

# ECE595 / STAT598: Machine Learning I

## Lecture 32 Validation

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



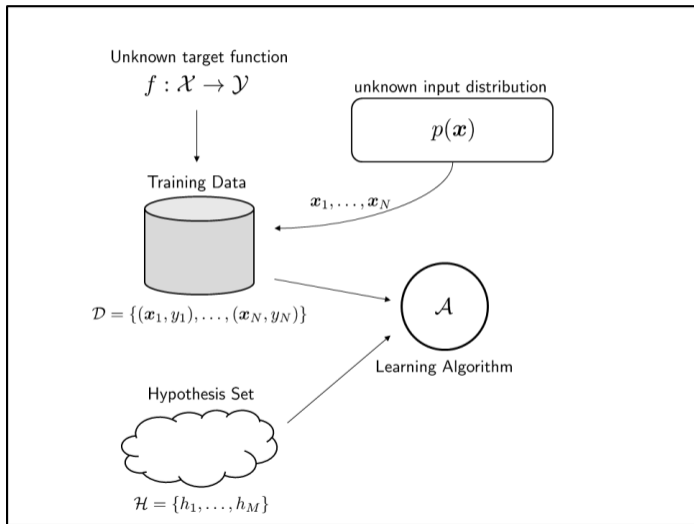
# Outline

- Lecture 31 Overfit
- Lecture 32 Regularization
- **Lecture 33 Validation**

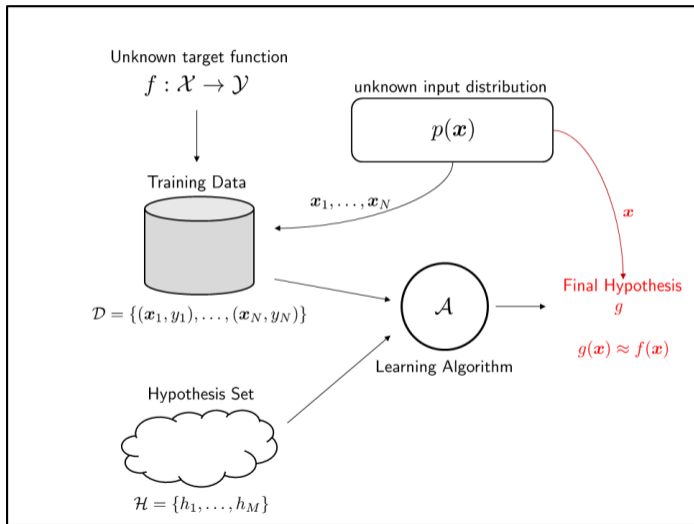
## Today's Lecture:

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

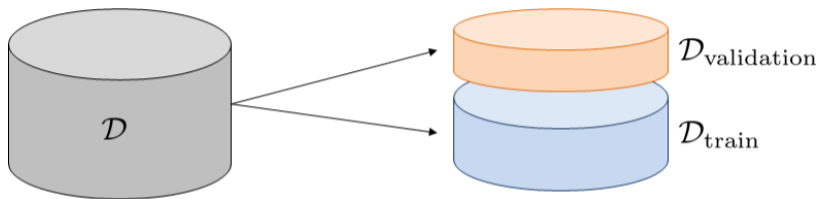
# Evaluating Your Model



# Evaluating Your Model



# Validation Set



- What does  $\mathcal{D}_{\text{val}}$  buy you?
- Generalization bound using  $\mathcal{D}_{\text{val}}$ ?
- How to use  $\mathcal{D}_{\text{val}}$ ?
- Validation vs Cheating
- Cross Validation

# The Role of Validation

- Recall the generalization error:

$$E_{\text{out}}(h) = E_{\text{in}}(h) + \underbrace{\text{overfitpenalty}}_{\text{regularization suppresses this term}}$$

- How about validation?

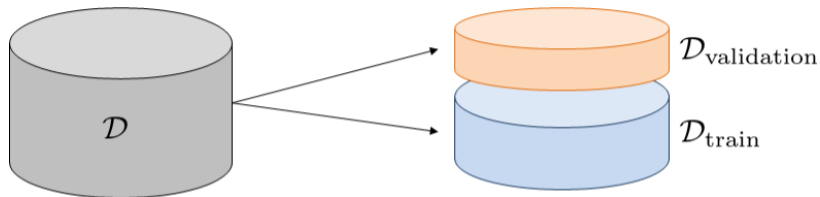
$$\underbrace{E_{\text{out}}(h)}_{\text{validation estimates this term}} = E_{\text{in}}(h) + \text{overfitpenalty}$$

- Is it the same as testing?

$$\underbrace{E_{\text{out}}(h)}_{\text{testing estimates this term}} = E_{\text{in}}(h) + \text{overfitpenalty}$$

- Testing: You cannot use testing set at any stage of training.
- Validation: You can use validation to make choices during training.

## Creating the Validation Set



- Data set:  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ .  $N$  samples.
- Validation set:  $\mathcal{D}_{\text{val}}$ .  $K$  samples.
- Training set:  $\mathcal{D}_{\text{training}}$ .  $N - K$  samples.
- If you run the learning algorithm on  $\mathcal{D}_{\text{train}}$ , you obtain

$$g^- \in \mathcal{H}$$

- $g^-$ : a hypothesis learned by “subtracting” some samples
- $g^-$  is not necessarily the final hypothesis you eventually report

## What does validation tell us?

Goal: Define the validation error  $E_{\text{val}}(g^-)$ , and analyze its statistical properties.

- The **validation error** is

$$E_{\text{val}}(g^-) = \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} e(g^-(\mathbf{x}_n), y_n)$$

- Average error over the *validation set*.  $e(g^-(\mathbf{x}_n), y_n)$ : Point-wise error.

- Classification:

$$e(g^-(\mathbf{x}), y) = \mathbb{I}[g^-(\mathbf{x}) \neq y]$$

- Regression:

$$e(g^-(\mathbf{x}), y) = (g^-(\mathbf{x}) - y)^2$$

- Want to analyze the **mean** and **variance** of  $E_{\text{val}}(g^-)$ .



## Property 1: Mean of $E_{\text{val}}(g^-)$

- Let us analyze the mean of  $E_{\text{val}}(g^-)$ .
- We want to show that the validation error  $E_{\text{val}}(g^-)$  is an **unbiased estimate** of  $E_{\text{out}}$
- That is, the expectation of  $E_{\text{val}}(g^-)$  is  $E_{\text{out}}$
- Here is why:

$$\mathbb{E}_{\mathcal{D}_{\text{val}}}[E_{\text{val}}(g^-)] = \mathbb{E}_{\mathcal{D}_{\text{val}}}\left[\frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} e(g^-(\mathbf{x}_n), y_n)\right] \quad \text{(definition)}$$

$$= \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} \mathbb{E}_{\mathcal{D}_{\text{val}}}[e(g^-(\mathbf{x}_n), y_n)] \quad \text{(linearity)}$$

$$= \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} \mathbb{E}_{\mathbf{x}_n}[e(g^-(\mathbf{x}_n), y_n)] \quad \mathcal{D}_{\text{val}} = (\mathbf{x}_n, f(\mathbf{x}_n))$$

$$= \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} E_{\text{out}}(g^-) = E_{\text{out}}(g^-) \quad \mathbf{x}_n \sim p(\mathbf{x})$$

## Property 2: Variance of $E_{\text{val}}(g^-)$

- Define  $\sigma_{\text{val}}^2 = \text{Var}_{\mathcal{D}_{\text{val}}}[E_{\text{val}}(g^-)]$ .
- How does  $\sigma_{\text{val}}^2$  depend on  $K$ ?
- Let's do some calculation

$$\begin{aligned}\sigma_{\text{val}}^2 &= \text{Var}_{\mathcal{D}_{\text{val}}} \left[ \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} e(g^-(\mathbf{x}_n), y_n) \right] && \text{(definition)} \\ &= \frac{1}{K^2} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} \underbrace{\text{Var}_{\mathcal{D}_{\text{val}}} [e(g^-(\mathbf{x}_n), y_n)]}_{\stackrel{\text{def}}{=} \sigma^2(g^-)} && \text{(independence)} \\ &= \frac{1}{K^2} \sum_{\mathbf{x}_n \in \mathcal{D}_{\text{val}}} \sigma^2(g^-) \\ &= \frac{1}{K} \sigma^2(g^-).\end{aligned}$$

## Property 2: Variance of $E_{\text{val}}(g^-)$

- If we consider a classification problem so that  $e(g^-(\mathbf{x}), y) = \mathbb{I}[g^-(\mathbf{x}) \neq y]$
- Then

$$\sigma_{\text{val}}^2 = \frac{1}{K} \sigma^2(g^-) = \frac{1}{K} \text{Var}_{\mathcal{D}_{\text{val}}} [e(g^-(\mathbf{x}), y)] \quad (\text{definition})$$

$$= \frac{1}{K} \text{Var}_{\mathcal{D}_{\text{val}}} [\mathbb{I}[g^-(\mathbf{x}) \neq y]] \quad (\text{classification})$$

$$= \frac{1}{K} \mathbb{P}[g^-(\mathbf{x}) \neq y](1 - \mathbb{P}[g^-(\mathbf{x}) \neq y]) \quad (\text{Bernoulli}).$$

- Remark: If  $X$  is Bernoulli, then  $\text{Var}[X] = p(1 - p) \leq \frac{1}{4}$ .
- Therefore, we can bound  $\sigma_{\text{val}}^2$  using

$$\sigma_{\text{val}}^2 \leq \frac{1}{4K}.$$

- So as  $K \rightarrow \infty$ ,  $\sigma_{\text{val}}^2 \rightarrow 0$ .

## Does $E_{\text{val}}(g^-)$ Generalize?

- $E_{\text{val}}(g^-)$  is a **random variable**. So it fluctuates.
- Mean:  $\mathbb{E}_{\mathcal{D}_{\text{val}}}[E_{\text{val}}(g^-)]$ .
- Variance:  $\text{Var}_{\mathcal{D}_{\text{val}}}[E_{\text{val}}(g^-)]$ .
- Previous slide:  $\mathbb{E}_{\mathcal{D}_{\text{val}}}[E_{\text{val}}(g^-)] = E_{\text{out}}(g^-)$ .
- So we should expect Hoeffding inequality to apply:

$$E_{\text{out}}(g^-) \leq E_{\text{val}}(g^-) + \mathcal{O}\left(\frac{1}{\sqrt{K}}\right).$$

- Why? Recall Hoeffding inequality for *one* hypothesis:

$$E_{\text{out}}(h) \leq E_{\text{in}}(h) + \mathcal{O}\left(\frac{1}{\sqrt{N}}\right).$$

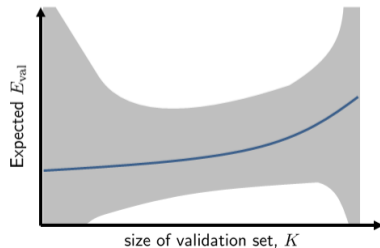
- So as  $K$  grows,  $E_{\text{val}}(g^-)$  actually generalizes  $E_{\text{out}}(g^-)$  very well.

## Large $K$ or Small $K$ ?

- No matter how you look at the result: Generalization bound or variance bound

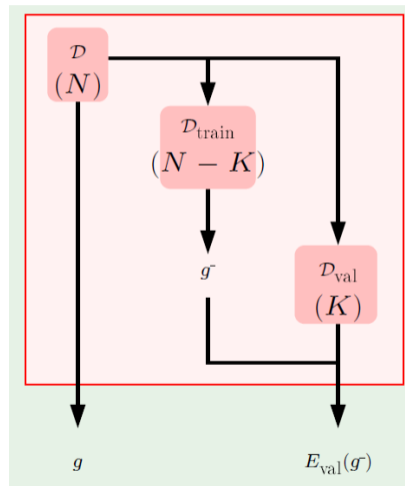
$$\sigma_{\text{val}}^2 \leq \frac{1}{4K}.$$

- If  $K \rightarrow \infty$ , then  $\sigma_{\text{val}}^2 \rightarrow 0$
- So large  $K$  is good.
- But can  $K$  be really really large?
- No.  $K$  for validation,  $N - K$  for training.



## Re-Using $K$

- Is it a waste if we can only use  $N - K$  samples for training?
- No. You are *allowed* to reuse the  $K$  samples
- Use  $\mathcal{D}_{\text{val}}$  to give an estimate of  $E_{\text{val}}(g^-)$
- Use  $E_{\text{val}}(g^-)$  as a guide to choose  $g$
- Here is a pictorial illustration
- Rule of Thumb:  $K = \frac{N}{5}$



# Outline

- Lecture 31 Overfit
- Lecture 32 Regularization
- **Lecture 33 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

## Validation for Model Selection

- Consider a set of  $M$  models:  $\mathcal{H}_1, \dots, \mathcal{H}_M$
- E.g., linear / quadratic / logistic, etc
- E.g., linear model with different regularization parameters, etc
- How to choose the model?
- Use  $\mathcal{D}_{\text{train}}$  to train  $g_1^-, \dots, g_M^-$ .
- Evaluate

$$E_m = E_{\text{val}}(g_m^-),$$

for  $m = 1, \dots, M$ .

- $E_m$  is an **unbiased estimate** of the out-sample error  $E_{\text{out}}(g_m^-)$ .
- Select the one with the minimum validation error:

$$m^* = \underset{m}{\operatorname{argmin}} E_m$$

- The model  $\mathcal{H}_{m^*}$  is the best model



# Generalization Bound for Model Selection

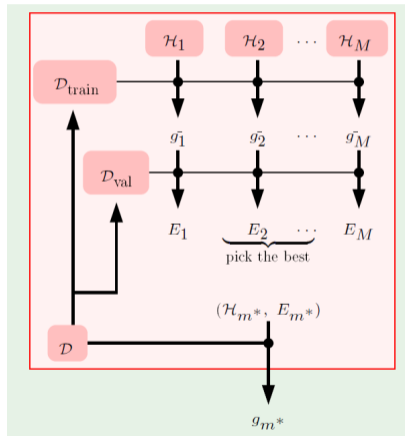
- If you choose  $g_{m^*}^-$  from  $g_1^-, \dots, g_M^-$
- You are effectively considering

$$\mathcal{H}_{\text{val}} = \{g_1^-, \dots, g_M^-\}.$$

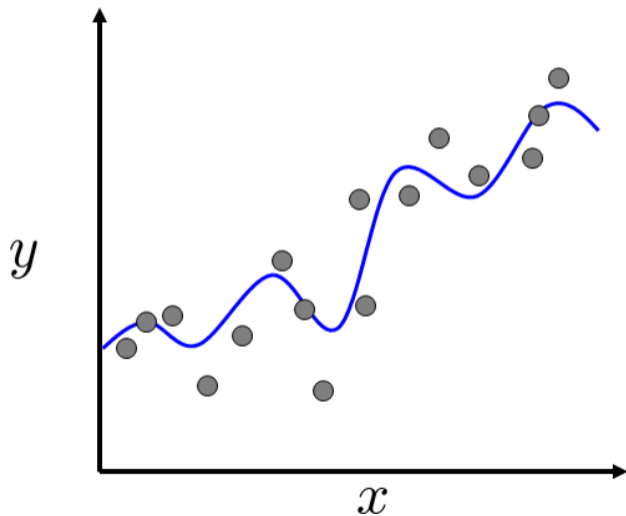
- So the price you need to pay in the generalization bound is

$$E_{\text{out}}(g_{m^*}^-) \leq E_{\text{val}}(g_{m^*}^-) + \mathcal{O}\left(\sqrt{\frac{\log M}{K}}\right).$$

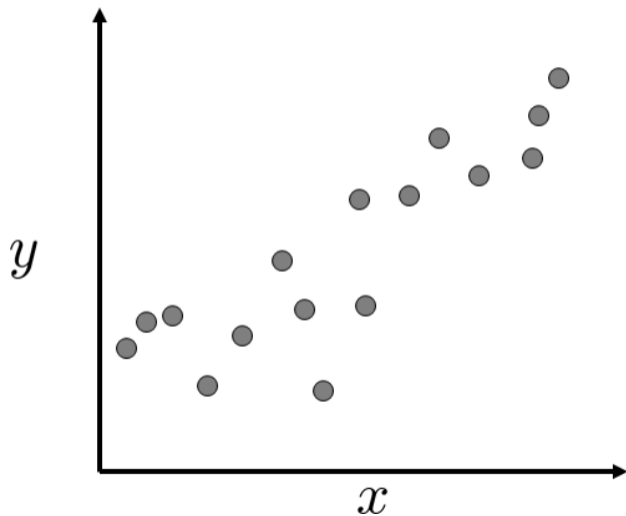
- Use  $g_{m^*}^-$  as the final hypothesis?
- No. Should choose  $\mathcal{H}_{m^*}$ , and train with  $N$  samples.



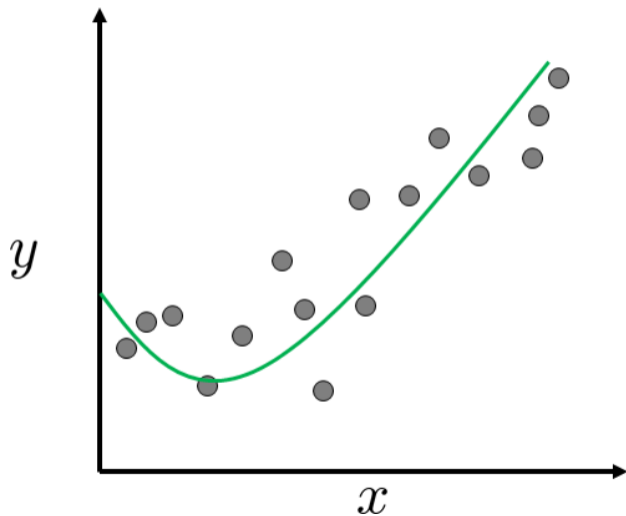
## Case Study: $\mathcal{H}_2$ vs $\mathcal{H}_5$



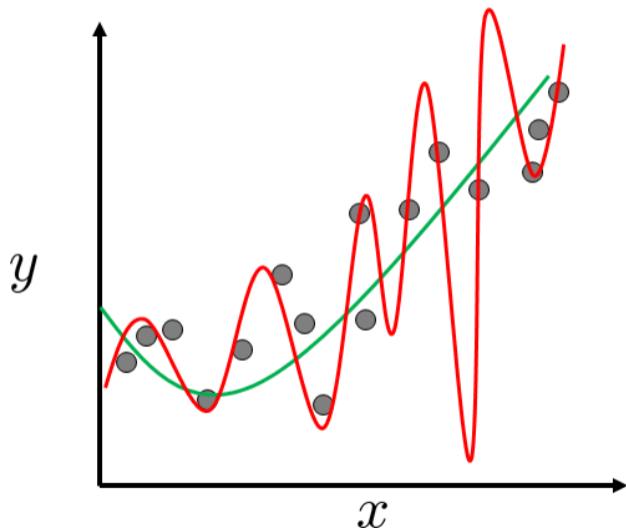
## Case Study: $\mathcal{H}_2$ vs $\mathcal{H}_5$



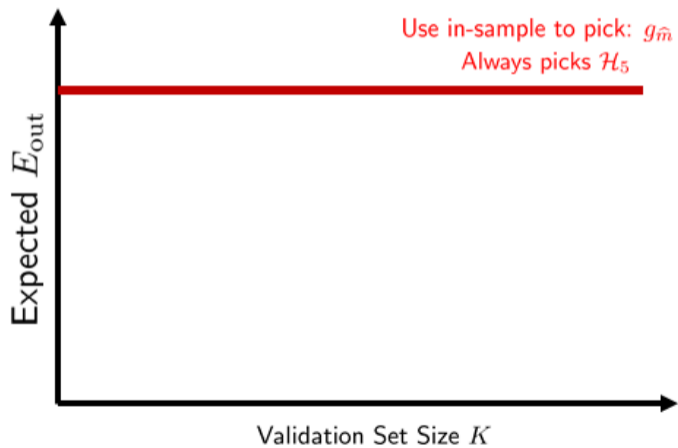
## Case Study: $\mathcal{H}_2$ vs $\mathcal{H}_5$



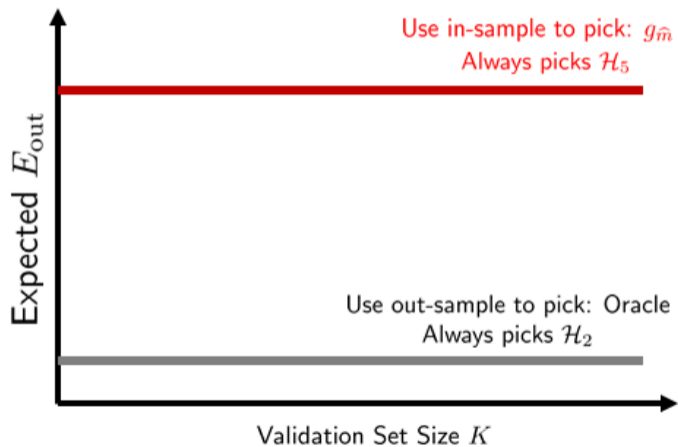
## Case Study: $\mathcal{H}_2$ vs $\mathcal{H}_5$



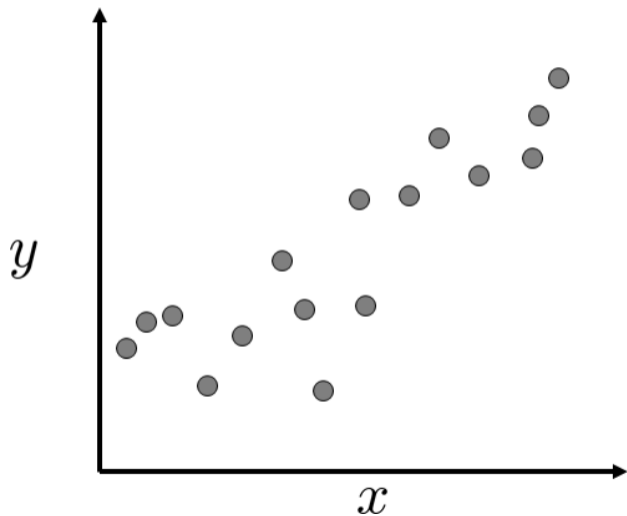
## Expected Error



## Expected Error

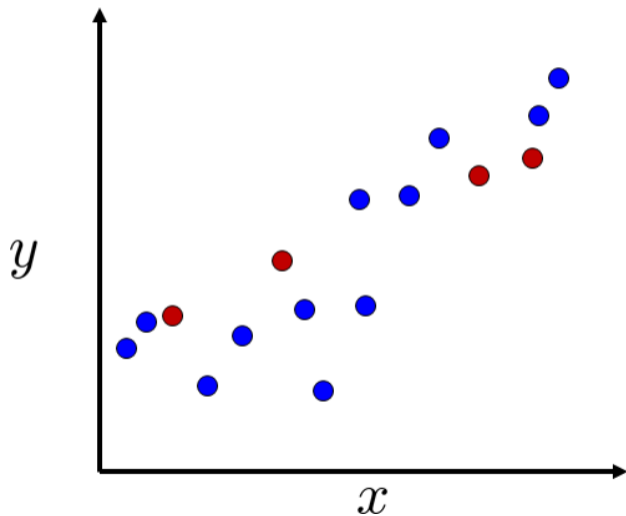


# Validation

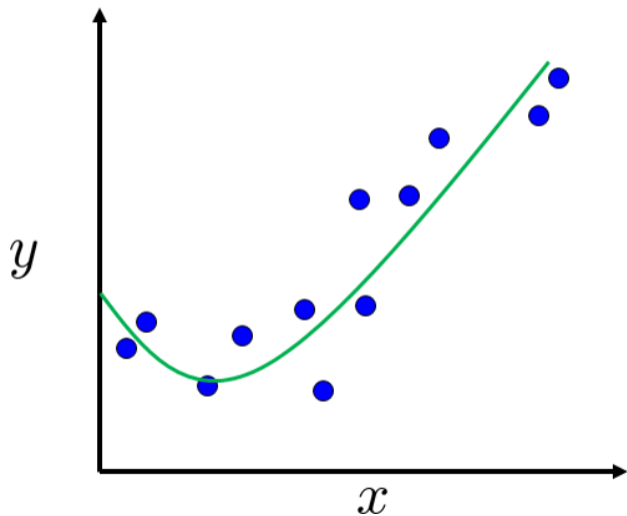




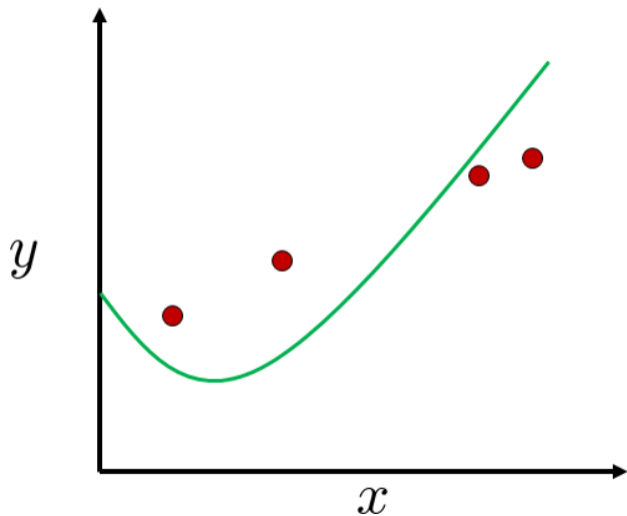
# Validation



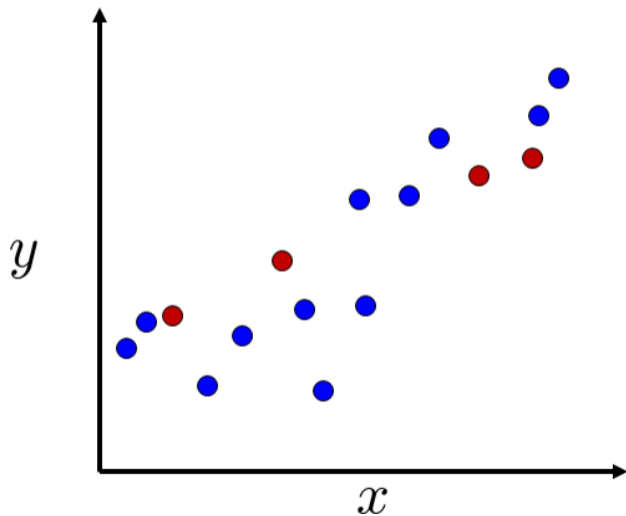
# Validation



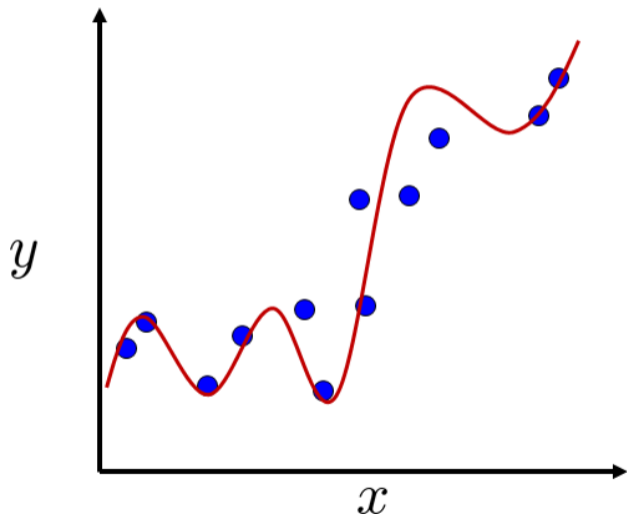
# Validation



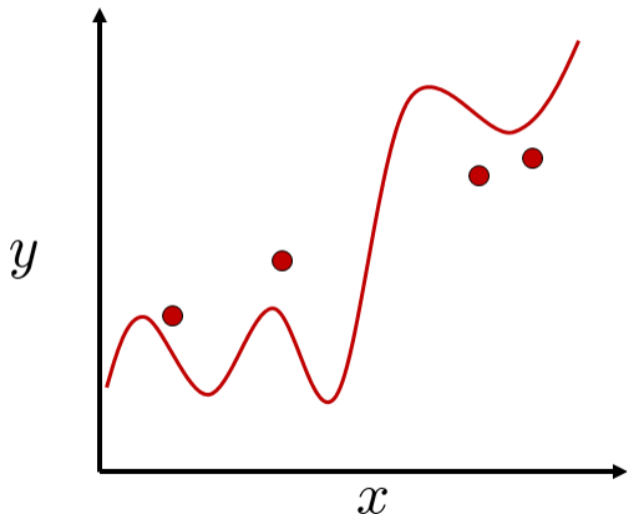
# Validation



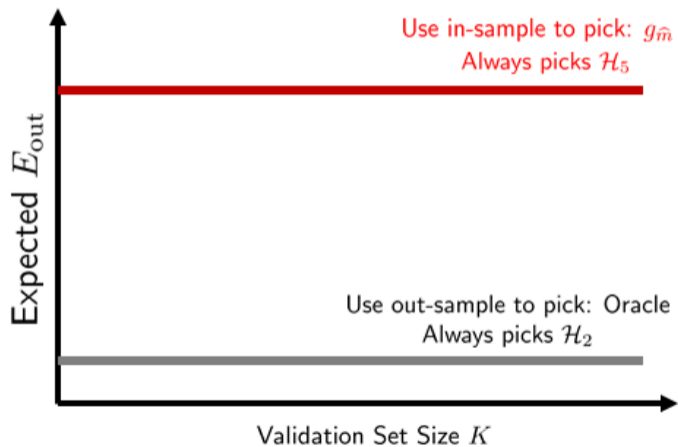
# Validation



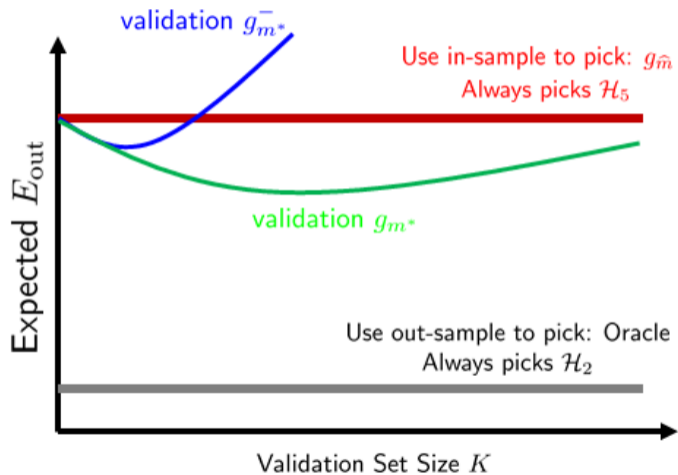
# Validation



## Expected Error



## Expected Error





# Observations

## Validation and $N - K$ samples for training:

- $\mathbb{E}[E_{out}(g_{m^*}^-)]$  drops and then rise.
- Compared to in-sample,  $\mathbb{E}[E_{out}(g_{m^*}^-)]$  uses a few samples to validate.
- This gives a good estimate of out-sample error.
- As  $K$  increases, the estimate improves. So  $\mathbb{E}[E_{out}(g_{m^*}^-)]$  drops.
- If  $K$  is too large, then only  $N - K$  samples for training.
- Poor training makes  $\mathbb{E}[E_{out}(g_{m^*}^-)]$  rise.

## Validation and $N$ samples for training:

- $\mathbb{E}[E_{out}(g_{m^*})]$  will be lower.
- Because you have chosen the best.

Therefore, you should always recycle the validation data for training the final hypothesis.

# Outline

- Lecture 31 Overfit
- Lecture 32 Regularization
- **Lecture 33 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  $(\mathbf{x}_n, y_n)$ .

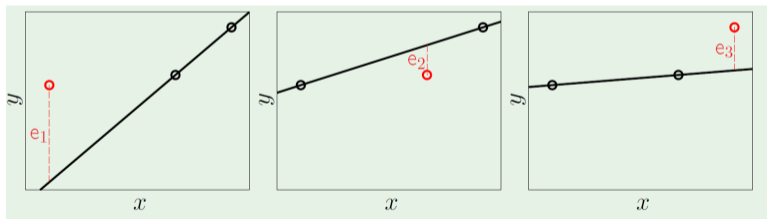
# Cross Validation

- 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  $\lambda$ ?
- Let

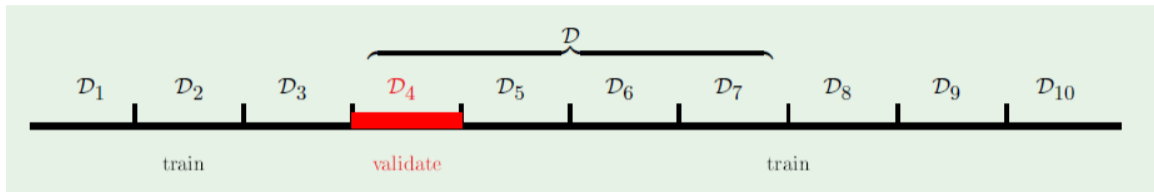
$$\begin{aligned} \mathbf{H}(\lambda) &= \mathbf{A}(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \\ \hat{\mathbf{y}} &= \mathbf{H} \mathbf{y} \end{aligned}$$

- 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  $\lambda$  that minimizes  $E_{\text{cv}}$

## 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  $E_{\text{val}}$   
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.