

Lecture 8 Data Splits, Models & Cross Validation

1. bias - variance trade-off

underfit : high bias

overfit : high variance

2. Regularization

examples: $\min_{\theta} \frac{1}{2} \sum_{i=1}^m \|y^{(i)} - \theta^T x^{(i)}\|^2 + \lambda \|\theta\|^2$

$$\operatorname{argmax}_{\theta} \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}; \theta) - \lambda \|\theta\|^2$$

SVM: $\min \|W\|^2$ has the same effect

3. Bayesian statistics and regularization

training set $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^m$

$$p(\theta | S) = \frac{p(S|\theta) \cdot p(\theta)}{p(S)}$$

logistic regression

$$\arg\max_{\theta} p(\theta | S) = \arg\max_{\theta} p(S|\theta) p(\theta) = \arg\max_{\theta} \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right) \cdot p(\theta)$$

assume: $p(\theta) \sim N(0, \tau^2 I)$ (Gaussian)

$$p(\theta) = \frac{1}{\sqrt{2\pi} (\tau^2 I)^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \theta^T (\tau^2 I)^{-1} \theta\right)$$

Frequentist:

$$\arg\max_{\theta} p(S|\theta) - \text{MLE}$$

Bayesian:

prior distribution $p(\theta)$; $\arg\max_{\theta} p(\theta | S) - \text{MAP estimation (maximum a posteriori)}$

4. Data Splits

$S_{\text{train}}, S_{\text{dev}}, S_{\text{test}}$ (dev = development / w data set)

Simple Hold-out Cross Validation:

- 1) Train each model i (option for degree of polynomial / λ / ...) on S_{train} . Get some hypothesis h_i .
- 2) Measure error on S_{dev} . Pick model with lowest error on S_{dev} .
- 3) Optional: apply algorithm to S_{test} and report the error

Traditionally: S_{train} 70% S_{dev} 30% or S_{train} 60% S_{dev} 20% S_{test} 20%

Now: if huge data set S_{train} 90% S_{dev} 5% S_{test} 5%

I. k -fold cross validation (small data set)

Leave-one-out CV (smaller data set, 20-50 rows)

b. Feature Selection

Start with $\mathcal{F} = \emptyset$

Repeat: $\{$

- 1) Try adding each feature to F , and see which single feature addition most improve dev set performance
 - 2) Add that feature to F
- }