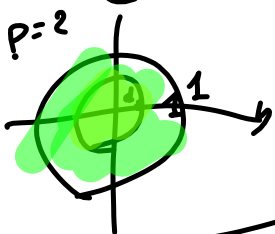


$$\Pr[\text{all } N \text{ points have distance} \geq d] = \frac{1}{2}$$

$t_i :=$  distance between  $x_i$  and origin

$$\Pr[t_i \geq d] = 1 - \Pr[t_i < d]$$



$$= 1 - \frac{\text{Volume of ball radius } d}{\text{Volume of ball radius } 1}$$

Volume of a ball with radius  $r = \frac{\pi^{p/2}}{(p/2)!} r^p$

$p=2 \quad \frac{\pi r^2}{(2/2)!}$

$$= 1 - \frac{\frac{\pi^{p/2}}{(p/2)!} d^p}{\frac{\pi^{p/2}}{(p/2)!} 1^p} = 1 - d^p$$

$$\Pr[\underbrace{\forall i}_{\text{"1/2"}}, t_i \geq d] = \prod_{i=1}^N \Pr[t_i \geq d] = (1 - d^p)^N$$

$$(1 - d^p)^N = \frac{1}{2}$$

$$1 - d^p = \left(\frac{1}{2}\right)^{1/N}$$

$$d^p = 1 - \left(\frac{1}{2}\right)^{1/N}$$

$$d = \left(1 - \left(\frac{1}{2}\right)^{1/N}\right)^{1/p}$$

qed

find  $\hat{f}(x)$  as a useful approximation of  $f(x)$   
 last week:  $L(y, f(x)) = (y - f(x))^2$  squared loss function  
 leads to  $\underline{f(x) = E[Y|X=x]}$

K-nearest neighbour

$$\underbrace{E[Y|X=x]}_{\text{average}} \rightarrow \text{neighbour}$$

issues

- curse of dimensionality
  - e.g. nearest point gets further to the point of interest as increasing
- balance between bias - variance

## Statistics vs Machine Learning

Statistical approach:

Starting from the model

$$Y = f(x) + \varepsilon, \quad E[\varepsilon] = 0, \quad \varepsilon \perp\!\!\!\perp X$$

additive model  $\rightarrow$  approximation of the truth  
 statistician  $\hat{f}(x)$  approx  $f(x)$

- we do not suppose  $Y = f(x)$  (deterministic)

BUT

we add an error term which captures:

- measurements errors;
- effects of non-measured variables;
- ...

often  $\varepsilon \sim \text{iid}$ ,  $\varepsilon \sim \underline{N}(0, \sigma^2)$

most natural approach, least square

## Machine learning approach

Assume  $Y = f(x) + \varepsilon$ → start from  $\hat{f}(x)$ , possibly simple  $\hat{f}(x) = c$  (initialization)→ evaluate  $\hat{f}(x)$  on our training set

e.g.  $\sum_{i=1}^N (y_i - \hat{f}(x_i))^{[k]}^2$

→ modify  $\hat{f}(x)$  to improve the prediction

$$\sum_{i=1}^N (y_i - \hat{f}(x_i))^{[k+1]}^2 < \sum_{i=1}^N (y_i - \hat{f}(x_i))^{[k]}^2$$

→ use training set  $(x_i, y_i)$ 

→ learning by examples

K is the step of the algorithm

→ focus is on the learner, that does not aim to approximate a true  $f(x)$ 

statistical/mathematical approach

 $(x_i, y_i)$  points of a  $p+1$  dimensional space

$$f(x): \underset{\substack{\text{in} \\ \mathbb{R}^p}}{X} \rightarrow \underset{\substack{\text{in} \\ \mathbb{R}}}{y}$$

goal: giving an approximation of  $f(x)$  working at points in  $X$  given  $T$

Parametric approach

$$\mathcal{F} = \{ p(y, x; \underline{\theta}), \underline{\theta} \in \Theta \subseteq \mathbb{R}^P \}$$

$$Y = f(x; \underline{\theta}) + \varepsilon$$

$$f_{\theta}(x) = f(x; \theta)$$

e.g.:

- linear model:  $f_{\theta}(x) = X^T \beta$   $\theta = \beta$

- linear basis expansion:  $f_{\theta}(x) = \sum_{k=1}^K h_k(x) \theta_k$

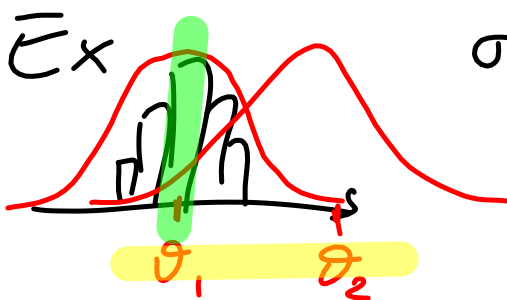
$$h_1(x) = x_1$$

$$h_2(x) = x_1^2$$

$$RSS(\theta) = \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

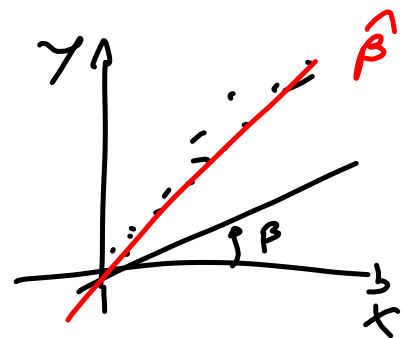
= function of  $\theta$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} RSS(\theta)$$



$$\sigma^2 = 1$$

$$\hat{\theta} = \theta_1$$



simplest cases

- least squares

more complicated frameworks

- maximum likelihood estimation

Likelihood estimation

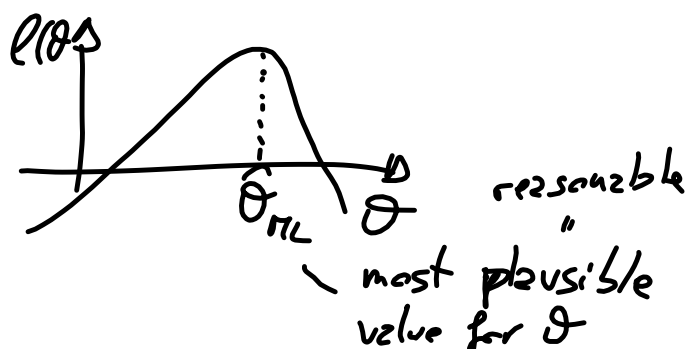
$Y_i \sim P$   $p(y)$  is indexed by  $\theta$   
 $Y_i$  iid.  $p(y; \theta)$

e.g. Gaussian distribution  $p(y; \underbrace{\mu, \sigma^2}_{\theta})$

$$L(\theta) = \prod_{i=1}^N p(y_i; \theta)$$

$$\ell(\theta) = \sum_{i=1}^N \log p(y_i; \theta)$$

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$



Note: when  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$   
 $\hat{\theta}_{ML} = \hat{\theta}_{LS}$

Restricted estimators

instead of estimating  $\theta$  as

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \text{RSS}(\theta)$$

we look for

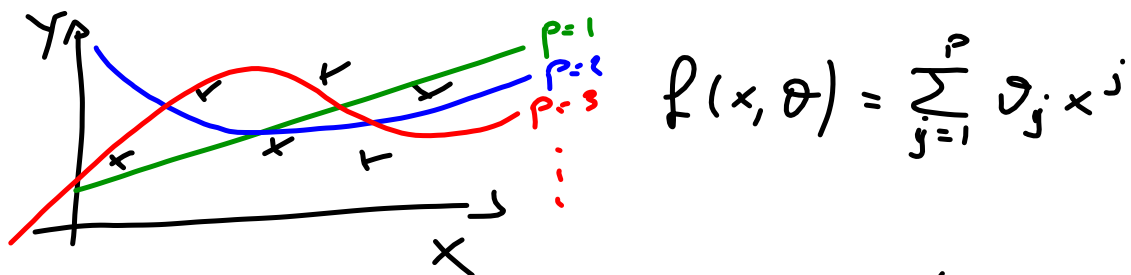
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \text{PRSS}(\theta)$$

where:  $\text{PRSS}(\theta) = \text{RSS}(\theta) + \lambda J(\theta)$

LASSO  $J(\theta) = \sum_{j=1}^p |\theta_j|$

RIDGE  $J(\theta) = \sum_{j=1}^p \theta_j^2$

Model selection and variance-bias trade-off



given enough parameters  $\theta_j$ , we will always be able to find a line which passes through all points (no bias)

$$\hat{\theta} = \arg \min_{\theta} \left\{ \text{RSS}(\theta) + \lambda J(\theta) \right\}$$

↗ penalization for the number of parameters

Bias-variance trade off

$$\text{error} = \sigma^2 + \underbrace{\text{variance} + \text{bias}^2}_{\text{MSE}}$$

$$y = f(x_0)$$

$$Y = f(x) + \epsilon$$

see figure 2.11

## Linear methods for regression

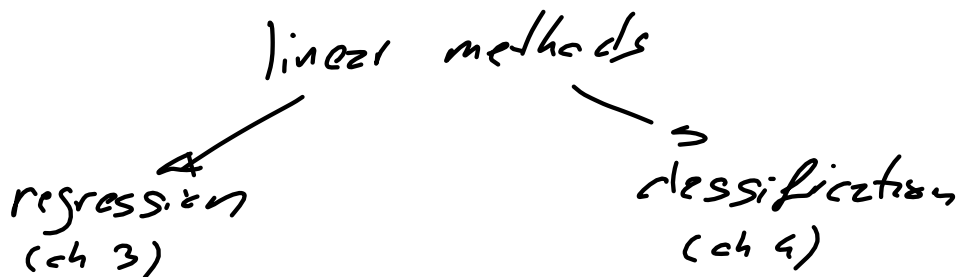
$E[Y|X]$  is linear in inputs

$$E[Y|X] = X^T \beta$$

- simple
- often adequate
- easy to interpret
- often outperforms fancier methods

↓  
simplicity!

- sample size is small
- sparse data
- low signal-to-noise ratio



## Regression

- consider continuous  $Y$
- linear regression  $f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \operatorname{RSS}(\beta)$$

$$\sum_{i=1}^n |y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j|$$

matrix form  $(y - X\beta)^T (y - X\beta)$

$$\frac{\partial \operatorname{RSS}(\beta)}{\partial \beta} = -2X^T (y - X\beta)$$

$$\frac{\partial^2 \operatorname{RSS}(\beta)}{\partial \beta^2} = 2X^T X$$

$$\frac{\partial RSS(\beta)}{\partial \beta} = 0 \quad \begin{aligned} X^T(y - X\beta) &= 0 \\ X^T y - X^T X \beta &= 0 \\ \hat{\beta} &= (X^T X)^{-1} X^T y \end{aligned}$$

$$\frac{\partial^2 RSS(\beta)}{\partial \beta^2} = 2 X^T X > 0 \Rightarrow \hat{\beta} \text{ minimum}$$

$$\hat{y} = X \hat{\beta} = \underbrace{X (X^T X)^{-1} X^T}_H y$$

hat matrix  
projection matrix

### Properties

$$\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2$$

$$\hat{\sigma}^2 = \frac{1}{N-p-1} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Focus on  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$   $\varepsilon$ : iid

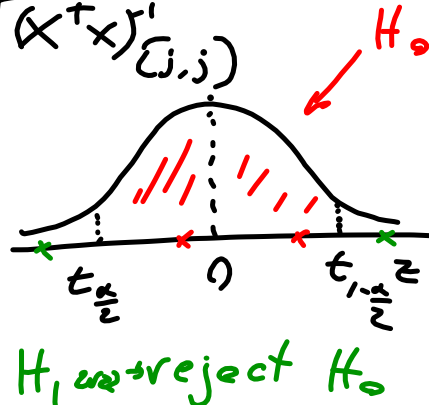
$$\hat{\beta} \sim \mathcal{N}(\beta, (X^T X)^{-1} \sigma^2)$$

$$(N-p-1) \hat{\sigma}^2 \sim \sigma^2 \chi_{N-p-1}^2$$



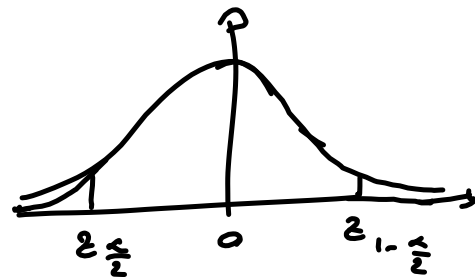
$$H_0: \beta_j = 0 \rightarrow z_j = \frac{\hat{\beta}_j - 0}{\text{sd}(\hat{\beta}_j)} = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{(X^T X)^{-1}_{(j,j)}}}$$

Under  $H_0$   $z_j \stackrel{H_0}{\sim} t_{N-p-1}$   
 $z_{obs}$



when  $\sigma^2$  is known

under  $H_0$   $z_j \stackrel{H_0}{\sim} \mathcal{N}(0, 1)$



---


$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$$

is  $x_1$  useful to predict  $y$ ?  
 $\beta_1 = 0 \rightarrow \begin{cases} \text{reject } H_0 & \text{Yes} \\ \text{do NOT reject } H_0 & \text{not } H_0 \end{cases}$

Are  $(x_2, x_3)$  useful to predict  $y$ ?

$H_0: \beta_2, \beta_3 = 0$

$$F = \frac{(RSS_0 - RSS_1) / (p_1 - p_0)}{RSS_1 / (N - p_1 - 1)}$$