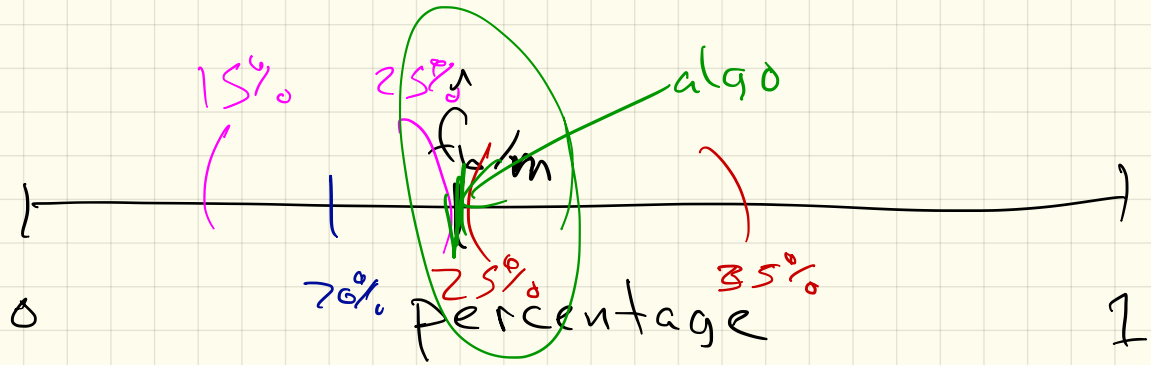


L13: Regression



might?  
must?

# Regression

Date  $(x, y)$

$x \in \mathbb{R}$  explanatorien

$y \in \mathbb{R}$  dependent

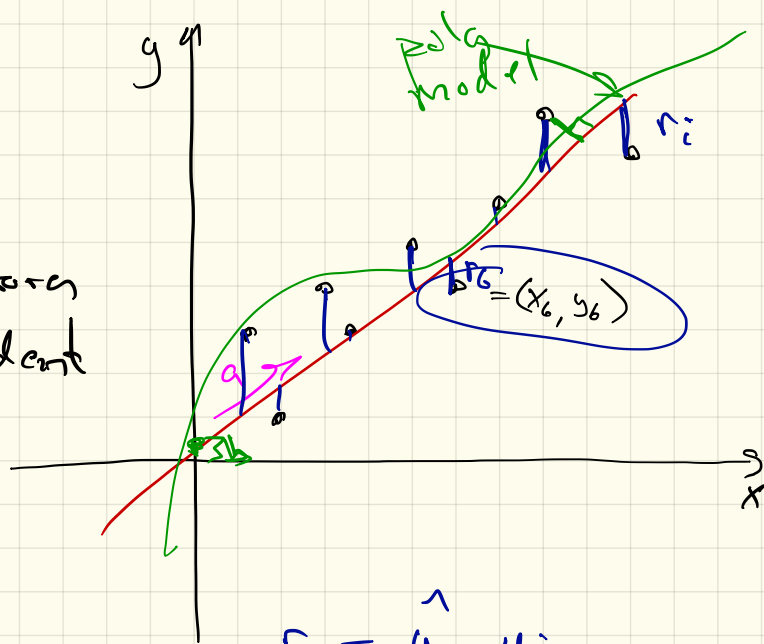
fit  $(x, y)$  w line

$$\hat{y}_i = f(x_i) = ax_i + b$$

$$0.8x_i + 0.2$$

Goal Find  $a, b$

$$\text{minimize}_{a, b} \sum_i (y_i - ax_i - b)^2$$



$$r_i = \hat{y}_i - y_i$$

$$\text{minimize} \sum_i r_i^2$$

$$\bar{x} = \frac{1}{n} \sum_i x_i \quad \bar{y} = \frac{1}{n} \sum_i y_i$$

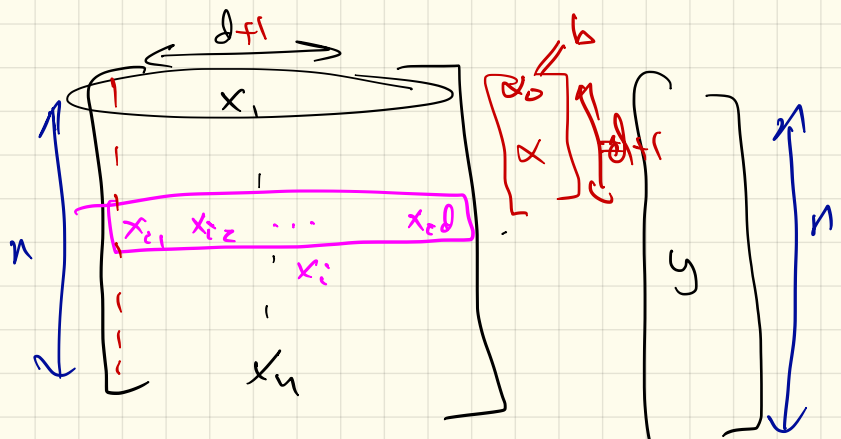
$$\text{Cov}(x, y) = \frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

$$a = \frac{\text{Cov}(x, y)}{\text{Cov}(x, x)} \quad b = \bar{y} - a\bar{x}$$

Data  $(X, y)$

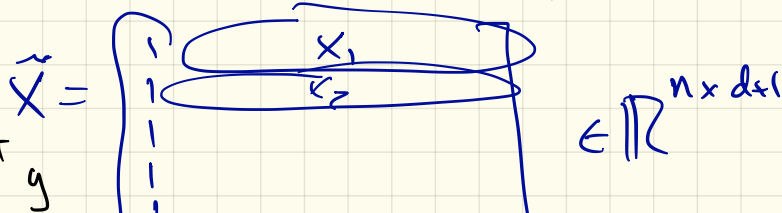
$$X \in \mathbb{R}^{n \times d}$$

$$y \in \mathbb{R}^n$$



Goal  $\alpha \in \mathbb{R}^{d+1}$   
 $\sum_i (y_i - \langle \alpha, \tilde{x}_i \rangle)^2$   
 $\downarrow$  solve

$$x_i = (x_{i1}, x_{i2}, \dots, x_{id}) \in \mathbb{R}^d$$



$$\alpha = (X^T X)^{-1} X^T y$$

$$\tilde{x}_i = [1, x_{i1}, x_{i2}, \dots, x_{id}]$$

$\uparrow$   
 $x_{i0}$

# Polynomial Regression

$$\text{Model } \hat{y}_i = M(x_i) = \alpha_0 x_i^0 + \alpha_1 x_i^1 + \alpha_2 x_i^2 + \alpha_3 x_i^3 \quad \leftarrow P$$

Input  $(x, y)$   $x_i \in \mathbb{R}$   
 $x \in \mathbb{R}$   
 $y \in \mathbb{R}$

$(+ \alpha_4 \sin(x_i))$

Goal  $\alpha \in \mathbb{R}^{P+1}$  minimize

$$\sum_i \left( y_i - \left( \sum_{j=0}^P \alpha_j x_i^j \right) \right)^2$$

map

$$X \in \mathbb{R}^{P+1}$$

$$\begin{matrix} \vec{x} \\ \downarrow \\ \uparrow \\ \vec{y} \end{matrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^P \\ 1 & x_2 & x_2^2 & \dots & x_2^P \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^P \end{bmatrix}$$

$\leftarrow P+1$

Soln

$$\alpha = (\tilde{X}_P^T \tilde{X}_P)^{-1} \tilde{X}_P^T y$$

# Gauss-Markov Thm

$\alpha = \underline{(X^T X)^{-1} X^T y}$  optimal

MLE Normal Noise

(1) Goal  $\sum (y_i - M(x_i))^2$   $\leftarrow$  good algo

(2) Assume residuals  $r_i = y_i - \hat{y}$  unbiased

(3) Want unbiased solution uncorrelated

Dim. Red

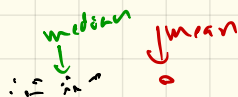


OLS: ok if units on  $x_i, y_i$  don't matter

$$\hat{y} = \alpha_0 + \alpha_1 x_i$$

## Robust Regression

mean vs. median



Theil-Sen Estimator  
 $a = \text{median all slopes}$   
 $b = \text{median } \{y_i - a x_i\}$

$$\left\{ \frac{(y_i - y_j)}{(x_i - x_j)} \right\}$$

# Ridge Regression (Tikhonov Regularization)

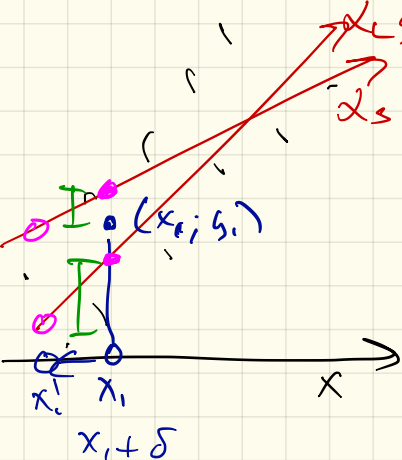
Goal  $\sum_i (y_i - \langle \alpha, x_i \rangle)^2 + s \|\alpha\|_2^2$

penalty on complex model  $\alpha$

$\exists s$  s.t. RR  $\alpha_s$  is better than  $\alpha_{LS}$  (unseen data)

$$\alpha_s = (X^T X + s^2 I)^{-1} X^T y$$

$d \times d$  identity matrix  $\begin{bmatrix} s^2 & & 0 \\ & \ddots & \\ 0 & & s^2 \end{bmatrix}$



smaller slope  
 $\|\alpha\|_2^2$  smaller

# Lasso (basis pursuit)

$$L_{1,s}(X, y, \alpha) = \sum_i (y_i - \langle x_i, \alpha \rangle)^2 + s \|\alpha\|_1$$

CONVEX  
not smooth

no simple linear algebra soln  $L_1$ -regularization

→ induce sparsity in  $\alpha$   
↳ many coefficients  $\alpha_i = 0$

Not perfect variable selection



only select variable  
iff  $\alpha_i \neq 0$