

ECE595 / STAT598: Machine Learning I

Lecture 32.3: Validation - Validation in Regularization

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline

- Lecture 30 Overfit
- Lecture 31 Regularization
- **Lecture 32 Validation**

Today's Lecture:

- Validation
 - Concept of validation
 - Properties of validation error
- Model Selection
 - Basic idea
 - Case study
- **Validation in Regularization**
 - **Cross validation**
 - **Parameter selection**

Cross Validation

- A principled way to estimate the out-sample error, without suffering from small K problem.
- Consider the **leave-one-out** approach.
- Let the data set be

$$\mathcal{D}_n = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{n-1}, y_{n-1}), \cancel{(\mathbf{x}_n, y_n)}, (\mathbf{x}_{n+1}, y_{n+1}), \dots, (\mathbf{x}_N, y_N)$$

- Remove the n -th training sample
- Learn the hypothesis function

$$g_n^- = \text{learn from } \mathcal{D}_n.$$

- Let error

$$e_n \stackrel{\text{def}}{=} E_{\text{val}}(g_n^-) = e(g_n^-(\mathbf{x}_n), y_n).$$

- Remark: e_n is based on a single data point (\mathbf{x}_n, y_n) .

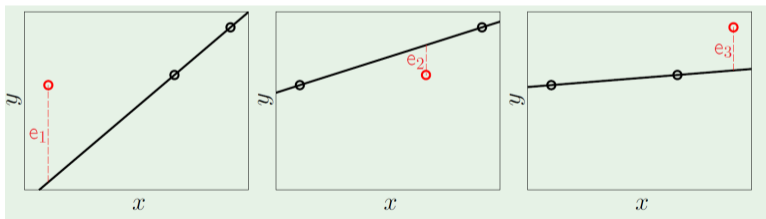
Cross Validation

- This will give you

$$e_1, e_2, \dots, e_N$$

- Let's compute the average

$$E_{cv} = \frac{1}{N} \sum_{n=1}^N e_n.$$



- Validation: Use K samples to validate
- Cross-Validation: Recycle the N samples to validate

Cross-Validation for Linear Regression

- Recall the linear regression model:

$$\mathbf{w}^* = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$$

- How to estimate the optimal λ ?
- Let

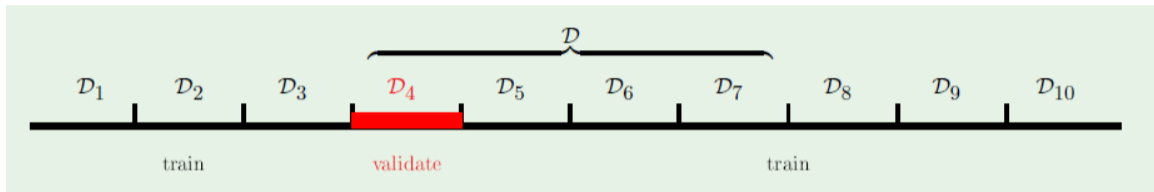
$$\begin{aligned} \mathbf{H}(\lambda) &= \mathbf{A}(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \\ \hat{\mathbf{y}} &= \mathbf{H} \mathbf{y} \end{aligned}$$

- Compute the cross validation score:

$$E_{\text{cv}} = \frac{1}{N} \sum_{n=1}^N \left(\frac{\hat{y}_n - y_n}{1 - H_{n,n}(\lambda)} \right)^2$$

- $H_{n,n}(\lambda) = \mathbf{x}_n^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{x}_n$. (See textbook Problem 4.26.)
- Pick λ that minimizes E_{cv}

V-fold validation



- Leave one out: N training sessions. Each session has $N - 1$ points.
- In practice: Partition the dataset into V sessions.
- Each session has N/V points.
- Train using $\mathcal{D} \setminus \mathcal{D}_V$.
- Test using \mathcal{D}_V .
- Rule of Thumb: $V = 10$. 10-fold cross-validation.

Summary

- Validation says: Break the dataset into testing and validation.
- Use validation set to help selecting models and parameters.
- Then reuse the data to report the final hypothesis.
- Can also use cross-validation to get a better estimate of E_{out} .
- Never use testing data for validation.

Reading List

- Yaser Abu-Mustafa, Learning from Data, Chapter 4.3

Appendix

Unbiasedness of E_{cv}

- Why care? If yes, then we can use E_{cv} to estimate E_{out}
- Recall $g^{(\mathcal{D})}$. The out-sample error for $g^{(\mathcal{D})}$ is

$$E_{out}(N) = \mathbb{E}_{\mathcal{D}} \left[E_{out}(g^{(\mathcal{D})}) \right].$$

- $E_{out}(N)$: Overall out-sample error average over all possible training sets
- $E_{out}(N)$: Function of N . If you have more training samples, then you have lower error
- We can show that

$$\begin{aligned} E_{out}(N) &\stackrel{?}{=} \mathbb{E}_{\mathcal{D}} [E_{cv}] \\ &= \mathbb{E}_{\mathcal{D}} \left[\frac{1}{N} \sum_{n=1}^N e_n \right] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\mathcal{D}} [e_n] = \mathbb{E}_{\mathcal{D}} [e_n]. \end{aligned}$$

Unbiasedness of E_{cv}

- So what is $\mathbb{E}_{\mathcal{D}}[e_n]$?

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[e_n] &= \mathbb{E}_{\mathcal{D}_n, (\mathbf{x}_n, y_n)}[e_n] \\ &= \mathbb{E}_{\mathcal{D}_n} \mathbb{E}_{(\mathbf{x}_n, y_n)}[e(g_n^-(\mathbf{x}_n), y_n)] \\ &= \mathbb{E}_{\mathcal{D}_n} \mathbb{E}_{(\mathbf{x}_n, y_n)}[E_{\text{val}}(g_n^-)] \\ &= \mathbb{E}_{\mathcal{D}_n} E_{\text{out}}(g_n^-) \\ &= E_{\text{out}}(N-1).\end{aligned}$$

decouple \mathcal{D}

unbiasedness of E_{val}
expectation of \mathcal{D}_n

- So,

$$\mathbb{E}_{\mathcal{D}}[E_{cv}] = E_{\text{out}}(N-1).$$

- That means: E_{cv} is an unbiased estimate of $E_{\text{out}}(N-1)$
- Remark: This gives us the mean of E_{cv} . The variance is a lot harder because \mathcal{D}_m and \mathcal{D}_n overlaps.