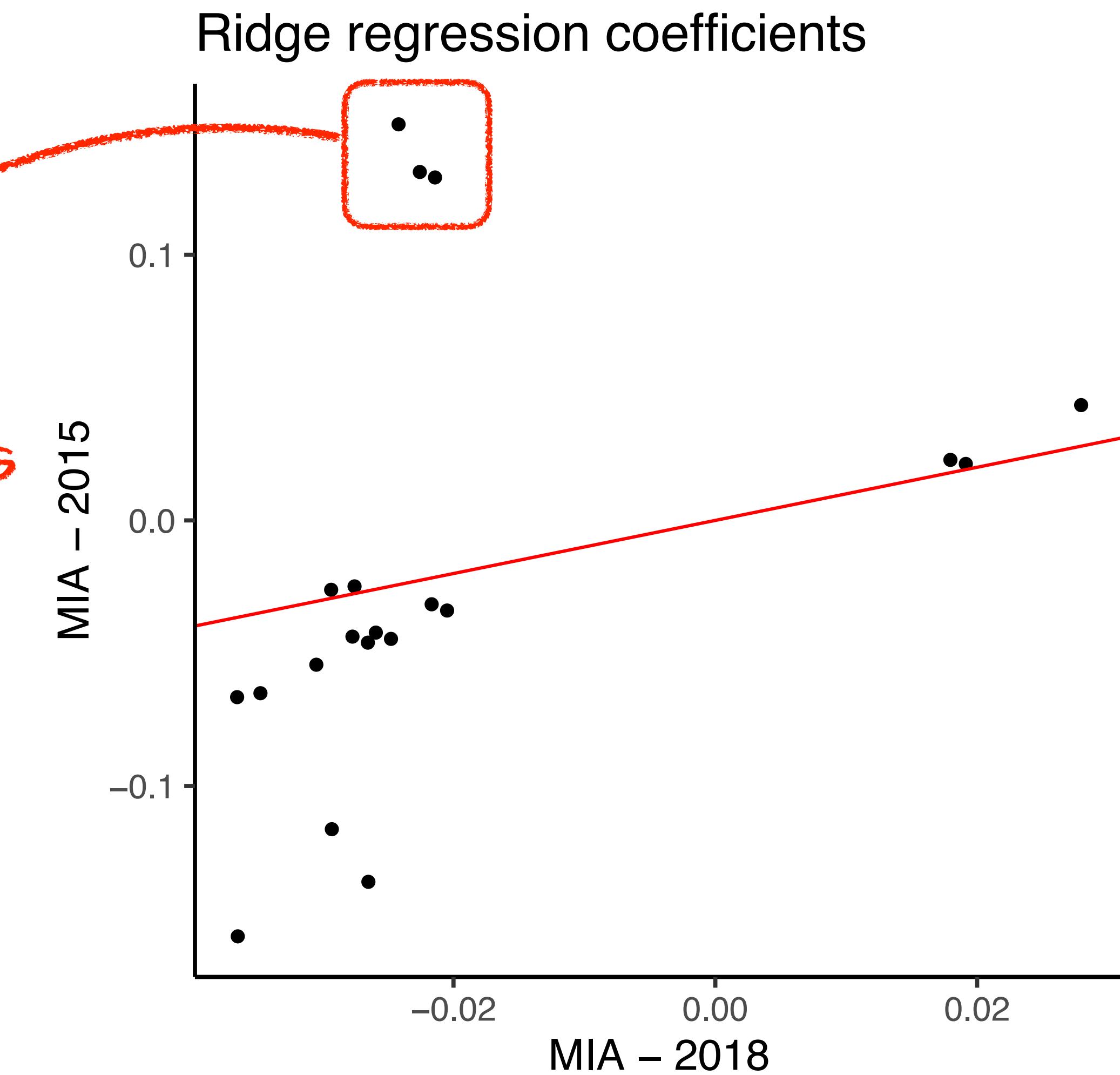
Step 2: feature estimation stability



Step 2: feature estimation stability

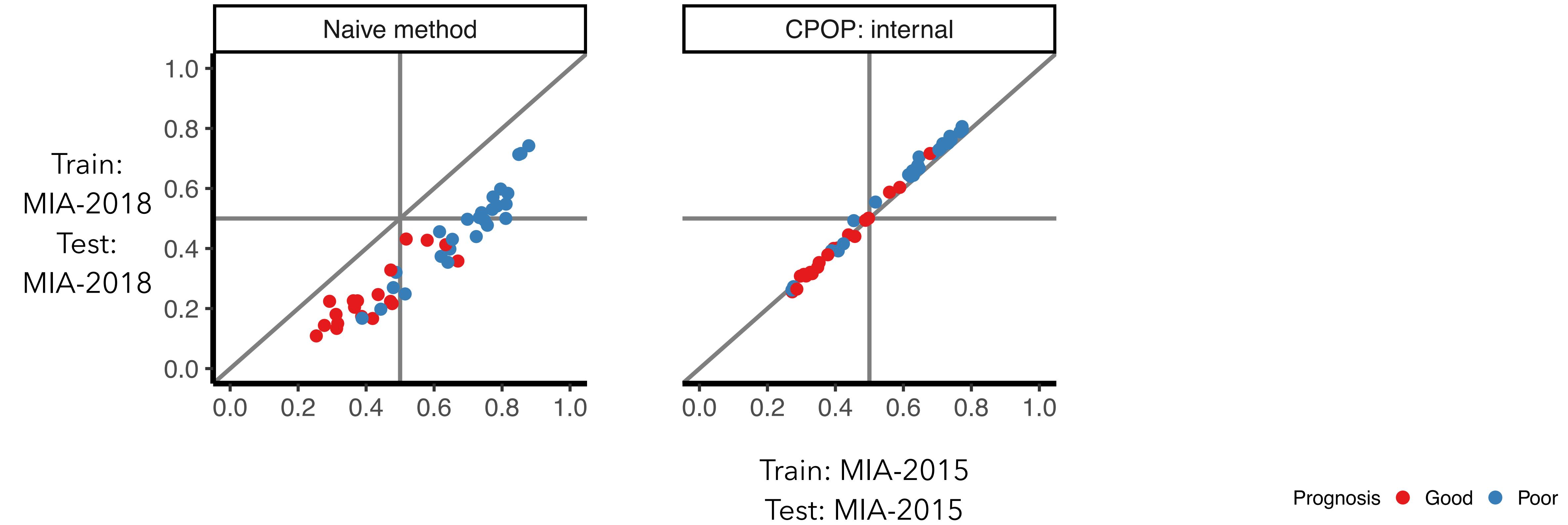


features with
unstable estimates
across data



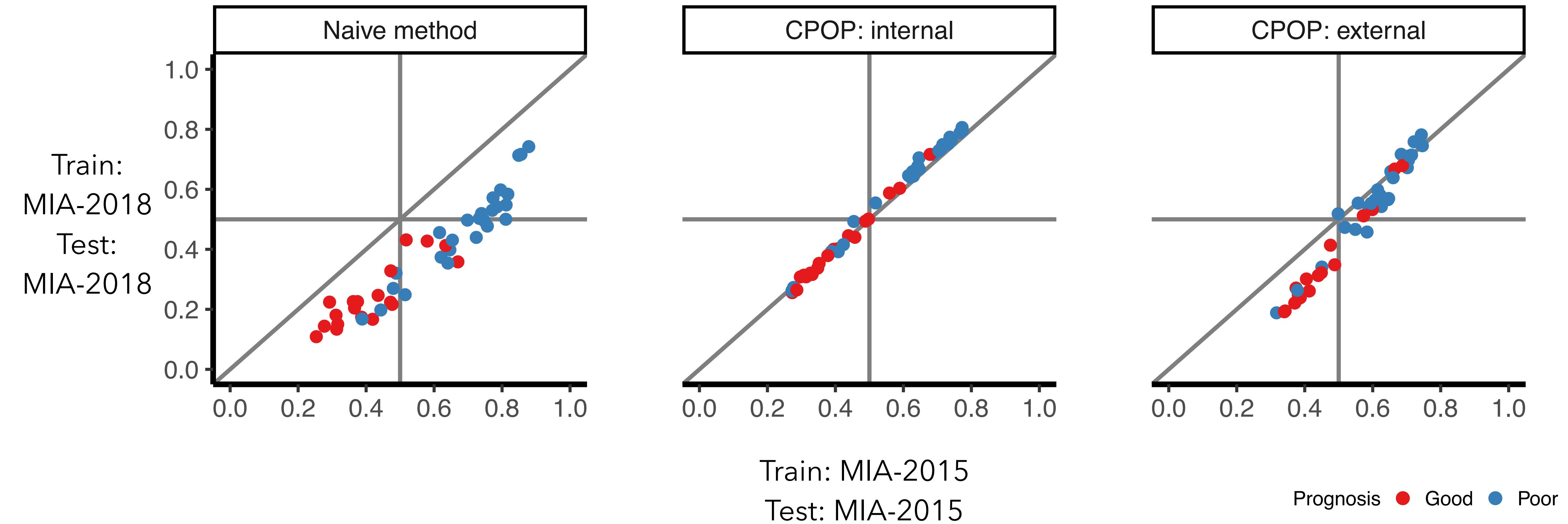
Results

CPOP results 1: MIA 2015 vs 2018 data



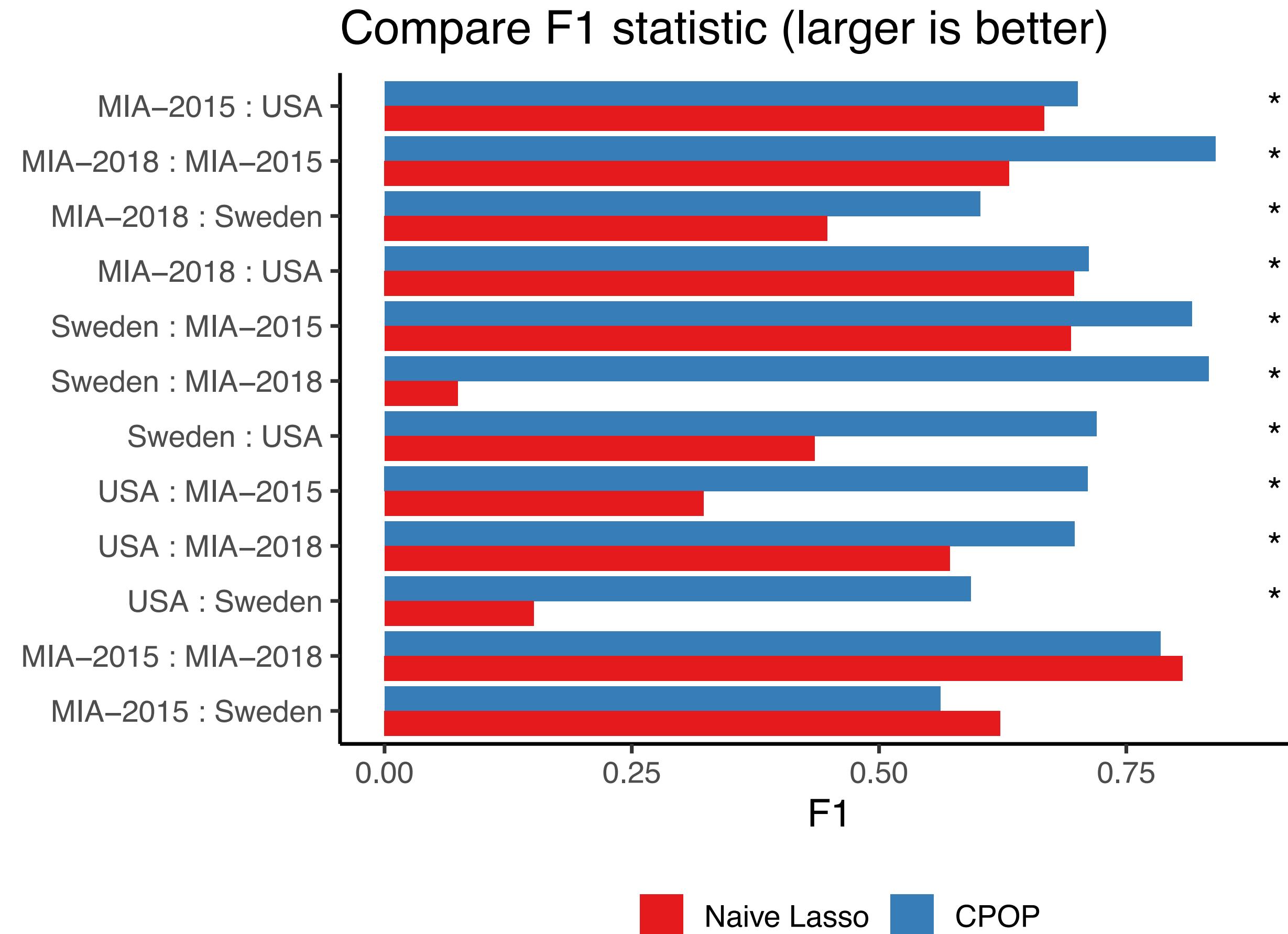
Small deviation in **predicted values** across datasets

CPOP results 1: MIA 2015 vs 2018 data



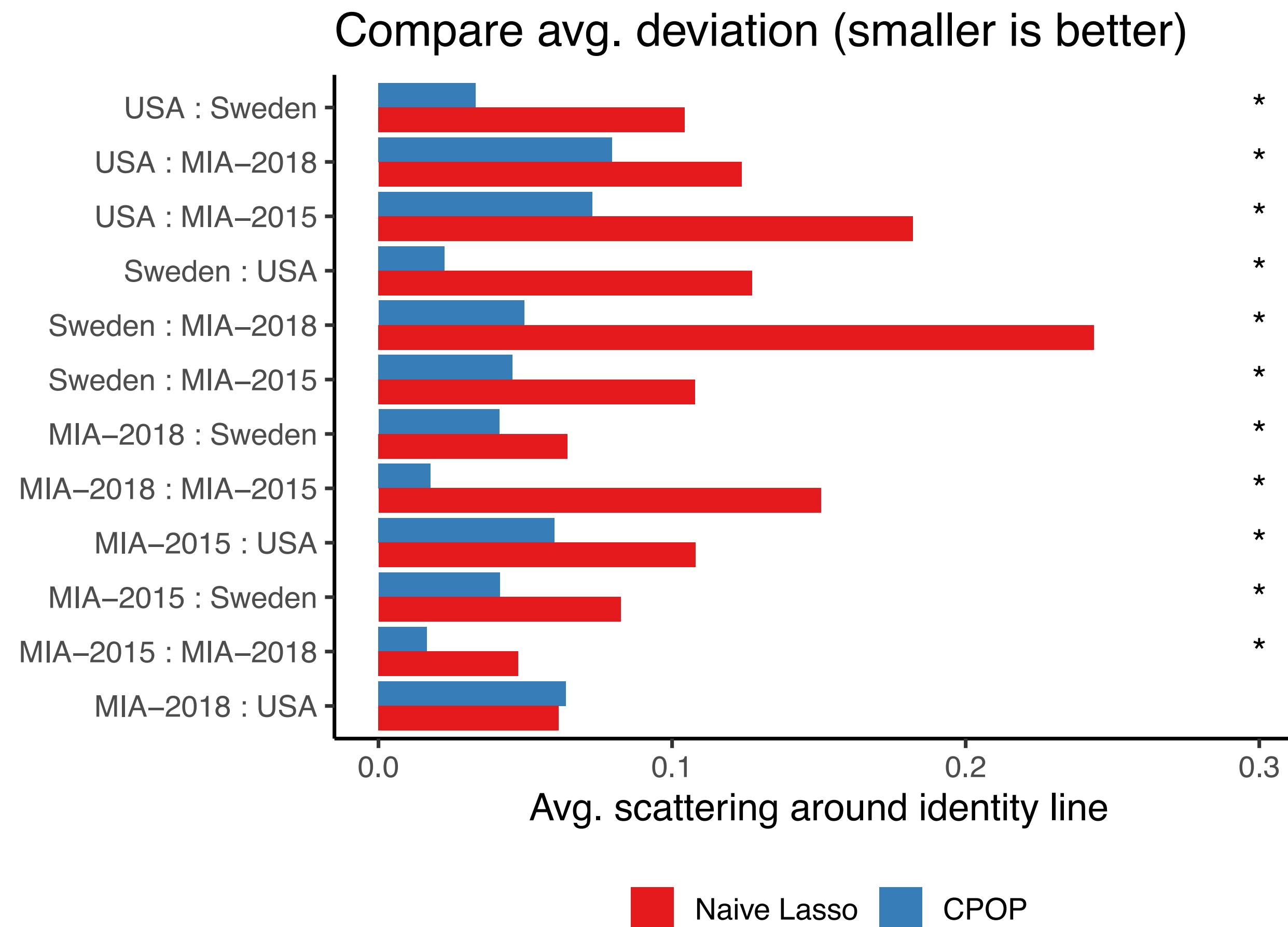
Small deviation in **predicted values** across datasets

CPOP results 2: four melanoma data



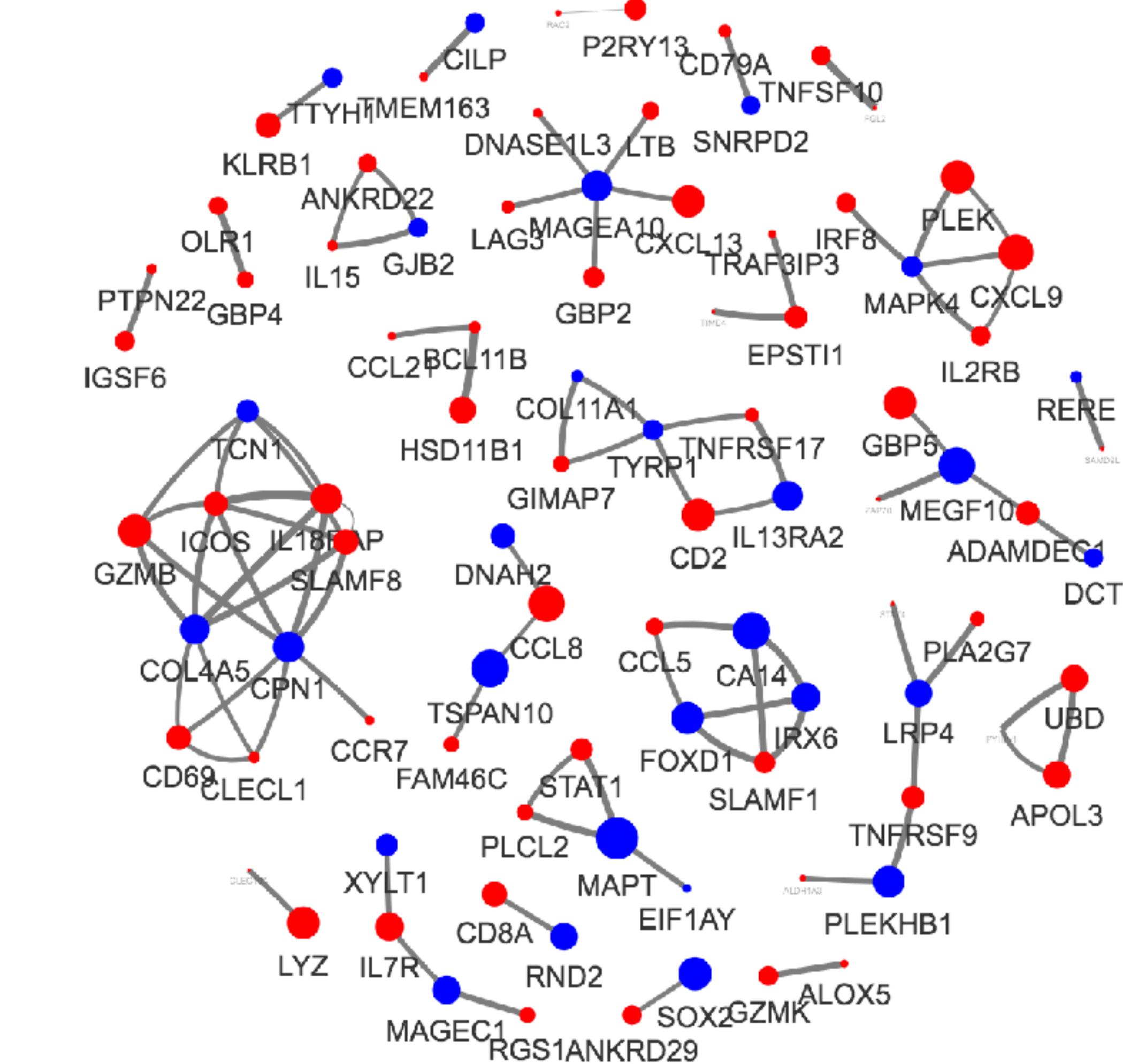
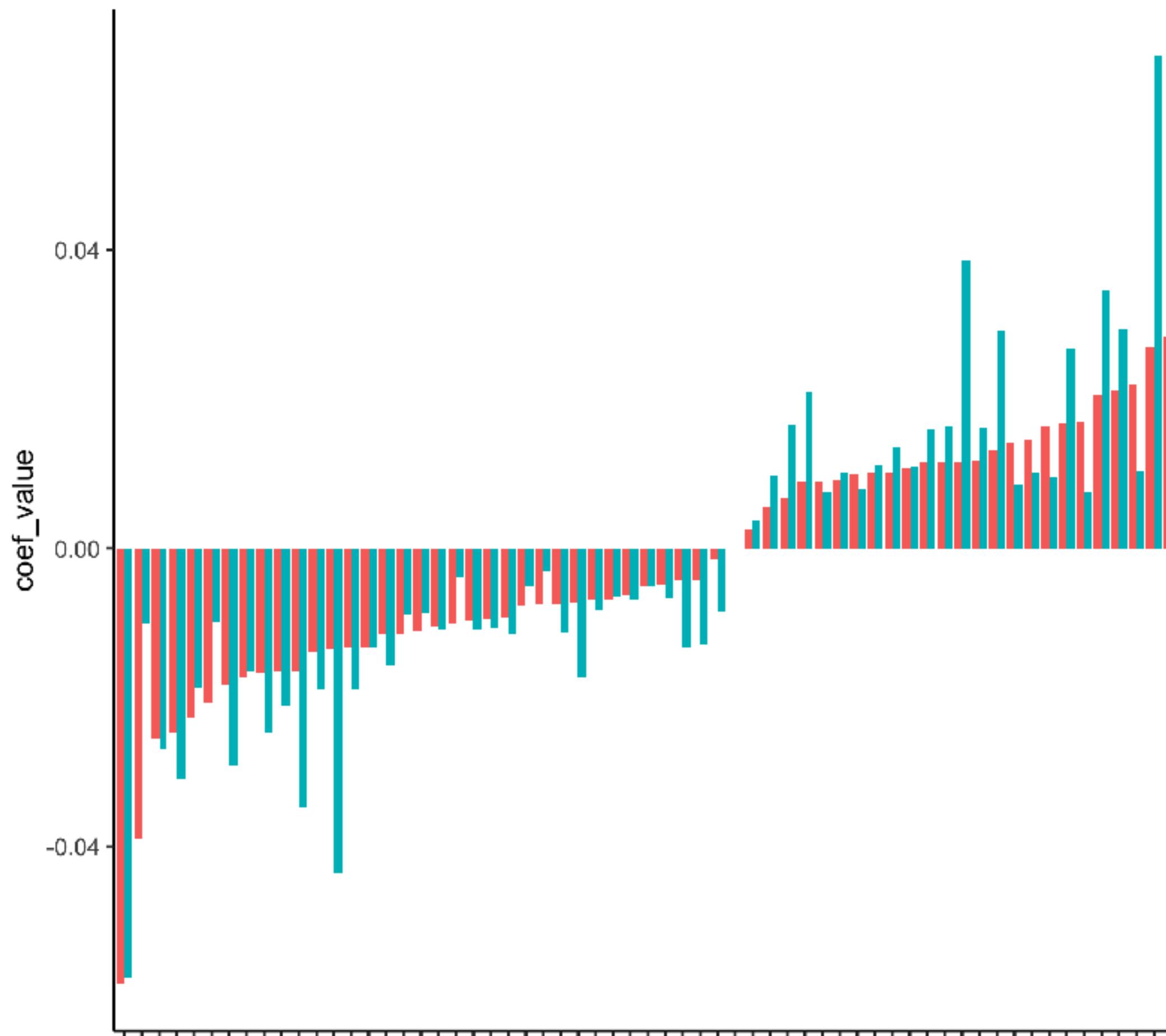
CPOP is highly predictive

CPOP results 2: four melanoma data

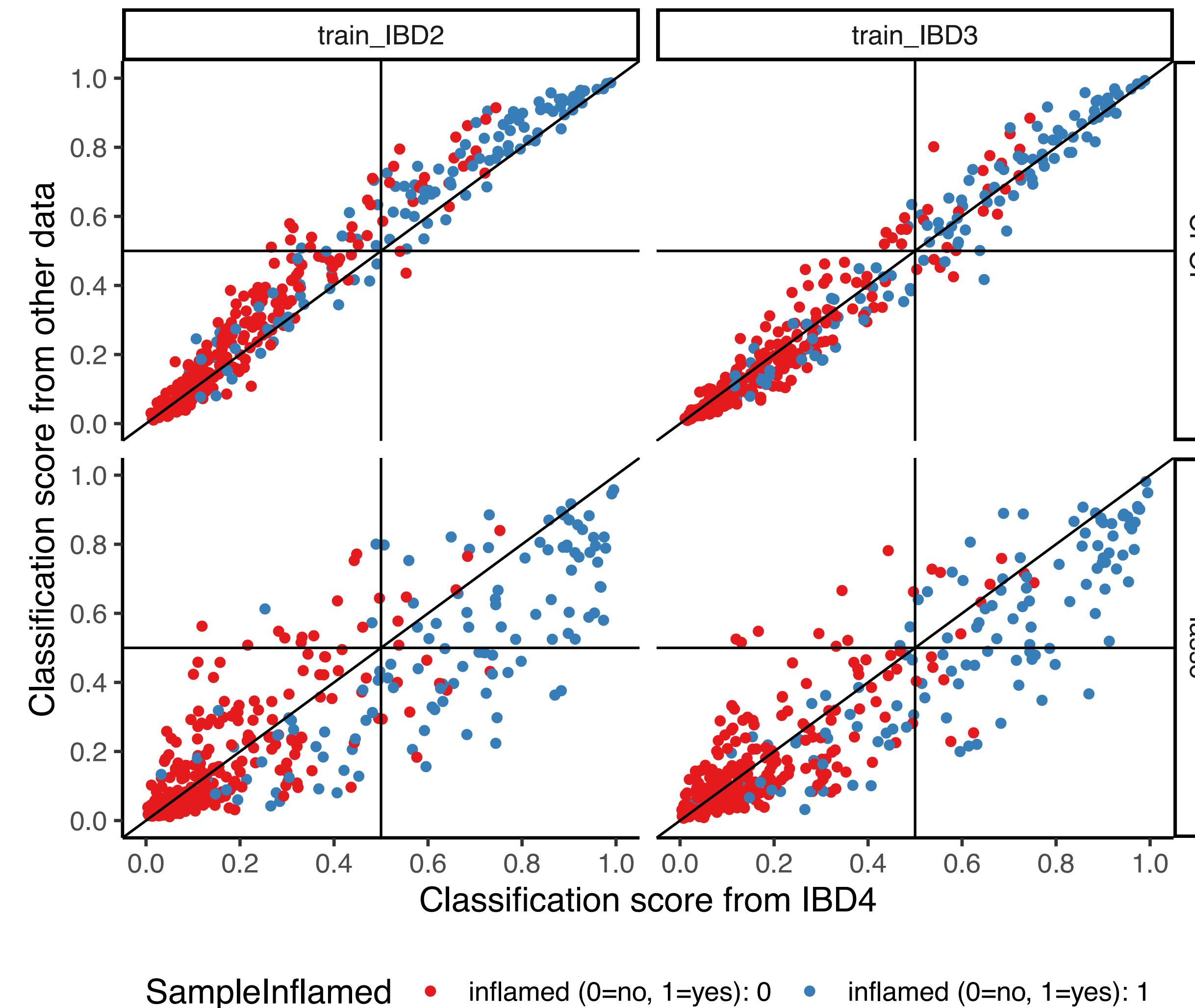


CPOP is highly stable

CPOP features



CPOP results 3: prospective prediction on inflammatory bowel disease



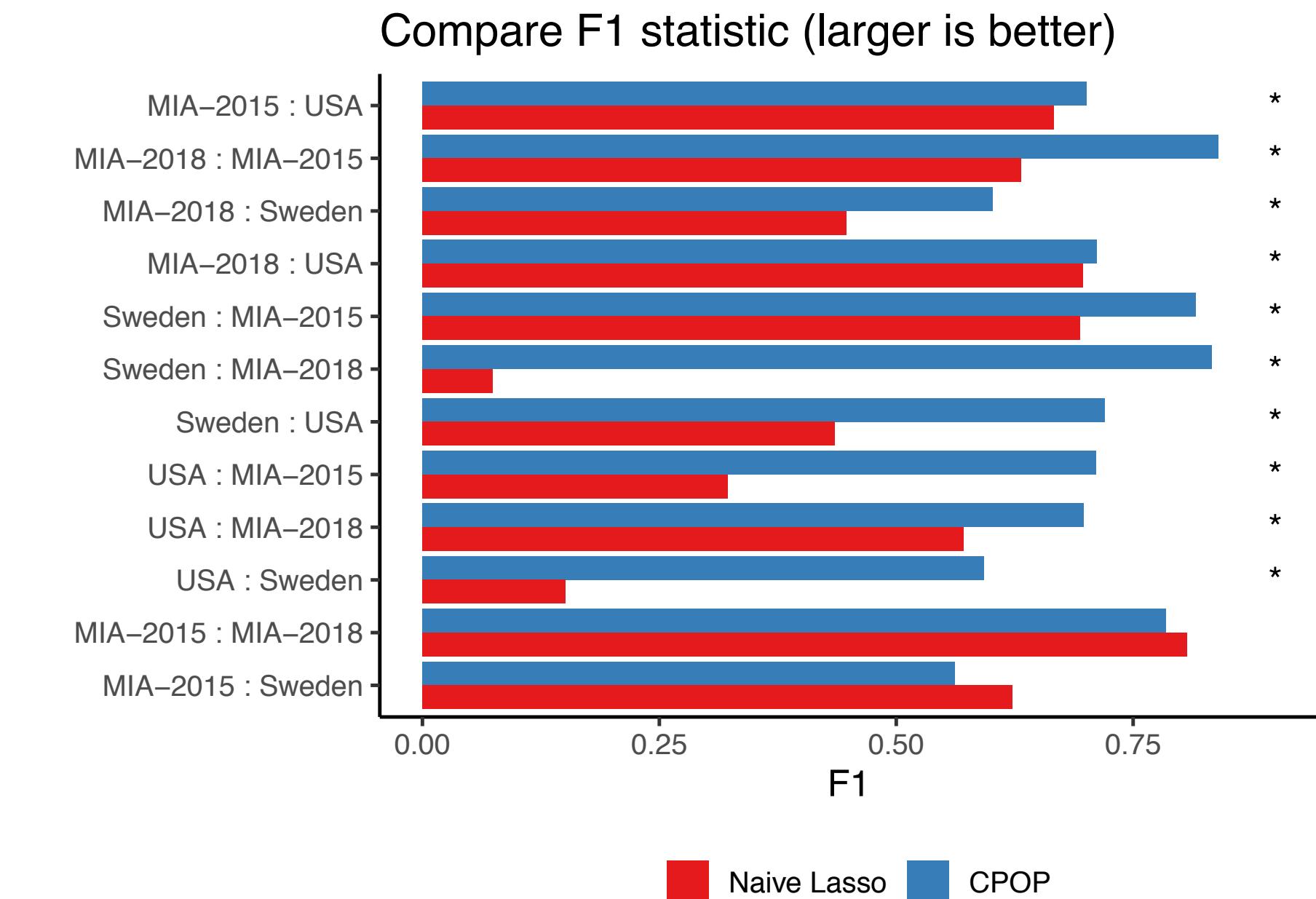
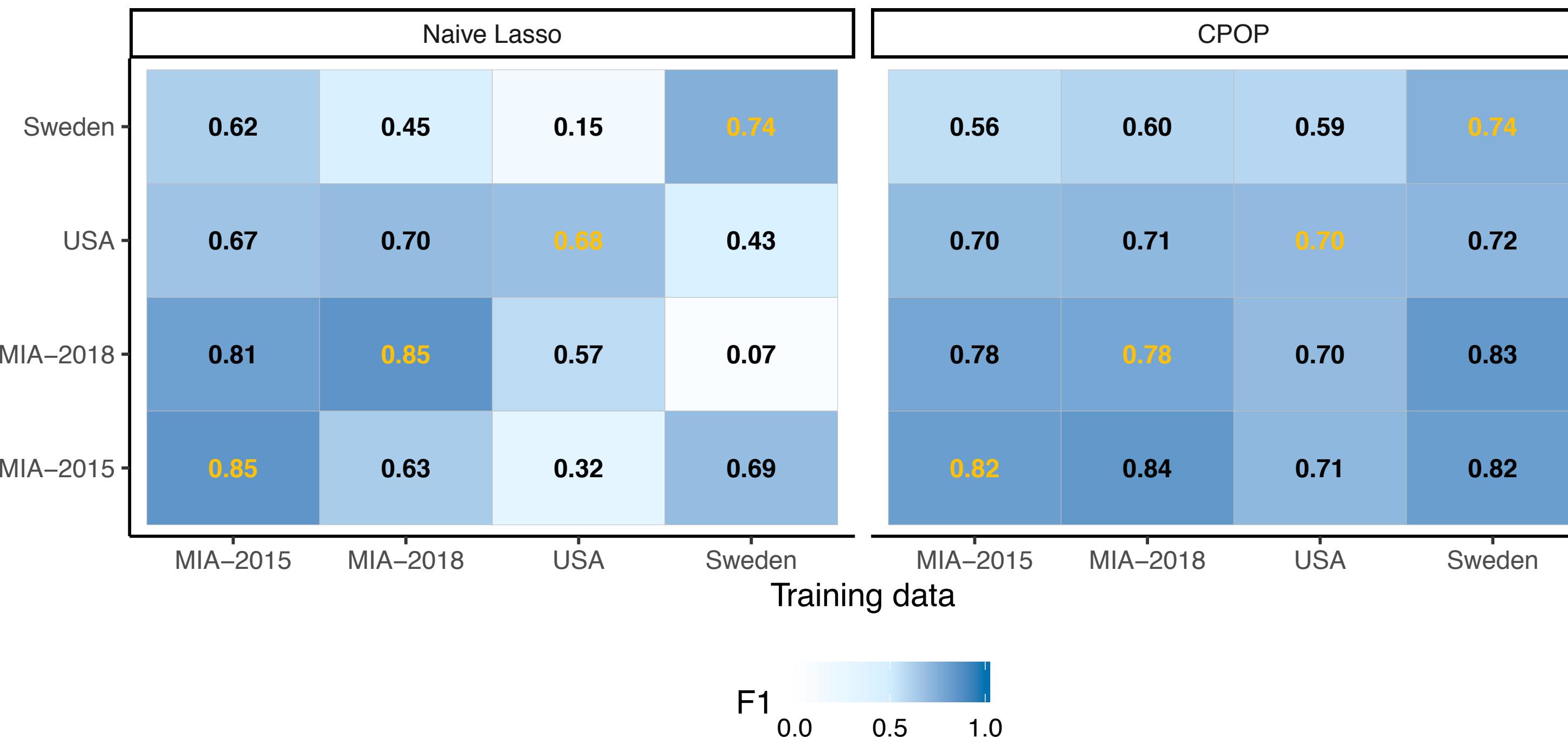
Concluding remarks

- ▶ Integrates clinical implementation constraints into the model
 1. Log-ratios enable prospective and multi-centres prediction
 2. Stable variable selection and estimation components
- ▶ A flexible framework with many adaptable components
- ▶ Potentials to handle data with higher relevance to precision medicine (e.g. drug sensitivity)

I think that is all.

CPOP results 2: four melanoma data

F1 classification statistic



CPOP is highly predictive

Motivation for CPOP: one patient cohort, two gene expression data

$$X_1 \hat{\beta}_1 \approx X_2 \hat{\beta}_2$$

loosely translate to

$$X_1 \approx X_2$$

column-wise

(feature distribution stability)

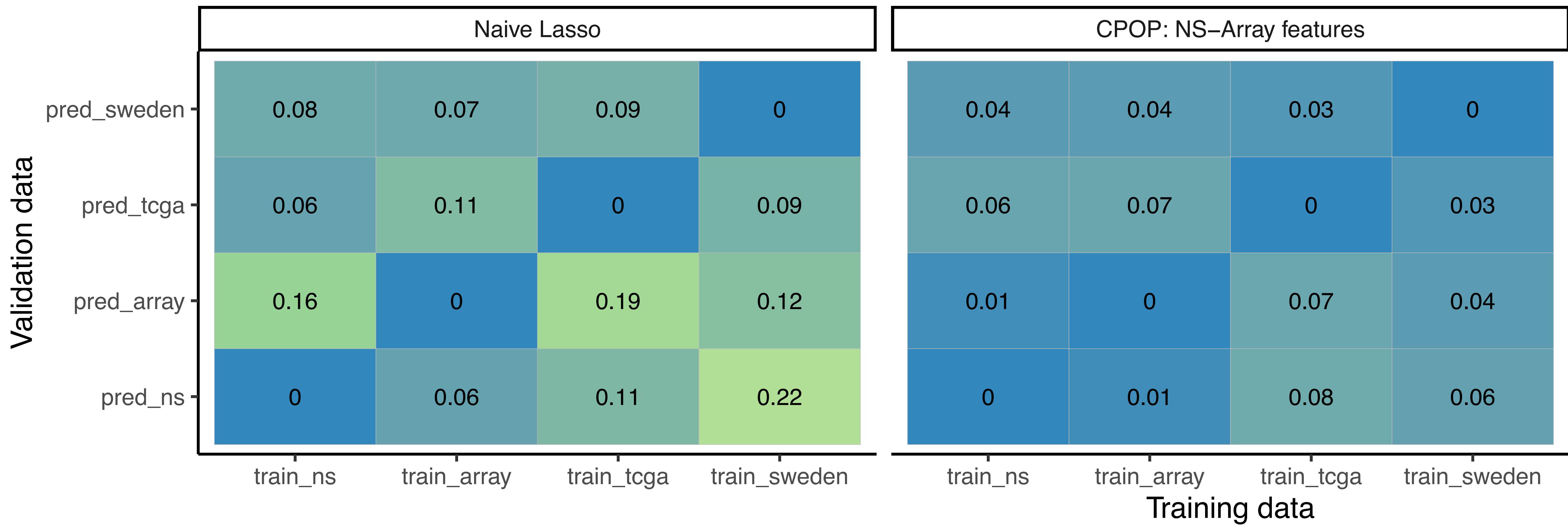
$$\hat{\beta}_1 \approx \hat{\beta}_2$$

element-wise

(mode estimation stability)

CPOP results 1: four melanoma data

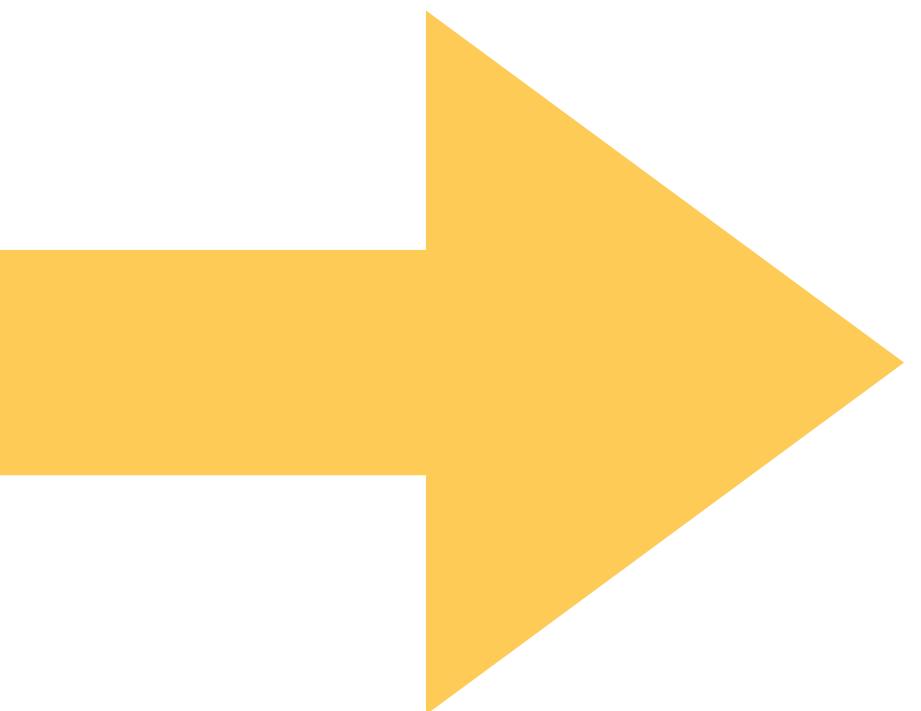
Identity distance between predicted values under various models



Small deviation in **predicted values** across datasets

Framingham heart disease risk score:

- ▶ Age (Years)
- ▶ Cholesterol (mg/dL)
- ▶ If smoker (Yes/No)
- ▶ HDL cholesterol (mg/dL)
- ▶ Systolic blood pressure (mm Hg)



$$\hat{y} = X \hat{\beta}$$

20 points model

Concluding remarks

- ▶ CPOP is a flexible procedure that allows for:
 - ▶ cross-platform omics prediction
 - ▶ stable single-patient prediction
- ▶ Not everyone can smooth-sailing through the PhD process, find your own way to deal with it (e.g. insert random pictures into your slides)

But what about breast cancer?

- ▶ Alvarado et. al. (2015) reported poor concordance in the prediction scores
- ▶ Hyeon et. al. (2017) considered NanoString as a viable alternative to RT-PCR

Name	Predictors	Targets	Prediction	Technology	Legit?
Oncotype DX	21 genes	ER +	Score	qRT-PCR	ASCO, NCCN
Prosigna	50 genes	Hormone receptor +	Score	NanoString	FDA 510k
MammaPrint	70 genes	Any ER status	Binary	DNA microarray	FDA

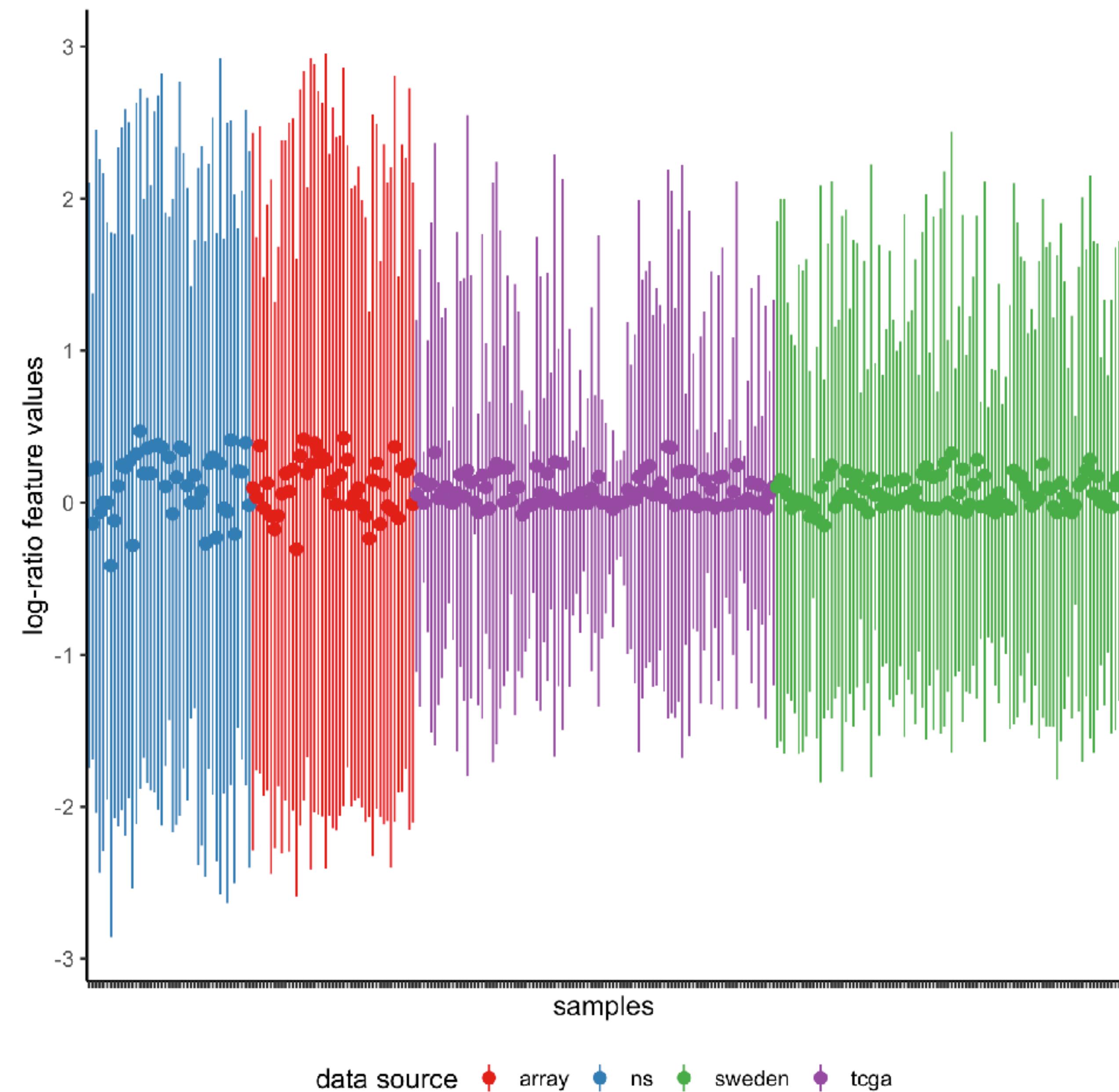
Is log-ratio really a new innovation?

Here is a list of papers that uses genes as predictors

Here is a list of papers that uses a single ratio for prediction

Our contribution is the advocacy using a whole collection of ratios for prediction

This has extra implications in terms of the statistics, but we are happy to tackle these.



Within-sample feature standardisation

