

# Evaluation the Learning Algorithm

Prof. Dr. Christina Bauer

[christina.bauer@th-deg.de](mailto:christina.bauer@th-deg.de)

Faculty of Computer Science

# DEBUGGING

---

- Example:
    - E.g. house price prediction
    - $J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$
  - Large errors in prediction → What to try next???
- Get more training examples
  - Try smaller set of features (to avoid overfitting)
  - Try getting additional features
  - Try adding polynomial features (e.g.  $x_1^2$ ,  $x_2^2$ ,  $x_1 * x_2$ )
  - Try decreasing  $\lambda$
  - Try increasing  $\lambda$

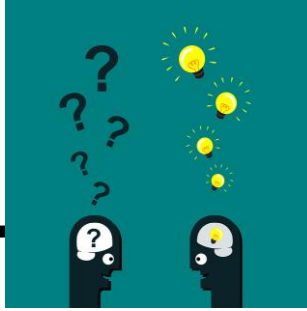
# DEBUGGING

---

- Machine Learning Diagnostics
- Diagnostic: A test that you can run to get insight what is or is not working with a learning algorithm, and gain guidance as to how best improve its performance.
- Diagnostic can take time to implement, but can help a lot.

# QUESTION

---



Which of the following statements about diagnostics are true? Check all that apply.

A: It's hard to tell what will work to improve a learning algorithm, so the best approach is to go with gut feeling and just see what works.

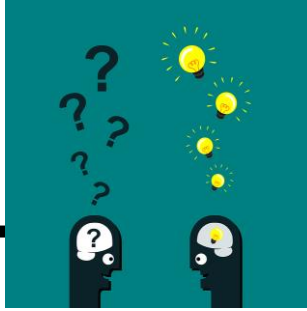
B: Diagnostics can give guidance as to what might be more fruitful things to try to improve a learning algorithm.

C: Diagnostics can be time-consuming to implement and try, but they can still be a very good use of your time.

D: A diagnostic can sometimes rule out certain courses of action (changes to your learning algorithm) as being unlikely to improve its performance significantly.

# QUESTION

---



Which of the following statements about diagnostics are true? Check all that apply.

A: It's hard to tell what will work to improve a learning algorithm, so the best approach is to go with gut feeling and just see what works.

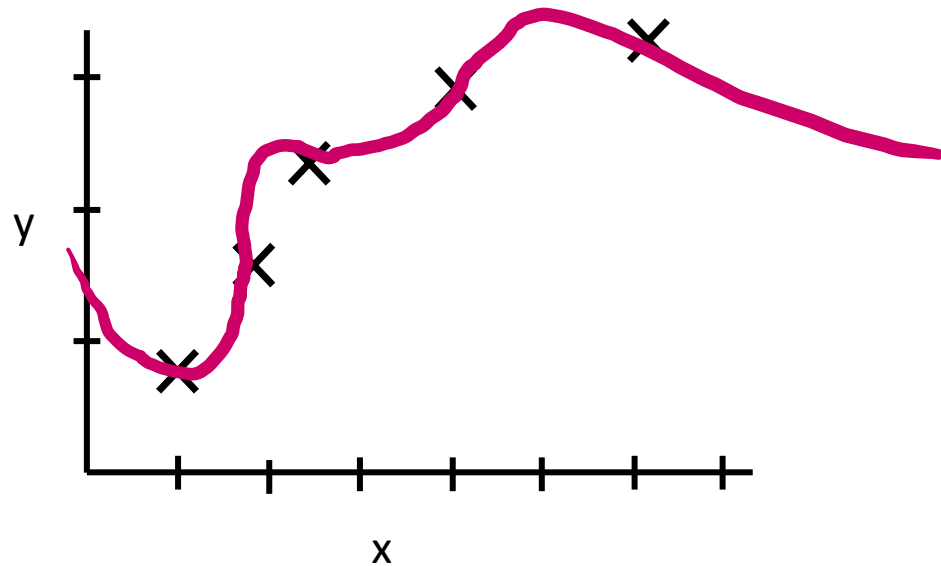
☒ B: Diagnostics can give guidance as to what might be more fruitful things to try to improve a learning algorithm.

☒ C: Diagnostics can be time-consuming to implement and try, but they can still be a very good use of your time.

☒ D: A diagnostic can sometimes rule out certain courses of action (changes to your learning algorithm) as being unlikely to improve its performance significantly.

# EVALUATING A HYPOTHESIS

---



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

→ Overfit + high variance

Fails to generalize to new examples that are not part of the training set.

$X_1$  = size of the house

$X_2$  = no. of bedrooms

$X_3$  = no. of floors

$X_4$  = age of the house

$X_5$  = average income in neighbourhood

$X_6$  = kitchen size

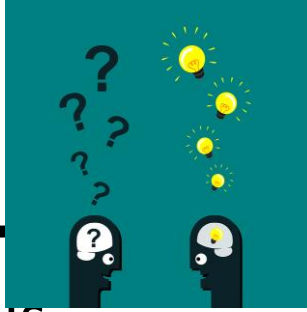
....  
 $X_{100}$

# EVALUATING A HYPOTHESIS

	Size in feet <sup>2</sup>	Price in K \$		
70 %	2104	400	Training set	$(x^{(1)}, y^{(1)})$
	1600	330		$(x^{(2)}, y^{(2)})$
	2400	369		....
	1416	232		$(x^{(m)}, y^{(m)})$
	3000	540		
	1985	300		
30 %	1534	315	Test set	$(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)})$
	1427	199		$(x_{\text{test}}^{(2)}, y_{\text{test}}^{(2)})$
	1380	212		....
	1494	243		$(x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})})$

# QUESTION

---



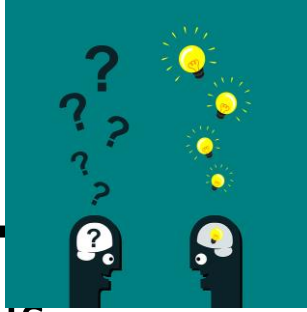
Suppose an implementation of linear regression (without regularization) is badly overfitting the training set. In this case, we would expect:

- A: The training error  $J(\theta)$  to be low and the test error  $J_{\text{test}}(\theta)$  to be high
- B: The training error  $J(\theta)$  to be low and the test error  $J_{\text{test}}(\theta)$  to be low
- C: The training error  $J(\theta)$  to be high and the test error  $J_{\text{test}}(\theta)$  to be low
- D: The training error  $J(\theta)$  to be high and the test error  $J_{\text{test}}(\theta)$  to be high



# QUESTION

---



Suppose an implementation of linear regression (without regularization) is badly overfitting the training set. In this case, we would expect:

- ~~A~~: The training error  $J(\theta)$  to be low and the test error  $J_{\text{test}}(\theta)$  to be high
- B: The training error  $J(\theta)$  to be low and the test error  $J_{\text{test}}(\theta)$  to be low
- C: The training error  $J(\theta)$  to be high and the test error  $J_{\text{test}}(\theta)$  to be low
- D: The training error  $J(\theta)$  to be high and the test error  $J_{\text{test}}(\theta)$  to be high

# EVALUATING A HYPOTHESIS

---

Example linear regression:

- Learn parameter  $\theta$  from training data (e.g. with 70% of data  $\rightarrow$  minimizing training error  $J(\theta)$ )
- Compute test set error

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^i) - y_{test}^i)^2$$

# EVALUATING A HYPOTHESIS

---

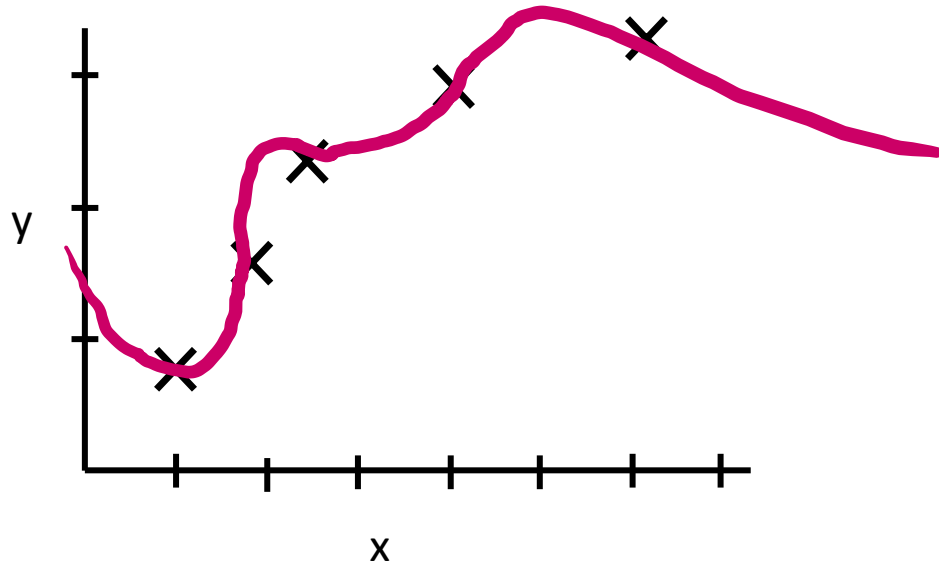
Example logistic regression:

- Learn parameter  $\theta$  from training data
- Compute test set error
- $J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \left[ \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log(h_{\theta}(x_{\text{test}}^{(i)})) + (1 - y_{\text{test}}^{(i)}) \log(1 - h_{\theta}(x_{\text{test}}^{(i)})) \right]$
- Alternative: Misclassification error (0/1 misclassification error)  
 $\text{err}(h_{\theta}(x), y) = 1$  if  $h_{\theta}(x) \geq 0.5$  and  $y = 0$  or  $h_{\theta}(x) < 0.5$  and  $y = 1$  (error case)  
0 otherwise

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

# MODEL SELECTION

---



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

→ Overfit

- Once the parameters were fit to some set of data (training data), the error of the parameters as measured on that data (the training error  $J(\theta)$ ) is likely to be (much) lower than the actual generalization error.

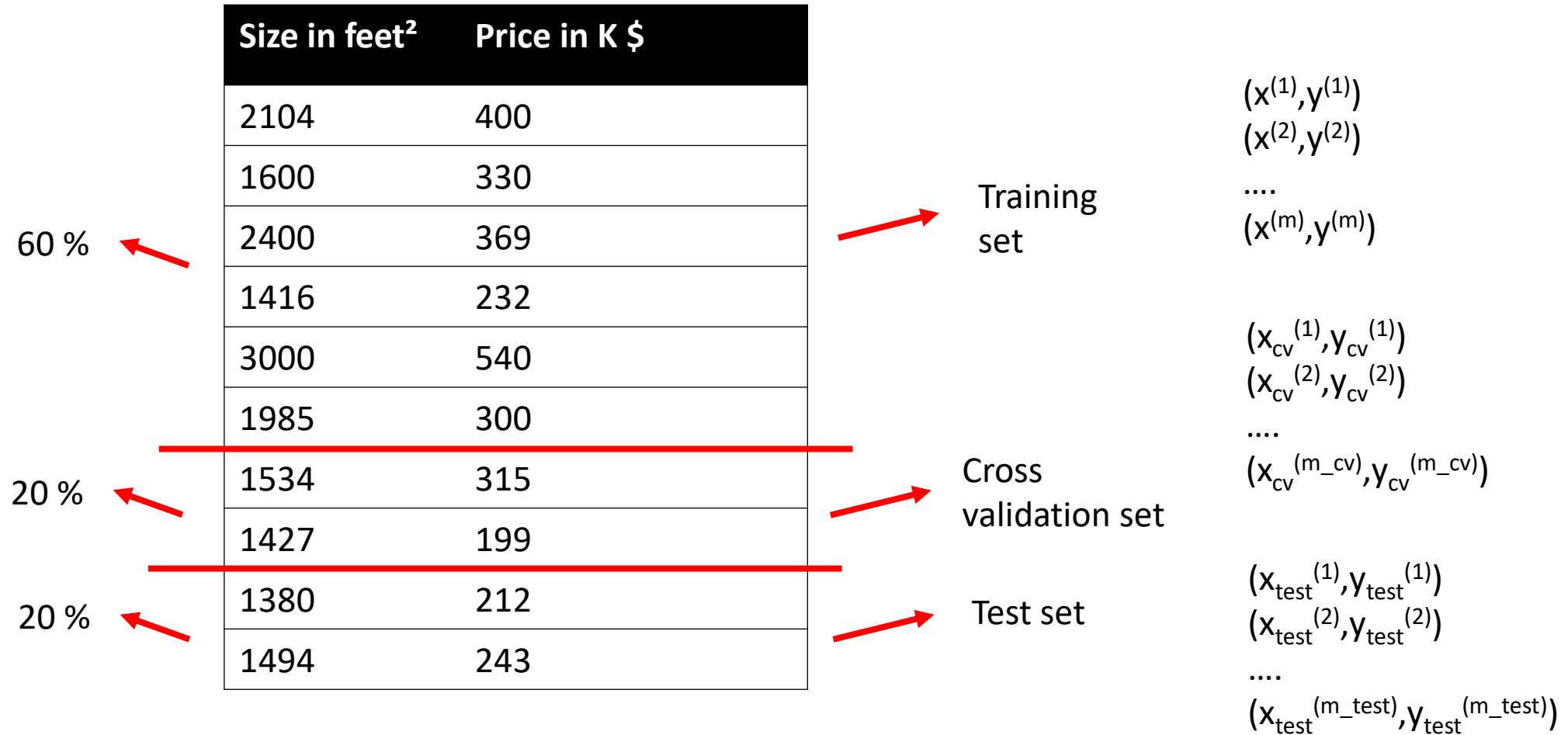
# MODEL SELECTION

---

d = degree of polynomial

- $d=1: h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \theta^1 \rightarrow J_{\text{test}}(\theta^1)$
- $d=2: h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^2 \rightarrow J_{\text{test}}(\theta^2)$
- $d=3: h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \rightarrow \theta^3 \rightarrow J_{\text{test}}(\theta^3)$
- ...
- $d=10: h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{10} \rightarrow J_{\text{test}}(\theta^{10})$
  
- Choose lowest test set error, e.g.  $h_{\theta}(x) = \theta_0 + \dots + \theta_5 x^5$
- Report test set error  $J_{\text{test}}(\theta^5)$
- Problem: d is fit to the test set  $\rightarrow J_{\text{test}}(\theta^5)$  is likely to be an overly optimistic estimate of generalization error. Hypothesis is likely to do better on test set than it would on new examples.

# EVALUATING A HYPOTHESIS



# EVALUATING A HYPOTHESIS

---

- Training error:

- $J_{train}(\theta) \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$

- Cross validation error:

- $J_{cv}(\theta) \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^i) - y_{cv}^i)^2$

- Test error:

- $J_{test}(\theta) \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^i) - y_{test}^i)^2$

# MODEL SELECTION

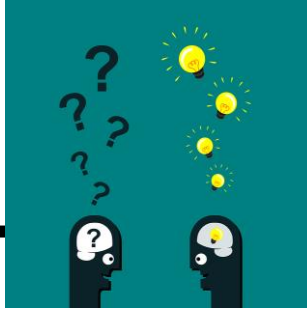
---

- $d=1: h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \theta^1 \rightarrow J_{cv}(\theta^1)$
- $d=2: h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^2 \rightarrow J_{cv}(\theta^2)$
- $d=3: h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \rightarrow \theta^3 \rightarrow J_{cv}(\theta^3)$
- ...
- $d=10: h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{10} \rightarrow J_{cv}(\theta^{10})$
  
- Choose lowest test cross validation error, e.g.  $h_{\theta}(x) = \theta_0 + \dots + \theta_4 x^4$
- Estimate generalization error for test set, e.g.  $J_{test}(\theta^4)$



# QUESTION

---

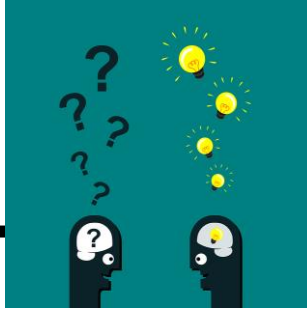


Consider the model selection procedure where we choose the degree of polynomial using a cross validation set. For the final model (with parameters  $\theta$ ), we might generally expect  $J_{cv}(\theta)$  to be lower than  $J_{test}(\theta)$  because:

- A: An extra parameter ( $d$ , the degree of the polynomial) has been fit to the cross validation set.
- B: An extra parameter ( $d$ , the degree of the polynomial) has been fit to the test set.
- C: The cross validation set is usually smaller than the test set.
- D: The cross validation set is usually larger than the test set.

# QUESTION

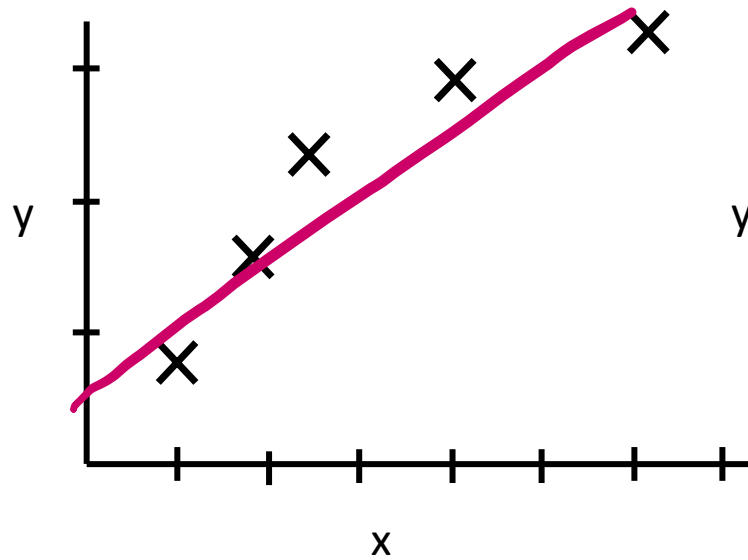
---



Consider the model selection procedure where we choose the degree of polynomial using a cross validation set. For the final model (with parameters  $\theta$ ), we might generally expect  $J_{cv}(\theta)$  to be lower than  $J_{test}(\theta)$  because:

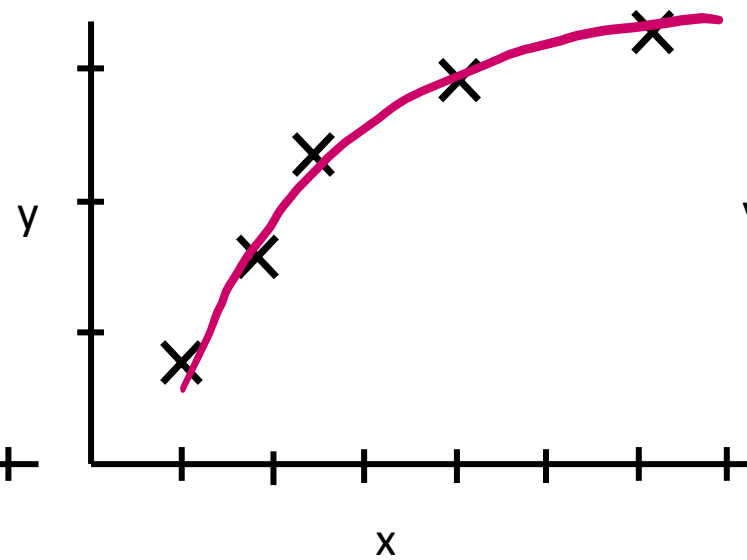
- ~~A~~: An extra parameter ( $d$ , the degree of the polynomial) has been fit to the cross validation set.
- B: An extra parameter ( $d$ , the degree of the polynomial) has been fit to the test set.
- C: The cross validation set is usually smaller than the test set.
- D: The cross validation set is usually larger than the test set.

# BIAS AND VARIANCE

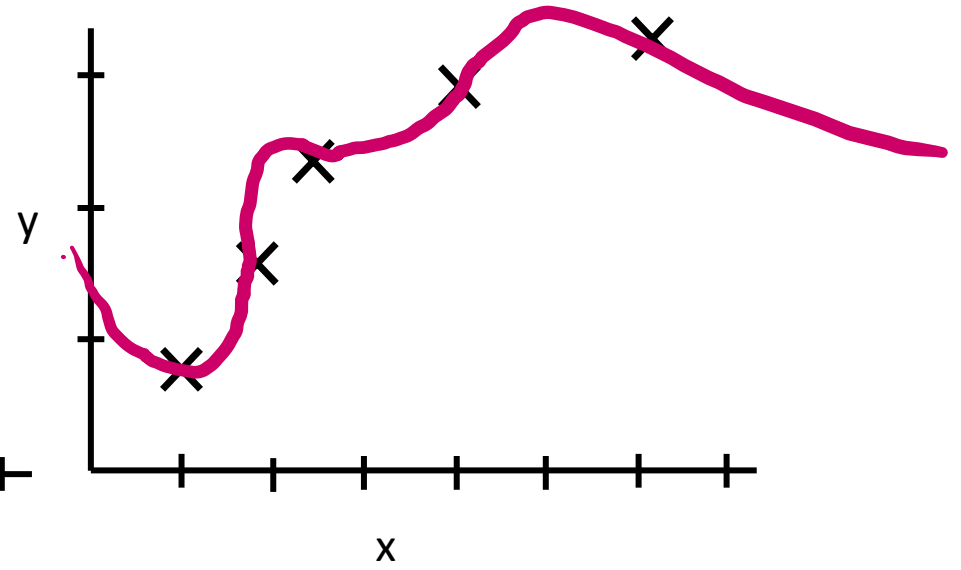


$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

→ Underfit + high bias



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

→ Overfit + high variance

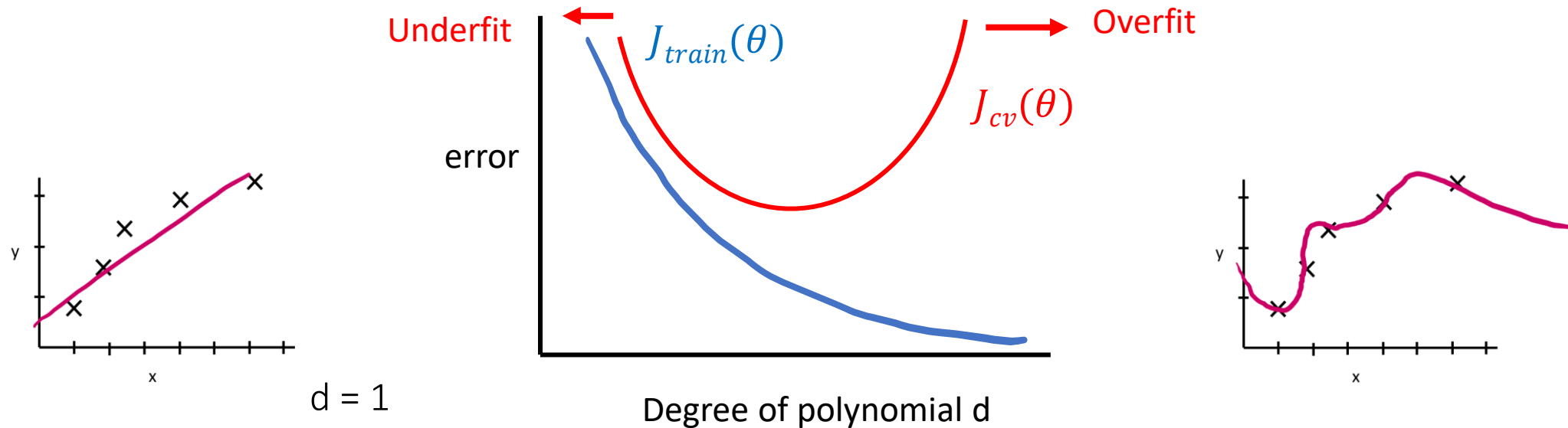
# BIAS AND VARIANCE

- Training error:

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$$

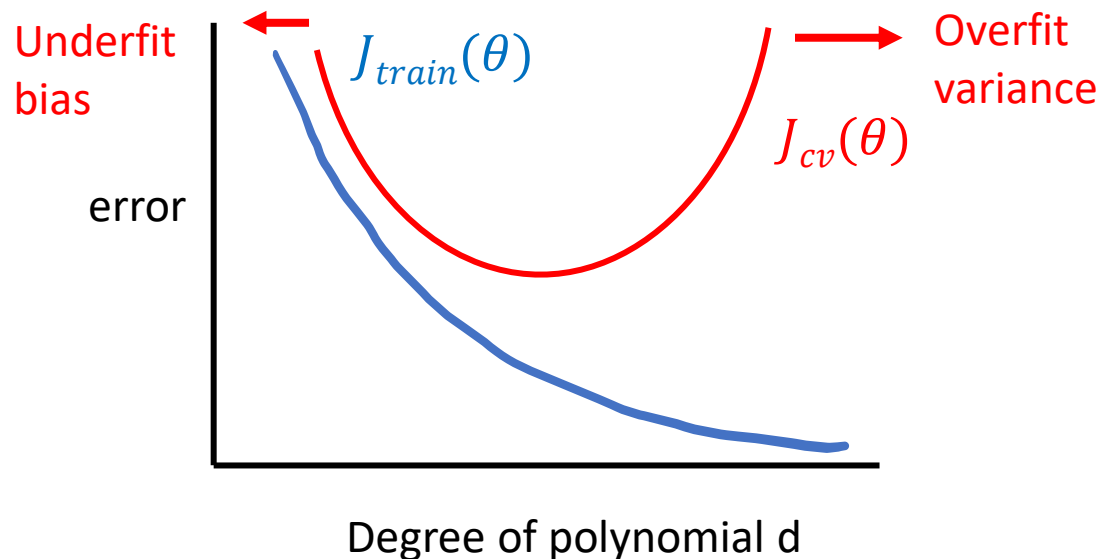
- Cross validation error:

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^i) - y_{cv}^i)^2 \quad (\text{similar for } J_{test}(\theta))$$



# DIAGNOSING BIAS AND VARIANCE

- Suppose your learning algorithm is performing less well than you are expecting (i. e.  $J_{cv}(\theta)$  or  $J_{test}(\theta)$  is high). Is it a bias problem or a variance problem?



**Bias (underfit):**

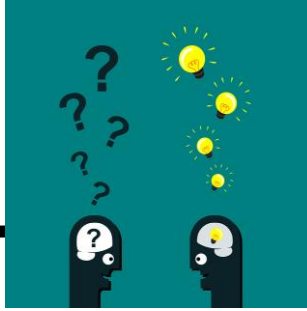
$J_{train}(\theta)$  and  $J_{cv}(\theta)$   
will be high

**Variance (overfit):**

$J_{train}(\theta)$  will be low  
 $J_{cv}(\theta) \gg J_{train}(\theta)$

# QUESTION

---



Suppose you have a classification problem. The (misclassification) error is defined as

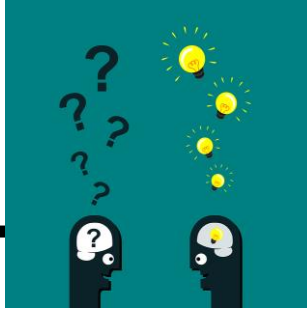
$\frac{1}{m} \sum_{i=1}^m \text{err}(h_{\theta}(x^{(i)}), y^{(i)})$ , and the cross validation (misclassification) error is similarly defined, using the cross validation examples  $(x_{\text{cv}}^{(1)}, y_{\text{cv}}^{(1)}), \dots, (x_{\text{cv}}^{(m_{\text{cv}})}, y_{\text{cv}}^{(m_{\text{cv}})})$

Suppose your training error is 0.10, and your cross validation error is 0.30. What problem is the algorithm most likely to be suffering from?

- A: High bias (overfitting)
- B: High bias (underfitting)
- C: High variance (overfitting)
- D: High variance (underfitting)

# QUESTION

---



Suppose you have a classification problem. The (misclassification) error is defined as

$\frac{1}{m} \sum_{i=1}^m \text{err}(h_{\theta}(x^{(i)}), y^{(i)})$ , and the cross validation (misclassification) error is similarly defined, using the cross validation examples  $(x_{\text{cv}}^{(1)}, y_{\text{cv}}^{(1)}), \dots, (x_{\text{cv}}^{(m_{\text{cv}})}, y_{\text{cv}}^{(m_{\text{cv}})})$

Suppose your training error is 0.10, and your cross validation error is 0.30. What problem is the algorithm most likely to be suffering from?

A: High bias (overfitting)

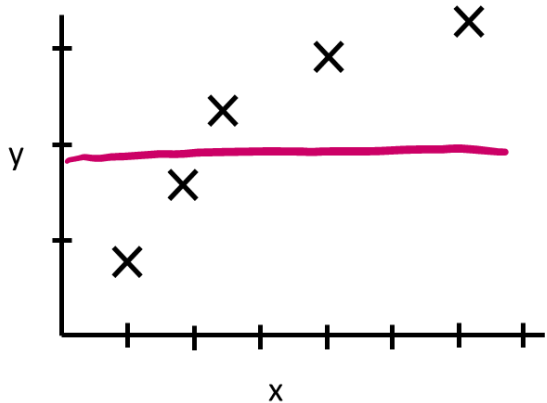
B: High bias (underfitting)

~~C: High variance (overfitting)~~

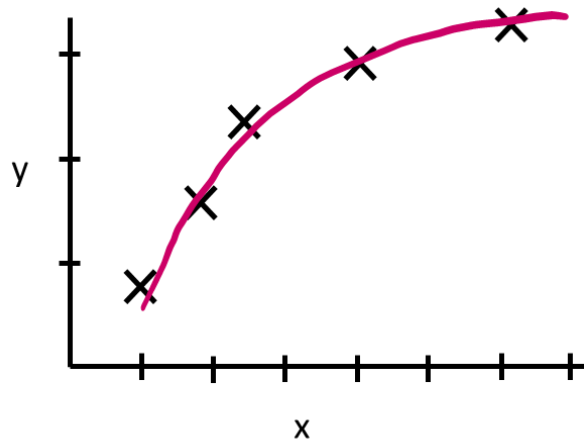
D: High variance (underfitting)

# REGULARIZATION

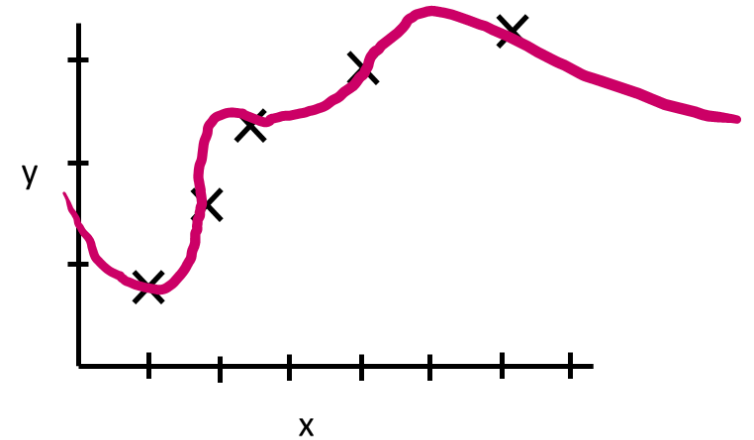
$$\text{Model: } h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \rightarrow J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$$



Large  $\lambda \rightarrow$  Underfit + high bias  
 $\rightarrow \lambda = 10000 \rightarrow \theta_1 \sim 0, \theta_2 \sim 0, ..$   
 $\rightarrow h_{\theta}(x) = \theta$



Intermediate  $\lambda$



Small  $\lambda \rightarrow$  high variance (overfit)  
 $\lambda \sim 0$



# CHOOSING $\lambda$

---

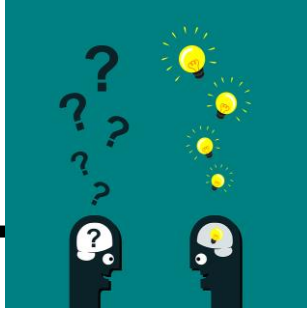
- Hypothesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
- Cost function:  $J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$
- Training error:
  - $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$
- Cross validation error:
  - $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^i) - y_{cv}^i)^2$
- Test error:
  - $J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^i) - y_{test}^i)^2$

# CHOOSING $\lambda$

---

- Hypothesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
  - Cost function:  $J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$
1. Try  $\lambda = 0 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^1 \rightarrow J_{cv}(\theta^1)$
  2. Try  $\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^2 \rightarrow J_{cv}(\theta^2)$
  3. Try  $\lambda = 0.02 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^3 \rightarrow J_{cv}(\theta^3)$
  4. Try  $\lambda = 0.04 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^4 \rightarrow J_{cv}(\theta^4)$
  5. Try  $\lambda = 0.08 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^5 \rightarrow J_{cv}(\theta^5)$
- ...
12. Try  $\lambda = 10 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{12} \rightarrow J_{cv}(\theta^{12})$
- Pick e.g.  $\theta^5$  (as the best results on the cv set)  $\rightarrow$  Test error:  $J_{test}(\theta^5)$

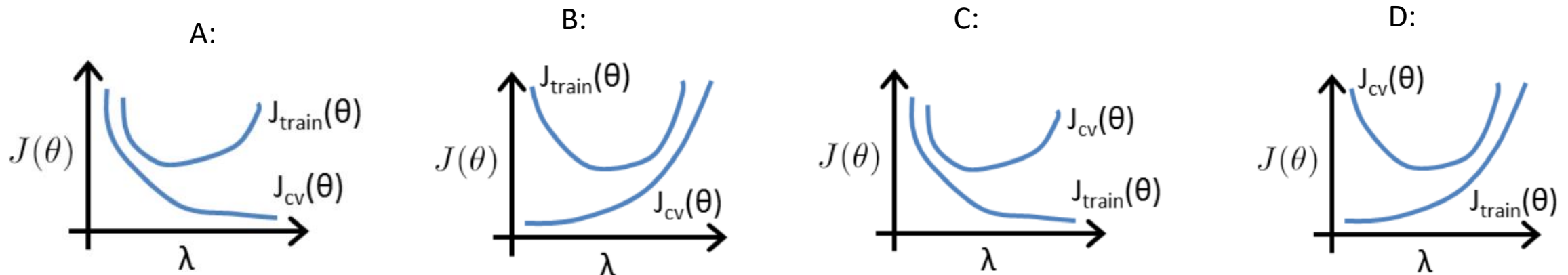
# QUESTION



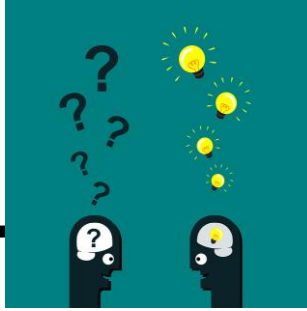
Consider regularized logistic regression. Let

- $J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$
- $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_{\text{train}}^i) - y_{\text{train}}^i)^2$
- $J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^i) - y_{\text{cv}}^i)^2$

Suppose you plot  $J_{\text{train}}$  and  $J_{\text{cv}}$  as a function of the regularization parameter  $\lambda$ . which of the following plots do you expect to get?



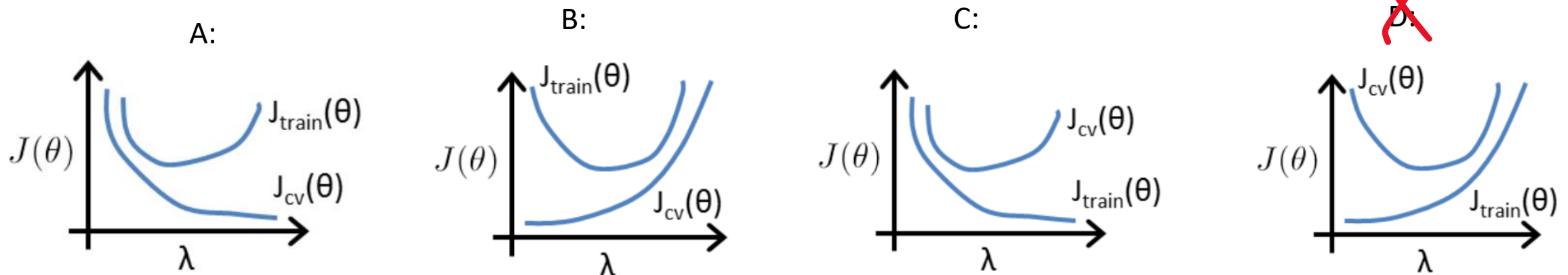
# QUESTION



Consider regularized logistic regression. Let

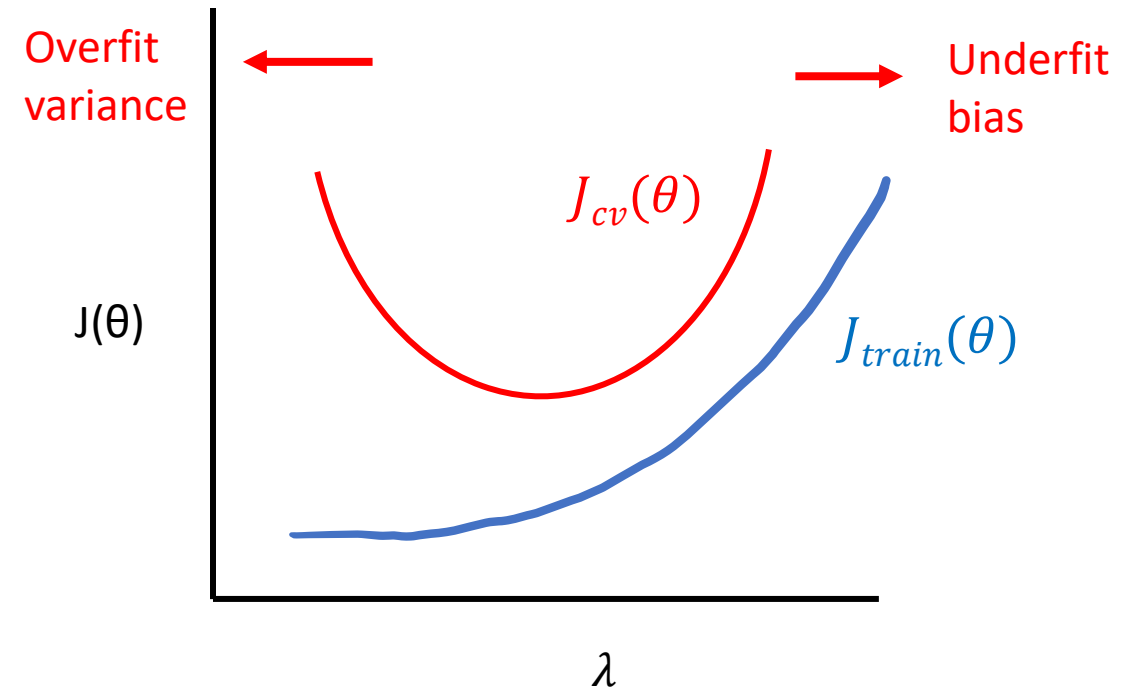
- $J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$
- $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_{train}^i) - y_{train}^i)^2$
- $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^i) - y_{cv}^i)^2$

Suppose you plot  $J_{train}$  and  $J_{cv}$  as a function of the regularization parameter  $\lambda$ . which of the following plots do you expect to get?



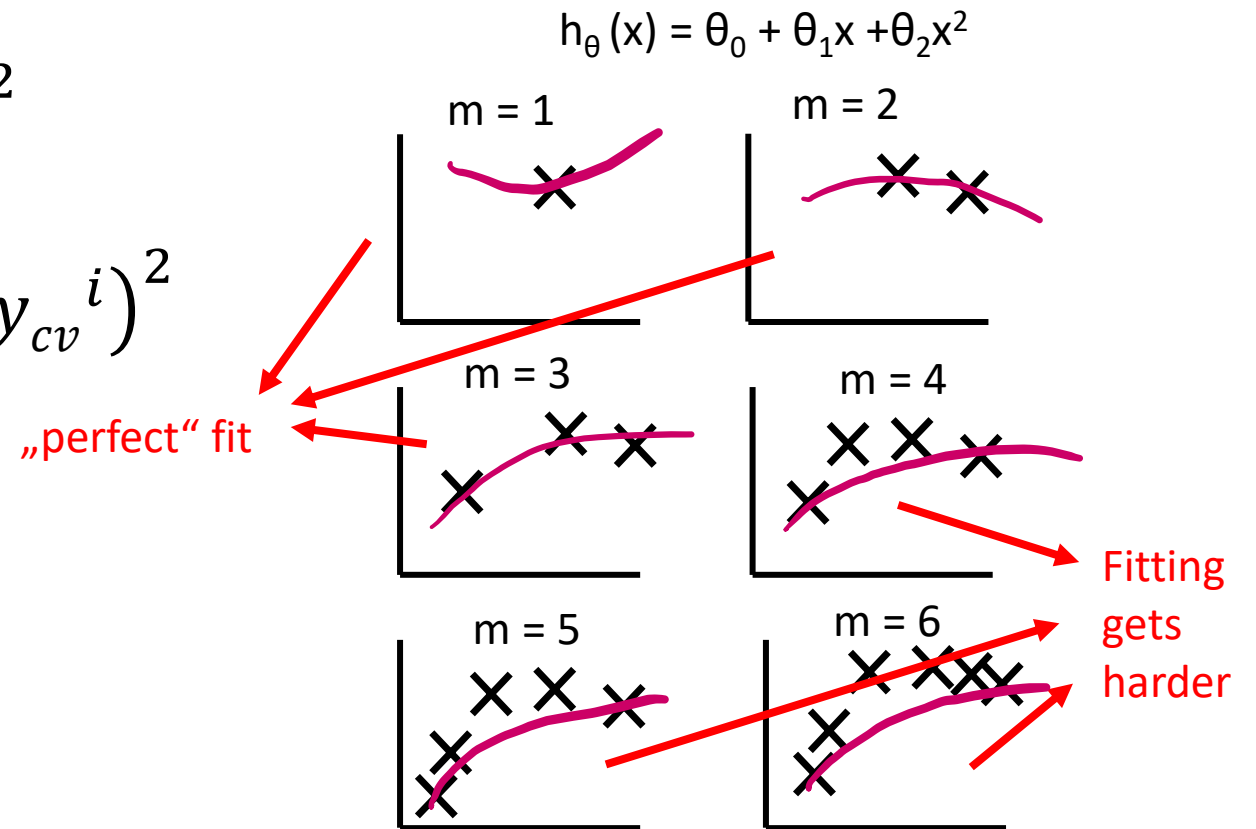
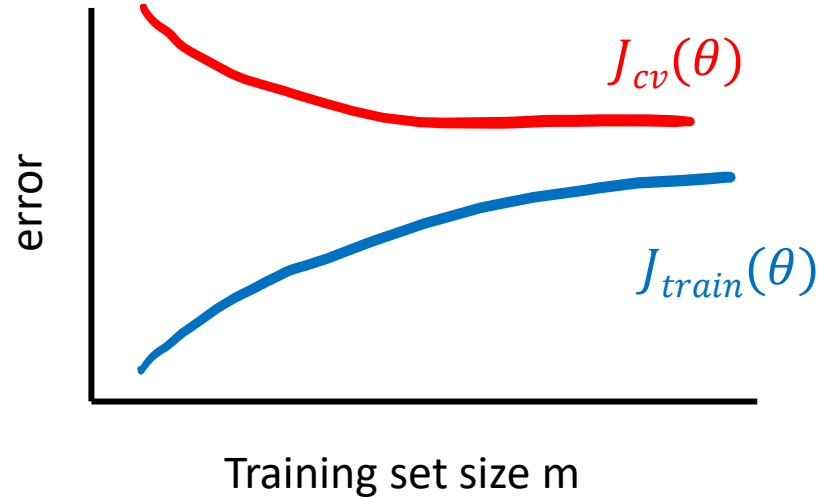
# BIAS AND VARIANCE AS A FUNCTION OF $\lambda$

- $J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$
- $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_{train}^i) - y_{train}^i)^2$
- $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^i) - y_{cv}^i)^2$

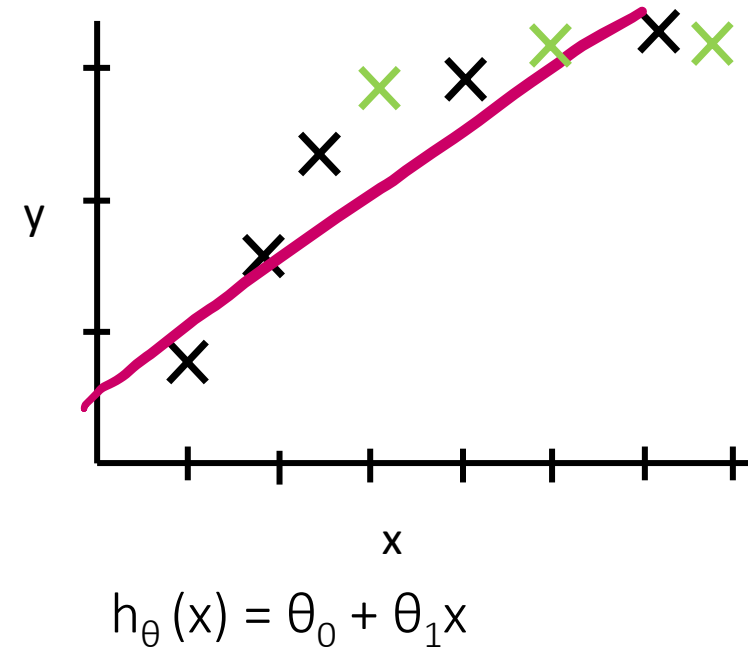
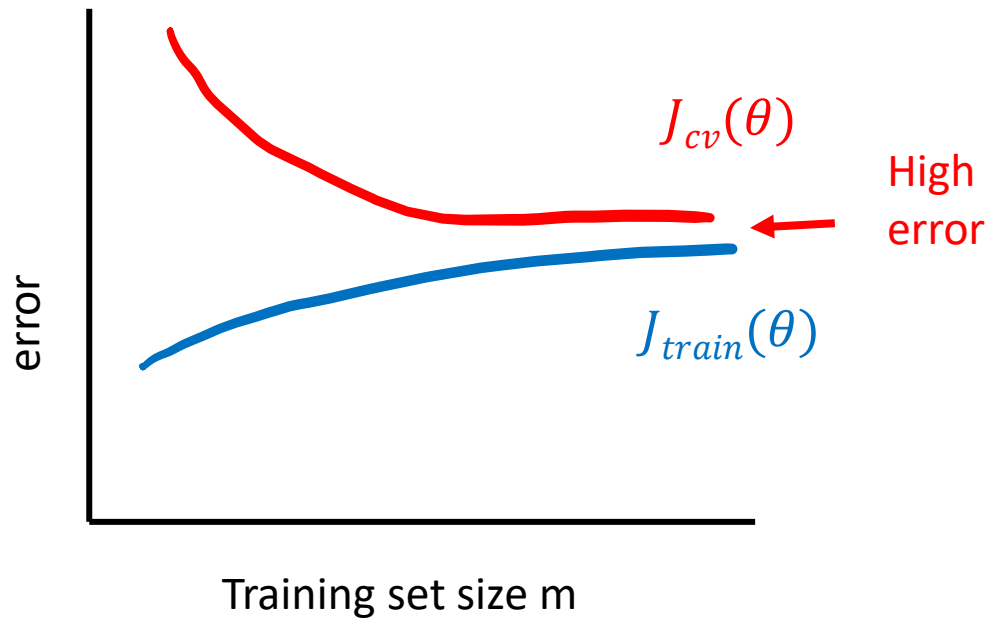


# LEARNING CURVES

- $J_{train}(\theta) \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$
- $J_{cv}(\theta) \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^i) - y_{cv}^i)^2$

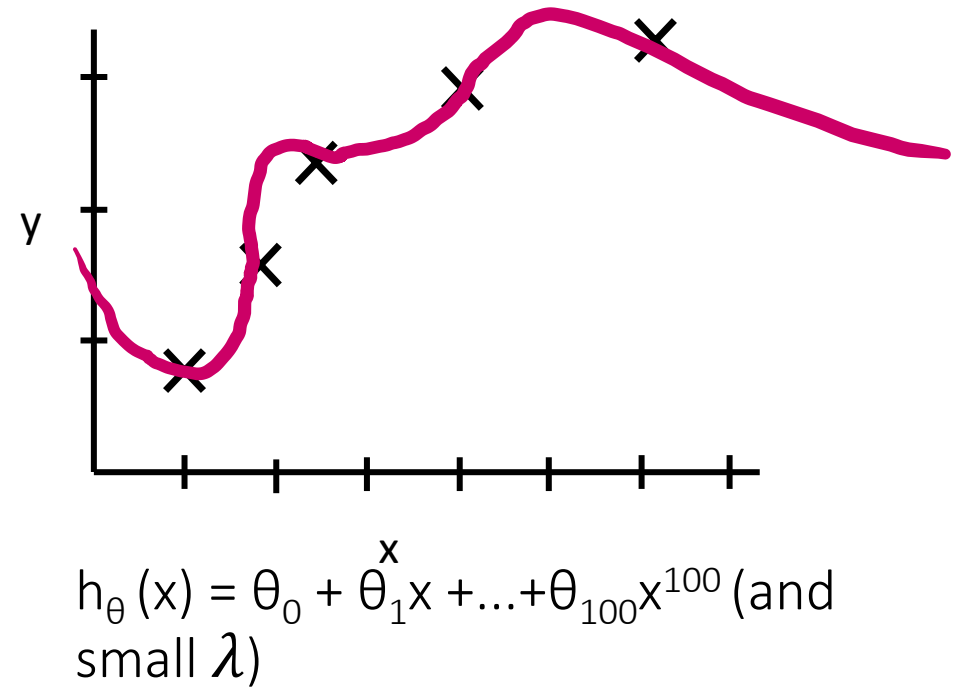
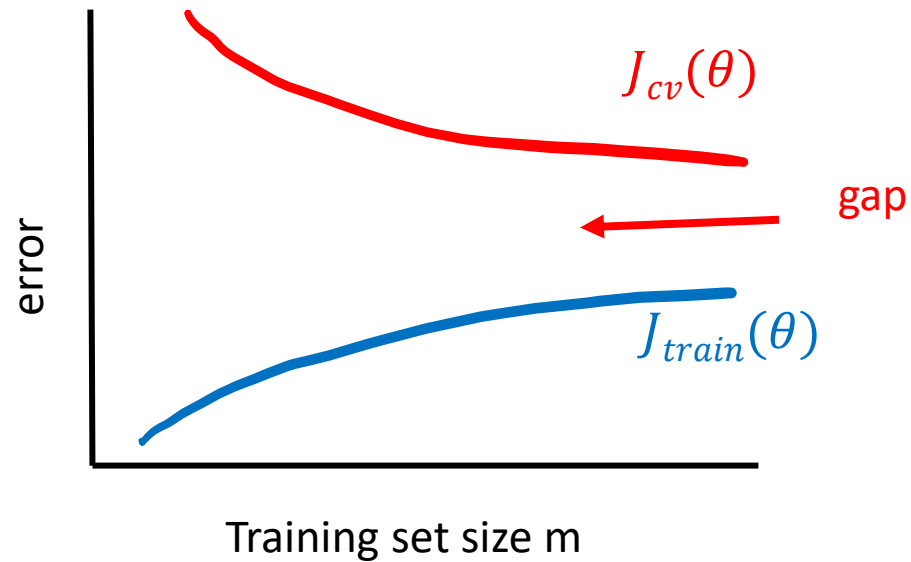


# LEARNING CURVES – HIGH BIAS - UNDERFIT



If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

# LEARNING CURVES – HIGH VARIANCE – OVERFIT

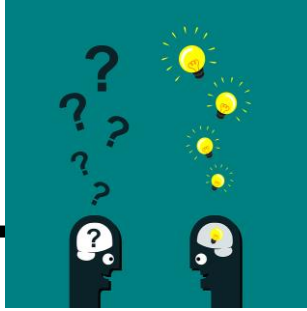


If a learning algorithm is suffering from high variance, getting more training data will likely help.



# QUESTION

---



In which of the following circumstances is getting more training data likely to significantly help a learning algorithm's performance?

A: Algorithm is suffering from high bias.

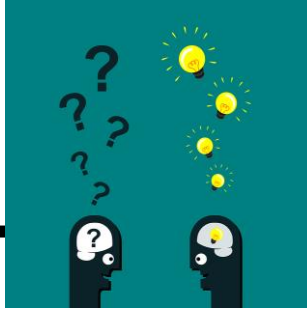
B: Algorithm is suffering from high variance.

C:  $J_{CV}(\theta)$  (cross validation error) is much larger than  $J_{train}(\theta)$  (training error).

D:  $J_{CV}(\theta)$  (cross validation error) is about the same as  $J_{train}(\theta)$  (training error).

# QUESTION

---



In which of the following circumstances is getting more training data likely to significantly help a learning algorithm's performance?

A: Algorithm is suffering from high bias.

~~B: Algorithm is suffering from high variance.~~

~~C:  $J_{cv}(\theta)$  (cross validation error) is much larger than  $J_{train}(\theta)$  (training error).~~

D:  $J_{cv}(\theta)$  (cross validation error) is about the same as  $J_{train}(\theta)$  (training error).

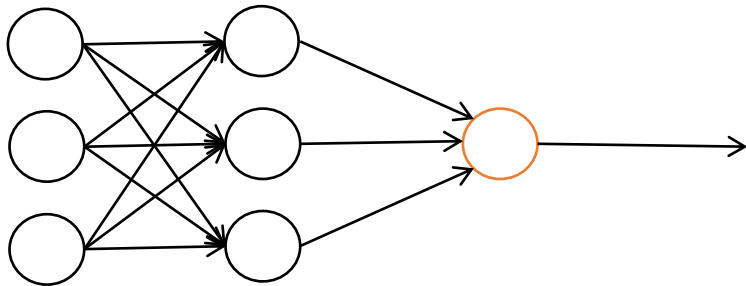
# DEBUGGING

---

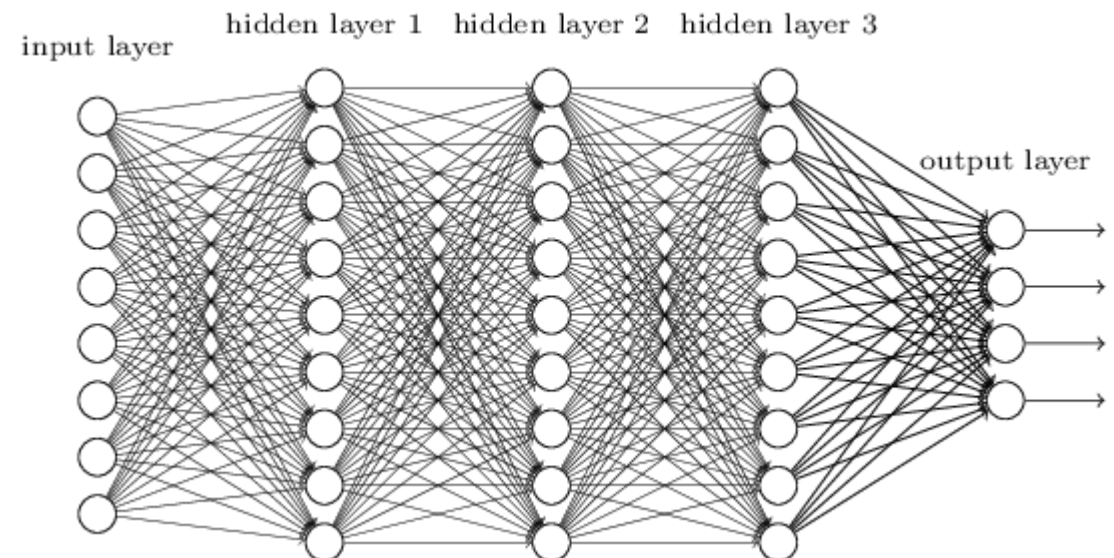
- Large errors in prediction → What to try next???
- Get more training examples → fixes high variance
- Try smaller set of features (to avoid overfitting) → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features (e.g.  $x_1^2$ ,  $x_2^2$ ,  $x_1 * x_2$ ) → fixes high bias
- Try decreasing  $\lambda$  → fixes high bias
- Try increasing  $\lambda$  → fixes high variance

# NEURAL NETWORKS AND FITTING PROBLEMS

- “Small” network
- fewer parameters
- more prone to underfitting
- Computationally cheaper

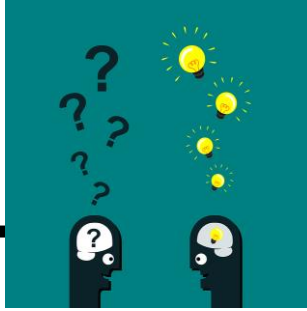


- “Large” network
  - more parameters
  - more prone to overfitting
  - Computationally more expensive
- User regularization to avoid overfitting



# QUESTION

---



Suppose you fit a neural network with one hidden layer to a training set. You find that the cross validation error  $J_{cv}(\theta)$  is much larger than the training error  $J_{train}(\theta)$ . Is increasing the number of hidden units likely to help?

A: Yes, because this increases the number of parameters and lets the network represent more complex functions.

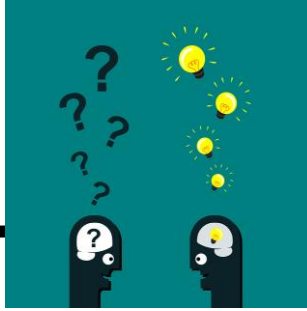
B: Yes, because it is currently suffering from high bias.

C: No, because it is currently suffering from high bias, so adding hidden units is unlikely to help.

D: No, because it is currently suffering from high variance, so adding hidden units is unlikely to help.

# QUESTION

---



Suppose you fit a neural network with one hidden layer to a training set. You find that the cross validation error  $J_{cv}(\theta)$  is much larger than the training error  $J_{train}(\theta)$ . Is increasing the number of hidden units likely to help?

A: Yes, because this increases the number of parameters and lets the network represent more complex functions.

B: Yes, because it is currently suffering from high bias.

C: No, because it is currently suffering from high bias, so adding hidden units is unlikely to help.

~~D: No, because it is currently suffering from high variance, so adding hidden units is unlikely to help.~~

# WRAP-UP

---

## Evaluating a Hypothesis

- A hypothesis may have a low error for the training examples but still be inaccurate (because of overfitting). Thus, to evaluate a hypothesis, given a dataset of training examples, we can split up the data into two sets: a training set and a test set. Typically, the training set consists of 70 % of your data and the test set is the remaining 30 %.
- The new procedure using these two sets is then
  - Learn  $\Theta$  and minimize  $J_{\text{train}}(\Theta)$  using the training set
  - Compute the test set error  $J_{\text{test}}(\Theta)$

# WRAP-UP

---

## Evaluating a Hypothesis

- The test set error

- For linear regression:  $J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^i) - y_{test}^i)^2$

- For classification ~ Misclassification error (aka 0/1 misclassification error):

- $err(h_{\theta}(x), y) = 1$  if  $h_{\theta}(x) \geq 0.5$  and  $y = 0$  or  $h_{\theta}(x) < 0.5$  and  $y = 1$  (error case)

0 otherwise

- This gives us a binary 0 or 1 error result based on a misclassification. The average test error for the test set is:

- Test error =  $\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err(h_{\theta}(x_{test}^{(i)}), y_{test}^{(i)})$

- This gives us the proportion of the test data that was misclassified.



# WRAP-UP

---

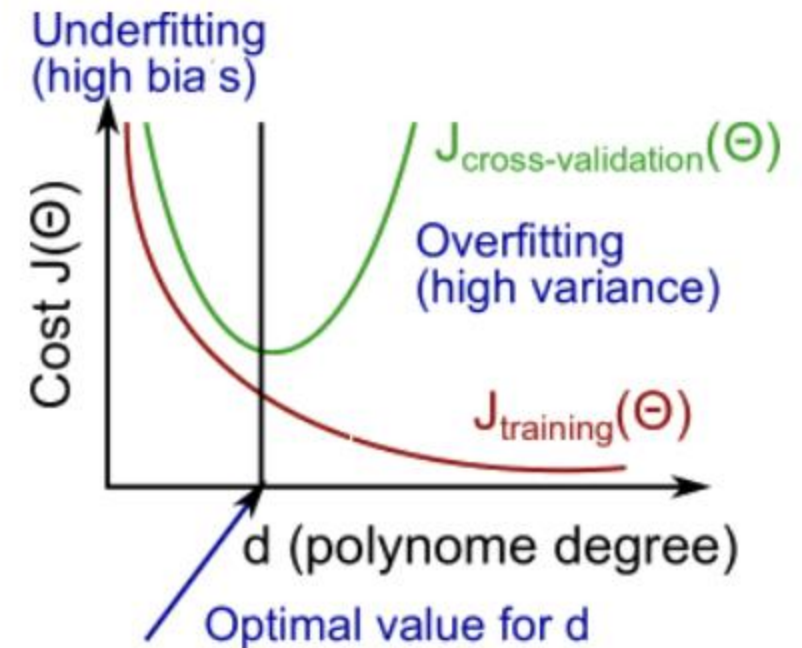
## Model Selection and Train/Validation/Test Sets

- Just because a learning algorithm fits a training set well, that does not mean it is a good hypothesis. It could overfit and as a result your predictions on the test set would be poor. The error of your hypothesis as measured on the data set with which you trained the parameters will be lower than the error on any other data set.
- Given many models with different polynomial degrees, we can use a systematic approach to identify the 'best' function. In order to choose the model of your hypothesis, you can test each degree of polynomial and look at the error result.
- One way to break down our dataset into the three sets is:
  - Training set: 60%
  - Cross validation set: 20%
  - test set: 20%
- We can now calculate three separate error values for the three different sets using the following method:
  - Optimize the parameters in  $\Theta$  using the training set for each polynomial degree.
  - Find the polynomial degree  $d$  with the least error using the cross validation set.
  - Estimate the generalization error using the test set with  $J_{\text{test}}(\Theta^{(d)})$ , ( $d$  = theta from polynomial with lower error);
  - This way, the degree of the polynomial  $d$  has not been trained using the test set.

# WRAP-UP

## Diagnosing Bias vs. Variance

- We have to examine the relationship between the degree of the polynomial  $d$  and the underfitting or overfitting of our hypothesis.
- We need to distinguish whether bias or variance is the problem contributing to bad predictions.
- High bias is underfitting and high variance is overfitting. Ideally, we need to find a golden mean between these two.
- The training error will tend to decrease as we increase the degree  $d$  of the polynomial.
- At the same time, the cross validation error will tend to decrease as we increase  $d$  up to a point, and then it will increase as  $d$  is increased, forming a convex curve.
- **High bias (underfitting):** both  $J_{\text{train}}(\Theta)$  and  $J_{\text{CV}}(\Theta)$  will be high. Also,  $J_{\text{CV}}(\Theta) \approx J_{\text{train}}(\Theta)$ .
- **High variance (overfitting):**  $J_{\text{train}}(\Theta)$  will be low and  $J_{\text{CV}}(\Theta)$  will be much greater than  $J_{\text{train}}(\Theta)$ .



# WRAP-UP

---

## Regularization and Bias/Variance

As  $\lambda$  increases, our fit becomes more rigid. On the other hand, as  $\lambda$  approaches 0, we tend to overfit the data. In order to choose the model and the regularization term  $\lambda$ , we need to:

- Create a list of lambdas (i.e.  $\in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$ )
- Create a set of models with different degrees or any other variants.
- Iterate through the  $\lambda$ s and for each  $\lambda$  go through all the models to learn some  $\Theta$ .
- Compute the cross validation error using the learned  $\Theta$  (computed with  $\lambda$ ) on the  $J_{cv}(\Theta)$  without regularization or  $\lambda = 0$ .
- Select the best combo that produces the lowest error on the cross validation set.
- Using the best combo  $\Theta$  and  $\lambda$ , apply it on  $J_{test}(\Theta)$  to see if it has a good generalization of the problem.

# WRAP-UP

---

## Learning Curves

Training an algorithm on a very few number of data points (such as 1, 2 or 3) will easily have 0 errors because we can always find a e. g. quadratic curve that touches exactly those number of points.

- As the training set gets larger, the error for a quadratic function increases.
- The error value will plateau out after a certain  $m$ , or training set size.

## Experiencing high bias:

- **Low training set size:** causes  $J_{\text{train}}(\Theta)$  to be low and  $J_{\text{CV}}(\Theta)$  to be high.
- **Large training set size:** causes both  $J_{\text{train}}(\Theta)$  and  $J_{\text{CV}}(\Theta)$  to be high with  $J_{\text{train}}(\Theta) \approx J_{\text{CV}}(\Theta)$ .
- If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

## Experiencing high variance:

- **Low training set size:**  $J_{\text{train}}(\Theta)$  will be low and  $J_{\text{CV}}(\Theta)$  will be high.
- **Large training set size:**  $J_{\text{train}}(\Theta)$  increases with training set size and  $J_{\text{CV}}(\Theta)$  continues to decrease without flattening. Also,  $J_{\text{train}}(\Theta) < J_{\text{CV}}(\Theta)$  and the difference between them remains significant.
- If a learning algorithm is suffering from high variance, getting more training data is likely to help.

# WRAP-UP

---

## **Deciding What to Do Next**

- Getting more training examples: Fixes high variance
- Trying smaller sets of features: Fixes high variance
- Adding features: Fixes high bias
- Adding polynomial features: Fixes high bias
- Decreasing  $\lambda$ : Fixes high bias
- Increasing  $\lambda$ : Fixes high variance

# WRAP-UP

---

## Diagnosing Neural Networks

- A neural network with fewer parameters is prone to underfitting. It is also computationally cheaper.
- A large neural network with more parameters is prone to overfitting. It is also computationally expensive. In this case you can use regularization (increase  $\lambda$ ) to address the overfitting.
- Using a single hidden layer is a good starting default. You can train your neural network on a number of hidden layers using your cross validation set. You can then select the one that performs best.

## Model Complexity Effects:

- Lower-order polynomials (low model complexity) have high bias and low variance. In this case, the model fits poorly consistently.
- Higher-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance.
- In reality, we would want to choose a model somewhere in between, that can generalize well but also fits the data reasonably well.

# QUIZ - QUESTION 1

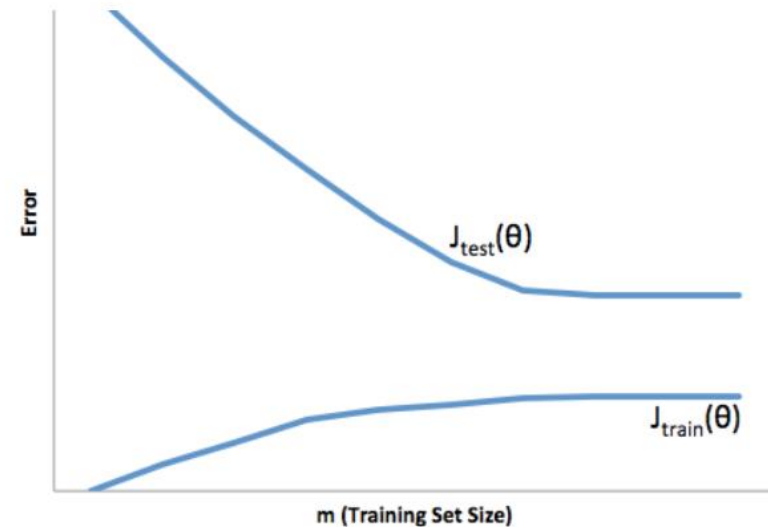
---

You train a learning algorithm and find that it has unacceptably high error on the test set. You plot the learning curve and obtain the figure below. Is the algorithm suffering from high bias, high variance, or neither?

A: High bias

B: Neither

C: High variance



# QUIZ - QUESTION 1

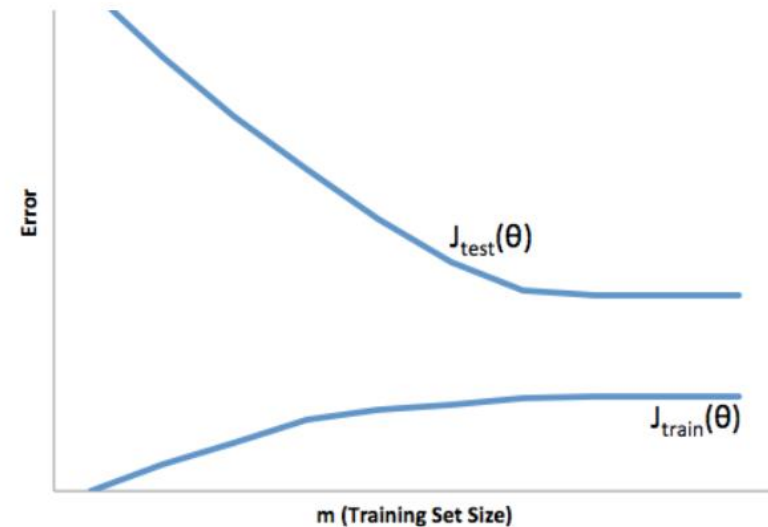
---

You train a learning algorithm and find that it has unacceptably high error on the test set. You plot the learning curve and obtain the figure below. Is the algorithm suffering from high bias, high variance, or neither?

A: High bias

B: Neither

~~C: High variance~~





# QUIZ - QUESTION 2

---

Suppose you have implemented regularized logistic regression to classify what object is in an image (i.e., to do object recognition). However, when you test your hypothesis on a new set of images, you find that it makes unacceptably large errors with its predictions on the new images. However, your hypothesis performs well (has low error) on the training set. Which of the following are promising steps to take? Check all that apply.

A: Try evaluating the hypothesis on a cross validation set rather than the test set.

B: Try using a smaller set of features.

C: Try decreasing the regularization parameter  $\lambda$ .

D: Try increasing the regularization parameter  $\lambda$ .

# QUIZ - QUESTION 2

---

Suppose you have implemented regularized logistic regression to classify what object is in an image (i.e., to do object recognition). However, when you test your hypothesis on a new set of images, you find that it makes unacceptably large errors with its predictions on the new images. However, your hypothesis performs well (has low error) on the training set. Which of the following are promising steps to take? Check all that apply.

A: Try evaluating the hypothesis on a cross validation set rather than the test set.

☒ B: Try using a smaller set of features.

C: Try decreasing the regularization parameter  $\lambda$ .

☒ D: Try increasing the regularization parameter  $\lambda$ .

# QUIZ - QUESTION 3

---

Suppose you have implemented regularized logistic regression to predict what items customers will purchase on a web shopping site. However, when you test your hypothesis on a new set of customers, you find that it makes unacceptably large errors in its predictions. Furthermore, the hypothesis performs poorly on the training set. Which of the following might be promising steps to take? Check all that apply.

- A: Try decreasing the regularization parameter  $\lambda$ .
- B: Try adding polynomial features.
- C: Try evaluating the hypothesis on a cross validation set rather than the test set.
- D: Use fewer training examples.

# QUIZ - QUESTION 3

---

Suppose you have implemented regularized logistic regression to predict what items customers will purchase on a web shopping site. However, when you test your hypothesis on a new set of customers, you find that it makes unacceptably large errors in its predictions. Furthermore, the hypothesis performs poorly on the training set. Which of the following might be promising steps to take? Check all that apply.

☒ A: Try decreasing the regularization parameter  $\lambda$ .

☒ B: Try adding polynomial features.

☐ C: Try evaluating the hypothesis on a cross validation set rather than the test set.

☐ D: Use fewer training examples.

# QUIZ - QUESTION 4

---

Which of the following statements are true? Check all that apply.

A: Suppose you are training a logistic regression classifier using polynomial features and want to select what degree polynomial to use. After training the classifier on the entire training set, you decide to use a subset of the training examples as a validation set. This will work just as well as having a validation set that is separate (disjoint) from the training set.

B: It is okay to use data from the test set to choose the regularization parameter  $\lambda$ , but not the model parameters ( $\theta$ ).

C: A typical split of a data set into training, validation and test sets might be 60% training set, 20% validation set, and 20% test set.

D: Suppose you are using linear regression to predict housing prices, and your dataset comes sorted in order of increasing sizes of houses. It is then important to randomly shuffle the dataset before splitting it into training, validation and test sets, so that we don't have all the smallest houses going into the training set, and all the largest houses going into the test set.

# QUIZ - QUESTION 4

---

Which of the following statements are true? Check all that apply.

A: Suppose you are training a logistic regression classifier using polynomial features and want to select what degree polynomial to use. After training the classifier on the entire training set, you decide to use a subset of the training examples as a validation set. This will work just as well as having a validation set that is separate (disjoint) from the training set.

B: It is okay to use data from the test set to choose the regularization parameter  $\lambda$ , but not the model parameters ( $\theta$ ).

☒ C: A typical split of a data set into training, validation and test sets might be 60% training set, 20% validation set, and 20% test set.

☒ D: Suppose you are using linear regression to predict housing prices, and your dataset comes sorted in order of increasing sizes of houses. It is then important to randomly shuffle the dataset before splitting it into training, validation and test sets, so that we don't have all the smallest houses going into the training set, and all the largest houses going into the test set.

# QUIZ - QUESTION 5

---

Which of the following statements are true? Check all that apply.

A: A model with more parameters is more prone to overfitting and typically has higher variance.

B: If the training and test errors are about the same (but high), adding more features will not help improve the results.

C: If a learning algorithm is suffering from high variance, adding more training examples is likely to improve the test error.

D: If a learning algorithm is suffering from high bias, only adding more training examples may not improve the test error significantly.

# QUIZ - QUESTION 5

---

Which of the following statements are true? Check all that apply.

- ☒ A: A model with more parameters is more prone to overfitting and typically has higher variance.
- ☒ B: If the training and test errors are about the same (but high), adding more features will not help improve the results.
- ☒ C: If a learning algorithm is suffering from high variance, adding more training examples is likely to improve the test error.
- ☒ D: If a learning algorithm is suffering from high bias, only adding more training examples may not improve the test error significantly.