

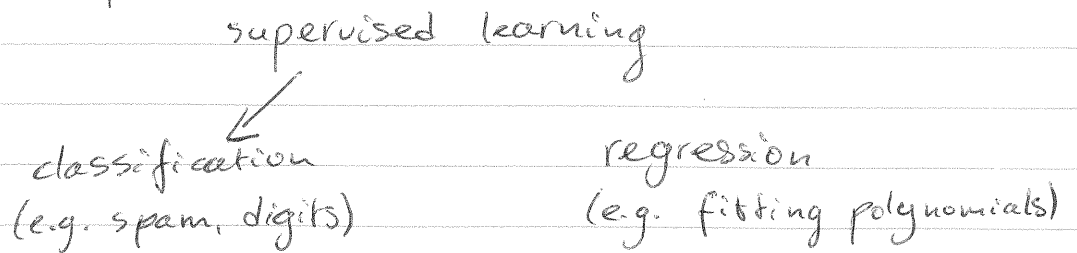
Statistical Learning II (5-11-2015)

1. Organization
2. Warm up:
 - a) statistical decision theory
 - b) ERM
3. Regression I: least squares
 - a) interpretations:
 - ERM, ML, projection
 - b) computation
 - c) bias-variance decomposition
4. k-Fold Cross-validation

1. Organization

- Anyone more time on exam?
- Next week: start 11:15, work on homework 1 after lecture
- Lecture notes on website (not Blackboard)

2. Warm up



$$\mathcal{T} = \left(\begin{matrix} y_1 \\ x_1 \end{matrix} \right), \dots, \left(\begin{matrix} y_N \\ x_N \end{matrix} \right)$$

↳ $\hat{f} \in \mathcal{F}$ and predict $\hat{y} = \hat{f}(x)$ for new x

2

2a Statistical Decision Theory

$EPE(f) = \mathbb{E}_{x,y} [L(Y, f(x))]$ measures quality of f

regression: $L(Y, f(x)) = (Y - f(x))^2$ "squared error"

classification:

$L(Y, f(x)) = \begin{cases} 0 & \text{if } f(x) = Y \\ 1 & \text{if } f(x) \neq Y \end{cases}$ "0/1-loss"

Bayes optimal predictor:

$$f_B = \operatorname{argmin}_f EPE(f)$$

regression: $f_B(x) = \mathbb{E}[Y|X]$

classification: $f_B(x) = \operatorname{argmax}_y \Pr(Y=y|X)$

Estimators:

\hat{f} depends on $T \Rightarrow EPE(\hat{f})$ depends on T
 \Rightarrow evaluate $\mathbb{E}[EPE(\hat{f})]$

2b) Empirical Risk Minimization

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(Y_i, f(X_i))$$

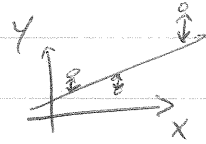
classification: - minimize #mistakes
- usually cannot compute efficiently for 0/1-loss.

regression: - minimize sum of squared errors
- can compute if \mathcal{F} is linear model

3. Regression I: least squares

Linear model: $\hat{y} = x^T \hat{\beta} = \hat{f}(x)$

$$RSS(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2$$



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

a) Interpretations

I. ERM for squared error

$$\mathcal{F} = \{f_{\beta}(x) = x^T \beta \mid \beta \in \mathbb{R}^p\}$$

$$\hat{f} = \hat{f}_{\hat{\beta}} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

II. Maximum Likelihood for Gaussian Errors

Model: $y = x^T \beta + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$

$$p_{\beta}(y_i | x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_i^T \beta)^2}{2\sigma^2}}$$

$\log(a \cdot b) = \log a + \log b$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmax}} \prod_{i=1}^N p_{\beta}(y_i | x_i)$$

$$= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N -\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y_i - x_i^T \beta)^2}{2\sigma^2}} \right)$$

$$= \underset{\beta}{\operatorname{argmin}} N \cdot \left(-\log \frac{1}{\sqrt{2\pi\sigma^2}}\right) + \sum_{i=1}^N -\log e^{-\frac{(y_i - x_i^T \beta)^2}{2\sigma^2}}$$

$$= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \frac{1}{2\sigma^2} (y_i - x_i^T \beta)^2$$

$$= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - x_i^T \beta)^2$$

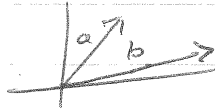
noise magnitude
σ² does not depend
on x_i

So least squares is often associated with assumption of Gaussian errors, but it can be useful more widely! (see I or II)

(4)

III L₂-Projection

$$a, b \in \mathbb{R}^N$$



L₂-norm of a : $\|a\| = \sqrt{\sum_{i=1}^N a_i^2}$ is length of a

Distance between a and b : $\|a-b\| = \sqrt{\sum_{i=1}^N (a_i-b_i)^2}$
 p features

$$V = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \quad U = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{N1} & \dots & x_{Np} \end{pmatrix} = \begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix}$$

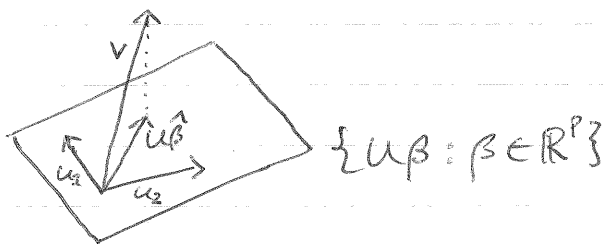
\uparrow u_1 \uparrow u_p

$$(U\beta)_i = x_i^T \beta$$

$$RSS(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2 = \|V - U\beta\|^2$$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} RSS(\beta) = \underset{\beta}{\operatorname{argmin}} \|V - U\beta\|^2 = \underset{\beta}{\operatorname{argmin}} \|V - U\beta\|$$

minimizes distance between V and $U\beta = u_1\beta_1 + \dots + u_p\beta_p$



Least squares projects V onto hyperplane
 $\{U\beta : \beta \in \mathbb{R}^p\}$

Can be sensible even for non-Gaussian errors!

($\hat{\beta}$ unique $\iff u_1, \dots, u_p$ linearly independent)

3b Computing the Least Squares Estimate $\hat{\beta}$

Before we start: minimizer $\hat{\beta}$ not unique if

x_1, \dots, x_p not linearly independent
(see projection picture) \mathcal{P}

handle by
pre-processing features

I. Gradient

$$a \in \mathbb{R}^p \quad \nabla h(a) = \begin{pmatrix} \frac{\partial h(a)}{\partial a_1} \\ \vdots \\ \frac{\partial h(a)}{\partial a_p} \end{pmatrix}$$

generalizes derivative
to functions of more
than one variable.

If h convex, then \hat{a} such that $\nabla h(\hat{a}) = \vec{0}$ is a minimum

II. Applied to $RSS(\beta)$

$$\begin{aligned} \frac{\partial RSS(\beta)}{\partial \beta_j} &= \frac{\partial}{\partial \beta_j} \sum_{i=1}^N (y_i - \sum_j x_{ij} \beta_j)^2 = \sum_{i=1}^N 2 (y_i - \sum_j x_{ij} \beta_j) (-x_{ij}) \\ &= 2 \sum_{i=1}^N x_{ij} (x_i^T \beta - y_i) \end{aligned}$$

$$\nabla RSS(\beta) = 2 \sum_{i=1}^N x_i (x_i^T \beta - y_i) = 2 \left(\sum_{i=1}^N x_i x_i^T \right) \beta - 2 \sum_{i=1}^N x_i y_i$$

$$\nabla RSS(\hat{\beta}) = \vec{0} : \left(\sum_{i=1}^N x_i x_i^T \right) \hat{\beta} = \sum_{i=1}^N x_i y_i$$

Suppose only one feature ($p=1$): $\hat{\beta} = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$

What is $\hat{\beta}$ if this feature is always 1? \leftarrow mean of y_i
~~But what if~~

6

For general p : $\left(\sum_{i=1}^N x_i x_i^T\right)_{j,k} = \sum_i x_{ij} x_{ik} = u_j^T u_k = (U^T U)_{j,k}$

$$\sum_i x_i x_i^T = U^T U$$

$$\sum_{i=1}^N x_i y_i = U^T v$$

Thus: $U^T U \hat{\beta} = U^T v$

$$\hat{\beta} = (U^T U)^{-1} U^T v$$

What if u_1, \dots, u_p are not linearly independent?
Then inverse does not exist!

3c. Bias-Variance Decomposition

Goal: understand how more features = larger model \mathcal{F} affects $\mathbb{E}[EPE(\hat{f})]$: (balance between under- and overfitting)

Thm. Let \hat{f} be any estimator. Let $\bar{f} = \mathbb{E}[\hat{f}]$. Then

$$\mathbb{E} EPE(\hat{f}) = EPE(f_B) \quad \leftarrow \text{Bayes error}$$

$$+ \mathbb{E}_x (f_B(x) - \bar{f}(x))^2 \leftarrow \text{bias (typically smaller with larger } \mathcal{F})$$

$$+ \mathbb{E}_{T, X, Y} (\hat{f}(x) - \bar{f}(x))^2 \leftarrow \text{variance (typically larger with larger } \mathcal{F})$$

Proof: see lecture notes week 1. \square

(show polynomial slides + bias-variance picture from book)

Suppose linear model is correct:

$$y = x^T \beta^* + \epsilon \quad \text{with independent noise s.t. } \mathbb{E}[\epsilon] = 0$$

Conditional on features u in training set:

$$\begin{aligned} \mathbb{E}_T[\hat{\beta} | u] &= \mathbb{E}_T[(u^T u)^{-1} u^T v | u] \\ &= (u^T u)^{-1} u^T \mathbb{E}[v | u] \\ &= (u^T u)^{-1} u^T u \beta^* \\ &= \beta^* \end{aligned}$$

$$\mathbb{E}_T[\hat{\beta}] = \mathbb{E}_{T|u}[\mathbb{E}_T[\hat{\beta} | u]] = \beta^*$$

Bias is 0, because $\bar{f}(x) = \mathbb{E}_T[x^T \hat{\beta}] = x^T \mathbb{E}_T[\hat{\beta}] = x^T \beta^* = f_B(x)$
 if model correct!!

Suppose model correct + Gaussian noise.

$$y = x^T \beta^* + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$EPE(f_B) = \sigma^2$$

$$\text{bias} = 0$$

$$\text{variance} \rightarrow \sigma^2 \cdot \frac{p}{N} \text{ as } N \rightarrow \infty \text{ (book p. 26)}$$

$$\mathbb{E}_T[EPE(\hat{f})] \rightarrow \sigma^2 + 0 + \sigma^2 \cdot \frac{p}{N} \text{ as } N \rightarrow \infty$$

↑ model correct + Gaussian noise! still OK for pretty large p .

(discuss bias-variance trade-off) (beat curse of dimensionality by imposing linear model?)

8

4. k-Fold Cross-Validation

Hyperparameter $m \in \{1, \dots, M\}$

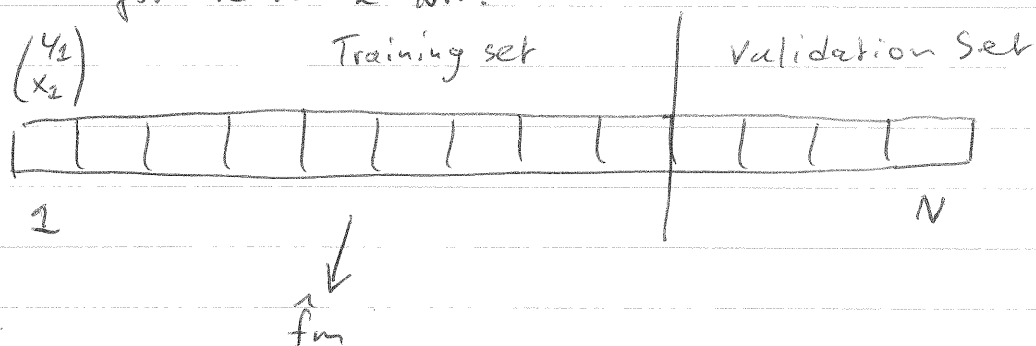
E.g. * m is degree of polynomial in linear regression
* m is k in k -nearest neighbour

Training set $\rightarrow \hat{f}_m$ for each m .

~~What happens if I look at the fit on training data
I construct least squares estimate \hat{f}_m for
each degree polynomial~~

What happens if I look at the fit on the training data to choose m ? (i.e., I use ERM)?

- for degree of polynomial.
- for k in k -NN.



Hold-out estimate:

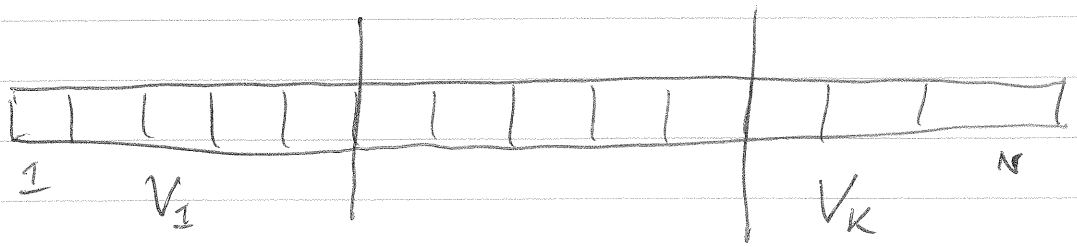
- Randomly set aside validation set V
- Construct \hat{f}_m on rest of training data for all m
- Use ERM on validation set:

$$\hat{m} = \operatorname{argmin}_m \frac{1}{|V|} \sum_{i \in V} L(y_i, \hat{f}_m(x_i))$$

- Works well if ~~subset~~ V is large enough, but that reduces training set a lot

K-Fold Cross validation:

Try to get away with smaller validation set by averaging multiple choices.



For $k = 1, \dots, K$:

- train \hat{f}_m on all data except V_k (for each m)

(use V_k as
validation set)

- evaluate on V_k :

$$\hat{L}_m^k = \frac{1}{|V_k|} \sum_{i \in V_k} L(y_i, \hat{f}_m(x_i))$$

$$\hat{m} = \underset{m}{\operatorname{argmin}} \frac{1}{K} \sum_{k=1}^K \hat{L}_m^k \quad (\text{average over all choices of validation set})$$

Often train final \hat{f}_m on all data, or take the ~~mean~~ average of the ~~estimates~~ K estimators found for \hat{m} during cross-validation.

$K = N$: "leave-one-out cross-validation"