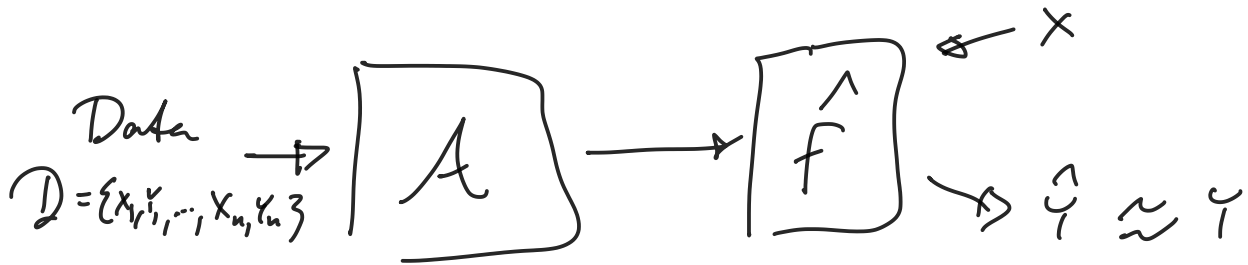


Lec 6

Tuesday, September 17, 2019 10:54

Recap: What makes a good supervised learning algorithm?



$$\begin{aligned}
 \mathcal{R}(A) &= \mathbb{E}_{\mathcal{D}} [\mathcal{R}(\hat{f})] \\
 &= \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{x, y} [(\hat{f}(x) - y)^2]] \\
 &= \mathbb{E}_x \text{Var}(y|x) + \mathbb{E}_x \text{Err}(x) \\
 \text{Err}(x) &= \mathbb{E}_{\mathcal{D}} [(\hat{f}(x) - f^*(x))^2] \\
 &= \text{Var} \hat{f}(x) + \underbrace{(\mathbb{E}_{\mathcal{D}} \hat{f}(x) - f^*(x))^2}_{\text{bias}}
 \end{aligned}$$

\uparrow variance \uparrow bias

Subset Selection

Find a subset of features that predict y well.

Why?

- Avoid overfitting
- Trade off bias/var
- Interpretability

Best subset

For the rest of lecture

$$\beta \in \mathbb{R}^{p+1} \quad \beta_0 \text{ intercept}$$

$$\beta_{1:p} \text{ coeffs}$$

$$\text{Define } \text{supp}(\beta) = \{j = 1, \dots, p : \beta_j \neq 0\}$$

$$\text{supp}\left(\begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}\right) = \{1, 2\}$$

for any set $S \subseteq \{1, \dots, p\}$

$$\hat{\beta}_S = \underset{\text{supp}(\beta) \subseteq S}{\text{argmin}} \|\mathbb{Y} - \mathbb{X} \beta\|_2^2 \leftarrow \text{OLS on the features } \{x_j : j \in S\} \text{ \& an intercept}$$

$$\hat{\beta}_{\text{best-}k} = \underset{|\text{supp}(\beta)| \leq k}{\text{argmin}} \|\mathbb{Y} - \mathbb{X} \beta\|_2^2$$

$$= \underset{\substack{|S| \leq k \\ S \subseteq \{1, \dots, p\}}}{\text{argmin}} \|\mathbb{Y} - \mathbb{X} \hat{\beta}_S\|_2^2$$

$$\hat{\beta}_{\text{best-}p} = \hat{\beta}_{\text{OLS}}$$

Computing $\hat{\beta}_{\text{best-}k}$ is hard for p, k large

$$\binom{p}{k} \sim \left(\frac{p}{k}\right)^k$$

Idea: sequentially add/remove the most "influential" features

Large coeff \rightarrow more influential?

But: subject to scaling
 \Rightarrow need to normalize

For $\hat{\beta}$ given by OLS, define the z-score

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma} \cdot \sqrt{v_j}}$$

$$\hat{\sigma}^2 = \frac{1}{n-k-1} \sum_i (y_i - \hat{y}_i)^2$$

$$v_j = ((X^T X)^{-1})_{jj}$$

General rule thumb: $|z_j| \geq 2 \Rightarrow$ "significant"
 $< 2 \Rightarrow$ "insignificant"

Backward - & Forward - stepwise Regression

BSR:

Start w/ $S = \{1, \dots, p\}$

While $|S| > k$:

 Compute $\hat{\beta}_S$

 Remove from S the feature w/
 smallest abs z-score

Repeat

FSR:

start w/ $S = \{\}$

While $|S| < k$:

 Find $j^* = \underset{j \notin S}{\operatorname{argmin}} \|Y - X \hat{\beta}_{S \cup \{j\}}\|_2^2$

 Add j^* to S

Repeat

At what k to stop?

The AML approach: use cross-validation.

Cross-Validation (CV)

CV is a technique to estimating $R(A)$

Split the data into \underline{K} folds:

$$S_1, \dots, S_K = S_1, \dots, S_K$$

$$\text{s.t. } S_i \cap S_j = \emptyset$$

$$|S_i| - |S_j| \leq 1$$

$$\hat{f}^{(j)} = A \left(\{(x_i, y_i) : i \notin S_j\} \right) \quad j=1, \dots, K$$

$$CV^{(j)} = \frac{1}{|S_j|} \sum_{i \in S_j} \ell(y_i, \hat{f}^{(j)}(x_i))$$

$$\hat{R}^{CV}(A) = \frac{1}{\underline{K}} \sum_{j=1}^{\underline{K}} CV^{(j)}$$

Collection of algorithms A_1, \dots, A_m

How to choose?

Naïve approach: choose A_j w/ smallest $\hat{R}^{CV}(A_j)$

The AML approach (AKA "one-std-err" rule of thumb)

$$\text{Std Err}(\hat{R}^{CV}(A)) = \frac{1}{\underline{K}} \sqrt{\sum_j (\hat{R}^{CV}(A) - CV^{(j)}(A))^2}$$

Pick the "simplest" algo w/ \hat{R}^{CV}

within one std err of the minimal one

Simplest:

- least # of variables
- least complexity
- least high order dependence

- least var