

## Regression.

Predict response  $Y$  using covariate  $X$ .

$$(X_i, Y_i) \stackrel{\text{iid}}{\sim} P.$$

Want to estimate conditional expectation

$$f^*(x) = E[Y | X=x].$$

Method: use a class  $\mathcal{F}$  of functions.

search  $f^*$  (or its approximation)

with  $\mathcal{F}$ .

( $f^*$  does not have to be in  $\mathcal{F}$ )

(We do not need to assume that)

$$Y_i = f^*(X_i) + \text{noise}$$

(no need to model the entire data-generating mechanism)

— Linear regression.

Feature vector  $X \mapsto \phi(X) \in \mathbb{R}^d$ .

$$\mathcal{F} = \left\{ X \mapsto \beta^T \phi(X) : \beta \in \mathbb{R}^d \right\}.$$

— Sparse linear regression.

where we consider a subset of  $\mathcal{F}$  in LR, with  $\beta$  having only a small subset of non-zero entries.

— Nonparametric smooth function classes.

$$\text{e.g. } \mathcal{F} = \left\{ f : [0,1] \rightarrow [0,1] \mid \begin{array}{l} |f(x) - f(y)| \leq |x-y| \\ \forall x, y \in [0,1] \end{array} \right\}$$

$$\text{e.g. } \mathcal{F} = \left\{ f : [0,1] \rightarrow [0,1] \mid |f^{(k)}(x)| \leq L, \quad \forall x \in [0,1] \right\}$$

— Neural nets.

# Linear regression.

$$\hat{\beta}_n := \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - \beta^T X_i)^2 \right\}$$

(for simplicity, we just let  $\phi(x) = x$ ).

In many cases, we may want to consider  
the model class  $x \mapsto \beta_0 + \beta^T x$ .

In such a case, we simply let  $\phi(x) = [x]$

"Least-square estimator"

(In general, we fit

$$\hat{f}_n = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2$$

Fact. Suppose that  $Y_i = f^*(X_i) + \varepsilon_i$

$$\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2) \quad (\sigma > 0)$$

then the least square estimator is MLE.

(and therefore consistent, asymptotically normal,  
asymptotically optimal, under certain assumptions)

Proof of the fact.

$$\text{log-likelihood} = \sum_{i=1}^n \left( -\frac{(Y_i - f(x_i))^2}{2\sigma^2} - \log(2\pi\sigma^2) \right)$$

note: this holds true for any regression problems.

For linear regression, let's start with models  
under strong assumptions.

$$Y_i = \beta^T X_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} P.$$

$$E[\varepsilon_i] = 0, \quad \text{var}(\varepsilon_i) = \sigma^2.$$

Fact.  $E[\hat{\beta}_n | X_1, \dots, X_n] = \beta^*$  ← well-specified  
linear model

$$\text{cov}(\hat{\beta}_n | X_1, \dots, X_n) = \sigma^2 \cdot (X^T X)^{-1}$$

↑ requires iid noises w/ finite var.

Where  $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{d \times n}$ .

Proof.

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n (y_i - \beta^T x_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n y_i^2 - \frac{2}{n} \sum_{i=1}^n y_i x_i^T \beta \\ & \quad + \beta^T \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) \beta. \end{aligned}$$

Solving the minimization problem

$$\hat{\beta}_n = \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right)^{-1} \cdot \left( \frac{1}{n} \sum_{i=1}^n y_i x_i \right)$$

(this holds true regardless of probabilistic assumptions)

Now using prob assumptions.

$$y_i = \beta_*^T x_i + \varepsilon_i$$

$$\sum_{i=1}^n y_i x_i = \sum_{i=1}^n x_i x_i^T \beta_* + \sum_{i=1}^n \varepsilon_i x_i$$

$$\hat{\beta}_n = \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \beta_* + \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right)$$

$$\mathbb{E}[\hat{\beta}_n | X_1, X_2, \dots, X_n] = \beta_*$$

(for this step, we only use

the fact that  $\mathbb{E}[Y_i | X_i = x] = \beta_* x$   
is actually linear, i.e.  $\mathbb{E}[\varepsilon_i | X_i = x] = 0$  ).

For the conditional variance.

$$\begin{aligned}\text{cov}(\hat{\beta}_n | X_1, \dots, X_n) &= \text{cov}\left(\left(\sum_{i=1}^n X_i X_i^\top\right)^{-1} \cdot \left(\sum_{i=1}^n \varepsilon_i X_i\right)\right) \\ &= \sum_{j=1}^n \text{cov}\left(\left(\sum_{i=1}^n X_i X_i^\top\right)^{-1} X_j \varepsilon_j\right) \\ &= \sigma^2 \cdot \sum_{j=1}^n (X X^\top)^{-1} X_j X_j^\top (X X^\top)^{-1} \\ &= \sigma^2 \cdot (X X^\top)^{-1}.\end{aligned}$$

(this step requires  $\varepsilon_i$ 's to be iid).

Asymptotics.

Fact. Under above assumptions.

summing that  $E[XX^T] \succ 0$ .

then we have

$$\sqrt{n}(\hat{\beta}_n - \beta^*) \xrightarrow{d} N(0, \sigma^2 \cdot (E[XX^T])^{-1}).$$

Proof:  $\hat{\beta}_n - \beta^* = \left( \frac{1}{n} \sum_{i=1}^n X_i X_i^T \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i X_i \right)$

(by LLN)  $\xrightarrow{P}$   $\xrightarrow{d}$   $N(0, \sigma^2 E[XX^T])$ .

(after  $\sqrt{n}$  scaling).

by Slutsky thm, we have

the convergence.

---

We can perform inference based on asymptotics.

• Straightforward approach.

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$$

$$\hat{\sigma}^2 = \frac{1}{n-d} \sum_{i=1}^n (Y_i - \hat{\beta}_n^T X_i)^2.$$

This  $\frac{1}{n-d}$  normalization makes  $\hat{\sigma}^2$   
 a conditionally unbiased estimator for  $\sigma^2$   
 (conditioned on  $X_1, X_2, \dots, X_n$ )  
 (when  $n \gg d$ , this does not make  
 much difference)

Putting them together

$$\sqrt{n} \cdot \frac{\hat{\Sigma}_n^{-\frac{1}{2}}}{\hat{\sigma}_n} (\hat{\beta}_n - \beta^*) \xrightarrow{d} N(0, I_d).$$

We can also use bootstrap etc.

Removing assumptions.

$\epsilon_i \overset{iid}{\sim}$

$$E(Y|X=x) = \beta^T x.$$

$$\text{var}(Y|X=x) = \sigma^2(x).$$

$$\text{Fact. } \sqrt{n}(\hat{\beta}_n - \beta^*) \xrightarrow{d} N(0, \Sigma_X^{-1} \Sigma^* \Sigma_X^{-1})$$

where  $\Sigma_X = \mathbb{E}[XX^T]$ .

$$\Sigma^* = \mathbb{E}[\sigma(X)^2 \cdot XX^T].$$

Proof.  $\sqrt{n}(\hat{\beta}_n - \beta^*) = \left( \frac{1}{n} \sum_i^n X_i X_i^T \right) \left( \frac{\sqrt{n} \sum_{i=1}^n \epsilon_i X_i}{\sqrt{n}} \right)$

$\downarrow P$   
 $\Sigma_X$   
 $\downarrow d$   
 $N(0, \Sigma^*)$ .

We can also perform inference

by estimating  $\Sigma_X$  and  $\Sigma^*$ .

$$\hat{\Sigma}_X = \frac{1}{n} \sum_i^n X_i X_i^T$$

$$\hat{\Sigma}^* = \frac{1}{n} \sum_i^n (Y_i - \hat{\beta}_n^T X_i)^2 X_i X_i^T.$$

Remark: - this method is more robust

, we can always use bootstrap.

Removing more assumptions.

~~$$\mathbb{E}[Y|X=x] = \beta_0 + \beta_1 x$$~~

Fact. Under only moment assumptions on  $X, Y$

We have that

$$\hat{\beta}_n \rightarrow \bar{\beta} = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}[(Y - X^\top \beta)^2].$$

$$(\underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}[|f(x) - \beta^\top x|^2])$$

(Best linear approximation of  
the conditional expectation function.

$$f^*(x) = \mathbb{E}[Y | X=x]$$

(we can also study  $\sqrt{n}(\hat{\beta}_n - \bar{\beta}) \rightarrow \text{stn}$ ).

Proof:  $\hat{\beta}_n = \underbrace{\left( \frac{1}{n} \sum_1^n X_i X_i^\top \right)^{-1}}_{\mathbb{E}[X X^\top]} \cdot \underbrace{\left( \frac{1}{n} \sum_1^n X_i Y_i \right)}_{\mathbb{E}[X Y]}$

$$\bar{\beta} = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \beta^\top \mathbb{E}[X X^\top] \beta - 2\beta^\top \mathbb{E}[X Y] \right\}$$

$$= \mathbb{E}[X X^\top]^{-1} \cdot \mathbb{E}[X Y].$$

Other extensions / procedures.

• Prediction interval.

— Learn a model using data  $(X_i, Y_i)_{i=1}^n$

— Given a new data  $x_0$

the goal is to construct an interval

$I$  (based on data)

s.t.  $P(Y \in I | X = x_0) \geq 1 - \alpha$ .

(both  $Y$  and  $I$  are random)

Naive idea: use CI for  $\beta_*^T x_0$   
(using  $\hat{\beta}_n^T x_i$ )

Problem:  $Y$  is also an r.v.

$Y$  is indep of  $I$ .

Assume  $E_i \stackrel{iid}{\sim} N(0, \sigma^2)$ .

then we can construct prediction interval.

$$\widehat{\beta}_n^T x_0 \pm \sqrt{\hat{\sigma}^2 + \frac{\hat{\sigma}^2 \cdot x_0^T \widehat{\Sigma}_x^{-1} x_0}{n}} \quad \text{Zach}$$

$$\sqrt{n} \left( \widehat{\beta}_n^T x_0 - \beta^T x_0 \right) \xrightarrow{d} N(0, \sigma^2 x_0^T \widehat{\Sigma}_x^{-1} x_0)$$

$(\widehat{\Sigma}_x = E[X X^T])$       ✓

$$Y - \beta^T x_0 \sim N(0, \sigma^2). \quad \xrightarrow{\text{independent}}$$

Taking difference

$$Y - \widehat{\beta}_n^T x_0 \sim N(0, \sigma^2 + \frac{\sigma^2 x_0^T \widehat{\Sigma}_x^{-1} x_0}{n})$$

---

Model selection in linear regression.

• Training error.

$$\widehat{R}_{tr} = \sum_{i=1}^n (Y_i - \widehat{f}_n(x_i))^2.$$

• Prediction risk (in sample)

$$R = \sum_1^n (Y_i^* - \widehat{f}_n(x_i))^2.$$

where  $\hat{Y}_i^*$  is a fresh indep sample  
from  $Y_i | X_i$ .

Fact.

$$E[\hat{R}_{tr} | X_1, \dots, X_n] \leq R$$

and  $R - E[\hat{R}_{tr} | X_1, \dots, X_n]$

$$= 2 + \sum_{i=1}^n \text{cov}(\hat{Y}_i, Y_i | X_1, \dots, X_n)$$

where  $\hat{Y}_i = \hat{f}_n(X_i)$ .  
(“overfitting” in ML).

, Mallow's Cp statistics.

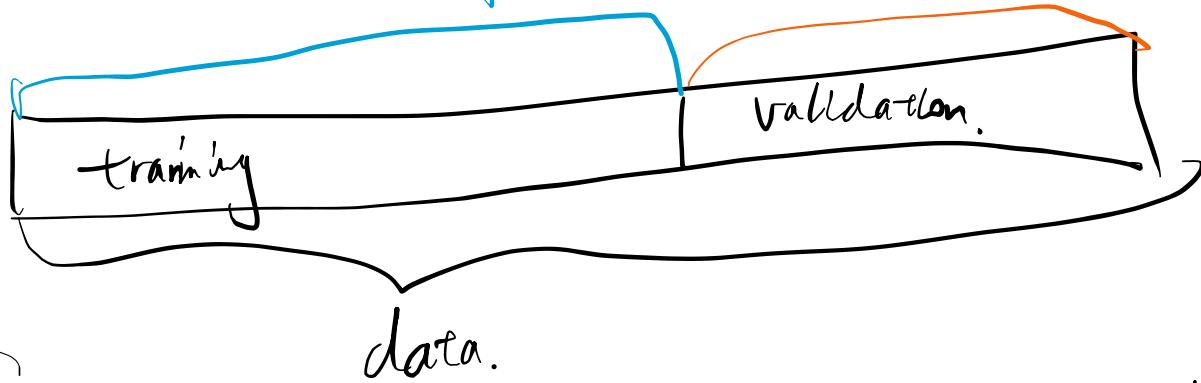
Consider linear regression using

subset  $S$  of the covariate variables.

$$\hat{R}(S) = \hat{R}_{tr}(S) + 2|S| \cdot \hat{\sigma}^2.$$

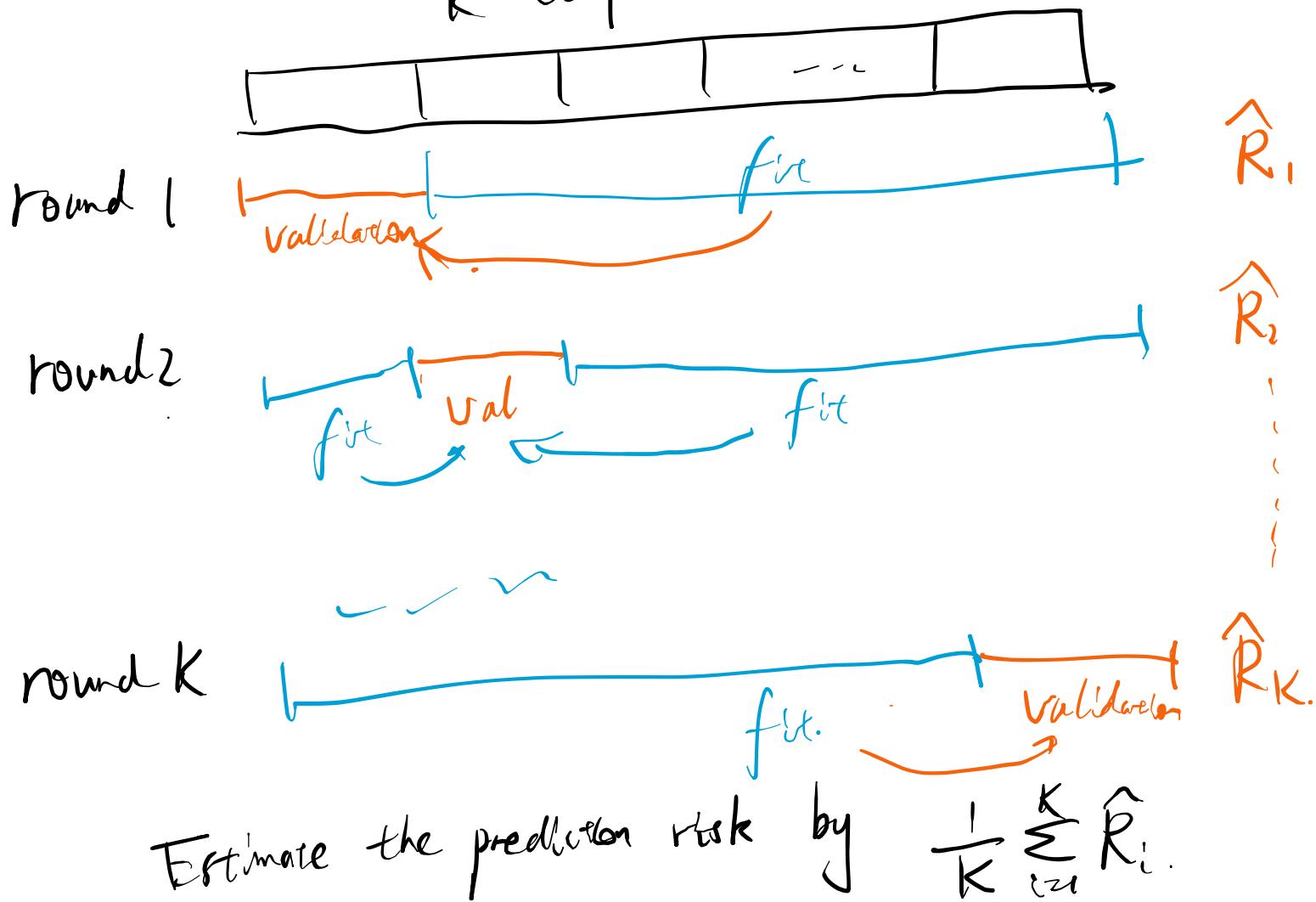
gives a better estimator for risk.

- cross Validation  $\hat{f}_n$   $\rightarrow$  estimate loss.



unbiased estimation for prediction risk.

- K-fold CV.  
K components.



- Leave-one-out CV  
( $K$ -fold with  $K=n$ ).

Logistic regression.

$$Y_i \in \{0, 1\}. \quad Y_i | X_i = x \sim \text{Ber}(p^*(x))$$

(Issue with linear model)

$$p^*(x) = x^T \beta^*$$

$x^T \beta^*$  may not always lie in  $[0, 1]$ .

Natural idea: put a link function  $\sigma(x^T \beta^*)$

Choice of  $\sigma$ :

- Smooth (for computation and convenience)

- Increasing
- mapping to  $[0, 1]$

Practical choice

$$\sigma(x) = \frac{e^x}{1+e^x} \quad (\text{sigmoid function})$$

Assuming  $Y_i | X_i \sim \text{Ber}(\sigma(X_i^\top \beta))$ .

$$\text{log likelihood} = \log \left( \prod_{i=1}^n p_i^{Y_i} (1-p_i)^{1-Y_i} \right)$$

$$= \sum_{i=1}^n (Y_i \log p_i + (1-Y_i) \log (1-p_i))$$

$$= \sum_1^n Y_i \log \sigma(X_i^\top \beta) + (1-Y_i) \log (1-\sigma(X_i^\top \beta))$$

$$= \sum_1^n Y_i \cdot (X_i^\top \beta - \log(1+e^{X_i^\top \beta}))$$

$$- (1-Y_i) \cdot \log(1+e^{X_i^\top \beta})$$

$$= \sum_{i=1}^n Y_i X_i^\top \beta - \log(1+e^{X_i^\top \beta})$$