

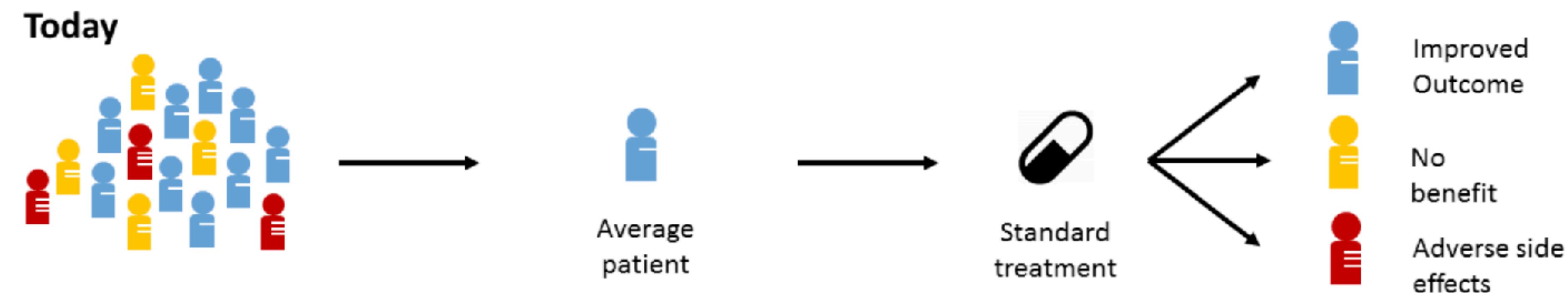
Kevin Wang

School of Mathematics and Statistics, University of Sydney

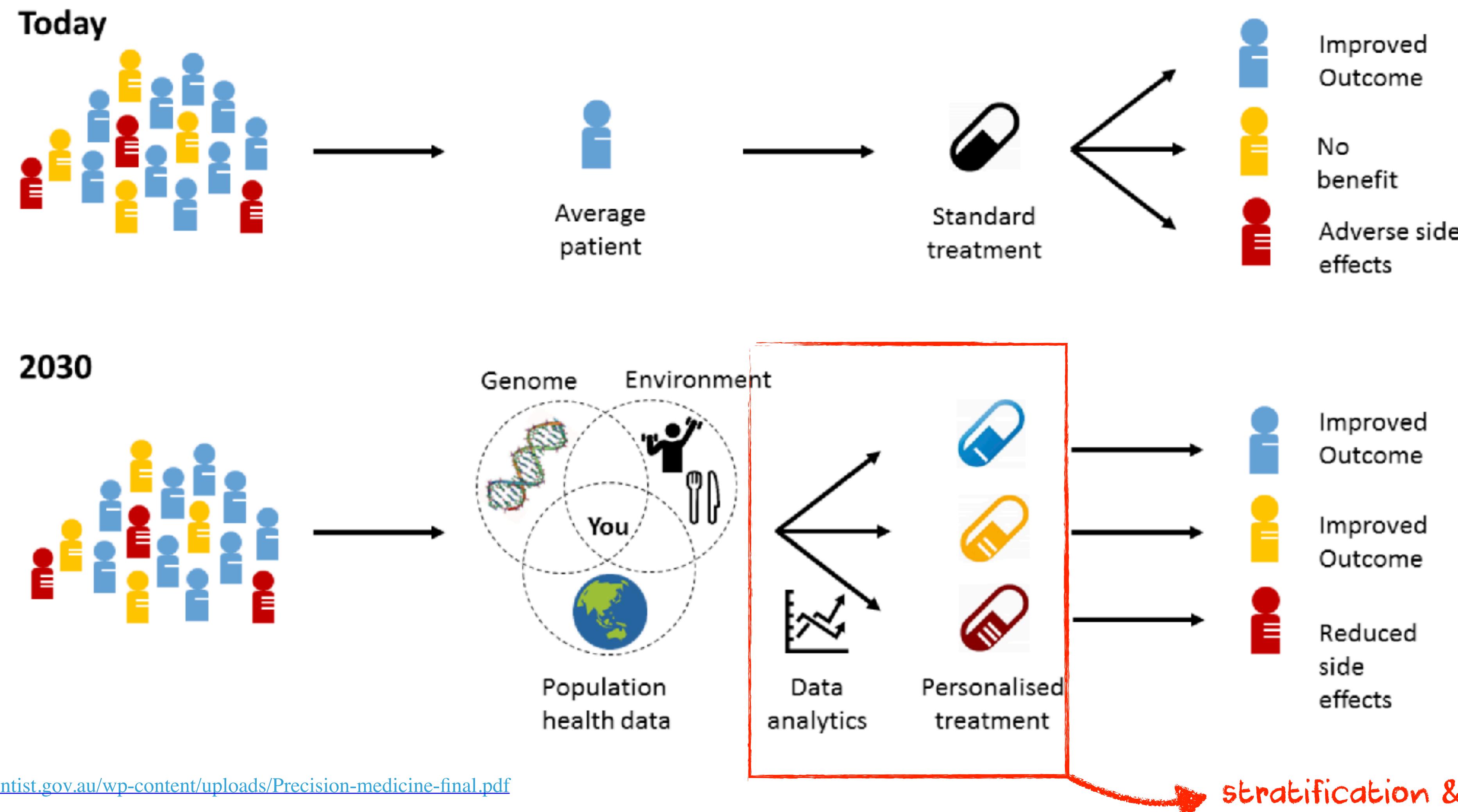
---

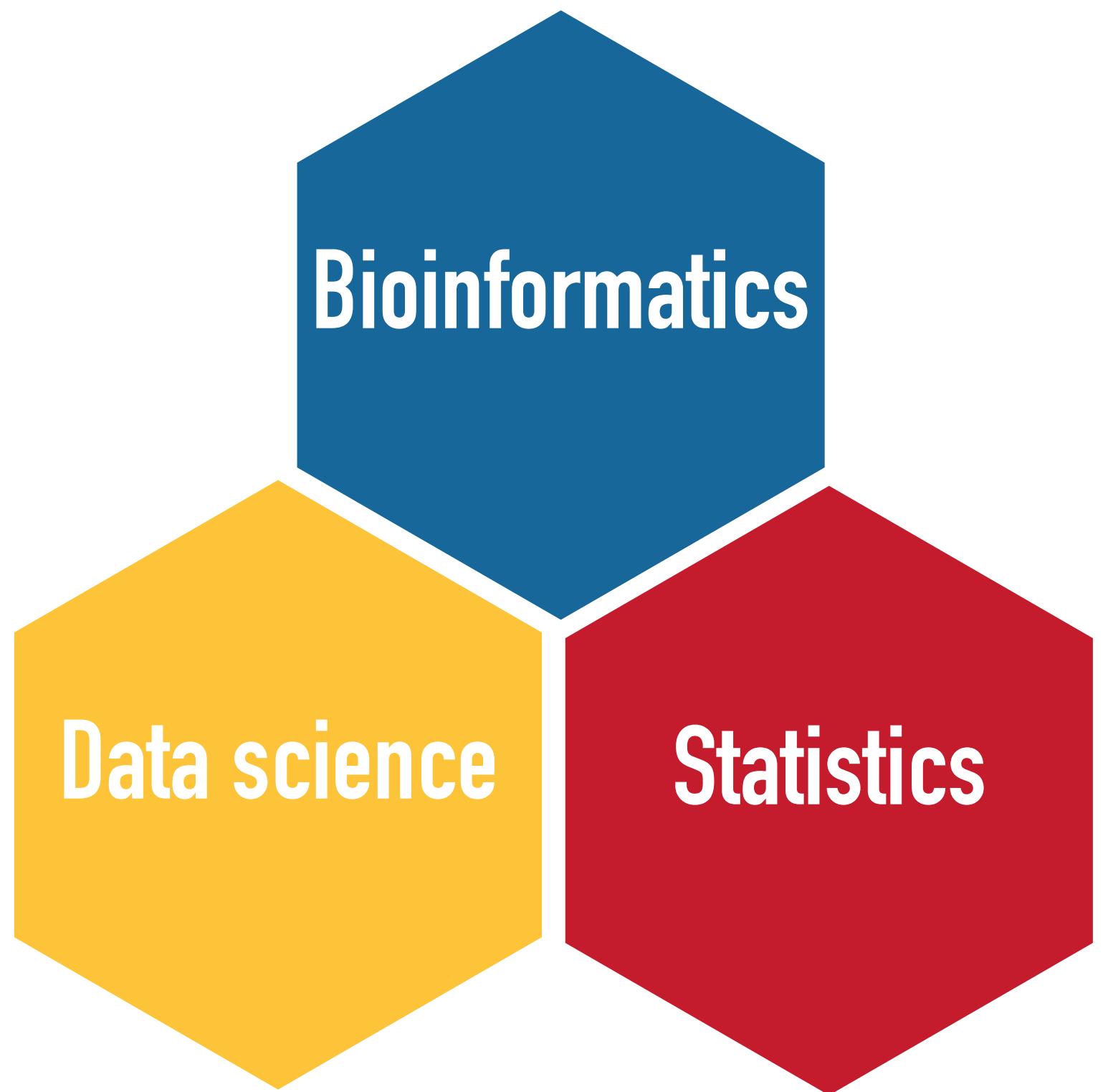
# Cross-Platform Omics Prediction: towards implementation of precision medicine

# Precision medicine: predicting best cause of action using omics data



# Precision medicine: predicting best cause of action using omics data





**BRAF-mutants**  
(CCR, 2019)

**scMerge**  
(PNAS 2019)

**bcGST**  
(Bioinformatics 2019)

**Bioinformatics**

**CPOP**  
(in prep)

**Melanoma  
Explorer**  
(Mel. Research 2019)

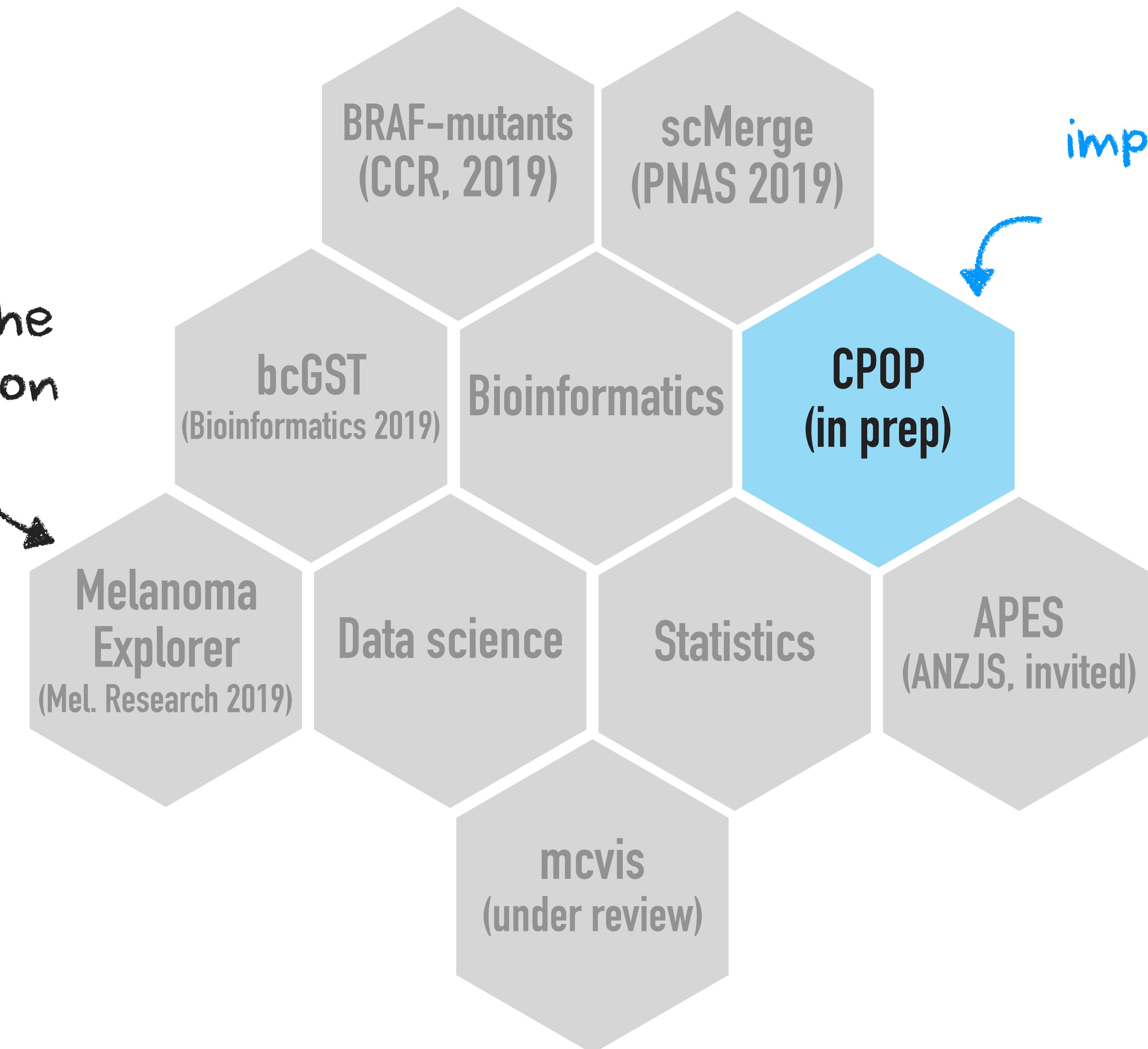
**Data science**

**Statistics**

**APES**  
(ANZJS, invited)

**mcvis**  
(under review)

Laying the  
foundation



Towards  
implementation

# Melanoma Institute Australia

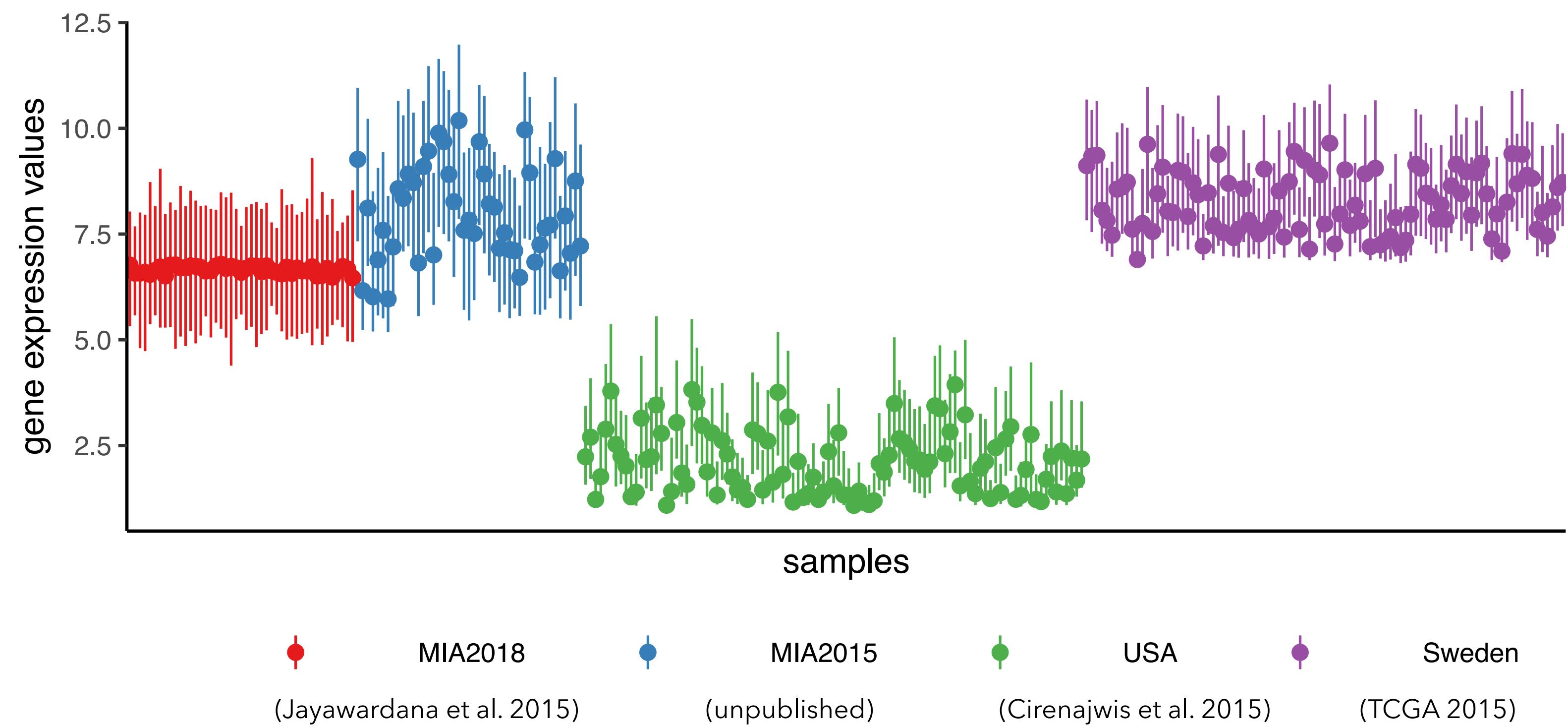
- ▶ Risk score using gene expression

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

- ▶ CRE grant will support implementation

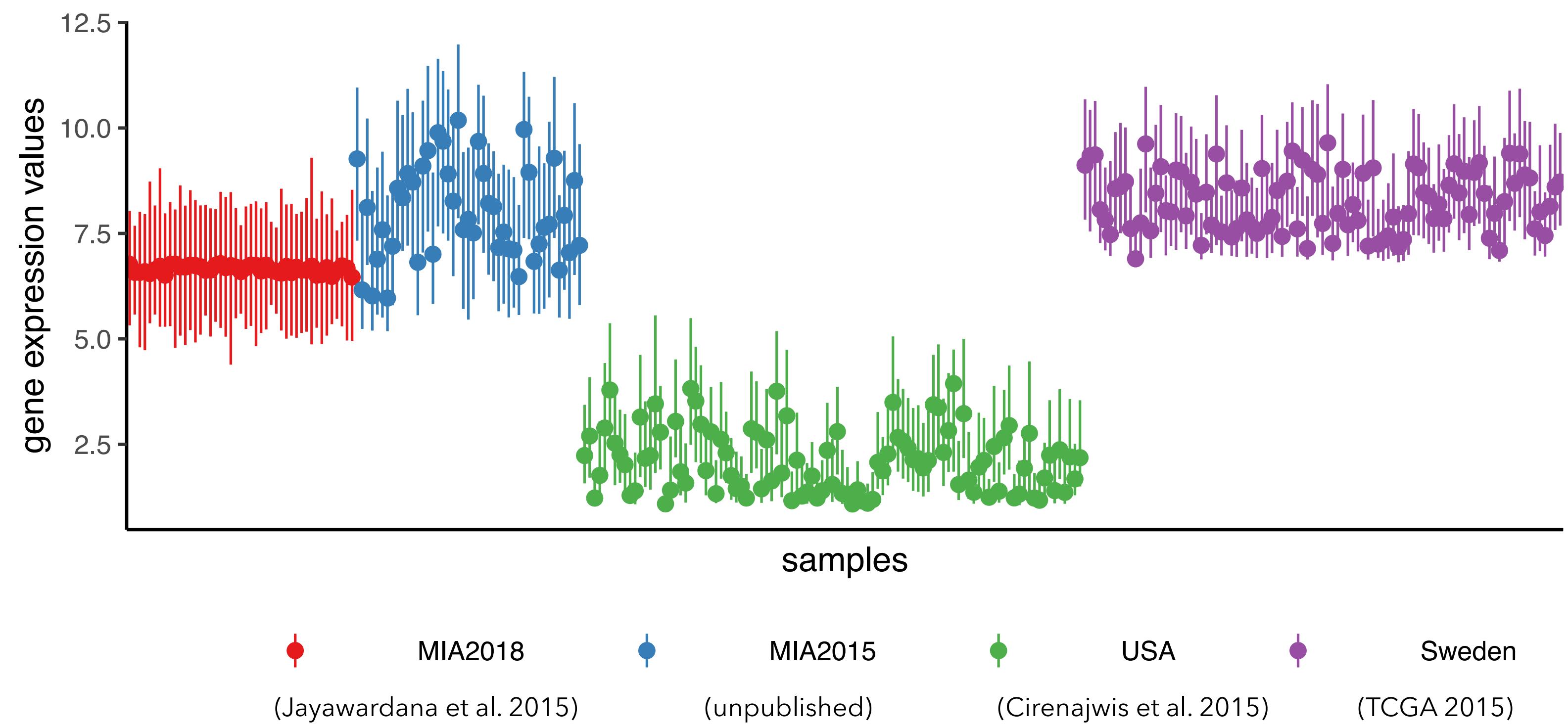
- ▶ prospective

- ▶ **multi-centres**



Prof. Graham Mann

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



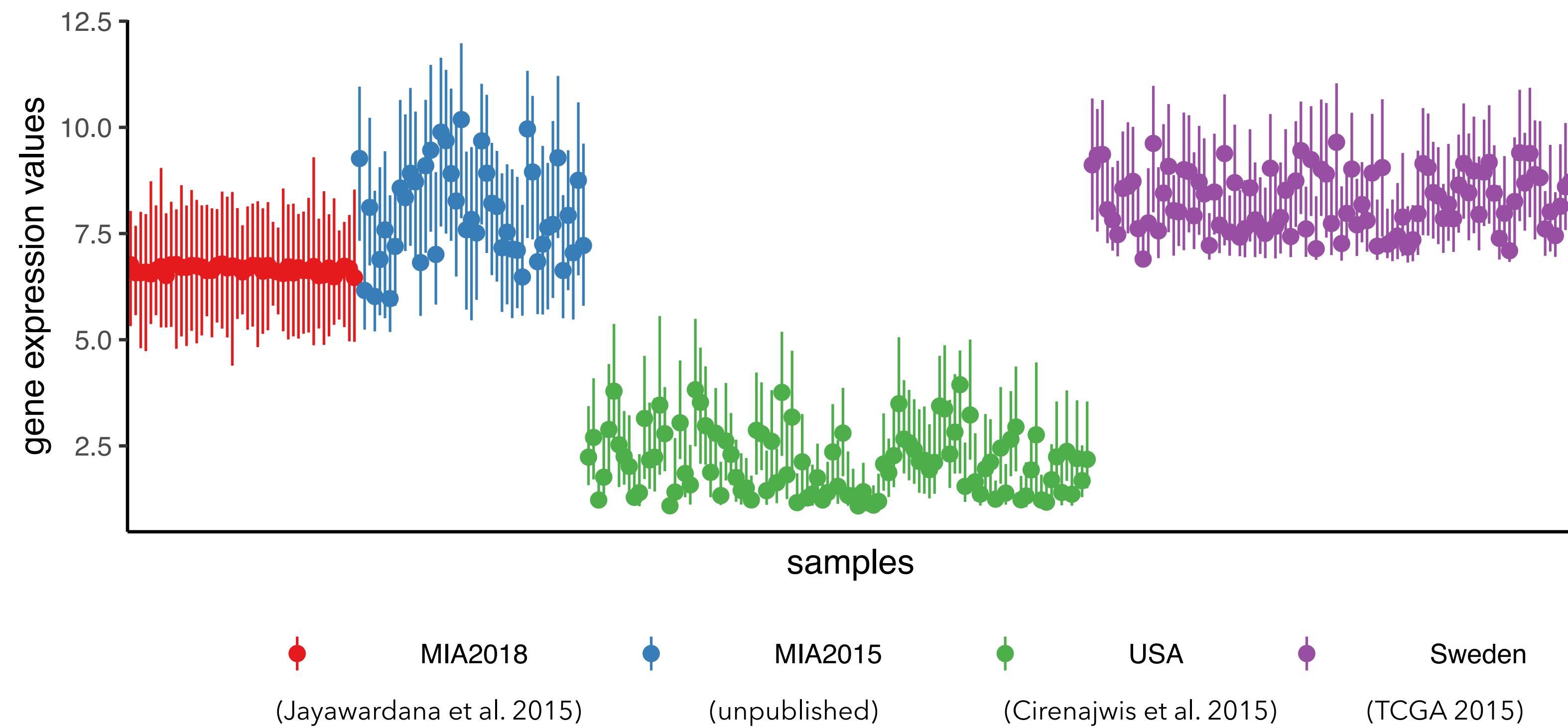
Gene expression features on a unit-less scale.

A typical value in one data can be an impossible value on another.

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

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

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



Gene expression features on a unit-less scale.

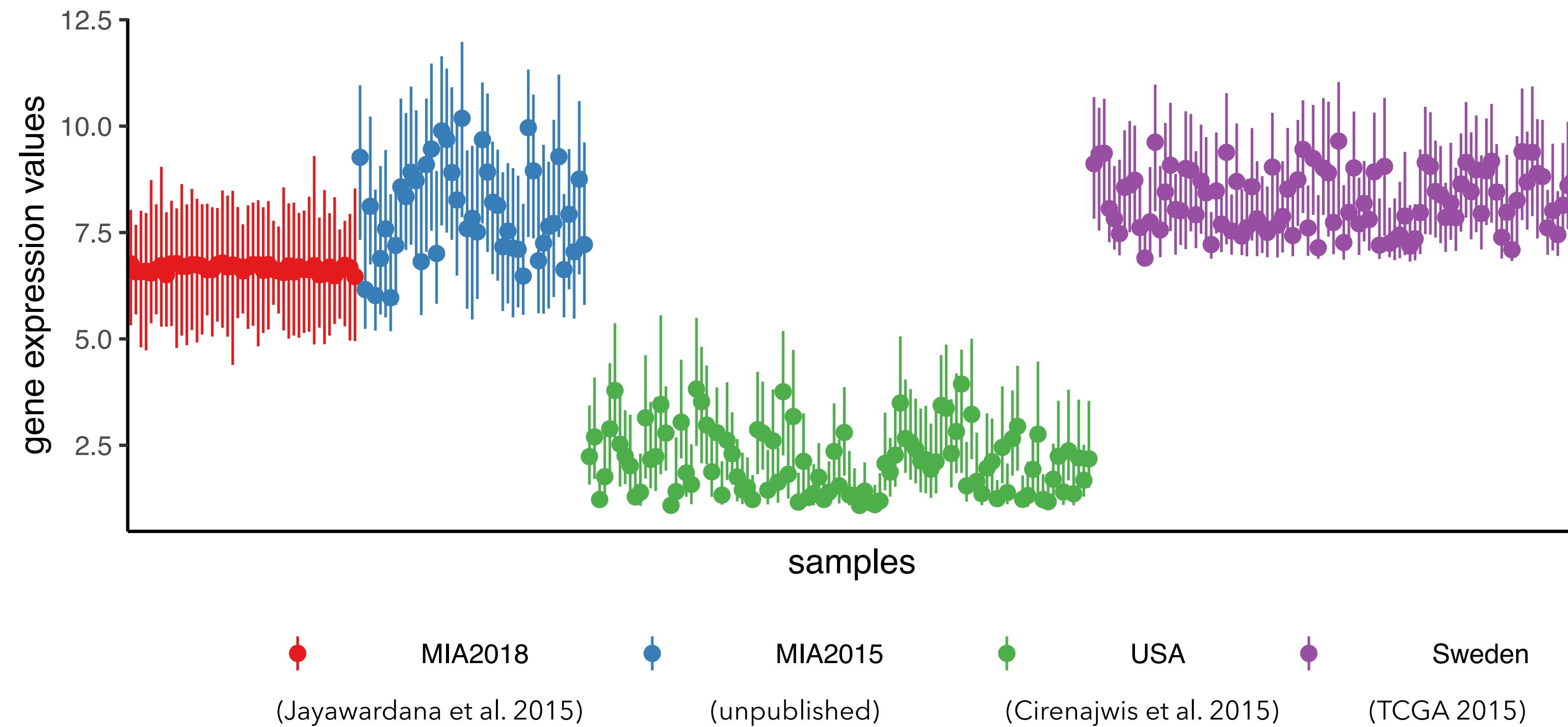
A typical value in one data can be an impossible value on another.

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

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

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

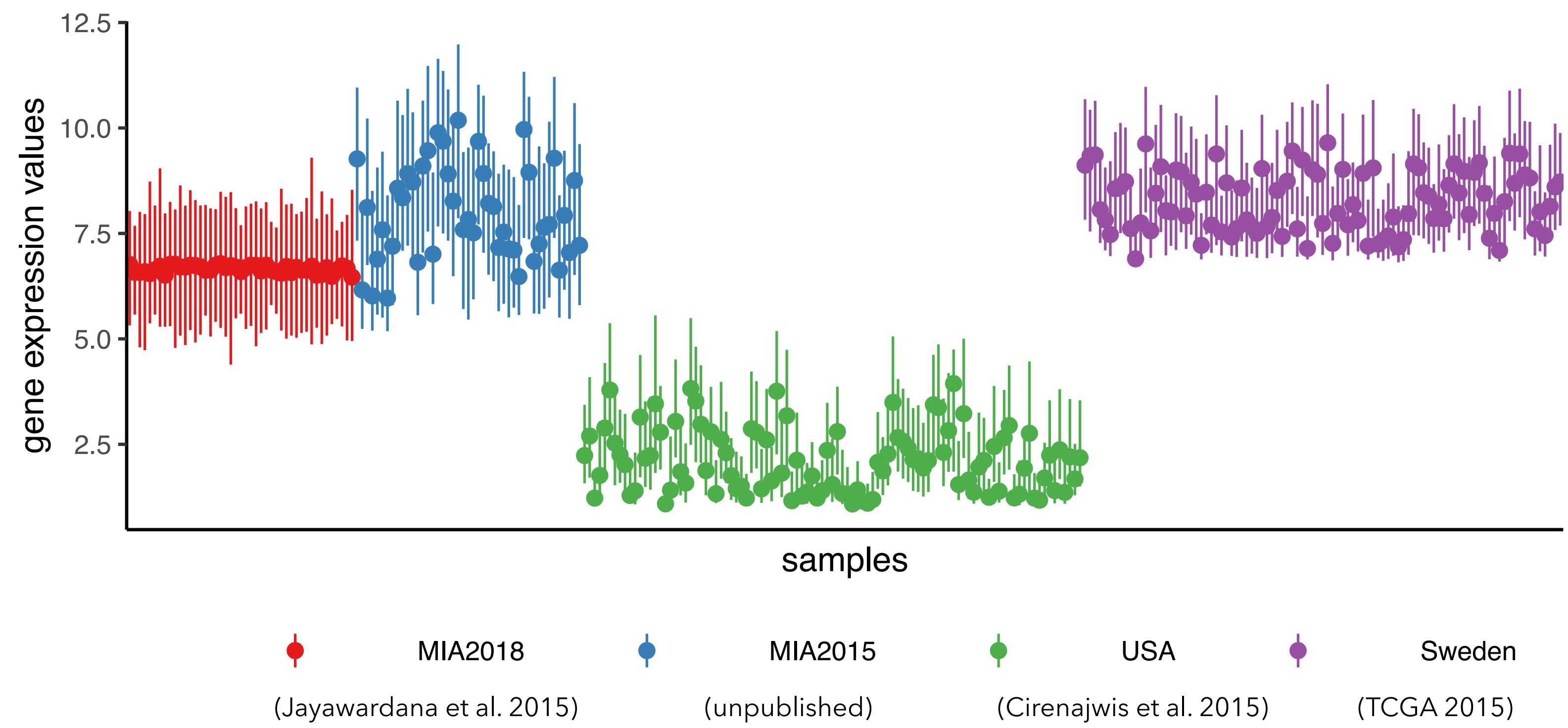
Assuming a  
noiseless shift



Gene expression features on a  
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?



**Transferability**  
A model trained from one  
data and should be predictive  
on another

# Existing approaches

# Components of a risk score

Data

$$(X_1, y_1)$$

$$(X_2, y_2)$$

Model

$$\hat{\beta}_1$$

Prediction

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

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

# Data-harmonisation

Data

$$(X_1, y_1)$$

$$(X_2, y_2)$$

Model

$$\hat{\beta}_1$$

Prediction

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

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

Harmonisation

Classical estimation  
methods

Classical predictions

# Data-harmonisation: normalisation or standardisation

Data

$(X_1, y_1)$

$(X_2, y_2)$

Harmonisation

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

**Prevented by ethics, privacy and data security**

# Clinical constraints in implementation

Data

$$(X_1, y_1)$$

$$(X_2, y_2)$$

Model

$$\hat{\beta}_1$$

No harmonisation

No re-training

Prediction

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

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

Scale-equivalent prediction

**CPOP: prediction model that  
respects implementation constraints**

## CPOP flowchart

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

# First component of CPOP: feature transform

---



## Log-ratio transformation: modelling relative gene expression

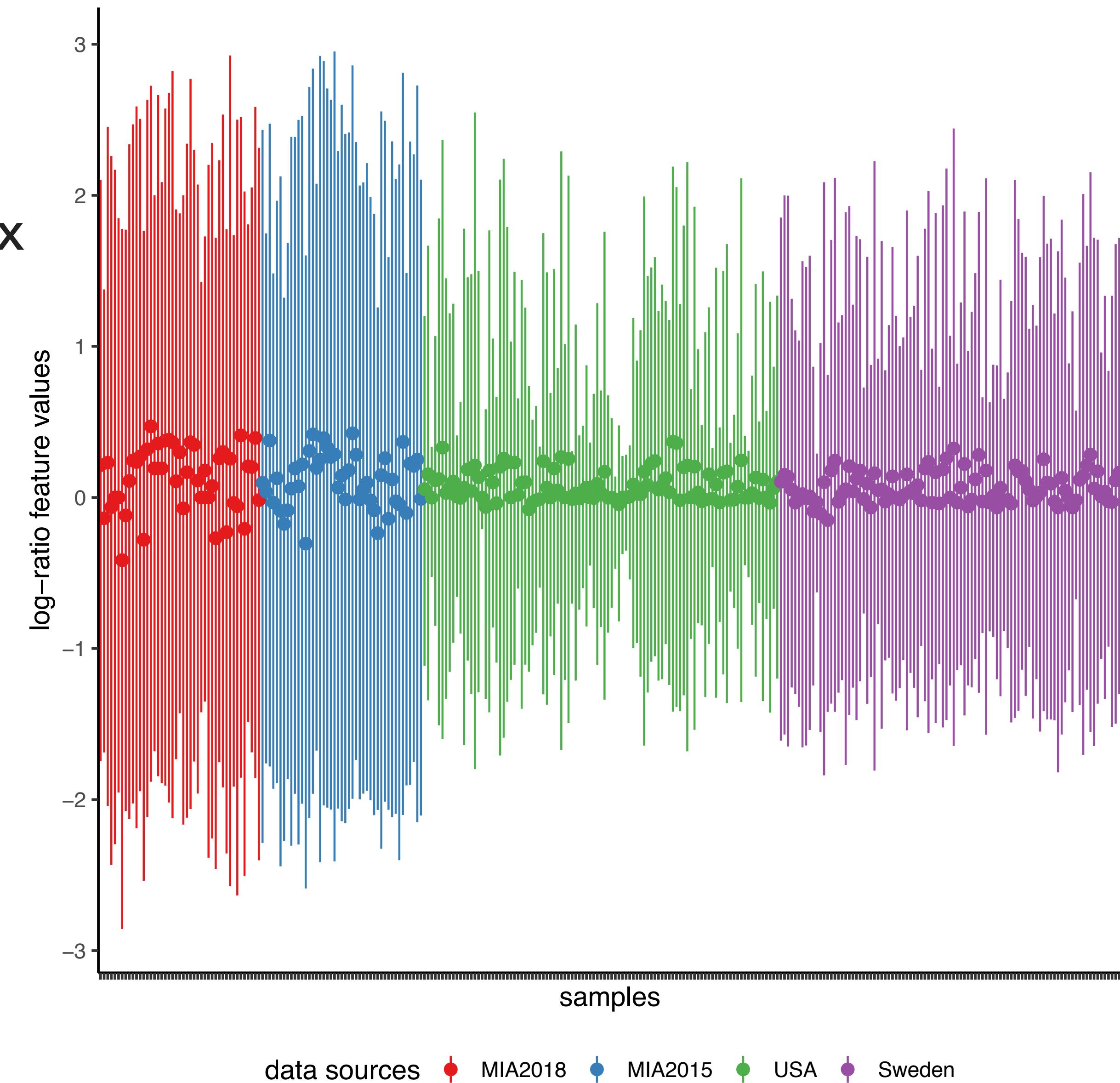
For each column in the gene expression matrix

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

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

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

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



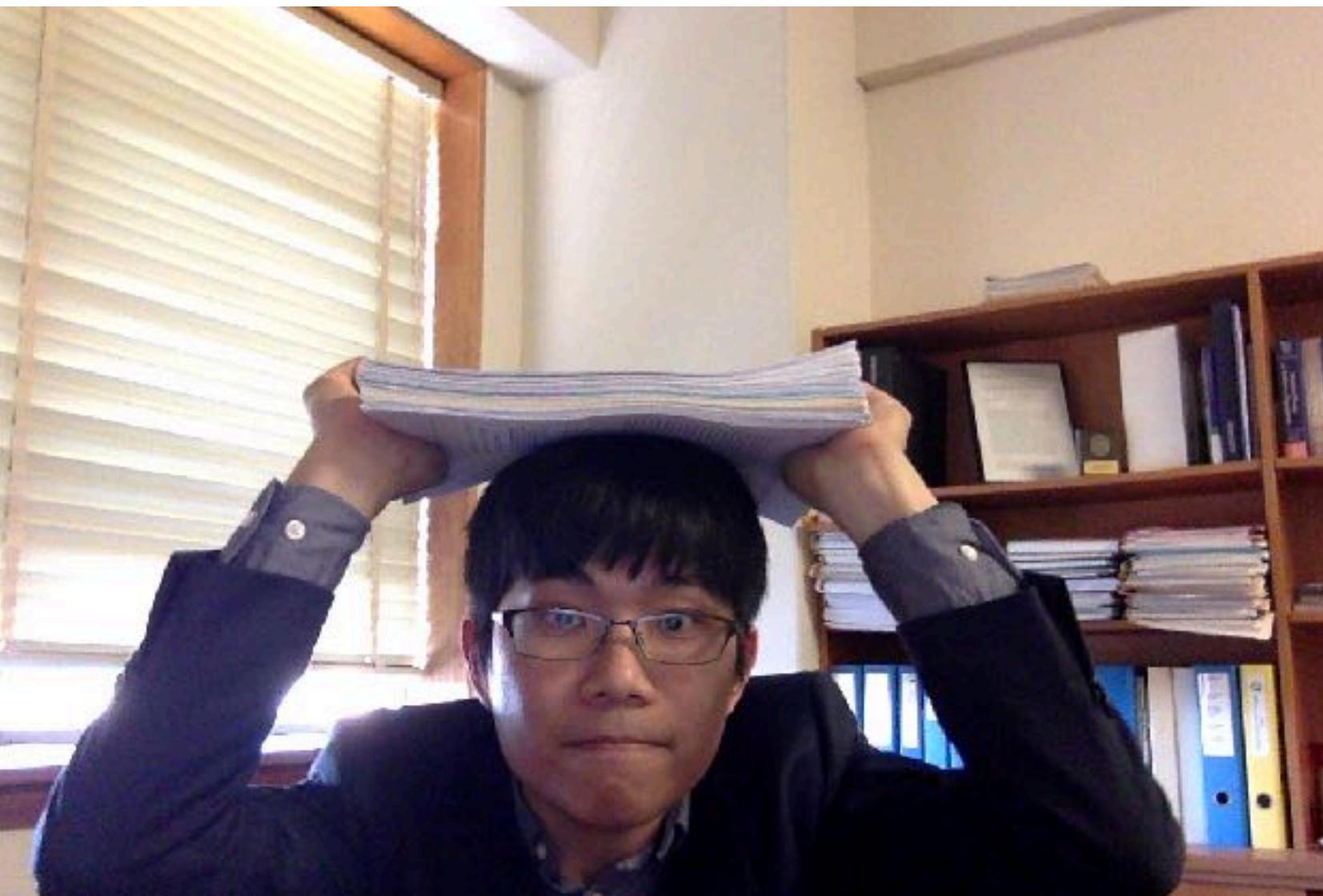
---

## Why log-ratios?

- ▶ Within-sample standardisation avoids re-normalisation and model re-training

# Why log-ratios?

- ▶ Within-sample standardisation, conforms to implementation constraints
- ▶ Modelling on relative expression of genes
- ▶ Opportunities for method developments

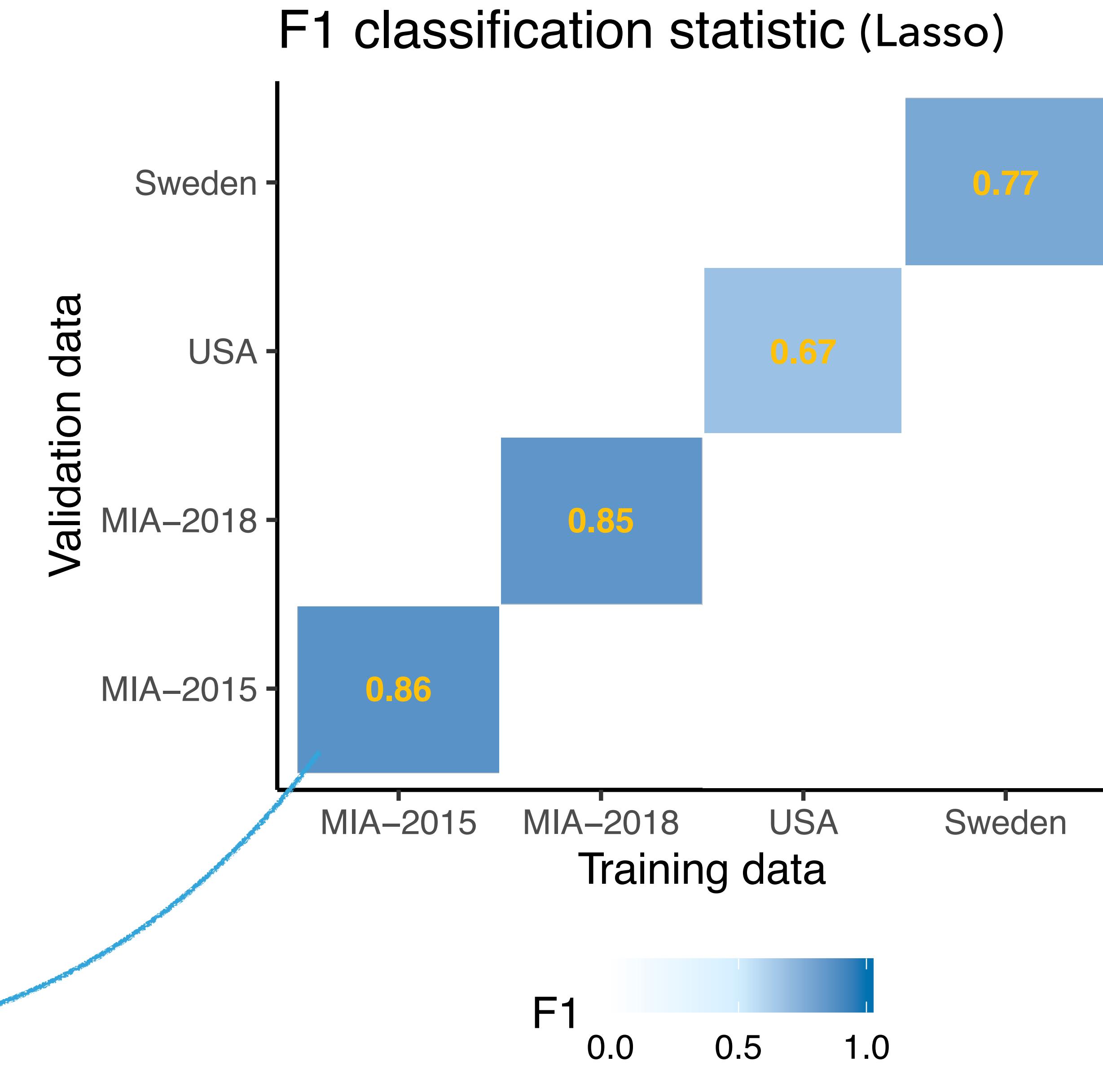
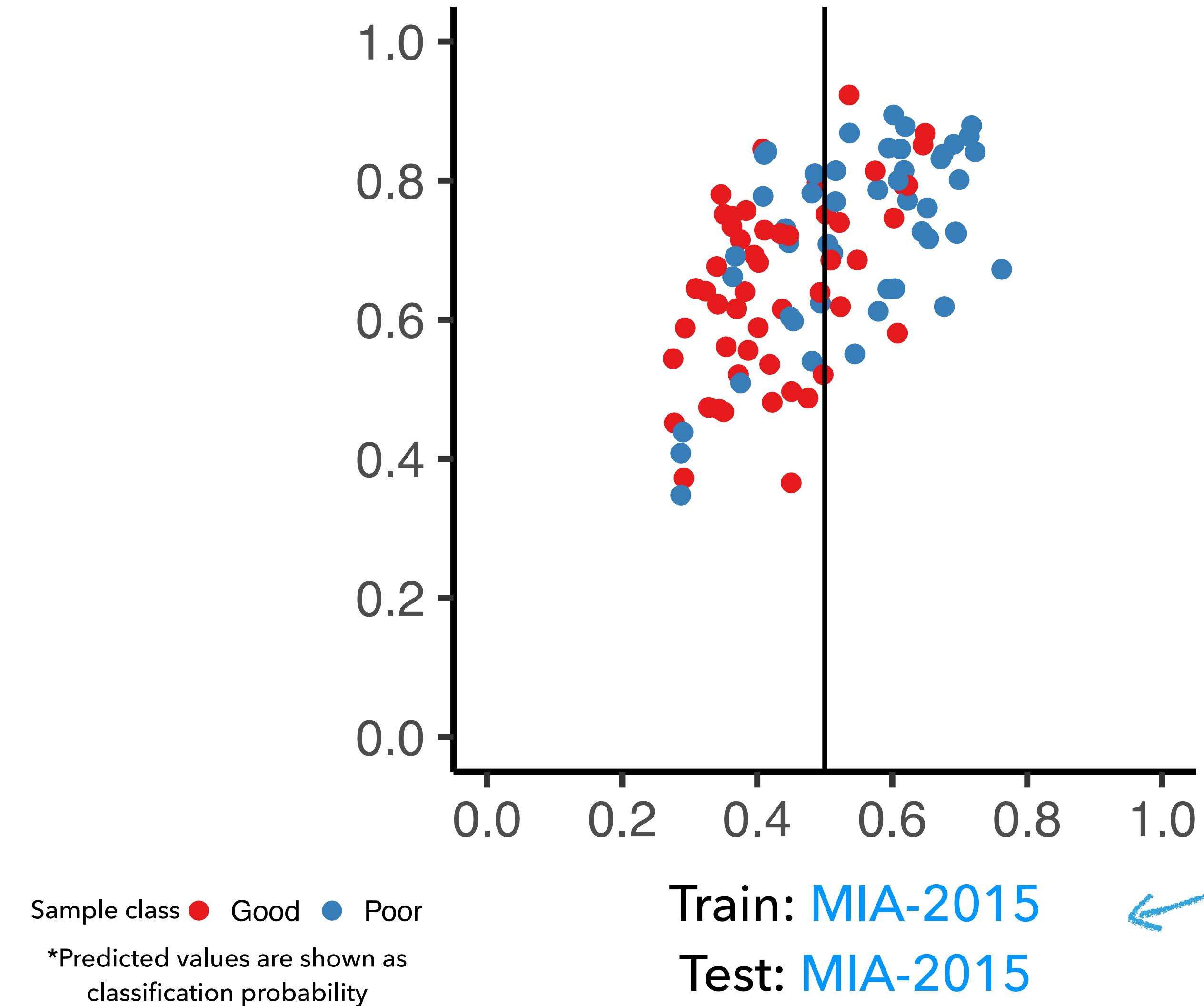


Papers that used gene expression

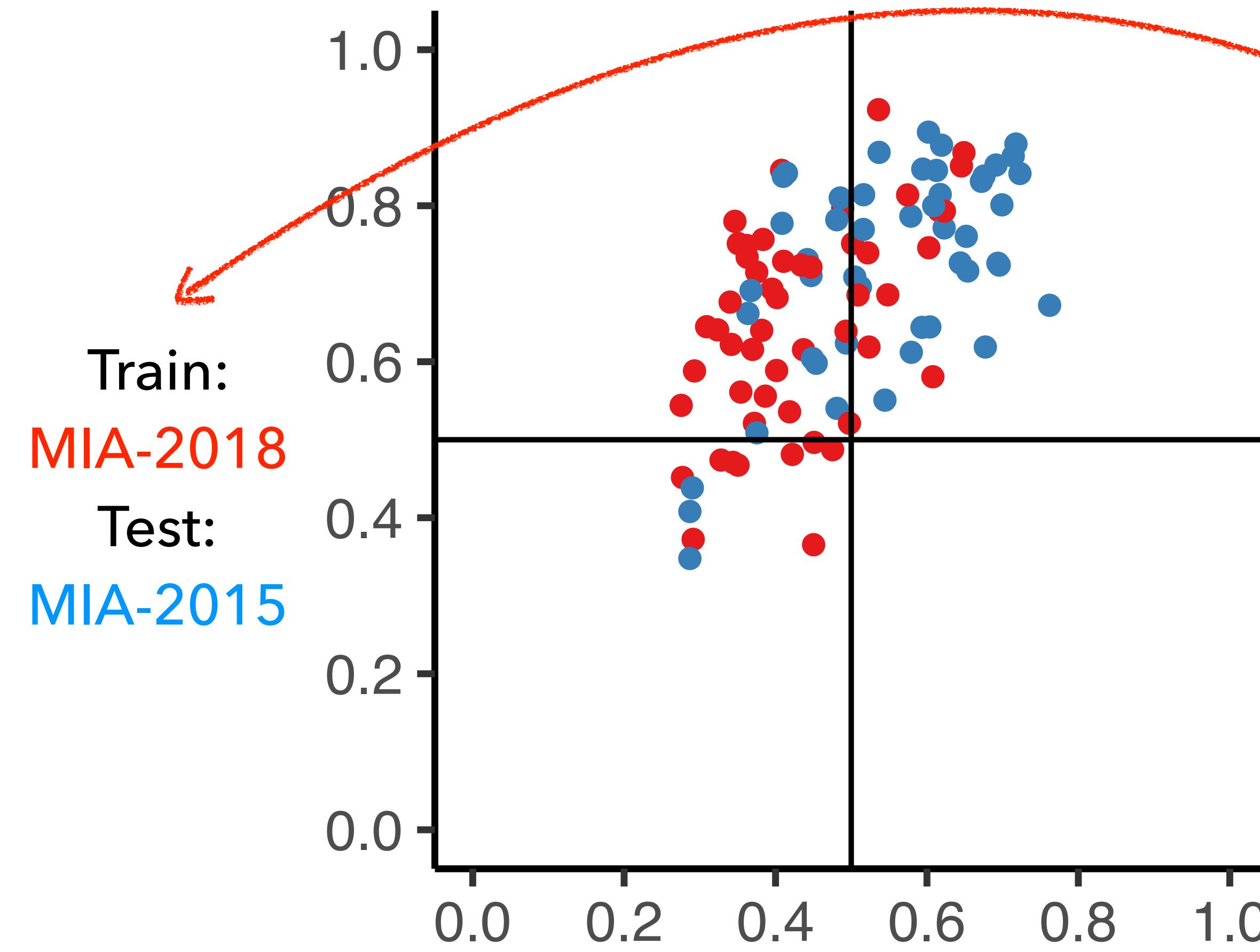


The paper that used all log-ratios

# Is log-ratio enough?

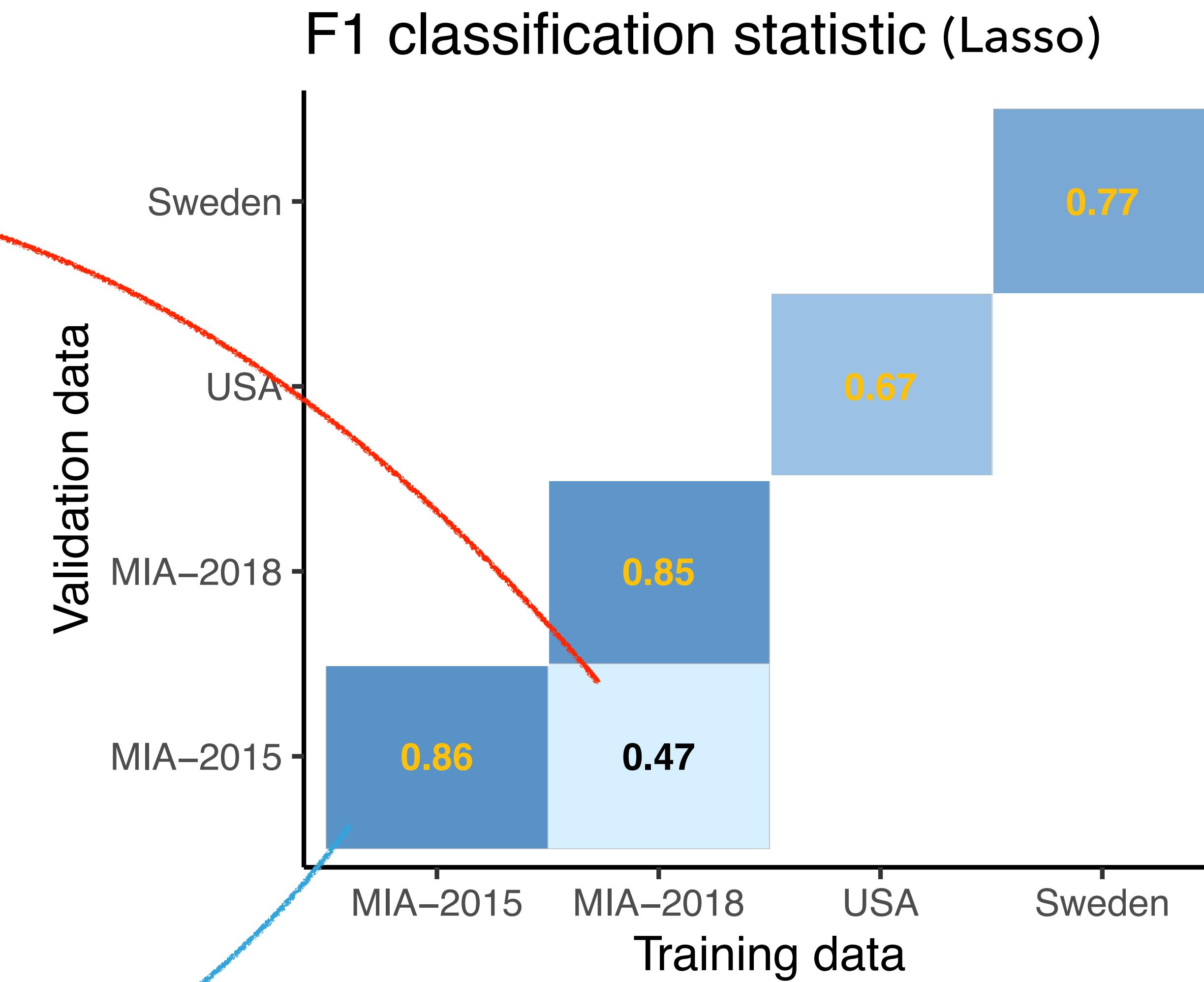


# Is log-ratio enough?

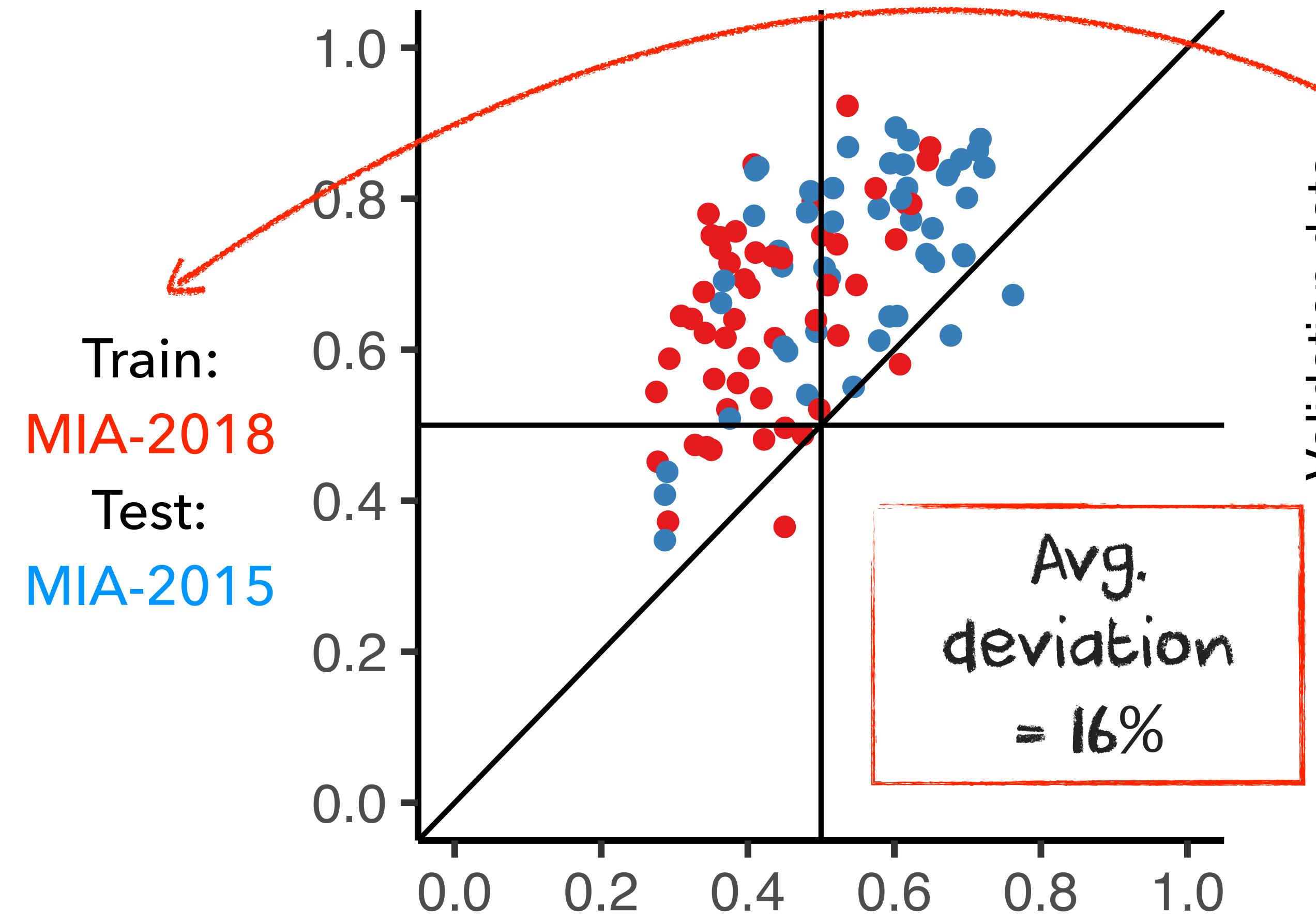


\*Predicted values are shown as  
classification probability

Train: MIA-2015  
Test: MIA-2015

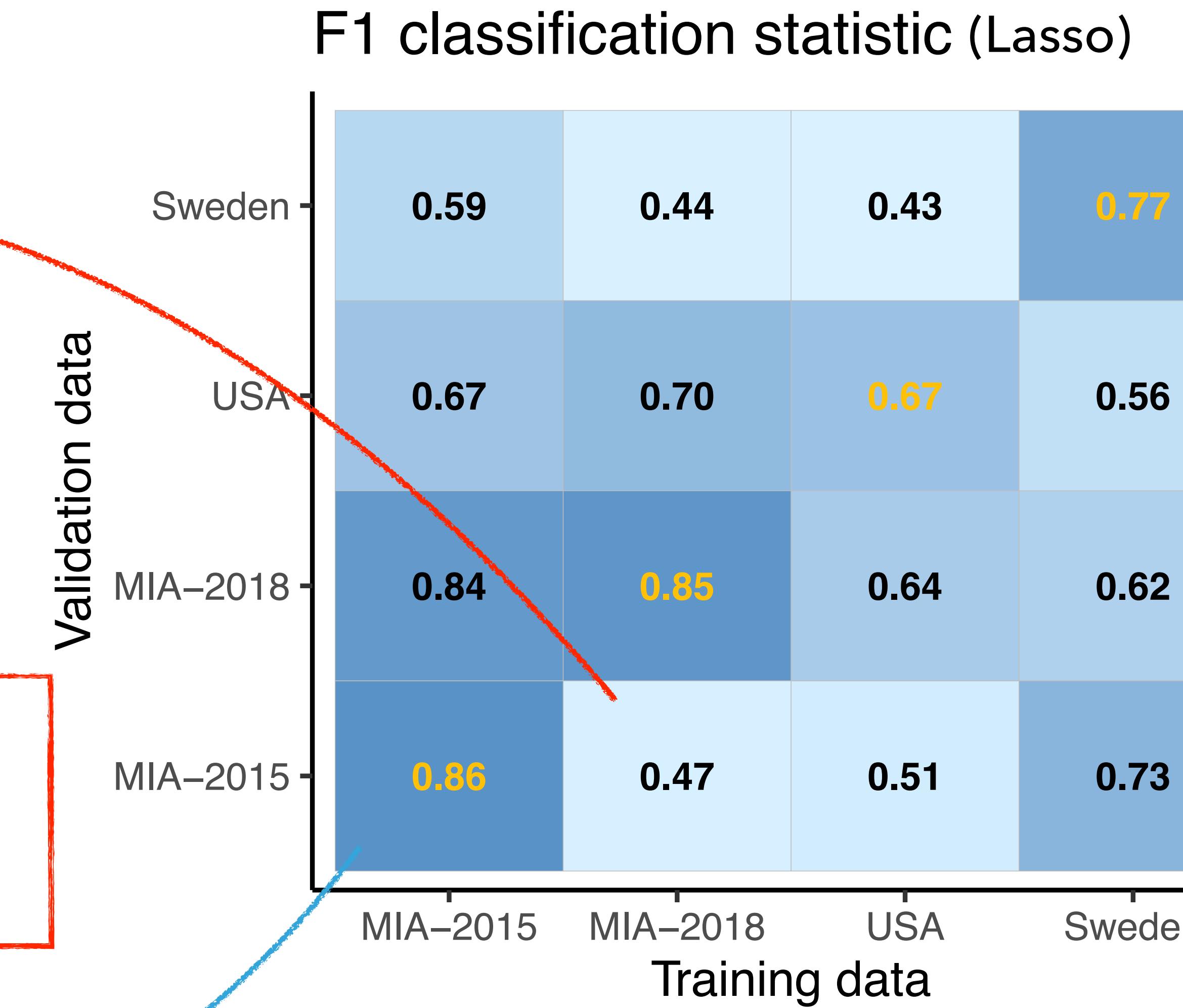


# Is log-ratio enough?



\*Predicted values are shown as classification probability

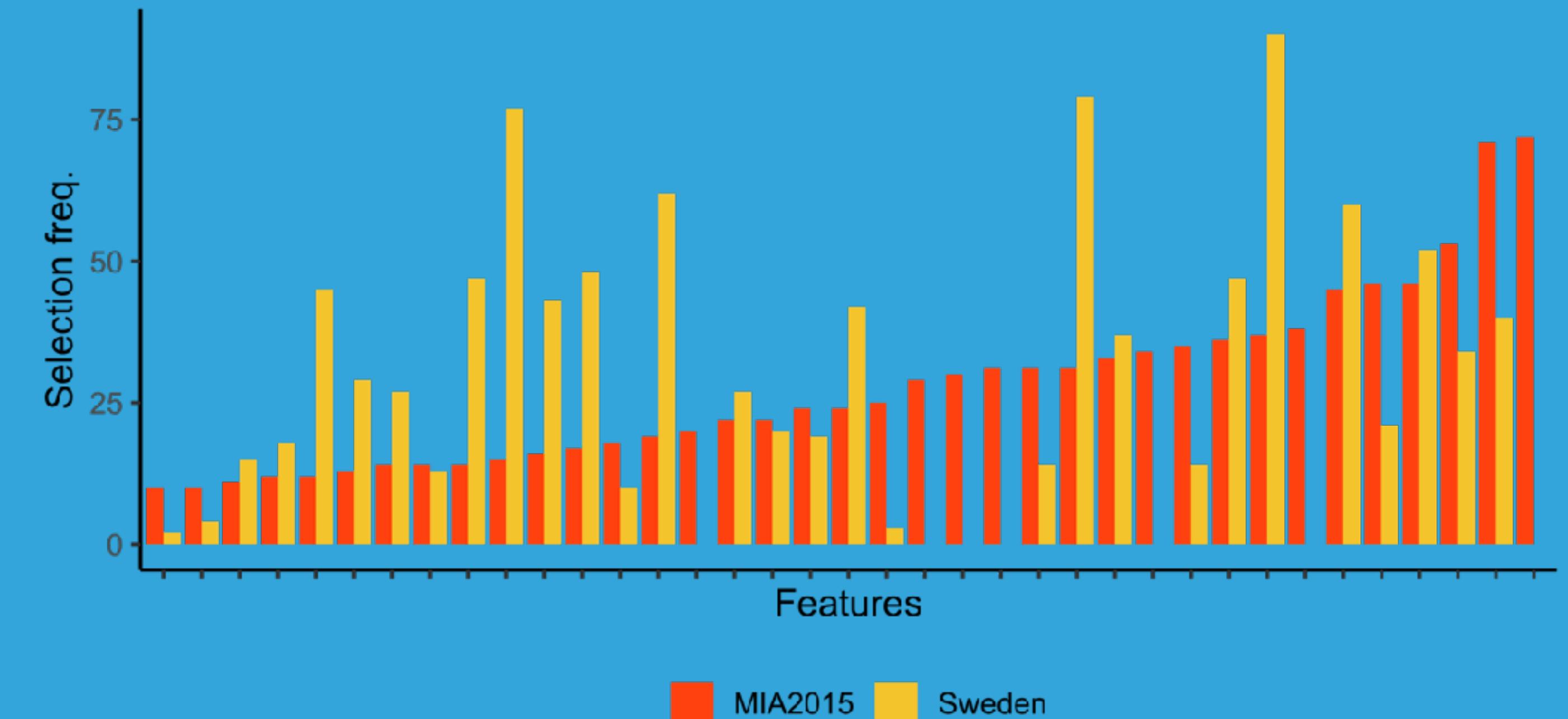
Train: MIA-2015  
Test: MIA-2015



---

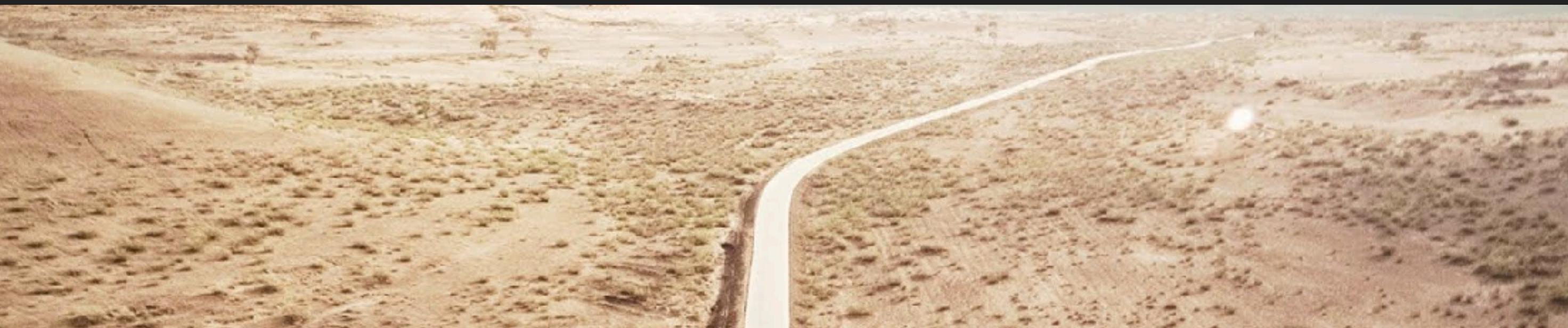
Transferability implies that predicted values should be evenly scattered around the identity line!

# Lasso variable selection is not stable!



# **Second component of CPOP: stable feature selection and estimation**

---



# 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

logistic loss  
function

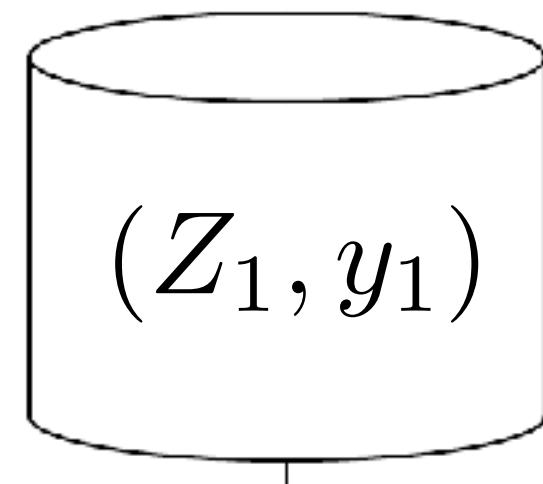
$$\hat{\beta}(y, Z) = \underset{\beta \in \mathbb{R}^{\binom{p}{2}}}{\operatorname{argmin}} \sum_{i=1}^n l(y_i, z_i^\top \beta) + \lambda \sum_{j=1}^q w_j \left[ \frac{(1-\alpha)}{2} \beta_j^2 + \alpha |\beta_j| \right]$$

modelling on log-ratio  
features

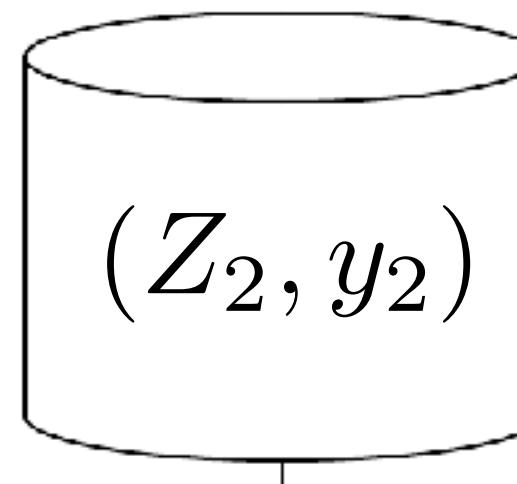
a mix of L1 and L2 penalties  
penalties on each feature,  
proportional to the stability of  
features

## Step 1: feature selection stability

Training data 1



Training data 2

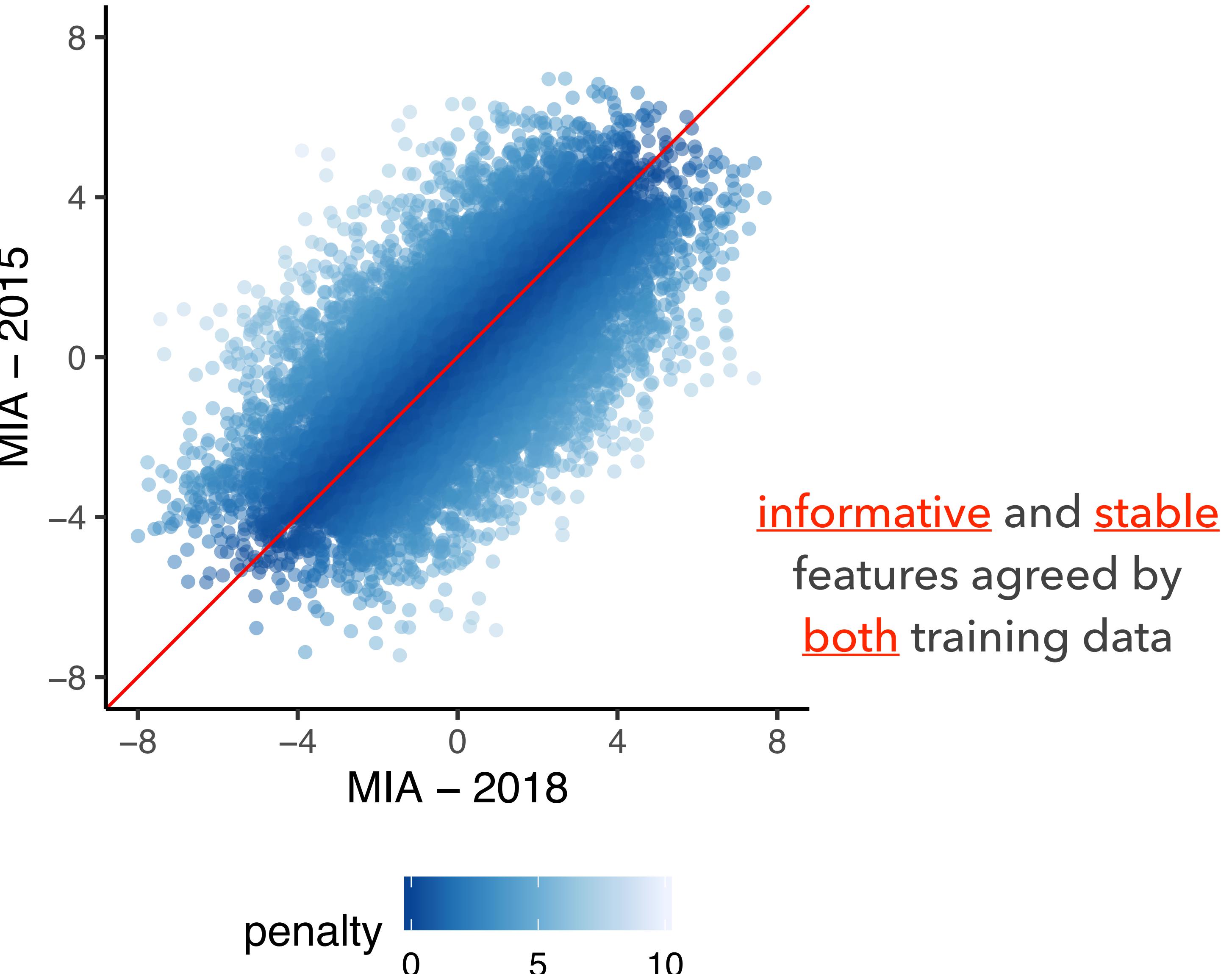


fit WEN using  $w_j$  to get  $\hat{\beta}_1^{(1)}$  and  $\hat{\beta}_2^{(1)}$

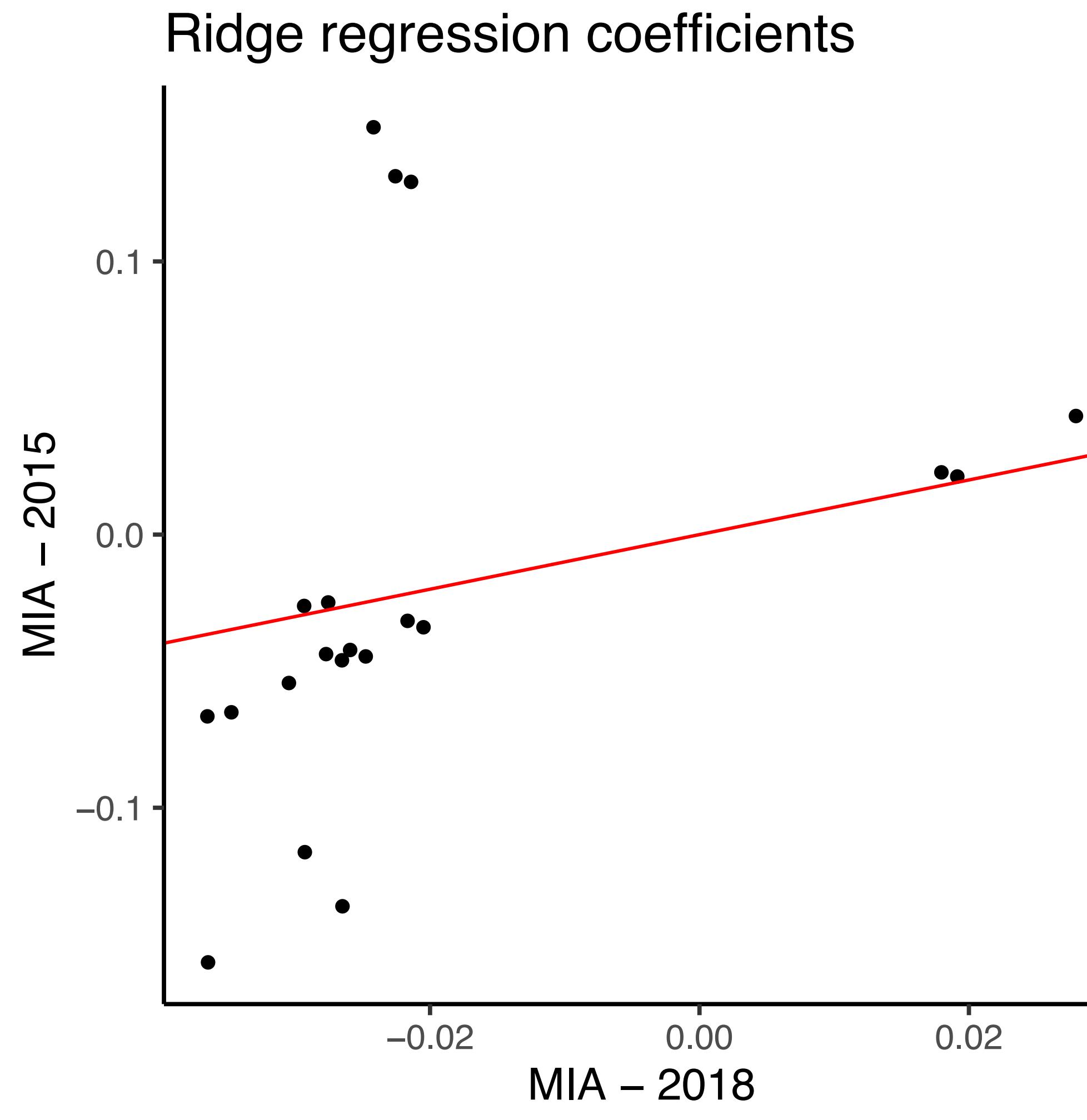
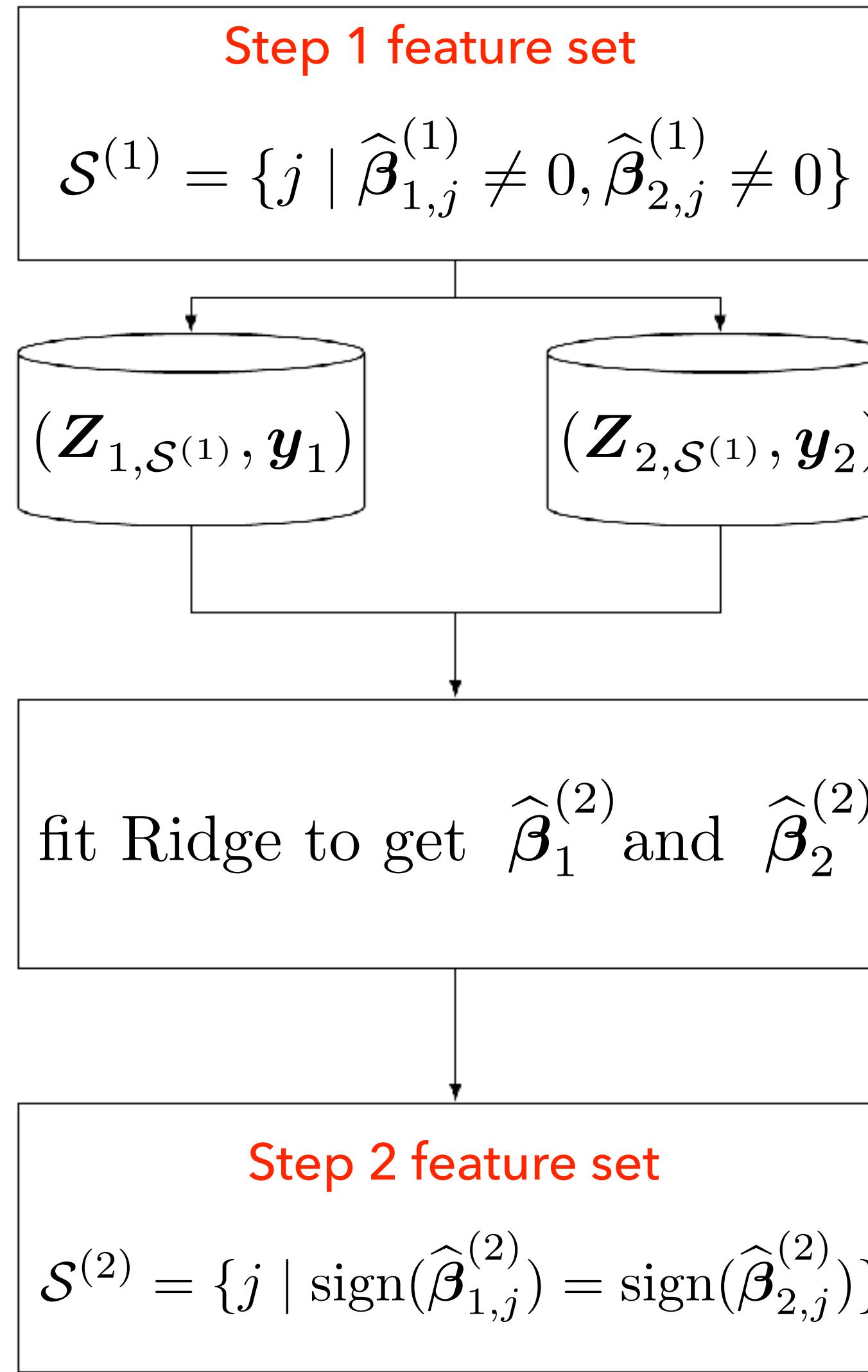
Step 1 feature set

$$\mathcal{S}^{(1)} = \{j \mid \hat{\beta}_{1,j}^{(1)} \neq 0, \hat{\beta}_{2,j}^{(1)} \neq 0\}$$

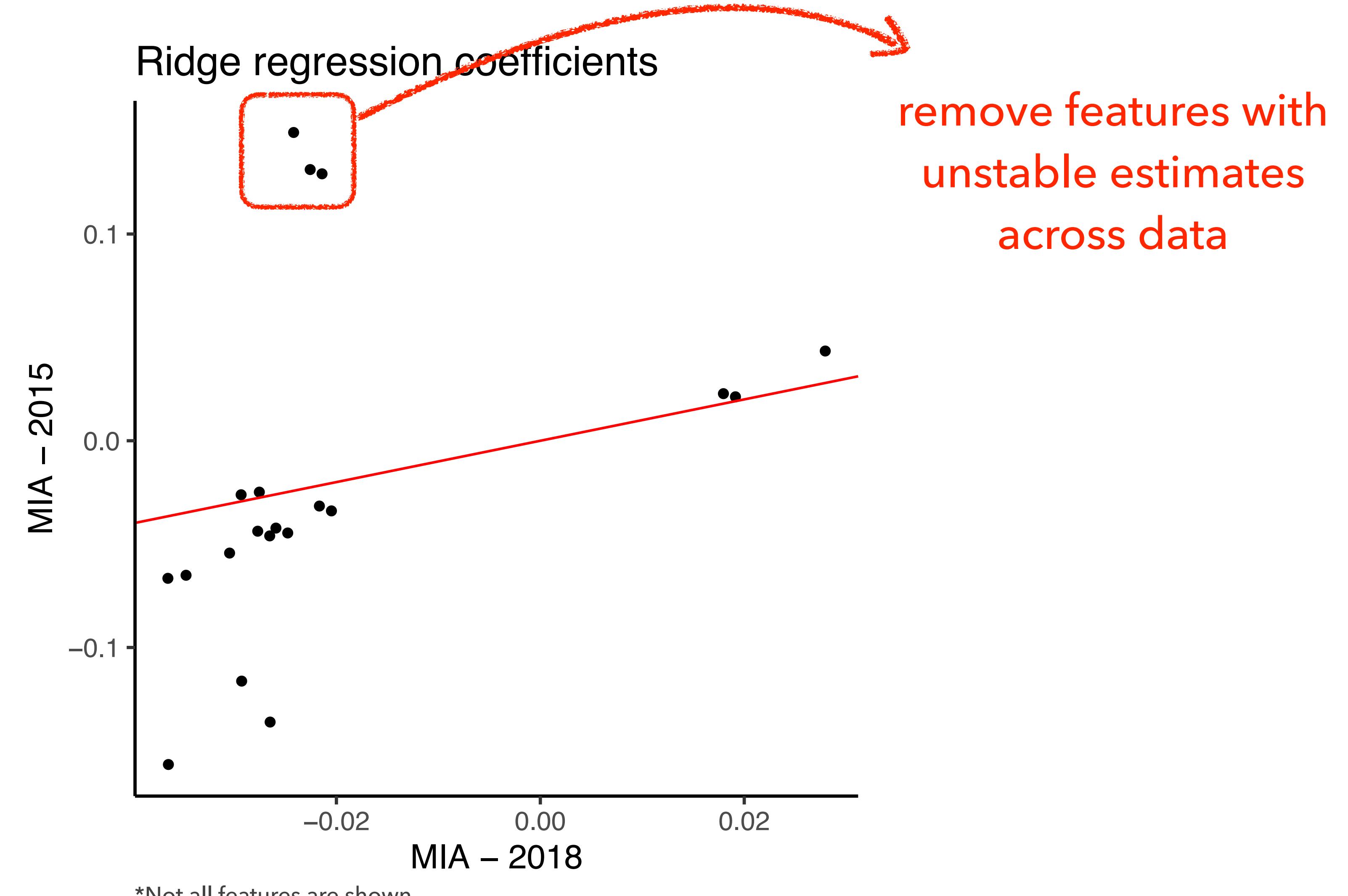
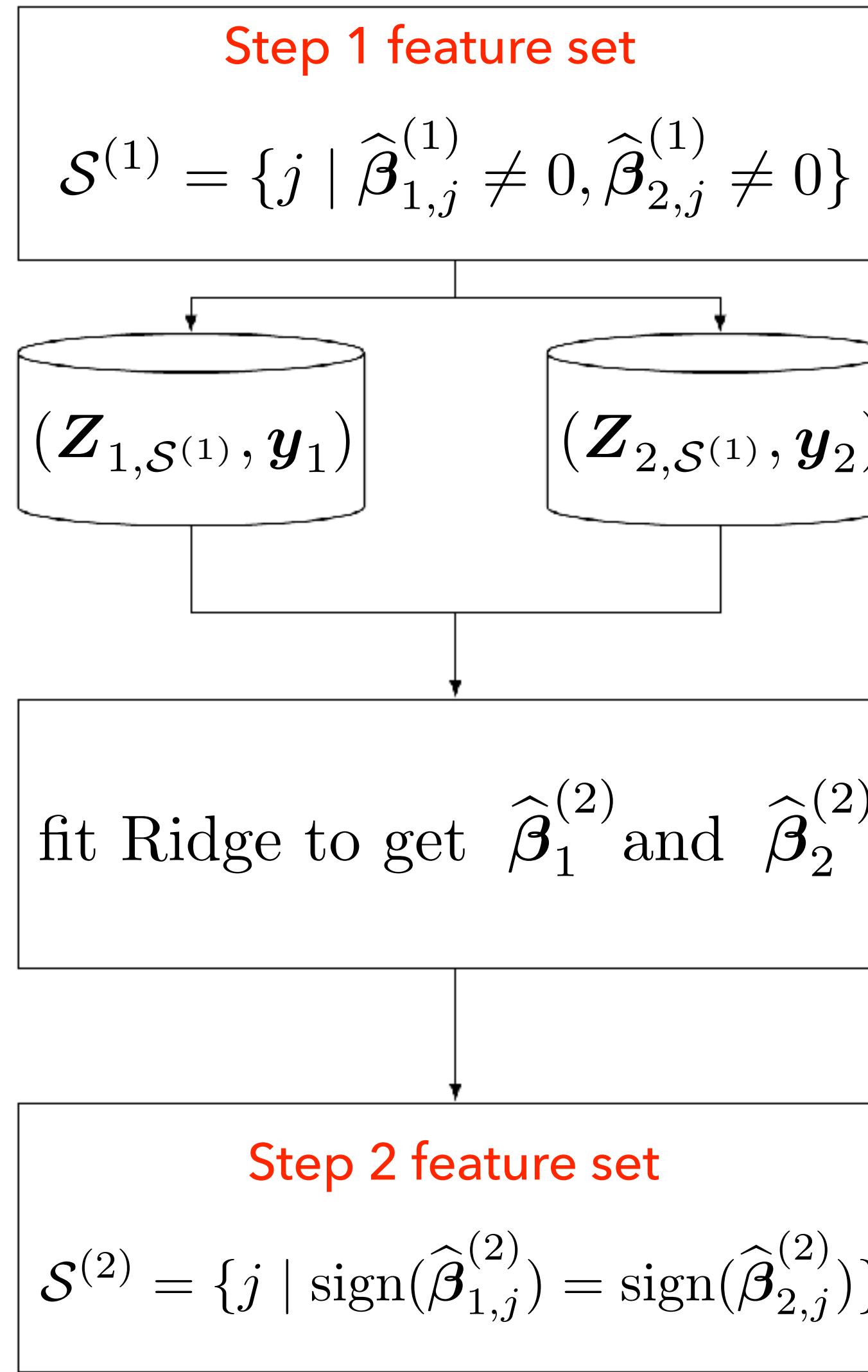
Mean of every log-ratio



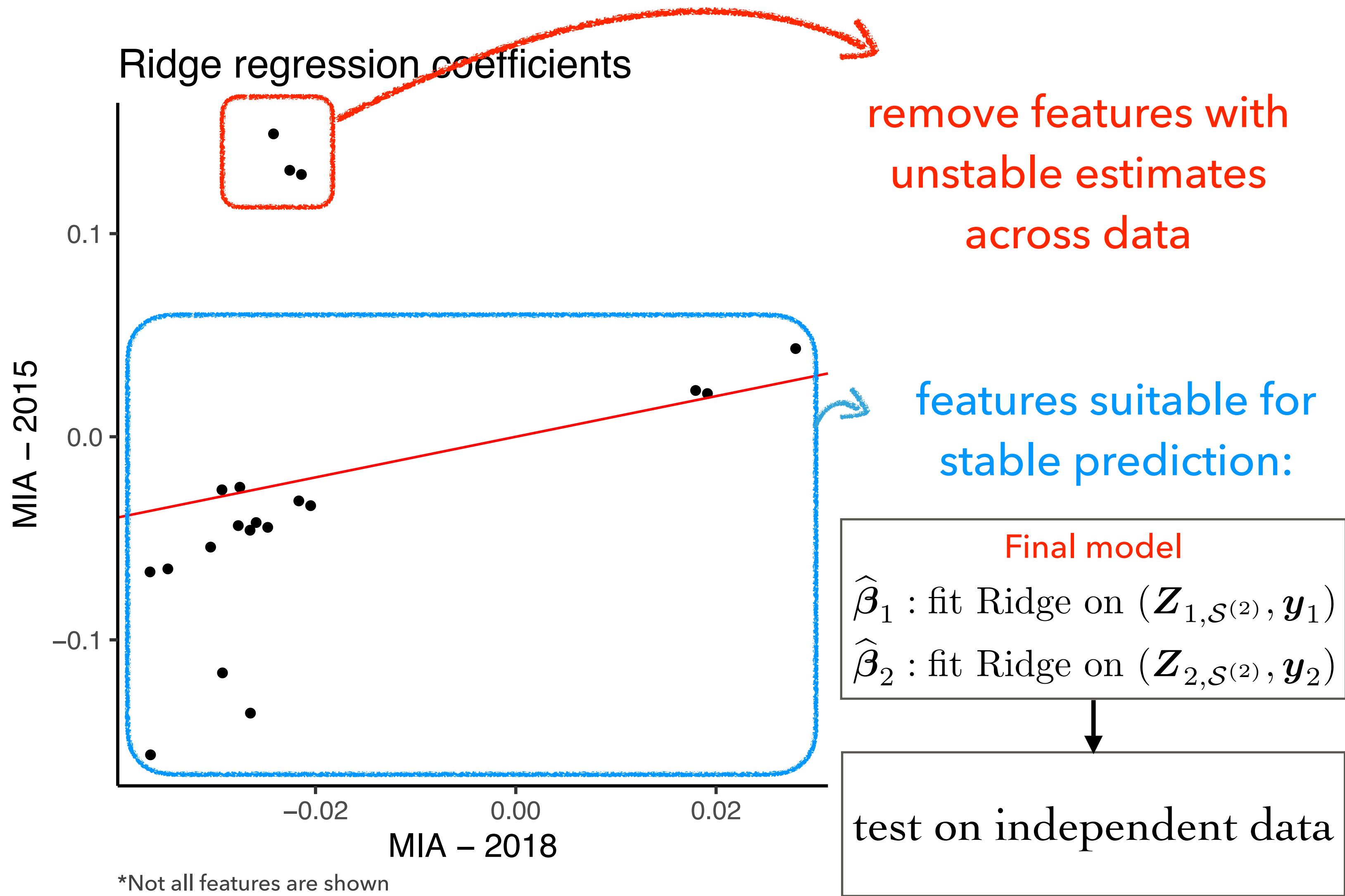
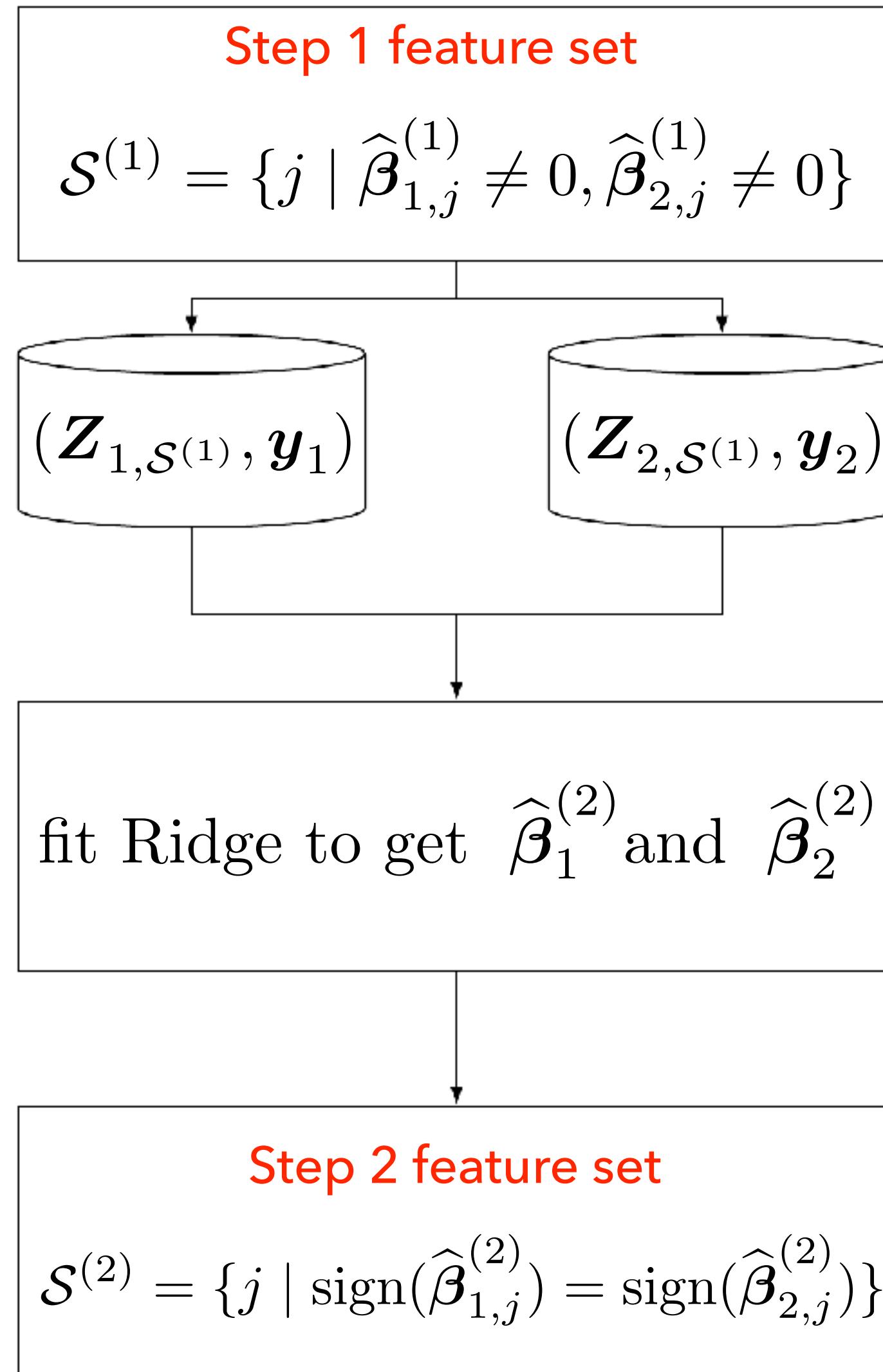
## Step 2: feature estimation stability



## Step 2: feature estimation stability

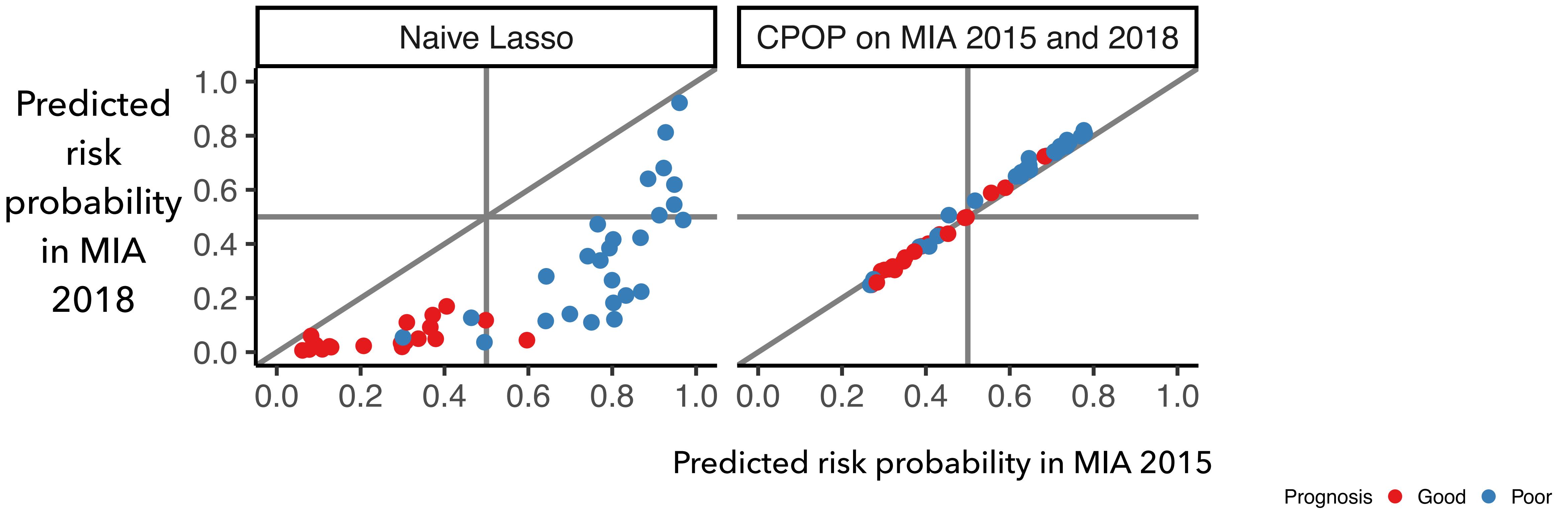


## Step 2: feature estimation stability



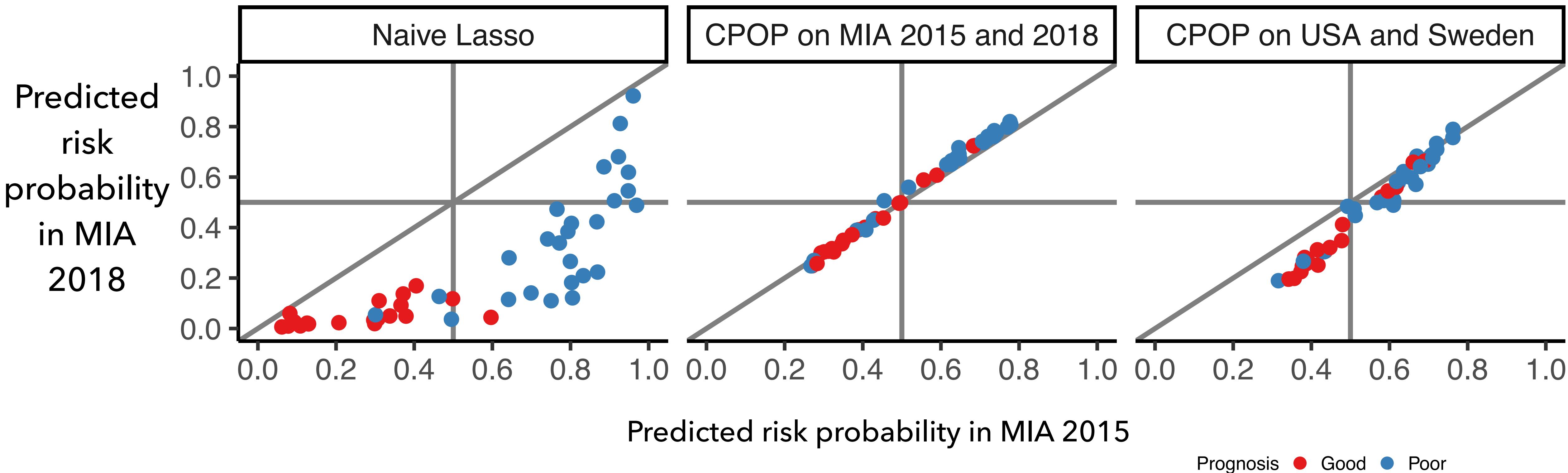
# Results

## Result 1: Predicted risk probability in MIA 2015 and 2018 data



Small deviation in **predicted values** across datasets

# Result 1: Predicted risk probability in MIA 2015 and 2018 data

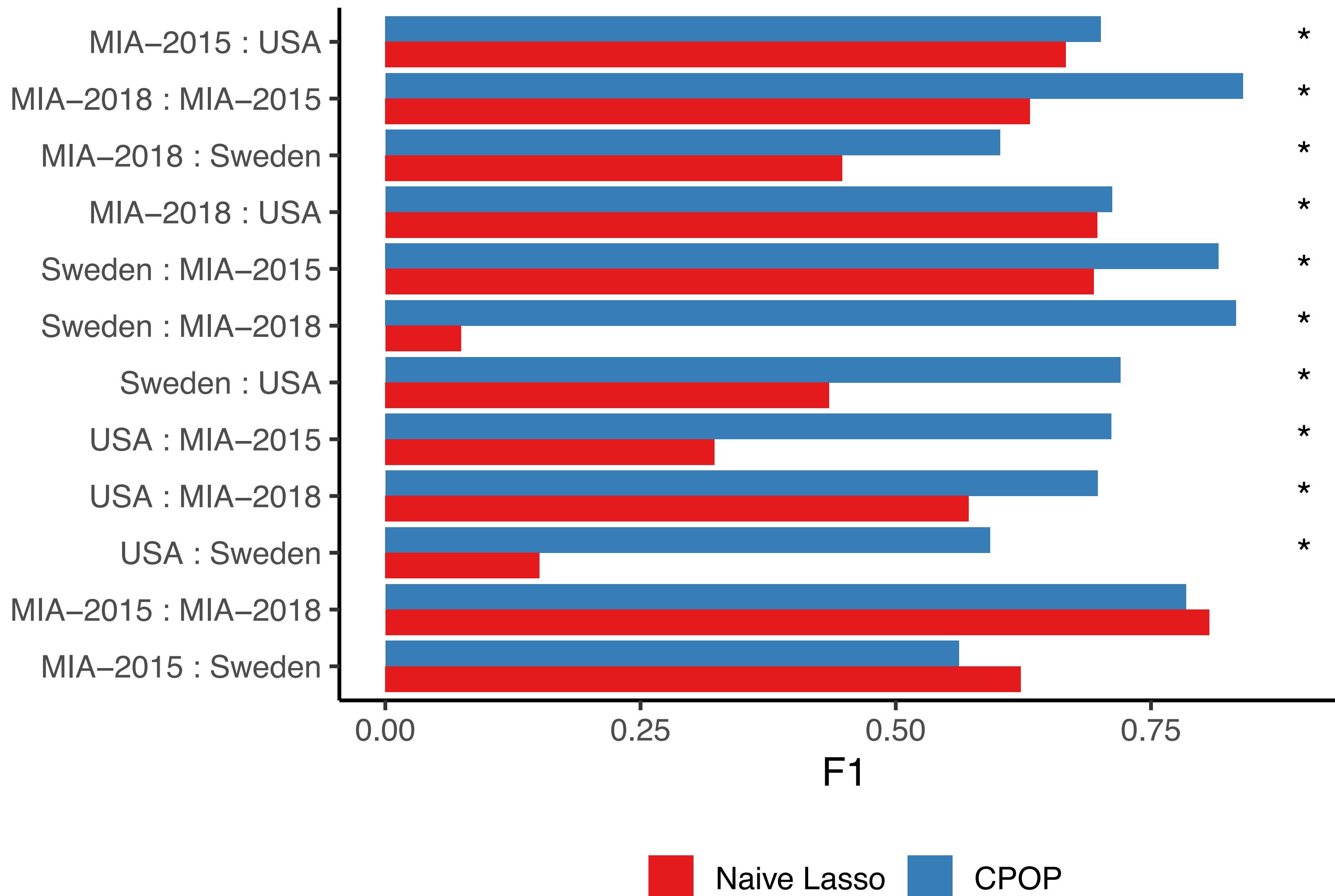


Small deviation in **predicted values** across datasets

## Result 2: four melanoma data

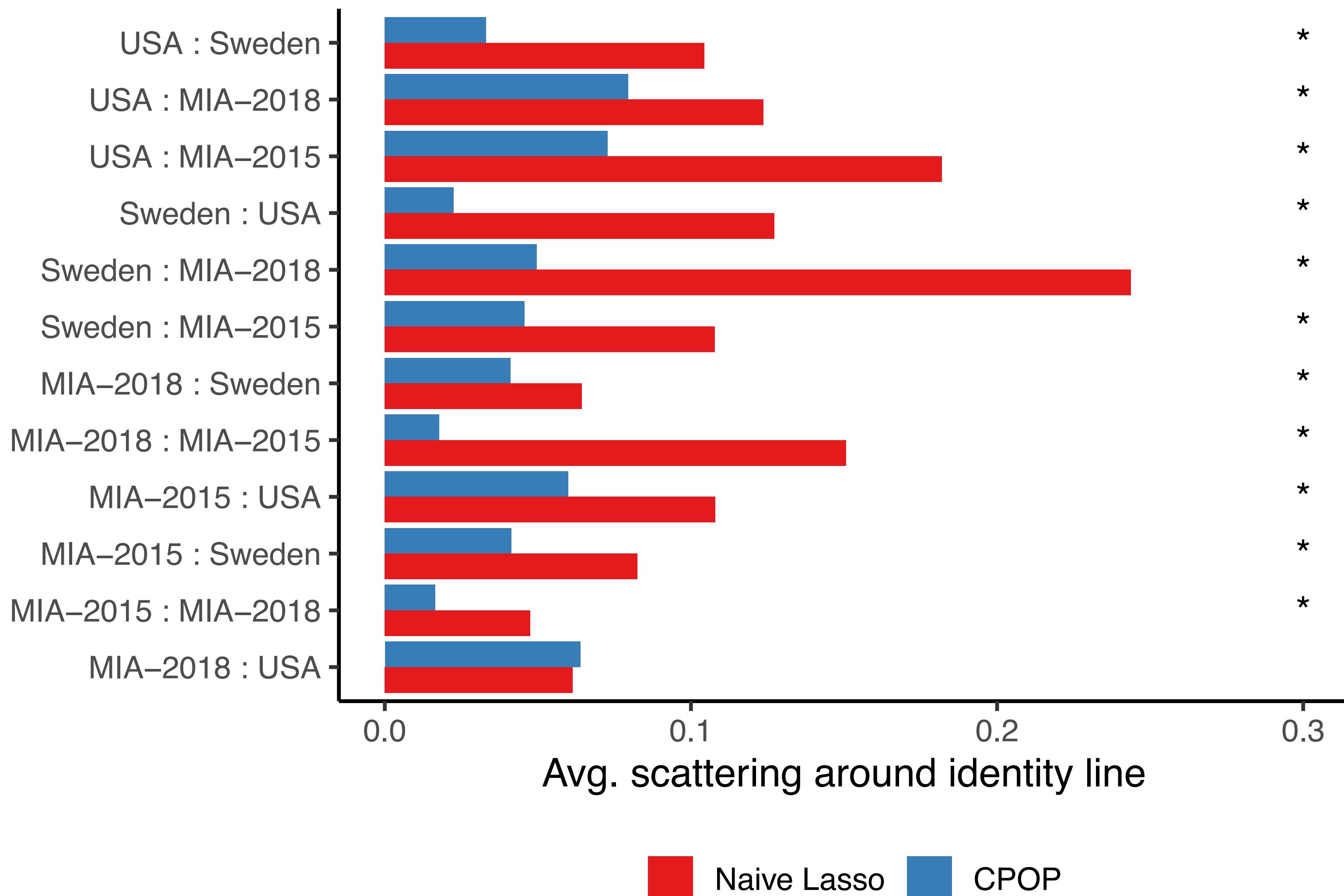
CPOP is highly predictive

Compare F1 statistic (larger is better)



## Result 2: four melanoma data

Compare avg. deviation (smaller is better)



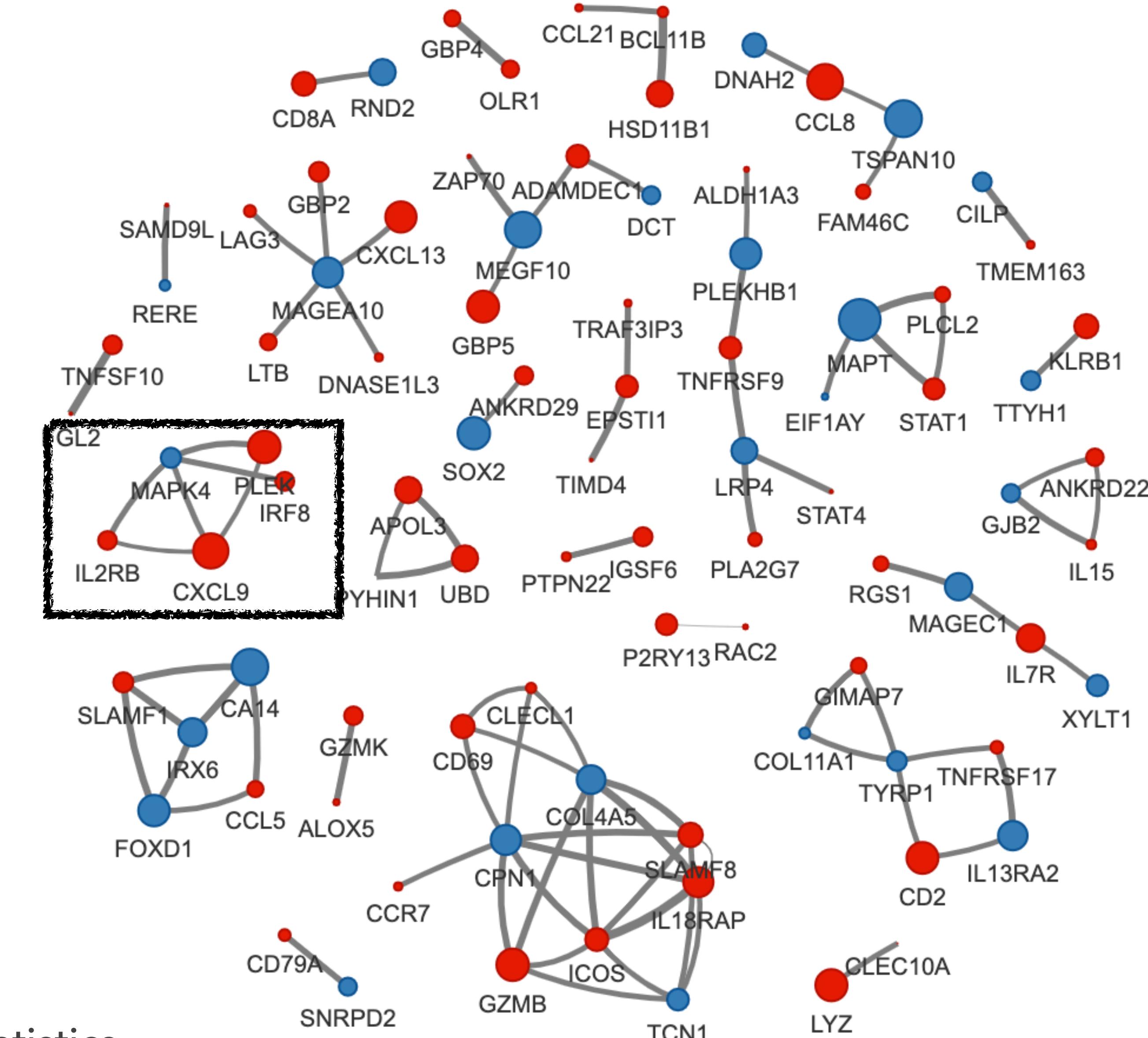
# Result 3: melanoma feature selection

# CPOP features offer new biological interpretation.

# Edges: log-ratio features selected by CPO

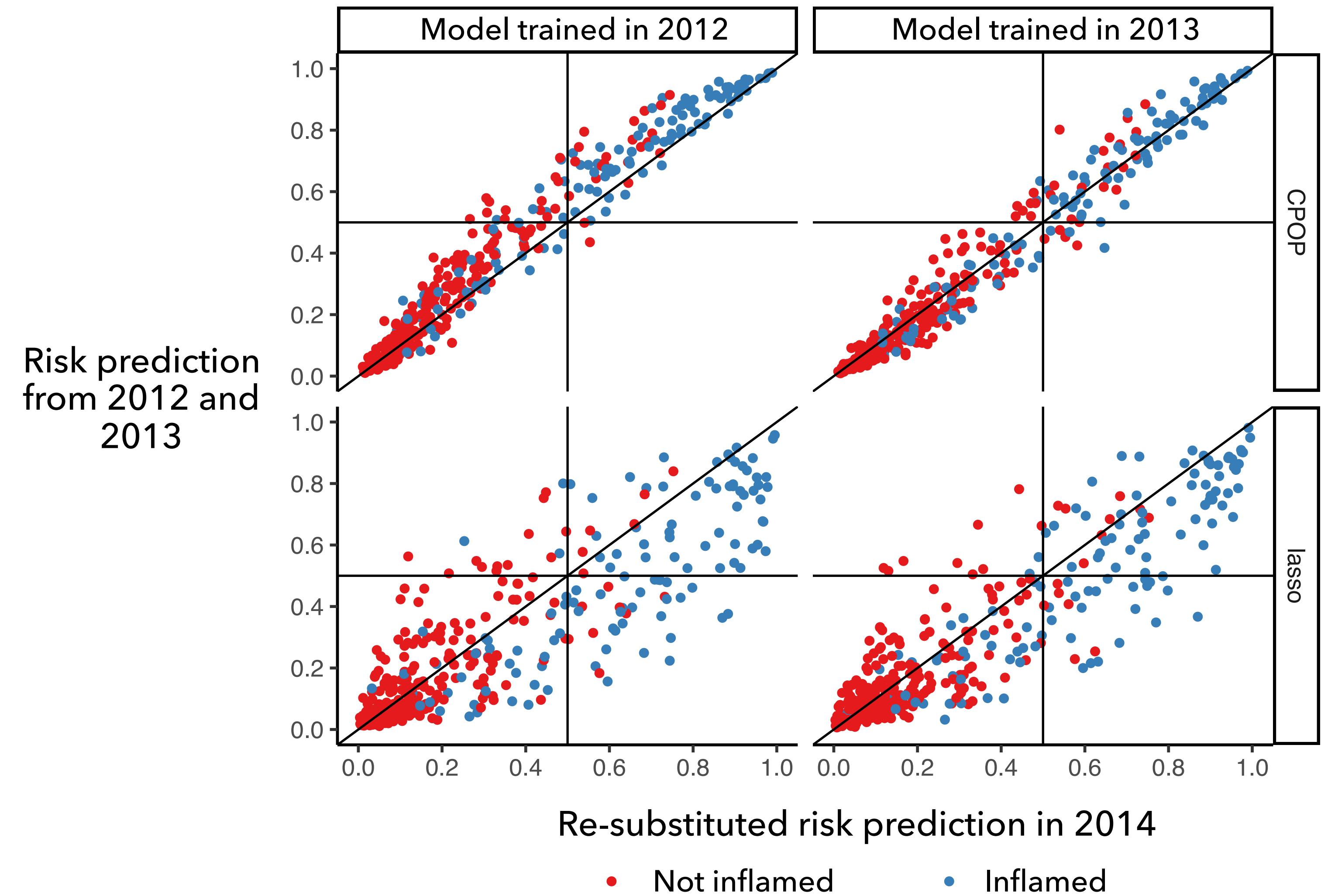
# Vertices: genes comprising the log-ratio feature

# Colour and size: sign and magnitude of a gene's t-statistic



## Result 4: prospective prediction on inflammatory bowel disease

CPOP works on  
prospective experiments



---

## Concluding remarks

1. Integrates clinical implementation constraints into the model
  2. Log-ratios enable prospective and multi-centres prediction
  3. Stable variable selection and estimation components
- ▶ A flexible and adaptable framework focuses on transferability
  - ▶ Potential to work under a CRE grant for implementation

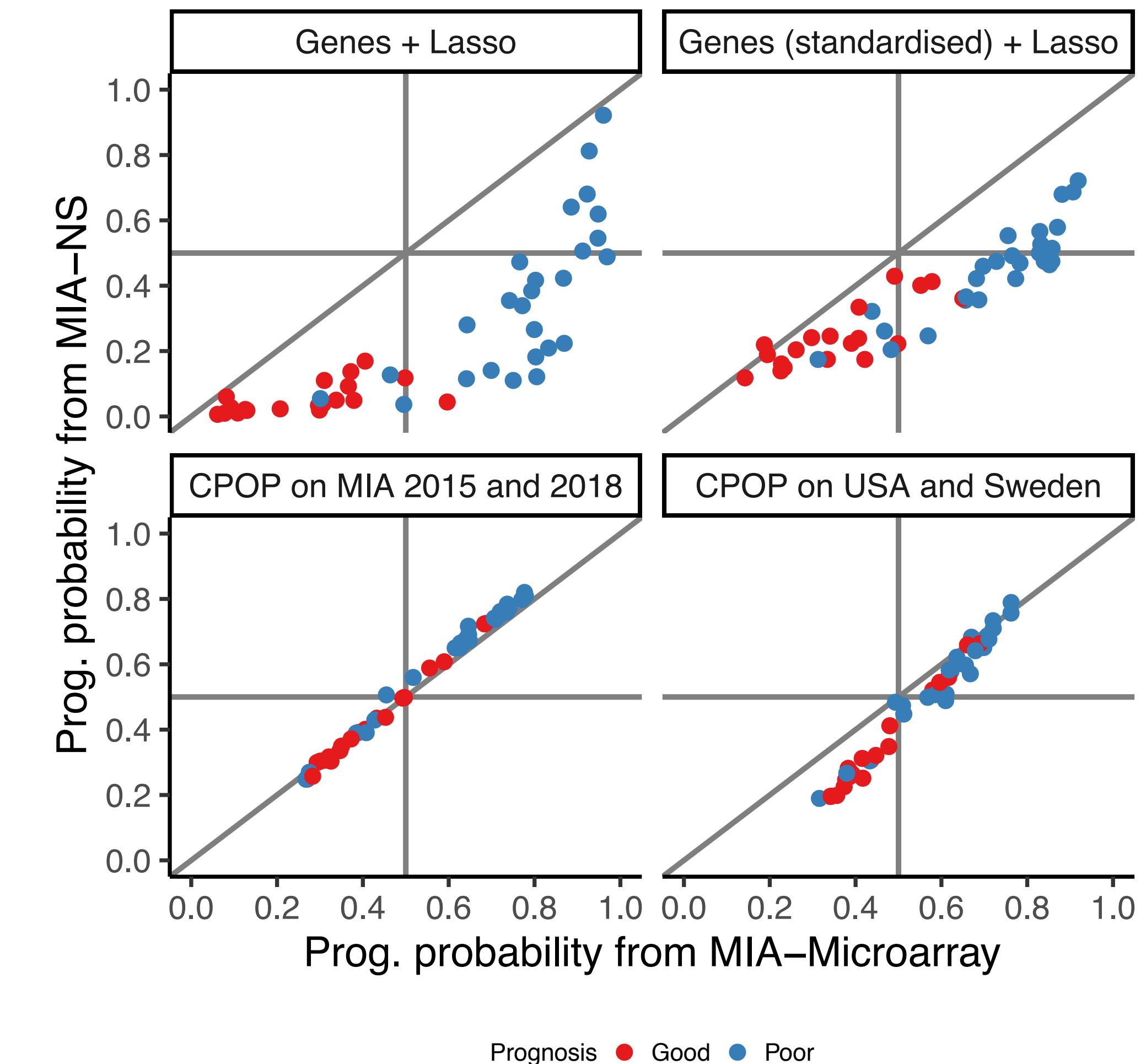
---

# Acknowledgement

- ▶ Supervision:
  - ▶ Jean Yang
  - ▶ Samuel Mueller
  - ▶ Garth Tarr
- ▶ Melanoma Institute Australia
- ▶ Sydney Precision Bioinformatics Group

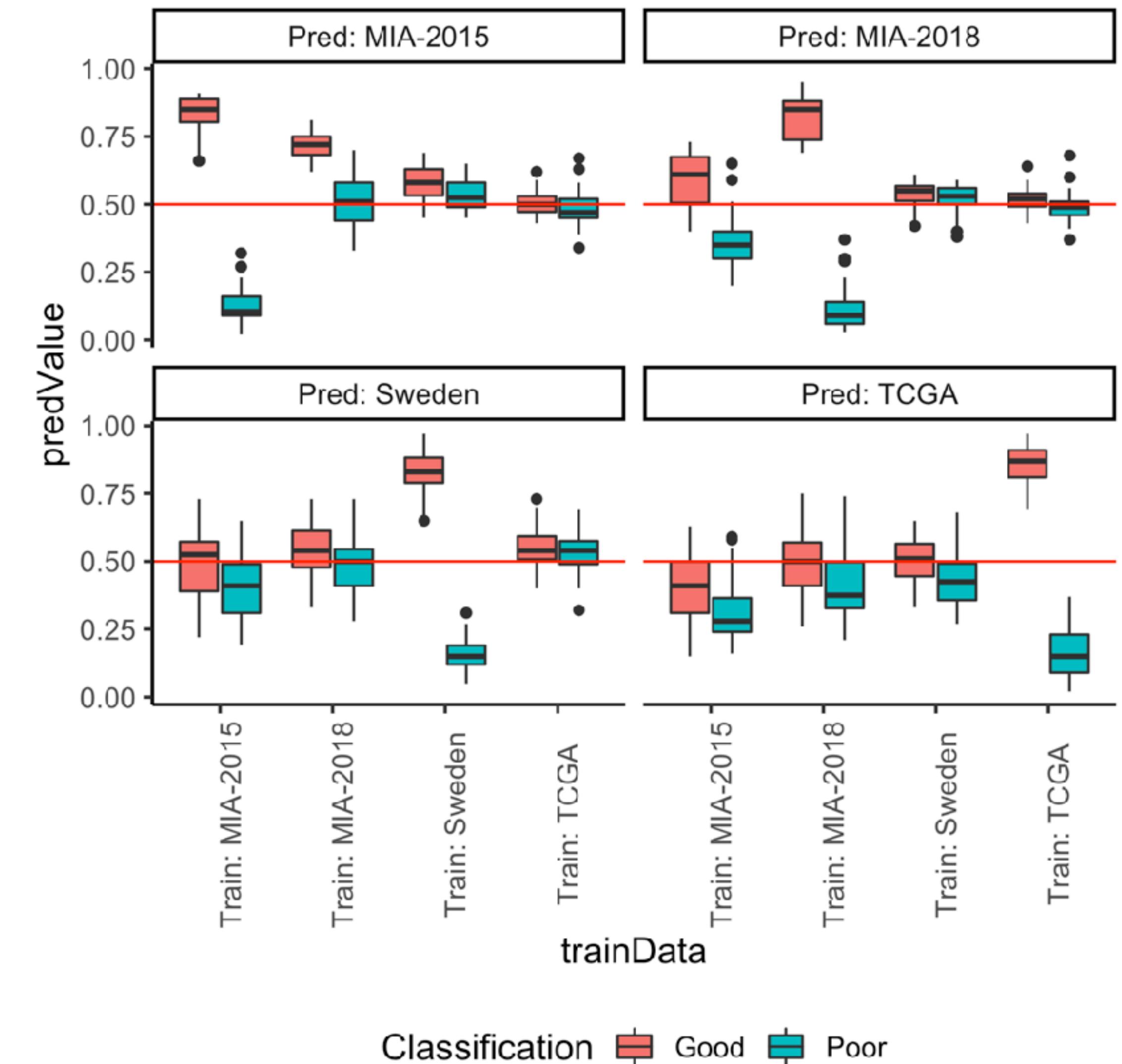
## FAQ 1: why don't you just standardise or scale the data by a reference gene?

- ▶ Contribution of CPOP is not restricted to these log-ratio
- ▶ Normalisation by a “housekeeping” gene doesn’t perform well empirically (left)
- ▶ Standardisation is not possible in prospective experiments
- ▶ Predicted outcome of a single sample depends on other samples that it just happens to be normalised with (McShane et al. 2013)



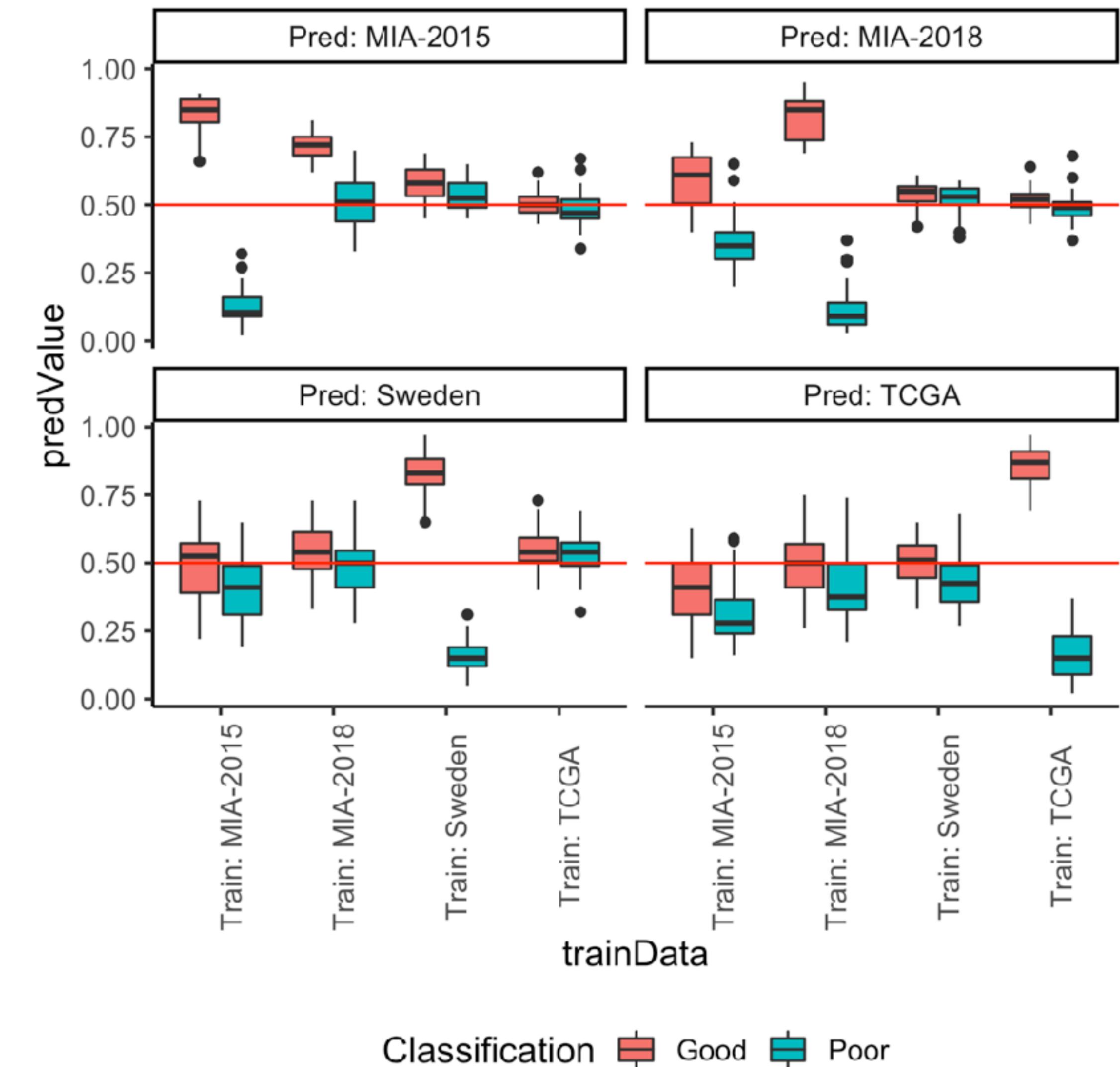
## FAQ 2: have you tried [another machine learning method]?

- ▶ Yes, but the stabilisation techniques in CPOP is not restricted to regression
- ▶ The regression methodologies already offer better transferability results in comparison to, for example, random forest (left)
- ▶ Developing CPOP under the EN regression framework allows us to tap into a much richer collection statistical literature that focuses on stability



## FAQ 2: have you tried [another machine learning method]?

- ▶ The regression model is usually more widely understood/accepted/established in the medical community
- ▶ In terms of future work, we want to extend CPOP to drug sensitivity prediction, which has been established using regression models, e.g. (Garnett et al, Nature, 2012) and (Lee et al, Nature Comm., 2018)



---

## FAQ 3: computational problems with log-ratios

- ▶ CPOP is not for biomarker discovery from thousands of genes. In biomedical research, prior to implementation, a gene set is restricted to about a few hundreds
- ▶ Variable pre-filtering could be applied before the construction of log-ratios (in the IBD data, 700 genes was converted to 250,000 features)

---

## FAQ 4: platform differences

MIA - 2015: Illumina HumanWG-6 v3.0 expression beadchip, 45 samples

MIA - 2018: Customised NanoString probes, 45 samples

USA: log2-FPKM Illumina RNA-Seq, 99 samples

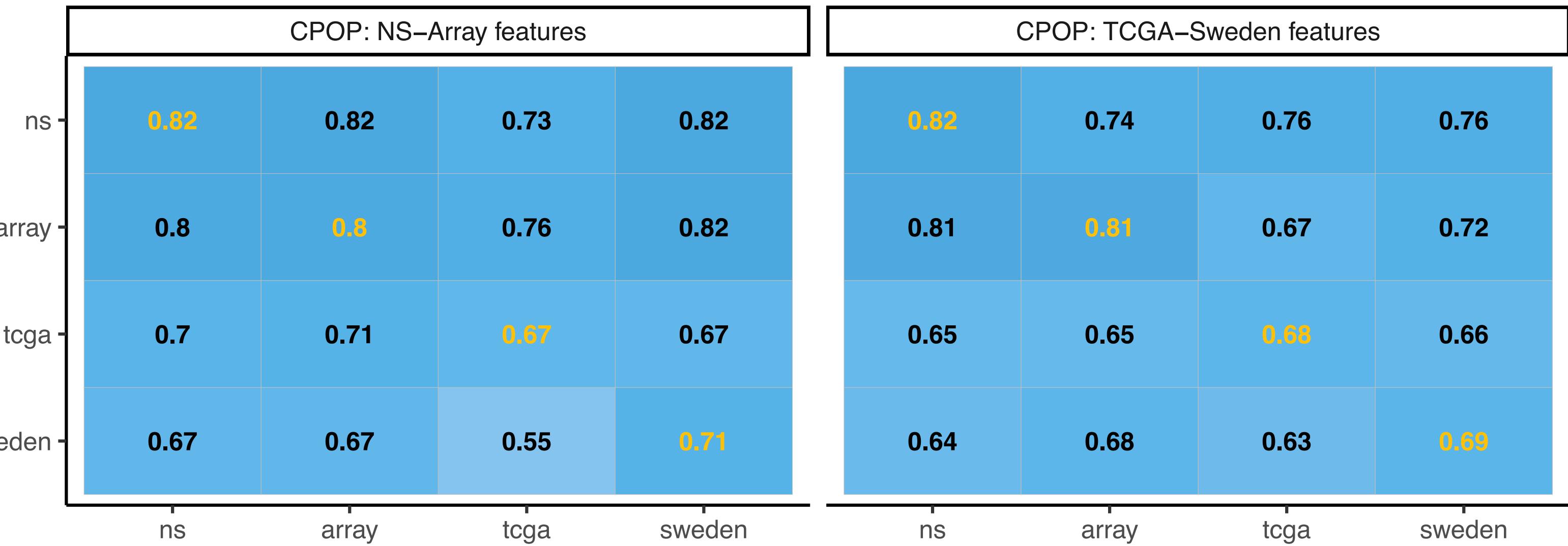
Sweden: Illumina HumanHT-12 version 4 microarray, 95 samples

- ▶ Concordance of log-ratios is much better compare to just gene expression

## FAQ 5: weights in CPOP

CPOP: absolute distance weights

F1 classification statistic under various models



CPOP: KS-statistics weights

