

Q&A SESSION 2



NEWS

① WELCOME FEM21045 PEOPLE!

② PROF. DR. DICK VAN DIJK
KAREL DE WIT

} Syllabus Update
(soon...)

③ NEW OPTIONAL MATERIAL FOR FEM3102
(readings, application discussions, videos, and so on.)

N 400 STUDENTS



KEEP
CALM
WE WILL
TAKE OVER EUR

OUTLINE

- OVERVIEW OF WEEK 2
- QUESTIONS FROM DISCUSSION FORUM 2
 - K-NN
 - K-FOLD CV
 - BIAS-VARIANCE
 - RIDGE-LASSO & CONSTRAINED OPTIMIZATION
 - INTEGER PROGRAMMING: CONVEX HULL
 - LEAST ANGLE REGRESSION
 - WHY REGULARIZATION?

OVERVIEW

Shrinkage Methods

$$Y \approx \beta_0 + \beta_1 X_1 + \underbrace{\beta_2 X_2}_{0} + \cdots + \underbrace{\beta_{p-1} X_{p-1}}_{0} + \beta_p X_p$$

Ridge Regression

$$\min_{\beta_0, \beta_1, \dots, \beta_p} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$$

Lasso

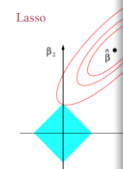
$$\min_{\beta_0, \beta_1, \dots, \beta_p} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 + \lambda \sum_{j=1}^p |\beta_j|$$

Ridge

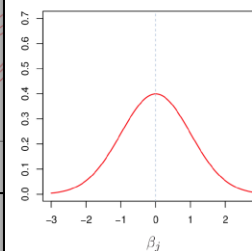
$$\min_{\beta} \{ (y - X\beta)^T (y - X\beta) : \|\beta\|_2^2 \leq \Delta \}$$

Lasso

$$\min_{\beta} \{ (y - X\beta)^T (y - X\beta) : \|\beta\|_1 \leq \Delta \}$$

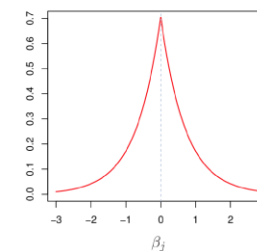


$$\beta_j \sim N(0, \psi^2)$$



$$\hat{\beta}_R = \arg \min_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$\beta_j \sim \text{Laplace}(0, \phi)$$



$$\hat{\beta}_L = \arg \min_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_1$$

Elastic Net

$$\hat{\beta}_{\text{EN}} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p (\alpha \beta_j^2 + (1 - \alpha) |\beta_j|)$$

Least Angle Regression

1. Scale the predictors: $\sum_{i=1}^n x_{ij} = 0$, $\sum_{i=1}^n x_{ij}^2 = 1$ for $j = 1, \dots, p$.
2. Set $\beta = 0$, $r = y - X\beta = y$.
3. Select the predictor that is most correlated with residual. Set r_j .
4. Move β_j towards its least squares coefficient of the current residual until some predictor r_k has the same correlation with the residual as r_j .
5. Move β_j and β_k in the direction defined by their joint least squares coefficient of the

Integer Programming Approach

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \|y - X\beta\|^2 \\ &\text{subject to} && \|\beta\|_0 \leq k. \end{aligned}$$

A simple reformulation as a mixed integer quadratic problem:

$$\begin{aligned} v_1 = & \text{minimize} && \frac{1}{2} \|y - X\beta\|^2 \\ & \text{subject to} && -Mz_j \leq \beta_j \leq Mz_j, \quad j = 1, \dots, p, \\ & && \sum_{j=1}^p z_j \leq k, \\ & && z_j \in \{0, 1\} \quad j = 1, \dots, p. \end{aligned}$$

convex hull of the feasible region

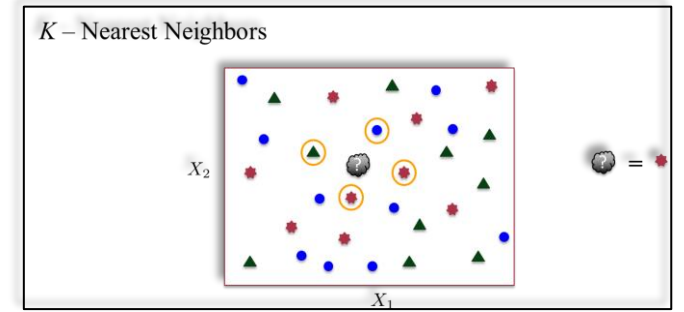
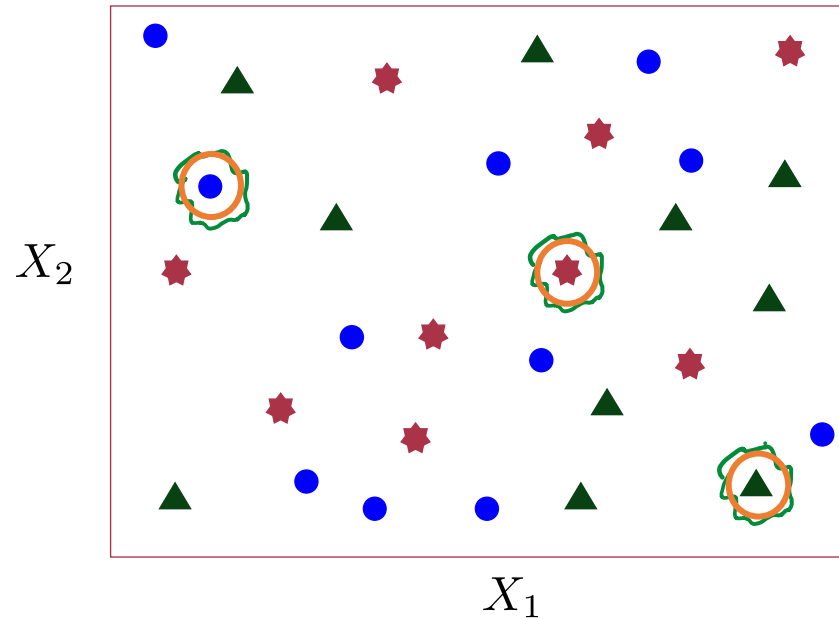
$$\{\beta : \|\beta\|_{\infty} \leq M, \|\beta\|_1 \leq Mk\} \subseteq \{\beta : \|\beta\|_1 \leq Mk\}$$

$$\begin{aligned} v_2 = & \text{minimize} && \frac{1}{2} \|y - X\beta\|^2 \\ & \text{subject to} && \|\beta\|_1 \leq Mk. \end{aligned}$$

$$v_2 \leq v_1$$

Lasso

K-NN



$K=1$
Training Error = 0 ?

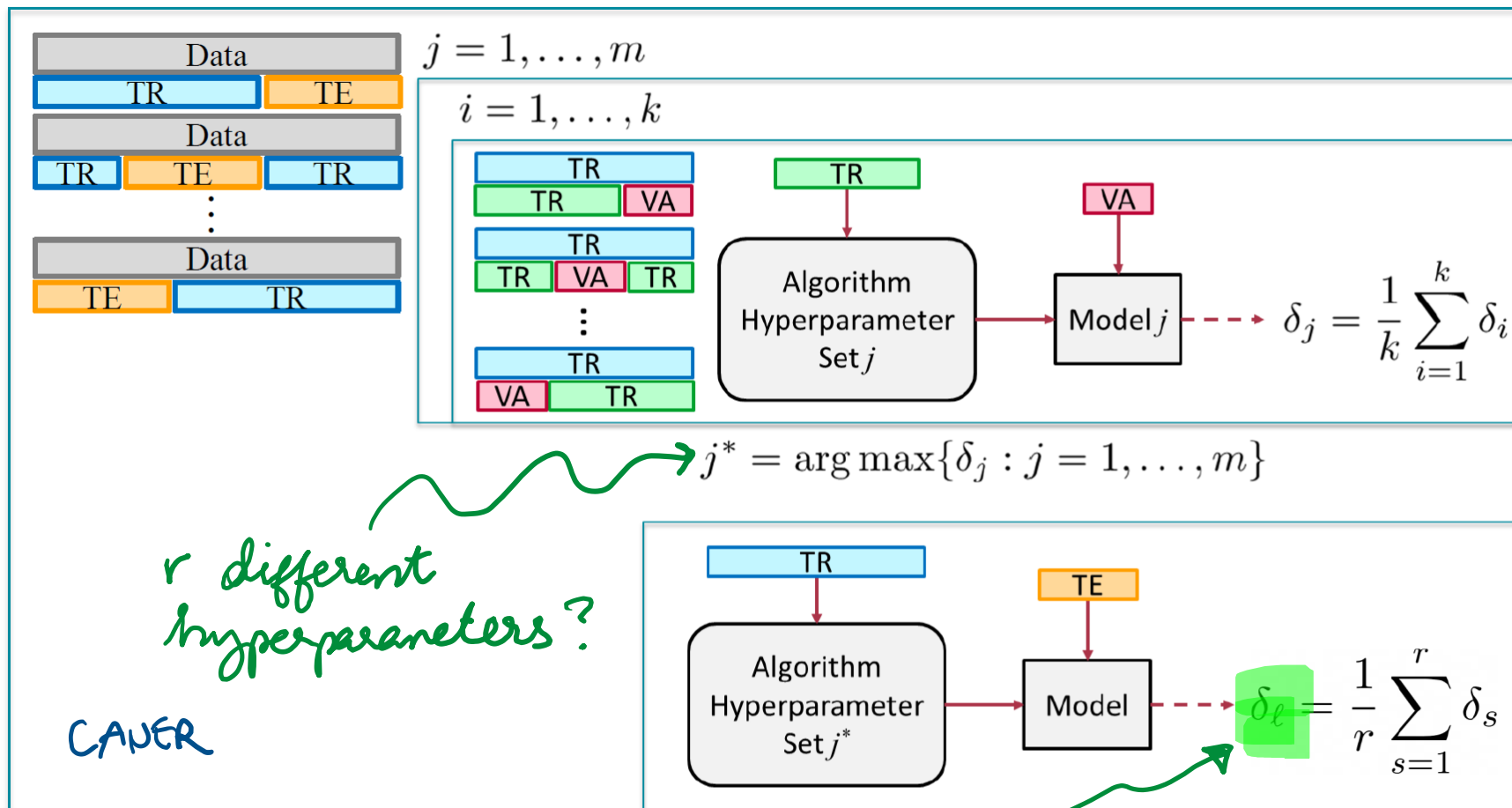
LIEVE
└ CANER
└ MARTEN

WAYNE
└ MARTEN

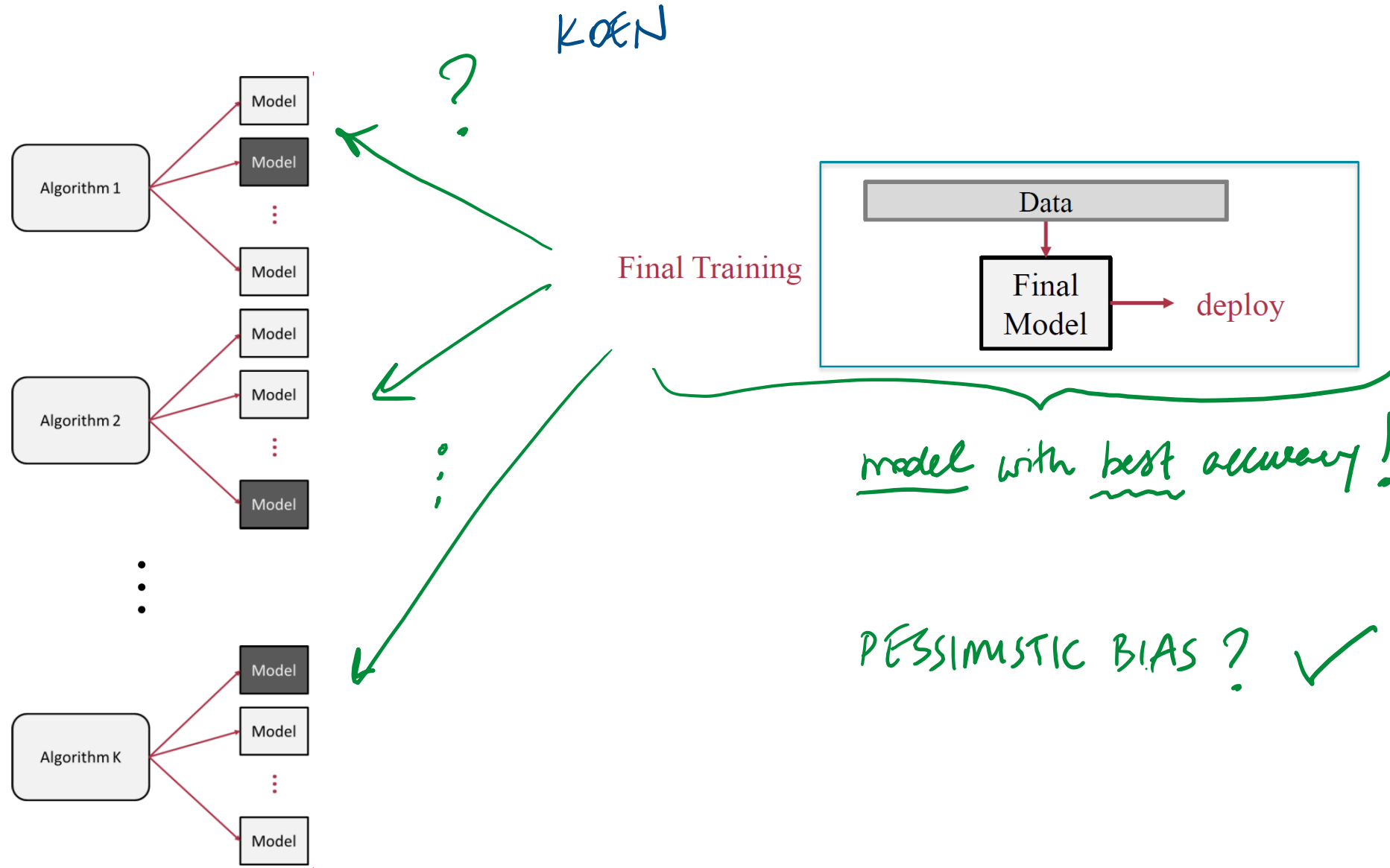
Ties ?
└ Random assignment
└ (next) add number for K

K-FOLD CV

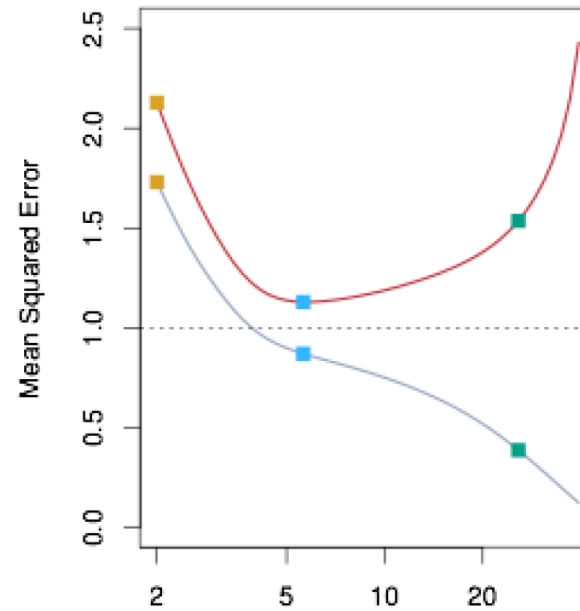
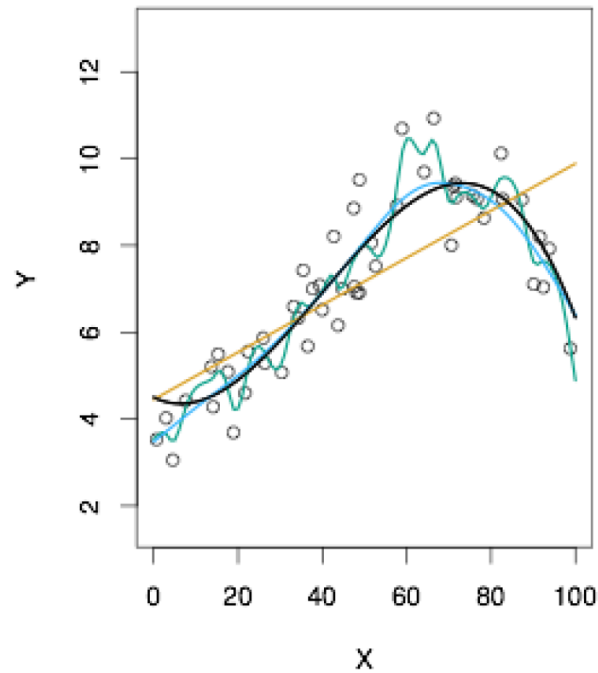
$s = 1, \dots, r$



K-FOLD CV



BIAS - VARIANCE

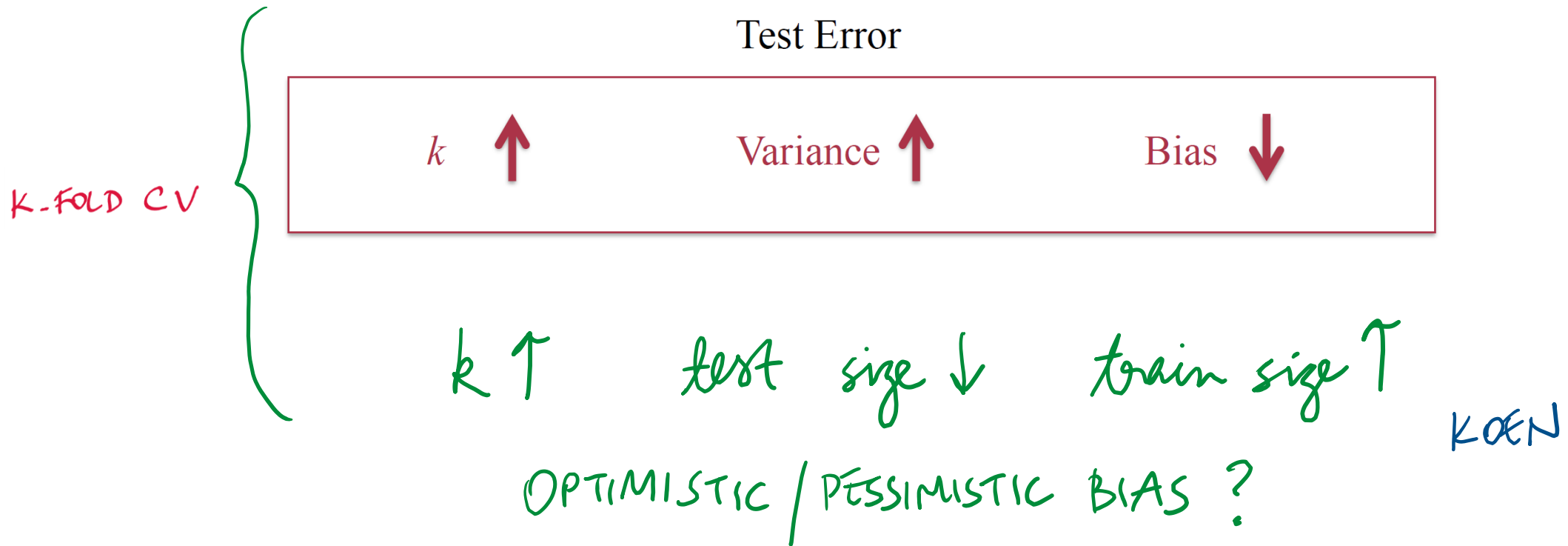


Flexibility ?
(continuous) parameter
of a smoothing spline

VIVI

\approx power of a polynomial
in curve fitting

BIAS - VARIANCE



Underestimates the test error rate \Rightarrow optimistic bias

LOCCV ($k=n$) gives almost the unbiased estimator of the test error.

RIDGE-LASSO & CONSTRAINED OPTIMIZATION

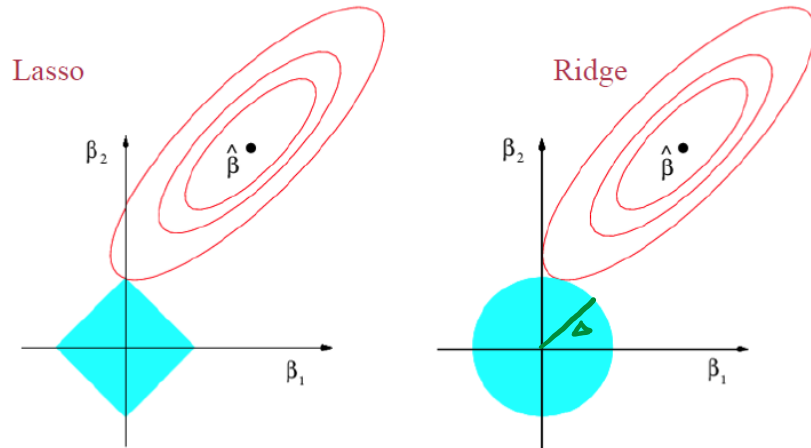
Ridge

(a) $\min_{\beta} \{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) : \|\beta\|_2^2 \leq \Delta \}$

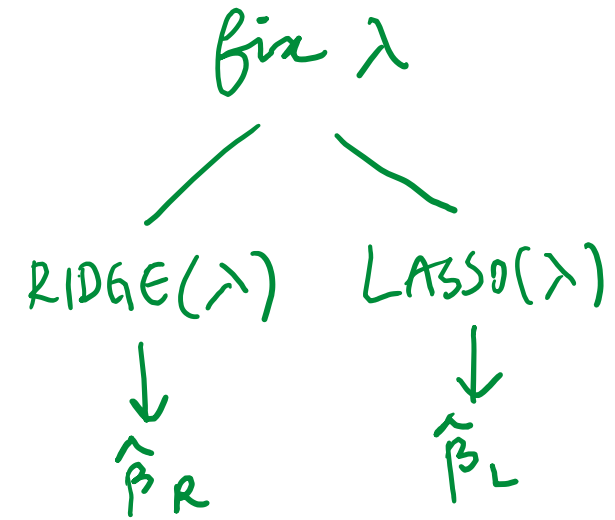
? LOURENS

Lasso

(b) $\min_{\beta} \{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) : \|\beta\|_1 \leq \Delta \}$



Least Squares Solution: $\hat{\beta} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}$



$\exists \Delta > 0$ such that

- solving (a) gives $\hat{\beta}_R$
- solving (b) gives $\hat{\beta}_L$

INTEGER PROGRAMMING: CONVEX HULL

$$\begin{aligned}
 v_1 = & \text{ minimize } \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\
 & \text{ subject to } \left. \begin{aligned} & -Mz_j \leq \beta_j \leq Mz_j, \quad j = 1, \dots, p, \\ & \sum_{j=1}^p z_j \leq k, \\ & z_j \in \{0, 1\} \quad j = 1, \dots, p. \end{aligned} \right\} \begin{aligned} & z_j = 1 \Rightarrow |\beta_j| \leq M \\ & \text{at most } k\text{-many} \\ & z_j \text{ can be } 1 \end{aligned}
 \end{aligned}$$

convex hull of the
feasible region

LOURENS ? $\{\boldsymbol{\beta} : \underbrace{\|\boldsymbol{\beta}\|_\infty \leq M}_{\text{max}\{|\beta_1|, \dots, |\beta_p|\} \leq M}, \underbrace{\|\boldsymbol{\beta}\|_1 \leq Mk}_{\sum_{j=1}^p |\beta_j| \leq Mk}\} \subseteq \{\boldsymbol{\beta} : \|\boldsymbol{\beta}\|_1 \leq Mk\}$

$$\max\{|\beta_1|, \dots, |\beta_p|\} \leq M$$

$$\sum_{j=1}^p |\beta_j| \leq Mk$$

LEAST ANGLE REGRESSION

CANER

β_1, β_2 ?

$\min \{n-1, p\}$ steps ?

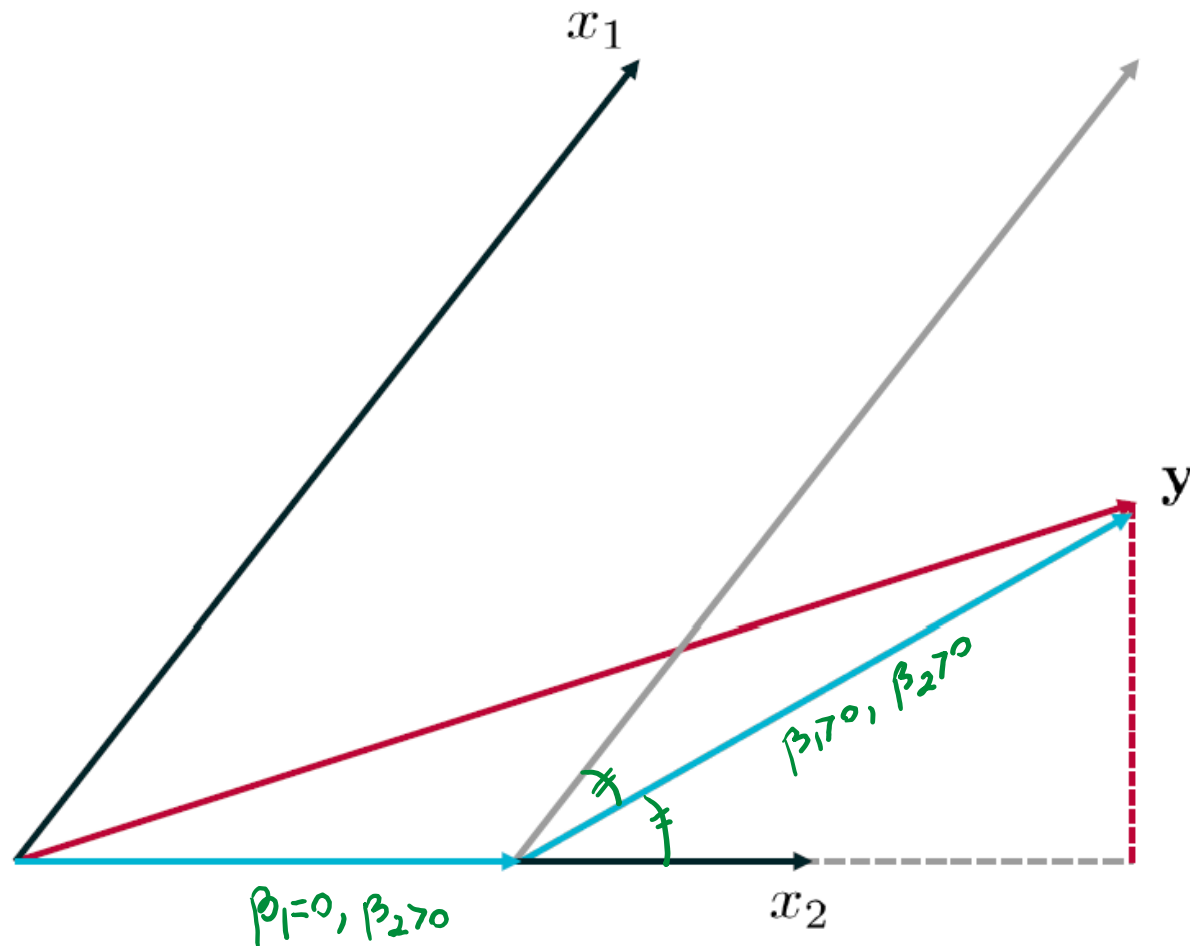
1. Scale the **predictors**: $\sum_{i=1}^n y_i = 0, \quad \sum_{i=1}^n x_{ij} = 0$

2. Set $\beta = 0, r = y - X\beta = y$



Columns are
centered

\Rightarrow row rank of $X = n-1$



WHY REGULARIZATION?

How about? { linear regression
+
significance tests
for parameters

MARTEN

OLS

→ Collinearity

→ $p > n$

Regularization:

- Dimension reduction (interpretability)
- Avoiding overfitting (e.g. NNs)

Forward Selection

→ Many tests

→ Problems with sequential testing

Sequential Selection Procedures and
False Discovery Rate Control

Max Grazioplene

Department of Statistics, Carnegie Mellon University, Pittsburgh, USA.

Stefan Wager

Department of Statistics, Stanford University, Stanford, USA.

Alexandra Chouldechova

Heinz College, Carnegie Mellon University, Pittsburgh, USA.

Robert Tibshirani

Departments of Health Research & Policy, and Statistics, Stanford University, Stanford, USA.

[link](#)

TILL NEXT WEEK !

* Assignment 2 (Due date: 14 Sep.)

* Videos for Lecture 3 (↓ 12 pm)