

ML : Week 4

Model selection (cv)

Bias - Var decomp

Motivation

- f is good?
→ generalization and

• Model selection.

$$\text{LS} \min_{W \in \mathbb{R}^D} \frac{1}{2N} \sum_{i=1}^N (y_i - w^T x_i)^2 + \lambda \|w\|_2^2$$

$$\underbrace{\{\lambda\}}_{\text{hyper para.}} \rightarrow W_f \text{ sol.} =$$

$$\lambda_1, \lambda_2, \dots, \lambda_k \rightarrow W_1, \dots, W_k$$

which λ should I use?

- Augmenting the features space: (x_1, x_2)

$$\rightarrow (x_1, x_2, x_1^2 + x_2^2)$$

$$(x_1, x_2, x_1^3 + 5x_2, \dots)$$

- Neural Network (Architecture, weight, depth...)

1st thing = consider Duval model

samples

$$S = \{(x_i, y_i)\}_{i=1}^N \sim D_{\text{unknown}}$$

nd iid

Learning Algorithms

$$A(S) = f_S$$

↑ ↑
input output

Question How good is f_S ?

→ True error:

$$L(f) = \mathbb{E}_{(x,y) \sim D} [\ell(y, f(x))]$$

where ℓ is the loss function

Loss function =

$$\bullet LS = L(y, y') = \frac{1}{2} (y - y')^2$$

- Logistic Function
- Cross entropy

Name in the literature

true risk
expected error loss

Problem: D is unknown!

Solutions: we have sample $S = \{(x_i, y_i)\}_{i=1}^N$ iid $\sim D$

• Empirical risk:

$$LS(f) = \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i))$$

$$\begin{aligned} \mathbb{E}_{S \sim D} [L_S(f)] &= \mathbb{E}_{S \sim D} \left[\frac{1}{N} \sum_i^N \ell(y_i, f(x_i)) \right] \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{S \sim D} [\ell(y_i, f(x_i))] \end{aligned}$$

$\hookrightarrow L_D(f)$

$$\mathbb{E}_{S \sim D} [L_S(f)] = L_D(f) \Rightarrow \text{(unbiased)}$$

Thus The empirical risk is an unbiased estimator!

But we can have fluctuations

In Application:

Given $S \sim D$

$$f_S = \hat{v}(S)$$

then $L_S(f_S) = \frac{1}{N} \sum_{i=1}^N \ell(y_i, f_S(x_i))$

Training error

In general, $L_S(f_S) \neq L_D(f_S)$

Basic idea:

Split data S :

$$S = S_{\text{train}} \cup S_{\text{test}}$$

• Learn $f_{S_{\text{train}}}$ on S_{train} (using training dat.)

• Validate it using $S_{\text{test}}: L_{\text{Test}}(f_{S_{\text{train}}})$

$$\Rightarrow L_{\text{Test}}(f_{S_{\text{train}}}) = \frac{1}{|S_{\text{test}}|} \sum_{i=1}^{|S_{\text{test}}|} \ell(y_i, f_{S_{\text{train}}}(x_i))$$

Test/validation risk (emp. risk)

$$\mathbb{E}_{S_{\text{train}} \sim D} [L_{\text{stat}}(f_{S_{\text{train}}})] = L_D(f_{S_{\text{train}}})$$

• Fluctuations:

• Fixed funct. f

S samples (S was not used to learn f)

$$|L_D(f) - L_S(f)| \leq \epsilon$$

→ can we bound this?

$$\rightarrow L_D(f) \leq \underbrace{L_S(f)}_{\text{we can compute}} + \epsilon$$

Assumptions:

• $\ell \in [0, 1]$

• $S \sim D$ iid $|S|$

• f

Claim:

$$\Pr [|L_D(f) - L_S(f)| \geq \sqrt{\frac{\log^2 S}{2|S|}}] \leq \delta$$

Proof we reduce it to the standard inequality.

$$\theta_n \sim \text{iid } \in [a, b]$$

$$\tilde{\theta} = \mathbb{E} [\theta_n]$$

$$\text{IP} \left[\left| \frac{1}{N} \sum_{i=1}^N \theta_i - \mathbb{E} [\theta_i] \right| \geq \varepsilon \right] \leq 2 \cdot e^{-\frac{2NE^2}{(a-b)^2}}$$

(concentration inequality).

Equivalence: $\theta_i = l(y_i, f(x_i))$

$$\varepsilon = \sqrt{\frac{\log \frac{1}{\delta}}{2N}}$$

How can we use it?

- [f given], Is f good?
S given

$$\text{IP} \left(L_D(f) \geq L_S(f) + \sqrt{\frac{(a-b)^2 \log \frac{1}{\delta}}{2|S|}} \right) \leq \delta$$

• S

we split $S = S_{\text{Train}} \cup S_{\text{Test}}$

$$V(S_{\text{Train}}) = \int S_{\text{Train}}$$

$$P\left[L_D(f_{S_{\text{Train}}}) > L_{S_{\text{Test}}}(f_{S_{\text{Train}}}) + \sqrt{\frac{(a-b)\log^2 S}{2|S_{\text{Test}}|}}\right] \leq S$$

→ bound on the generalization error

$$\text{Generalization error} = L_D - L_S$$

2nd

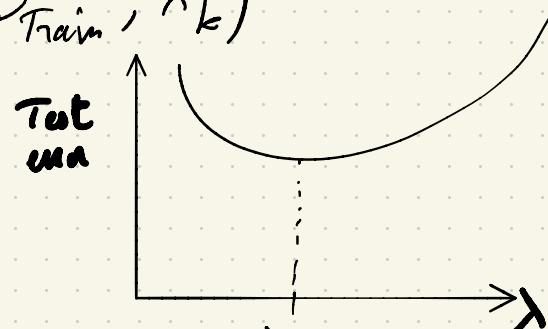
Model Selection

We have K hypo. $(\lambda_1, \lambda_2, \dots, \lambda_K)$

$$S = S_{\text{Train}} \cup S_{\text{Test}}$$

$$\cdot f_{S_{\text{Train}}, \lambda_k} = A(S_{\text{Train}}, \lambda_k)$$

$$\cdot L_{S_{\text{Test}}} (f_{S_{\text{Train}}, \lambda_k})$$



Union bound:

$$\delta_k = \delta_{S_{\text{Train}}, \lambda_k}$$

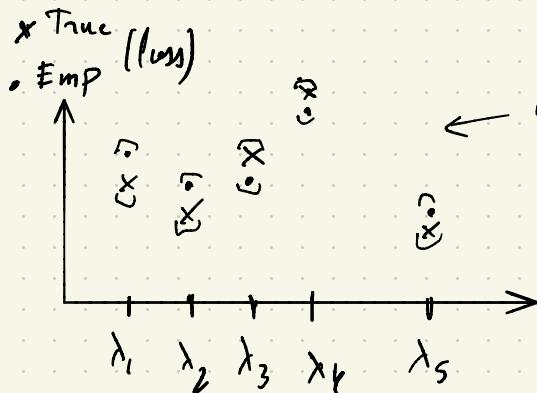
λ_{\min}
[Take This λ !]

$$P[\max_k |L_D(f_k) - L_{S_{\text{Test}}}(f_k)| > \varepsilon] \leq \delta$$

$$\leq \sum_k P[|L_D(f_k) - L_{S_{\text{Test}}}(f_k)| > \varepsilon]$$

$$\leq 2K e^{-\frac{2N\varepsilon^2}{(\alpha-\beta)^2}}$$

$$P\left[\max_k |L_D(f_k) - L_{S_{\text{Test}}}(f_k)| > \sqrt{\frac{(b-a)^2 \log \frac{2K}{\delta}}{2|S_{\text{Test}}|}}\right] \leq \delta$$



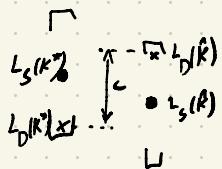
← all The width are The same

$$\hat{K} = \arg \min_k L_{S_{\text{Test}}} (f_{S_{\text{Train}}}, \lambda_k)$$

$$K^* = \arg \min_k L_D (f_{S_{\text{Train}}}, \lambda_k)$$

$$P \left[|L_D(f_{K^*}) - L_D(f_K)| \geq 2 \sqrt{\frac{(b-a)^2 \log \frac{2K}{\delta}}{2|S_{\text{Test}}|}} \right] \leq \delta$$

Worst case:



$$C \leq 2 \sqrt{\frac{(b-a)^2 \log \frac{2K_S}{\delta}}{2|S_{\text{Test}}|}}$$

$$K^* \hat{K}$$

- Extension to infinitely many models ($K=\infty$) (lecture notes)

Proof of Chernoff bound:

$$\alpha < \theta < b$$

Without loss of generality; assume $E[\theta] = 0$

$$\text{we will prove that } P\left[\frac{1}{N} \sum_{i=1}^N \theta_i > \varepsilon\right] \leq e^{-\frac{2N\varepsilon^2}{(b-a)^2}}$$

$$\text{together with } P\left[\frac{1}{N} \sum_{i=1}^N \theta_i \leq -\varepsilon\right] \leq e^{-\frac{2N\varepsilon^2}{(b-a)^2}}$$

$$\rightarrow P\left[\left|\frac{1}{N} \sum_{i=1}^N \theta_i\right| > \varepsilon\right] \leq 2e^{-\frac{2N\varepsilon^2}{(b-a)^2}}$$

$$P\left[\frac{1}{N} \sum_{i=1}^N \theta_i > \varepsilon\right] \stackrel{S \geq 0}{=} P\left[\frac{S}{N} I(\theta_i) > S\varepsilon\right]$$

$$= P\left[e^{\frac{S}{N} I(\theta_i)} > e^{S\varepsilon}\right]$$

$$(\text{markov}) \leq \min_S E\left[e^{\frac{S}{N} I(\theta_i)}\right] \cdot e^{S\varepsilon}$$

$$\text{Hoeffding lemma: } E[e^{Sx}] \leq e^{\frac{1}{2} S^2(b-a)^2}$$

$$\Rightarrow \leq \min_S e^{S^2(b-a)^2/2N} \cdot e^{-S\varepsilon}$$

$$\leq e^{-\frac{2N\varepsilon^2}{(b-a)^2}} \text{ for } S = \frac{4N\varepsilon}{(b-a)}$$

Proof of Hoeffding lemma: by convexity:

$$e^{Sx} \leq \frac{x-\alpha}{b-\alpha} e^{Sb} + \frac{b-x}{b-\alpha} e^{Sa}$$

$$E[x] = 0$$

$$E[e^{Sx}] \leq \frac{bc^{Sa} - ac^{Sb}}{b-a} \leq \dots \leq c^{S^2(b-a)^2}$$