

## Lecture 5. Regularization & Model Selection.

### ◦ Introduction:

(1) different model  $\rightarrow$  learning problem.

(2) finite sets of models.

e.g. linear model, SVM, CNN.

### ◦ Cross validation:

(1) empirical risk minimization.

(2) simple cross validations:

① training set  $S$ :

split into  $S_{\text{train}} / S_{\text{cv}}$

② train on  $S_{\text{train}}$

③ select  $M_i$ :

$\min \hat{\mathcal{E}}_{S_{\text{cv}}}(h_i)$  (retrain  $M_i$ )

\*: sensitive model: no retrain.

(3) drawback: wasting data (30%)

data is scarce (e.g.  $m=30$ )

promoted idea: **k-fold**.

①  $S$ : randomly split into disjoint subset

$S_1, S_2, S_3, \dots, S_k$ .

② For  $j=1, 2, \dots, k$ .

train  $M_i$  on  $S / S_j \Rightarrow h_{ij}$

test on  $S_j \Rightarrow \hat{\epsilon}_{S_j}(h_{ij})$

③ pick  $\min_{M_i} \sum_{j=1}^k \hat{\epsilon}_{S_j}(h_{ij})$

\*: retrain on the whole set.

(computational expensive)

(4) leave-one-out:

only hold out one training example.

• Feature Selection:

(1) small number of features  $\rightarrow$  relevant.

$n$  features  $\rightarrow 2^n$  subsets.  $\rightarrow$  heuristic search.

(2) forward search.

①  $F = \Phi$  (feature set)

② Repeat:

$i = 1, 2, \dots, n$  and  $i \notin F$

find  $F_i = F \cup \{i\}$  optimal.

③ Output  $F$ . ( $O(n^2)$ )

Similarly: backward.

(3) Filter feature selection:

score  $S(i) \rightarrow$  how informative.

mutual information:

$$MI(X_i, y) = \sum_{x_i} \sum_{y_j} P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(x_i) P(y_j)}$$

$$KL(P(X_i, y) \parallel P(X_i)P(y))$$

chow different distributions are)

rank  $\rightarrow$  choose  $k$  (using CV)

• Bayesian statistics & regularization:

(1) maximum likelihood:

$$\theta_{ML} = \arg \max_{\theta} \prod_{i=1}^m P(y_i | x_i, \theta)$$

view  $\theta$  as an unknown parameter.

constant value but unknown.

$\theta$  not random.

task: statistical procedure

$\rightarrow$  estimate parameter  $\theta$ .

} frequentist.

(2) Bayesian View:

$\theta$  is random

task: prior belief.  $\rightarrow$  posterior

$$\begin{aligned} P(\theta | S) &= \frac{P(S | \theta) P(\theta)}{P(S)} \\ &= \frac{\left( \prod_{i=1}^m P(y_i | x_i, \theta) \right) P(\theta)}{P(S)} \end{aligned}$$

usually:  $\theta \sim N(0, \tau^2 I)$

(3) Maximum a Posteriori.

$$\theta_{MAP} = \arg \max_{\theta} \prod_{i=1}^m P(y_i | x_i, \theta) P(\theta)$$

smaller norm  $\Rightarrow$  less likely to overfit.