

Oslo Bioinformatics Workshop Week 2025

# Statistical principles in machine learning for small biomedical data

Manuela Zucknick

Oslo Centre for Biostatistics and Epidemiology, University of Oslo  
[manuela.zucknick@medisin.uio.no](mailto:manuela.zucknick@medisin.uio.no)

December 12, 2025

Some of the figures in this presentation are taken from "An Introduction to Statistical Learning, with applications in R" (Springer, 2021) with permission from the authors.

# Schedule for Today

## Schedule

Time	Topic	Presenter
Now	<u>Preparations</u>	
09:00 - 10:00	<u>(Supervised) machine learning with small data</u>	Manuela Zucknick
	<u>R lab 1</u>	Manuela Zucknick
10:15 - 11:15	<u>Overfitting, regularisation and all that</u>	Manuela Zucknick
	<u>R lab 2</u>	Manuela Zucknick
11:30 - 12:00	<u>The potential of Bayesian modelling for complex biomedical data</u>	Manuela Zucknick

# Github, Workshop webpage and Posit Cloud project

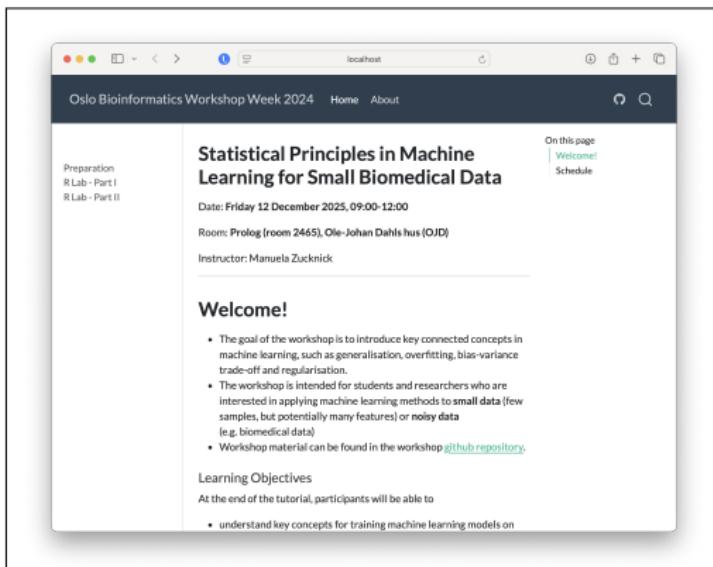
- Github:

<https://github.com/ocbe-uio/workshop-stat-higdim/>

- Workshop webpage:

<https://ocbe-uio.github.io/workshop-stat-higdim/>

- Posit Cloud project: <https://posit.cloud/content/5131383/>



# Some topics for this morning

## Part 1

- What is supervised machine learning?
- What do we mean by small data?
- What can we do to improve ML with small data?
  - Restrict the model space → Regularisation
  - Borrow information → Include known structure in the model

## Part 2

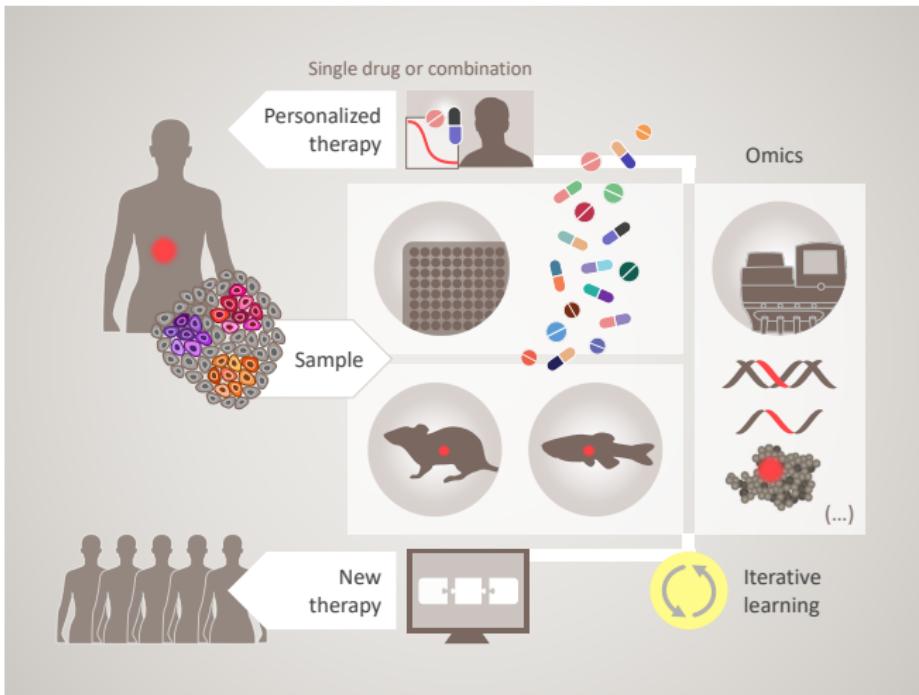
- Overfitting
- Variance vs bias
- Model selection, assessment & validation
- Prediction performance
- Resampling: Cross-validation

## Introductory example:

Integrative omics for personalized cancer therapy

# Personalized cancer therapy

...aims to find the best therapy for each patient based on data about the patient and tumor (e.g. genomic data).



# Predict sensitivity to multiple drugs $\mathbf{Y}$ from multi-omics $\mathbf{X}$

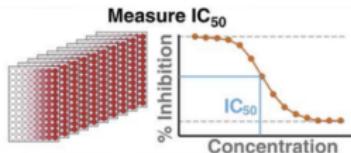
$$\mathbf{Y} = \mathbf{XB} + \epsilon$$

- **Multivariate  $\mathbf{Y}$ :**

Drug dose response

*drug sensitivity*

$$n \text{ cell lines} \left[ \begin{array}{c|c|c} & \dots & \\ \text{y}_{\bullet 1} & \dots & \text{y}_{\bullet m} \\ & \dots & \end{array} \right] = \mathbf{Y}$$

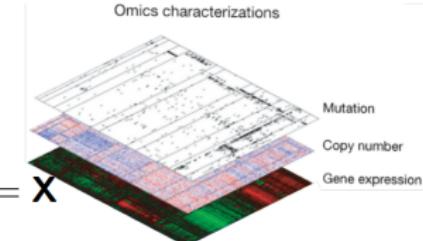


Source: Yang, et al. 2017

- **Heterogeneous  $\mathbf{X}$ :**

Integrative omics

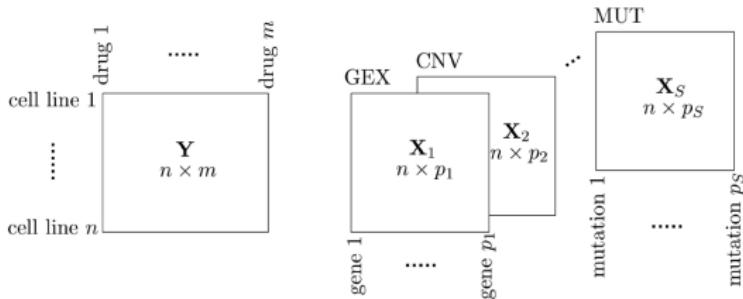
$$n \text{ cell lines} \left[ \begin{array}{c|c|c} \text{gene expression} & \text{copy number} & \text{mutation} \\ \hline \text{X}_1 & \text{X}_2 & \text{X}_3 \\ | & | & | \\ \end{array} \right] = \mathbf{X}$$



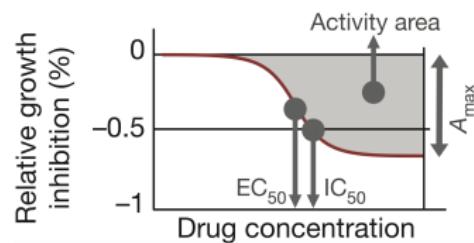
Source: TCGA, 2013

# Challenges and opportunities (1)

- Small sample size
- Several types of input data  $\mathbf{X}$ :  
E.g., gene expression, copy number, mutation
- Multivariate response  $\mathbf{Y}$



- Unclear how to define  $\mathbf{Y}$



## Challenges and opportunities (2)

The data are highly **structured**:

- In Y:** relationships between drugs, e.g. due to similar chemical drug composition, same target genes/pathways
- In X:** relationships between molecular data sources

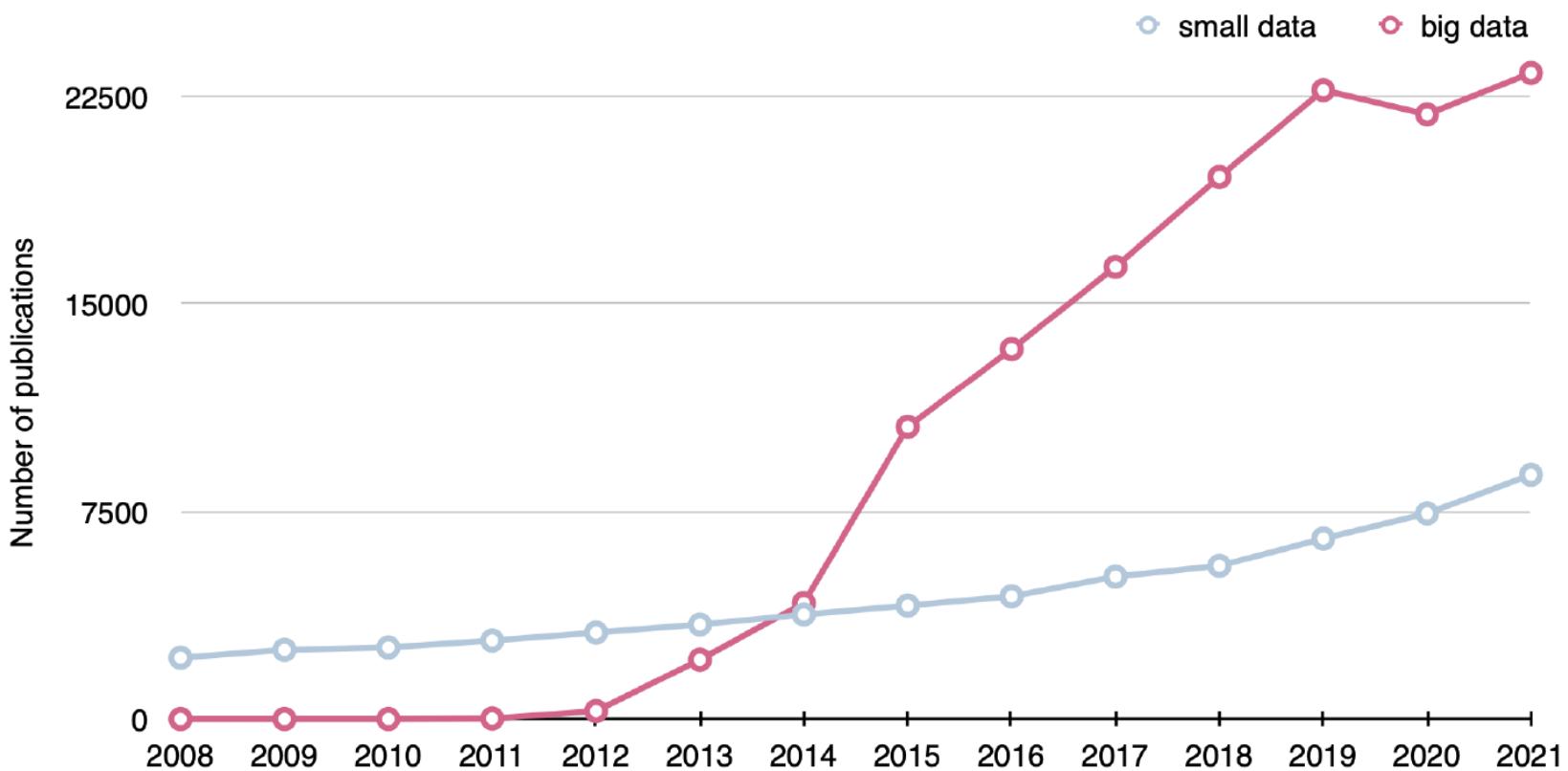
a	Function	Memory	Environment	Message	Product	Result
b	Central dogma of molecular biology	Genome (DNA)	Epigenome and other regulatory elements (e.g. chromatin modifications, miRNA, TFs) 	Transcriptome (mRNA) 	Proteome (protein) 	Phenome (cell, tissue, organism) 
c	Data types	CN, SNPs, LOH 	Histone modification TF binding, miRNA, methylation 	GE 	Protein expression 	Phenotype, clinical characteristics 

Ickstadt et al. (2018)

# (Supervised) Machine Learning with Small Data

Manuela Zucknick (with slides from Maren Hackenberg)

# Machine learning with small data



# Machine learning with small data

- What do we mean by “small data”?
  - Implications for machine learning?
  - Aspects when building (multi-omic) machine learning predictors of drug response (e.g. Sammut et al. Nature 2022):
    1. Biological knowledge +
    2. Feature selection +
    3. Prioritisation of accessible data types +
    4. Machine learning algorithms
- Develop ML methods that allow us to consider aspects 1 to 3.

# What is supervised machine learning?

## Supervised learning

refers to the task of inferring a functional relationship between **input data matrix  $\mathbf{X}$**  (e.g. gene expression array measurements) and **output data vector  $Y$**  (= response/ outcome).

The input data are used for **predicting** the outcome.

$$Y = f_{\beta}(\mathbf{X}) + \epsilon,$$

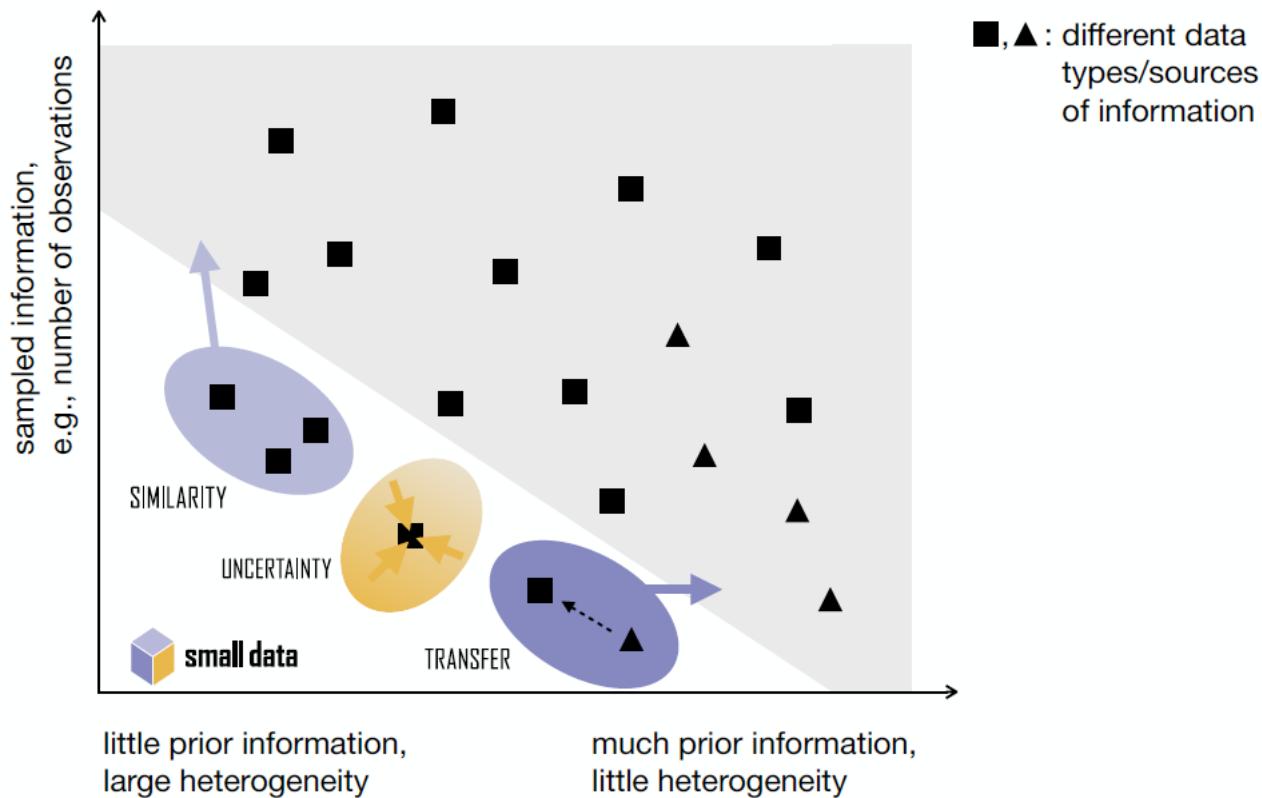
where  $\epsilon$  captures measurement errors and other discrepancies, e.g. by  $\epsilon \sim N(0, \sigma^2 I_n)$ .

In classical statistics, this task is usually performed by **(generalised) linear regression models**.

# What do we mean by small data?

- Large p, small n ( $p > n$ )
- Potentially, more variables in the model than we have samples
- Classical statistical methods (e.g. linear regression) do not work:
- More parameters (e.g. regression coefficients) to estimate than observations for estimating them
- Even if all parameters can be estimated: Danger of over-fitting
- Example: Predict treatment response using gene expression data  
( $n \sim 100$ ,  $p \sim 20000$ )

# What do we mean by small data?



What can we do?

- (1) Restrict the model space
- (2) Borrow information across observations
- (3) Increase sample size ☺

# What can we do?

## (1) Restrict the model space

- (A) Careful feature engineering:
  - Preselect variables by biological relevance
  - Non-specific filtering, e.g. keep only variables with variance across observations larger than a threshold
- (B) Make use of known structure in the data (biological knowledge)
- (C) Use of regularisation techniques:
  - L1 and L2 penalisation
    - add a penalty term to the loss function to reduce the complexity of the model
    - Bayesian equivalents: restrictions on the prior (Bayesian variable selection)
  - Early stopping
    - train a model iteratively only until the validation error starts to decrease (boosting, neural networks)
  - Dropout regularisation
    - randomly dropping out neurons while training (neural networks) or
    - randomly dropping features when building a regression tree (random forest)

## Penalised regression

- Standard regression cannot deal with  $p \gg n$ :
  - The maximum-likelihood estimate  $\hat{\beta} = \arg \max_{\beta} \ell(\beta)$  does not exist ( $\ell = \log\text{-likelihood}$ ).
- **Solution:**  
Penalise the likelihood function by subtracting a penalty term and maximise penalised log-likelihood instead:

$$\hat{\beta} = \arg \max_{\beta} (\ell(\beta) - \lambda \|\beta\|)$$

- $\lambda$  is a **penalty parameter**,
- $\|\beta\|$  represents the size of the regression coefficient vector,
- The larger  $\lambda$  is chosen, the more the algorithm is encouraged to find a solution where  $\|\beta\|$  is small  $\rightarrow$  **shrinkage**.

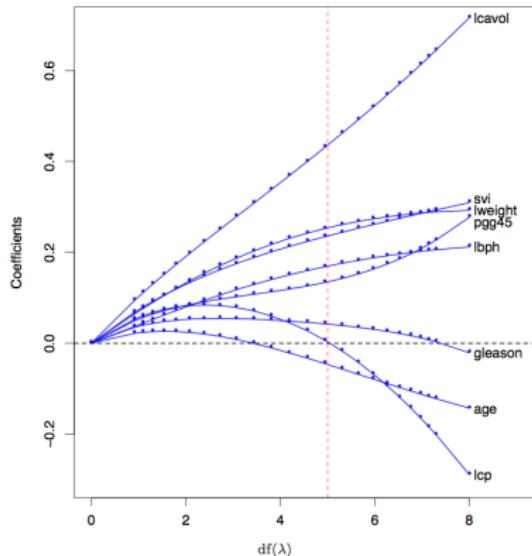
## Penalised regression

- Examples for penalty terms:
  - Ridge regression (Hoerl and Kennard 1970):
$$\lambda \|\beta\| := \lambda \sum_{g=1}^p \beta_g^2 \quad \rightarrow \mathbf{L}_2 \text{ penalty}$$
  - Lasso regression (Tibshirani 1996):
$$\lambda \|\beta\| := \lambda \sum_{g=1}^p |\beta_g| \quad \rightarrow \mathbf{L}_1 \text{ penalty}$$
  - Elastic net (Zou and Hastie 2005):  
Combination of both ridge and lasso penalty:
$$\lambda_1 \sum_{g=1}^p |\beta_g| + \lambda_2 \sum_{g=1}^p \beta_g^2$$
- Advantage of lasso and elastic net:  
Both will produce a sparse solution, where only a few genes have estimate  $\hat{\beta}_g \neq 0$ .

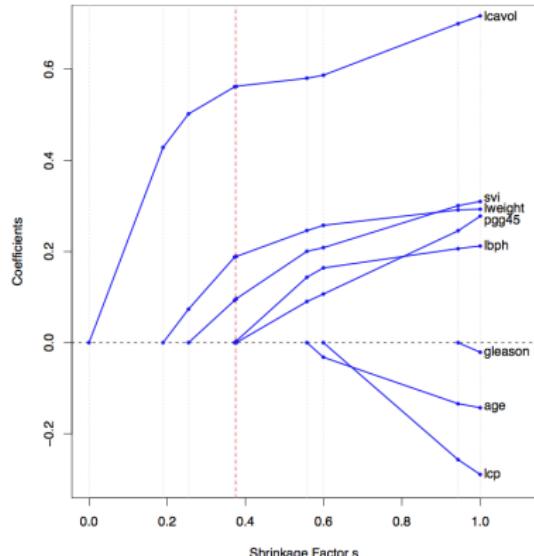
# Penalised regression

Examples for coefficient paths relative to penalty  $\lambda$ :

Ridge regression



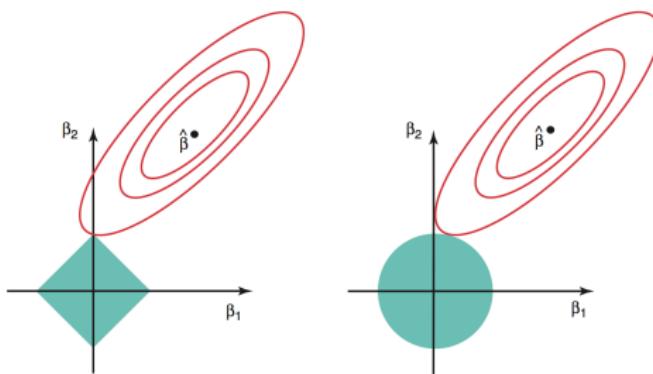
Lasso regression



Hastie et al. (2009), Figures 3.8 and 3.10

## Penalised regression

- Ridge regression  $L_2$ : shrinks all coefficients to small, but non-zero values.
- Lasso regression  $L_1$ : shrinks some coefficients to exactly zero.
- Elastic net: mixture of the two: does shrink some coefficients to exactly zero. Keeps more variables if there is correlation.



**FIGURE 6.7.** Contours of the error and constraint functions for the lasso (left) and ridge regression (right). The solid blue areas are the constraint regions,  $|\beta_1| + |\beta_2| \leq s$  and  $\beta_1^2 + \beta_2^2 \leq s$ , while the red ellipses are the contours of the RSS.

## Different penalties for different types of data

Assume two data matrices  $\mathbf{X}$  and  $\mathbf{Z}$ :

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \epsilon$$

- **Mandatory covariates:** Do not penalise the parameters  $\gamma$ :

$$\ell_{\text{pen}}(\beta, \gamma) = \ell(\beta, \gamma) - \lambda \|\beta\|$$

e.g. with R packages `glmnet` or `penalized`

- **Several types of molecular data sets:**

Allow different penalties for  $\beta$  and  $\gamma$ :

$$\ell_{\text{pen}}(\beta, \gamma) = \ell(\beta, \gamma) - \lambda_\beta \|\beta\| - \lambda_\gamma \|\gamma\|$$

e.g. with R packages `GRridge` (Van de Wiel *et al.*, 2016)  
<http://www.few.vu.nl/~mavdwiel/grridge.html>)

## Different penalties for different types of data

Assume two data matrices  $\mathbf{X}$  and  $\mathbf{Z}$ :

$$Y = \mathbf{X}\beta + \mathbf{Z}\gamma + \epsilon$$

- Several types of molecular data sets:
- Alternative: **Combine all data and use one penalty**, after scaling all features to unit variance to ensure that the data sources are treated equally.
- Example: Elastic Net models in Barretina et al. (2012)

## Bayesian interpretation of penalised regression

- Ridge regression as a penalised log-likelihood problem ...

$$\hat{\beta} = \arg \max_{\beta} (\ell - \lambda \sum_{i=1}^p \beta_i^2)$$

- ... is equivalent to *maximum a posteriori* solution of Bayesian linear regression with Gaussian prior

$$p(\beta|\tau) = N(0, \tau I_p), \text{ where } \tau = 1/(2\lambda) :$$

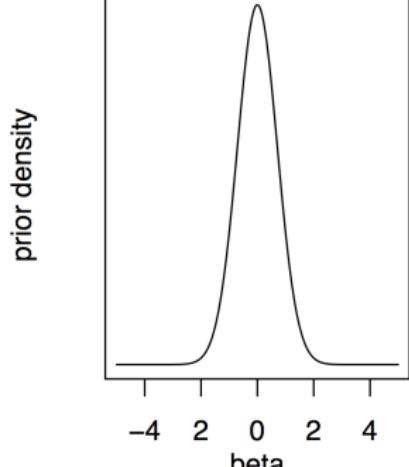
$$\begin{aligned} \text{posterior} &\propto \text{likelihood} \times \text{prior} \\ p(\beta|X, Y, \tau) &\propto p(Y|\beta, X)p(\beta|\tau) \\ \Leftrightarrow \log(\beta|X, Y, \tau) &\propto \ell + \log p(\beta|\tau) \\ &\propto \ell - \lambda \sum_{i=1}^p \beta_i^2 - C \end{aligned}$$

# Bayesian interpretation of penalised regression

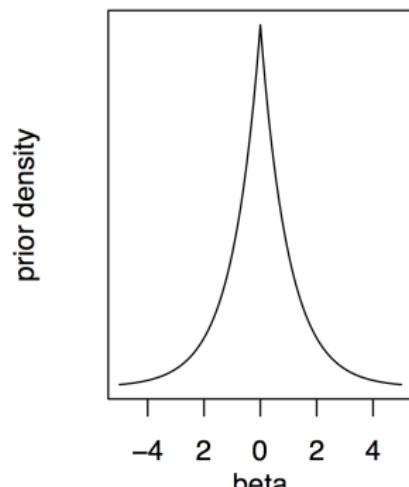
- Generalisation: Bridge regression (e. g. Frank and Friedman 1993)

$$\hat{\beta} = \arg \max_{\beta} (\ell - \lambda \sum_{i=1}^p |\beta_i|^q) \quad (q > 0)$$

Ridge ( $q = 2$ ): Gaussian prior



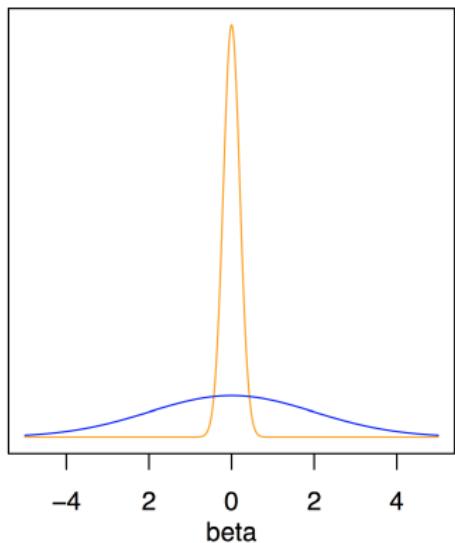
Lasso ( $q = 1$ ): Laplace prior



# Bayesian hierarchical model for variable selection (BVS)

- Bayesian variable selection model with indicator variable

$$\gamma_i = \begin{cases} 1 & \text{variable } i \text{ is included} \\ 0 & \text{variable } i \text{ is excluded} \end{cases}$$



E. g. normal mixture prior  
(George and McCulloch 1993):

$$\beta_i | \gamma_i \sim (1-\gamma_i)N(0, \sigma^2) + \gamma_i N(0, g\sigma^2)$$

where  $\sigma^2 > 0, g > 0$

# Bayesian hierarchical model for variable selection (BVS)

## Advantages:

- Flexibility in penalization through large variety of possible prior distributions for  $\beta$  (hierarchical models)
- Full posterior distributions, including probabilities for the selection of variables and posterior distributions for  $\beta$ ;
- Improve prediction performance by Bayesian Model Averaging

Bayesian Model Averaging for Linear Regression Models



Adrian E. Raftery; David Madigan; Jennifer A. Hoeting

*Journal of the American Statistical Association*, Vol. 92, No. 437 (Mar., 1997), 179-191.

Stable URL:  
<http://links.jstor.org/sici?&sici=0162-1459%28199703%2992%3A437%3C179%3ABMAFLR%3E2.0.CO%3B2-9>

*Journal of the American Statistical Association* is currently published by American Statistical Association.

# What can we do?

## (2) Borrow information

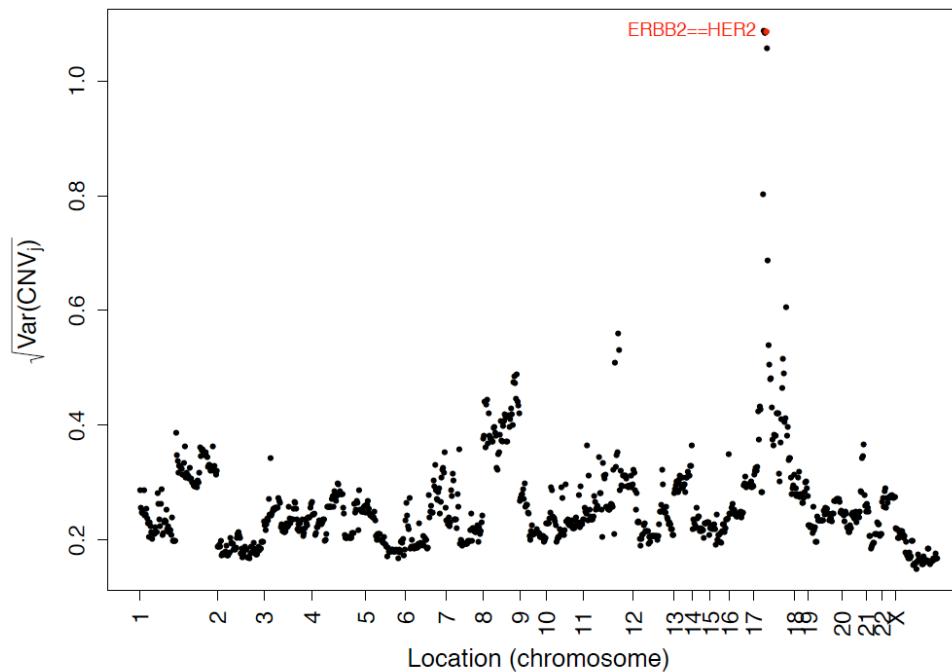
- Borrow information across observations in the data set
- If there is correlation, include this in your model
  - between variables (e.g. MRF prior for defining which variables to include together)
  - between samples (covariance matrix)
- Borrow information from external knowledge
  - E.g., use pathways to determine which genes should be included together
- Borrow information across data sets: transfer learning

# Make use of external (biological) knowledge

- (1) Use known relationships with one data source (CNV) to guide the variable selection in another (gene expression)
- (2) Combine the data-driven ML approach with knowledge-driven mechanistic modelling
- (3) Make use of correlations in the data
  - between input variables - to restrict the model space
  - between response variables - to borrow information

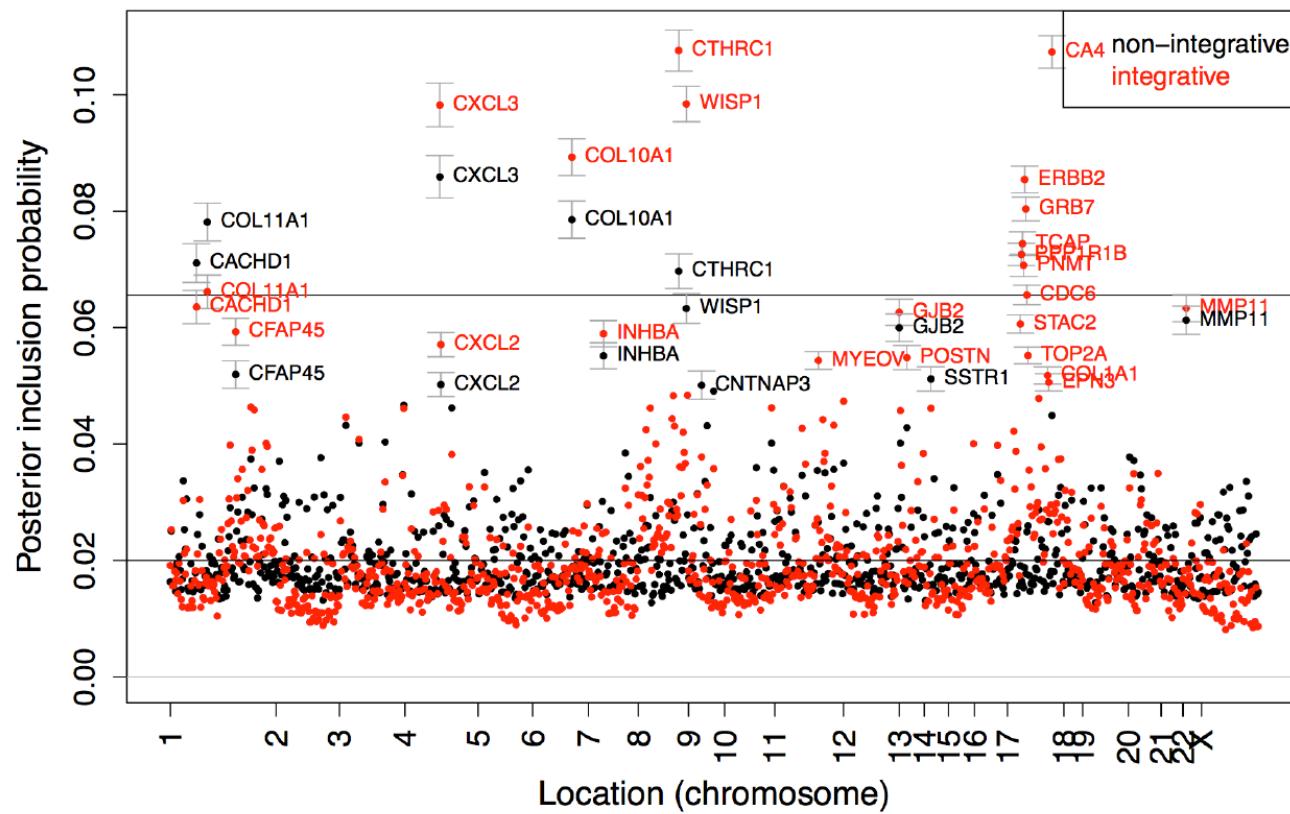
# (1) Use known relationships with one data source to guide the variable selection in another

Std. dev. of CNV data of HER2-pos. breast cancer and healthy tissue samples



**Idea:** Use CNV information to weigh prior inclusion probabilities of gene expression variables in Bayesian variable selection

# (1) Use known relationships with one data source to guide the variable selection in another



HER2 (= ERBB2) only selected in integrative analysis

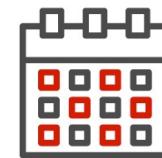
## (2) Combine the data-driven ML approach with knowledge-driven mechanistic modelling

An exemplary small data challenge: Learn disease trajectories of patients with spinal muscular atrophy



**Baseline characterisation**

- age
- SMA subtype
- ...



**Different motor function tests over time**

- RULM
- HFMSE
- ...



Latent health status

$$\frac{d}{dt} \mu(t) = ?$$

Explicit model



Subgroup-specific local models



Heterogeneity



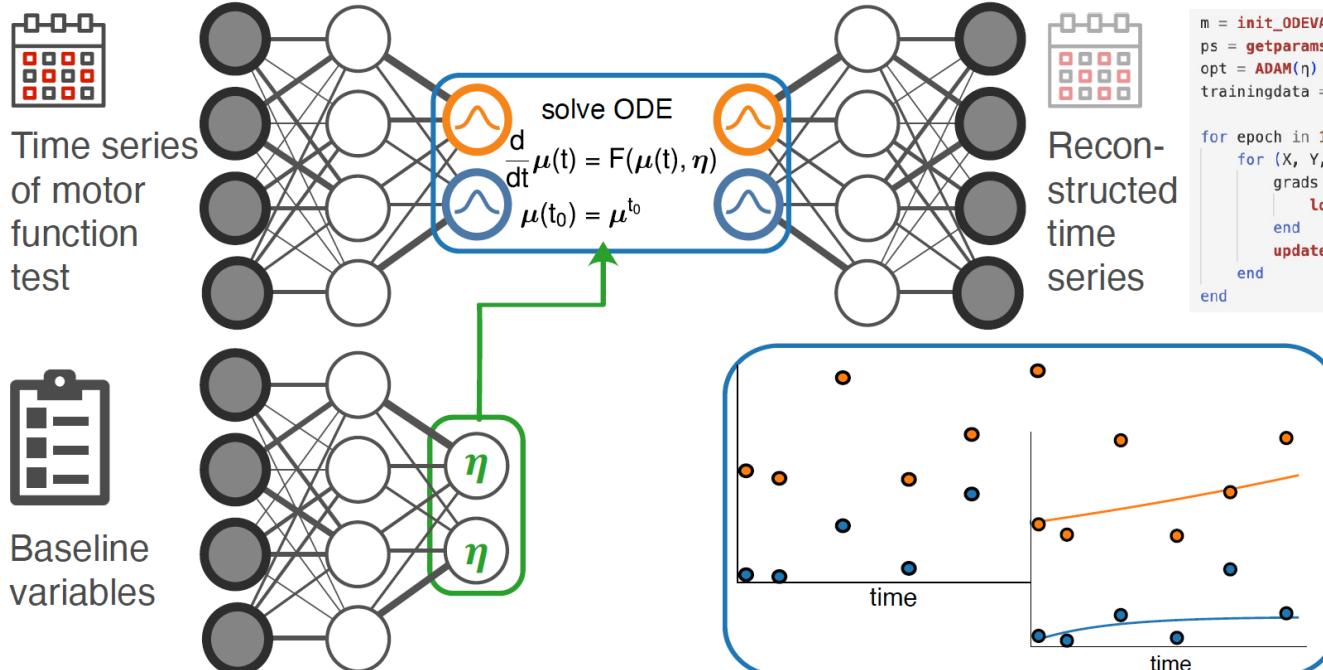
Irregular time points



RULM  
HFMSE  
**Different motor function tests**

## (2) Combine the data-driven ML approach with knowledge-driven mechanistic modelling

Describe individual SMA trajectories as ODEs in the latent space of a deep learning model



```
m = init_ODEVAE()
ps = getparams(m)
opt = ADAM(η)
trainingdata = zip(xs, xs_baseline, tvals)

for epoch in 1:epochs
    for (X, Y, t) in trainingdata
        grads = gradient(ps) do
            loss(X, Y, t, m, args=args)
        end
        update!(opt, ps, grads)
    end
end
```

# (3) Make use of correlations in the data: between input variables - to restrict the model space

BayesSUR: An R Package for High-Dimensional Multivariate Bayesian Variable and Covariance Selection in Linear Regression

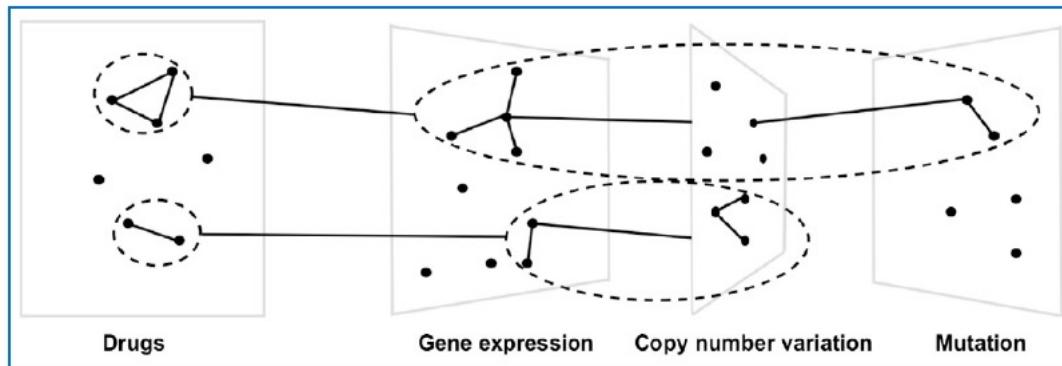
Zhi Zhao, Marco Banterle, Leonardo Bottolo, Sylvia Richardson, Alex Lewin, Manuela Zucknick

Vol. 100, Issue 11

 Paper

 R package (BayesSUR)

 R replication code



**-> See later**

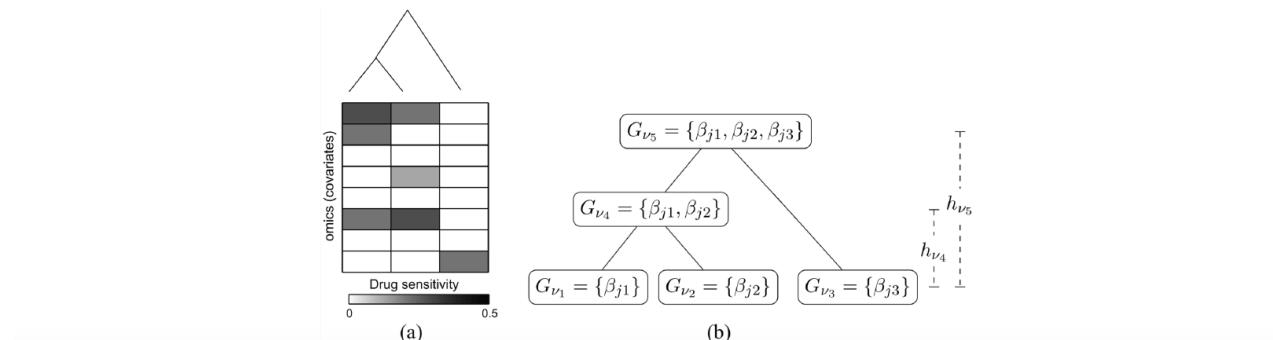
# (3) Make use of correlations in the data: between response variables - to borrow information

(Multi-response) Tree-guided group lasso (Kim & Xing 2012)

- Include dependencies between columns of  $\mathbf{Y}$  in a group lasso
- Extension to IPF-tree lasso

$$\text{Tree lasso: } \text{pen}(\mathbf{B}) = \lambda \sum_{j=1}^p \sum_{\nu \in \{V_{\text{int}}, V_{\text{leaf}}\}} \omega_\nu \|\beta_j^{G_\nu}\|_{\ell_2}$$

$$\text{IPF-tree lasso: } \text{pen}(\mathbf{B}) = \sum_s \lambda_s \left( \sum_{j_s} \sum_{\nu \in \{V_{\text{int}}, V_{\text{leaf}}\}} \omega_\nu \|\beta_{j_s}^{G_\nu}\|_{\ell_2} \right)$$



# (3) Make use of correlations in the data: between response variables - to borrow information

## Drug screens for precision cancer medicine: How to predict the drugs' effect with data on drugs and tumour?

ROYAL STATISTICAL SOCIETY  
DATA | EVIDENCE | DECISIONS

Journal of the Royal Statistical Society  
Applied Statistics  
Series C

Original Article | Open Access | CC BY SA

Structured penalized regression for drug sensitivity prediction

Zhi Zhao, Manuela Zucknick

