

Stanford ML

WEEK 6:

DEBUGGING

Suppose you have implemented regularized linear regression to predict housing prices.

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

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

Things to try:

- More training examples
- Try smaller set of features
- Try getting additional features
- Adding polynomial features
- Changing λ

How to decide
what to do?

ML DIAGNOSTIC: A test that you can run to identify what is/ isn't working with a learning algorithm, and gain guidance on how best to improve its performance

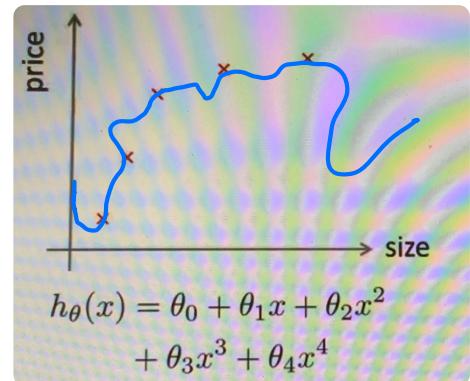
↳ MAY TAKE A LONG TIME BUT WORTH IT!

DIAGNOSTIC 1: Evaluating your Hypothesis

- Low training error may not be a good indication due to **OVERTFITTING** \Rightarrow FAILS TO GENERALIZE

SOLN: For a low # of features we could simply plot $h(\theta)$

→ NOT FEASIBLE w/ LARGER DATASETS



ALTERNATIVE:

- Split dataset into a TRAINING set and TEST set split ~ 70:30%.
- If there is an order within the set, RANDOMIZE test and training set

PROCEDURE: LIN REG.

- 1.) Learn parameters θ from training data (by minimizing $J(\theta)$)
- 2.) Compute test set error

$$J_{\text{test}}(\theta) = \frac{1}{m_t} \sum_{i=1}^{m_t} \{ h_\theta(x_{\text{test}}^{(i)}) - y^{(i)} \}^2$$

LIN.
REG.

PROCEDURE: LOG. REG.

- 1.) Learn parameters θ from training data (by minimizing $J(\theta)$)
- 2.) Compute test set error

$$J_x(\theta) = -\frac{1}{m} \sum_{i=1}^{m_x} \{ y_x^{(i)} \log(h_\theta(x_x^{(i)})) + (1-y_x^{(i)}) \log(1-h_\theta(x_x^{(i)})) \}$$

Dataset:		
Size	Price	
2104	400	Training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	
1427	199	Test set
1380	212	
1494	243	

OR we can use the MISCLASSIFICATION ERROR

$$\text{err}(h_\theta(x), y) = \begin{cases} 1 & \text{if } h_\theta(x) \geq 0.5, y = 0 \\ & \text{or if } h_\theta(x) \leq 0.5, y = 1 \\ 0 & \text{otherwise} \end{cases}$$

1 MISCLASSIFIED

$$\Rightarrow \text{Test error} = \frac{1}{m_x} \sum_{i=1}^{m_x} \text{err}\{h_\theta(x_t^{(i)}), y^{(i)}\}$$

Model Selection

1. $h_\theta(x) = \theta_0 + \theta_1 x \quad d=1$
2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \quad d=2$ etc
3. $h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3$
- ⋮
10. $h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$

- Almost as if there is another parameter to pick:

d = degree of polynomial

- We could find the parameters using each model generating $\theta(d=1), \theta(d=2), \dots, \theta(d=10)$
 - Calculate $J_t(\theta)$ for $\forall d = 1, \dots, 10$
- \Rightarrow select model w/ lowest J_t e.g. $d=5$

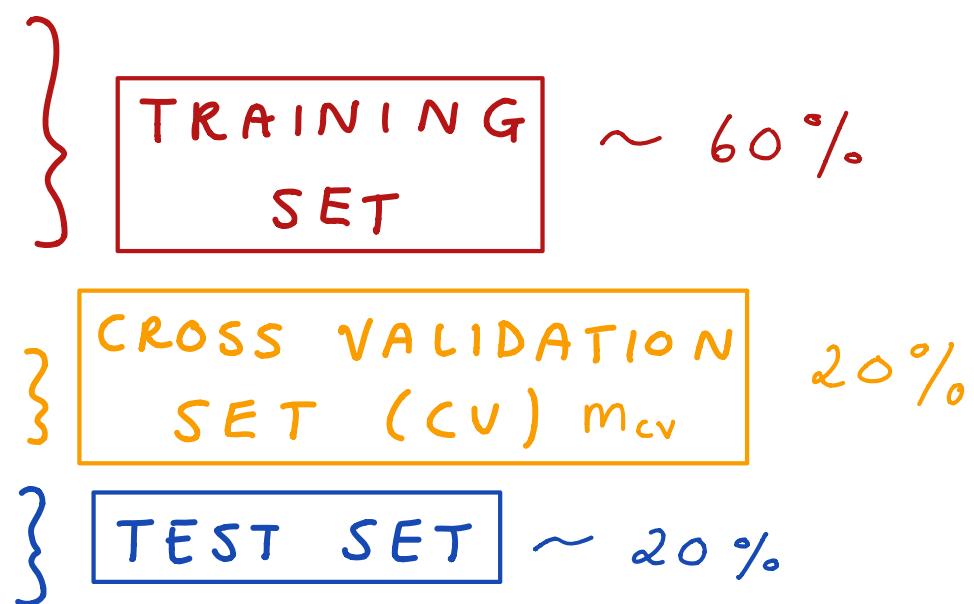
PROBLEM: Even though $J_t(\theta^{(5)})$ is the smallest, this still MAY NOT generalize well and will likely be an **OPTIMISTIC** estimate of generalization error.
i.e. OUR EXTRA PARAMETER WAS FITTED USING THE TEST SET

\Rightarrow BIASED IN FAVOUR OF TEST SET!

ALTERNATIVE: split the data in 3!

Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
1380	212
1494	243



- Instead of using the test set to pick the model we use the CV set

1.) Obtain $\theta^{(1)}, \theta^{(2)}, \dots$

$\theta^{(n)}, \theta^{(10)}$ via $\min_{\theta} J(\theta)$

2.) Compute $J_{cv}(\theta^{(i)})$ for $\forall i = 1, \dots, d$

3.) Pick lowest $J_{cv}(\theta)$ i.e. $\theta^{(4)}$; $\theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4$

4.) Estimate generalization error for test set $J_t(\theta^{(4)})$

Training error:

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

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:

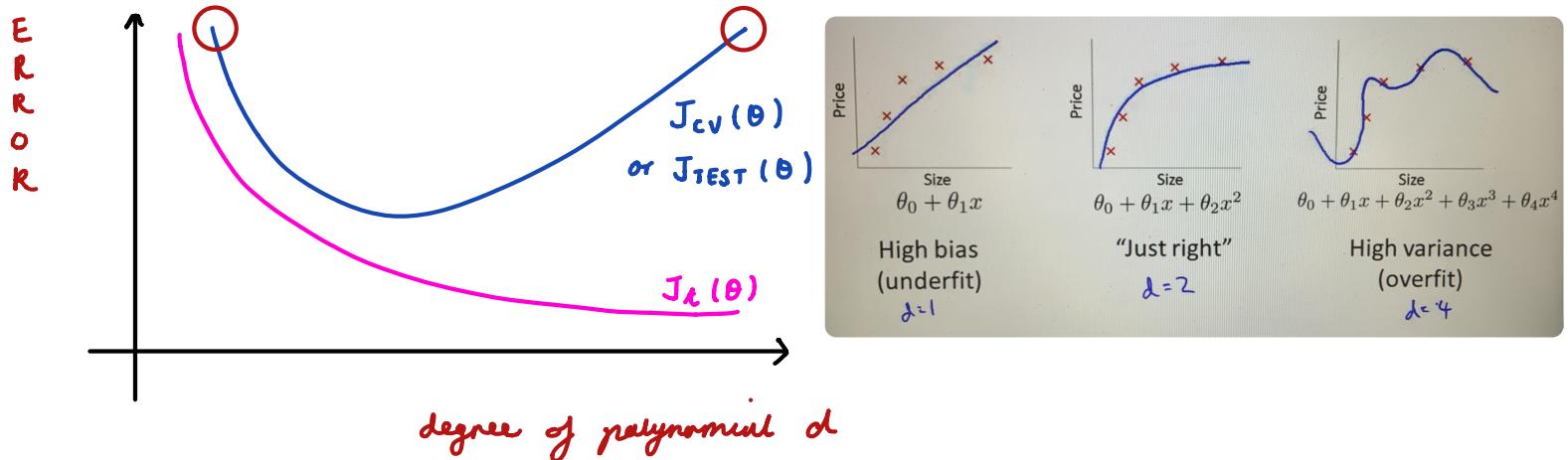
$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Bias vs. Variance

- Problems usually arise due to one or both aspects in play

TRAINING ERROR: $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2$

CROSS-VALIDATION ERROR: $J_{\text{cv}}(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)}]^2$



degree of polynomial d

- Large $J_{\text{cv}}(\theta)$ for $d = 1$ and $d = 10$

(UNDER)

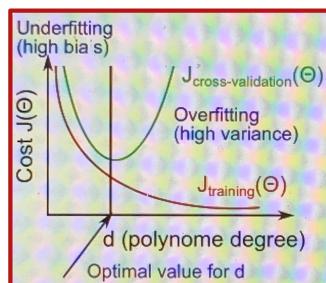
HIGH BIAS!

(OVER)

HIGH VARIANCE!



- BOTH J_{cv} and J_{train} ARE BOTH HIGH
 $J_{\text{cv}}(\theta) \sim J_{\text{train}}(\theta)$



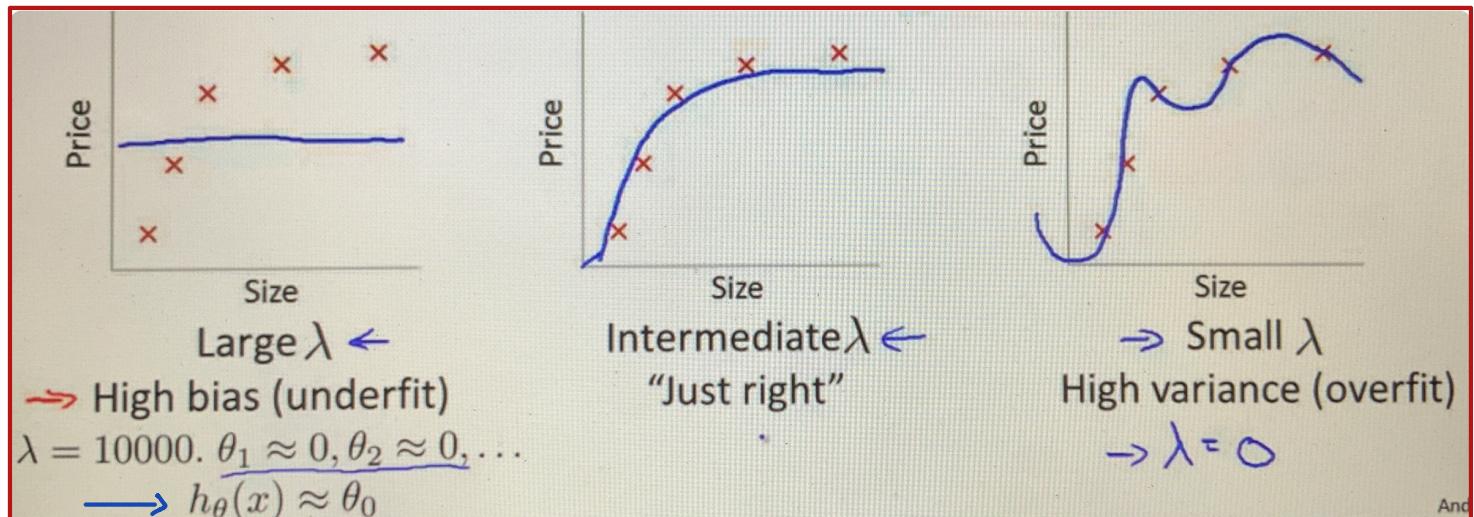
- J_{train} will be LOW but J_{cv} HIGH

$$J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$$

Regularization : BIAS / VARIANCE

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2 + \boxed{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2}$$



Choosing the right λ :

uit

* $J_t(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x_t^{(i)}) - y_t^{(i)}]^2$

same thing!

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

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

* $J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2 + \boxed{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2}$

PROCEDURE:

- 1.) Consider range of λ values (increase in MULTIPLES of 2)
- 2.) $\min_{\theta} J(\theta) \rightarrow \theta$ for each λ value
- 3.) Plug in each θ into $J_{cv}(\theta)$

1. Try $\lambda = 0 \leftarrow$
2. Try $\lambda = 0.01$
3. Try $\lambda = 0.02$
4. Try $\lambda = 0.04$
5. Try $\lambda = 0.08$
- ⋮
12. Try $\lambda = 10 \frac{10}{10.24}$

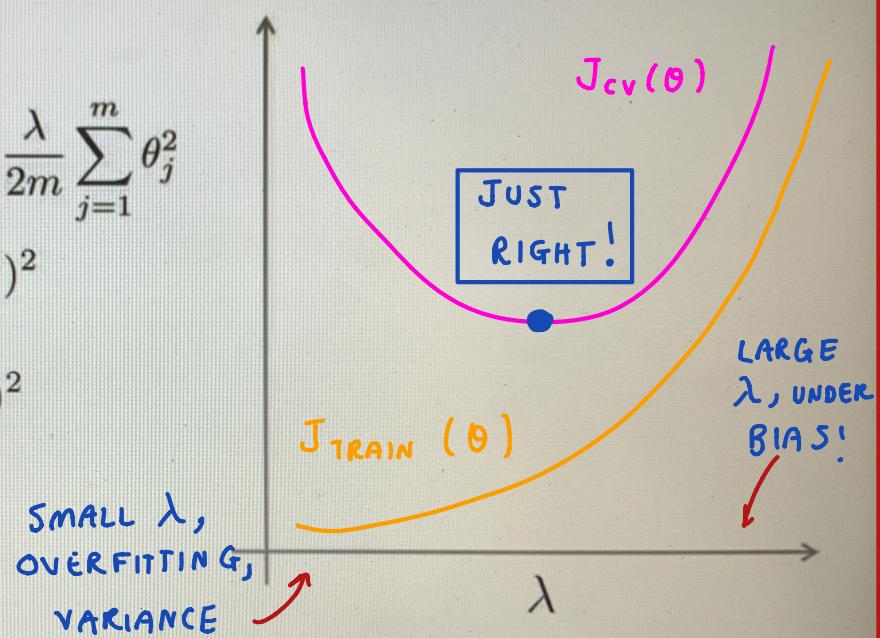
1.) Pick smallest $J_{cv}(\theta)$ i.e $J_{cv}(\theta^{(s)})$

5.) Test error : $J_{test}(\theta^{(s)})$

$$\Rightarrow J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

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

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



VARIANCE

- small λ , OVERFITTING

- $J_{cv}(\theta) \gg J_{train}(\theta)$

BIAS

- large λ , UNDERFITTING

- $J_{cv}(\theta) \approx J_{train}(\theta)$

GENERAL ALGORITHM

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

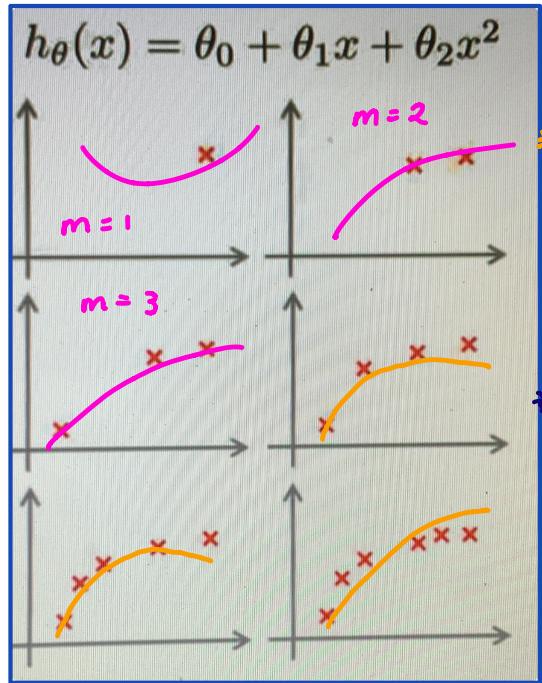
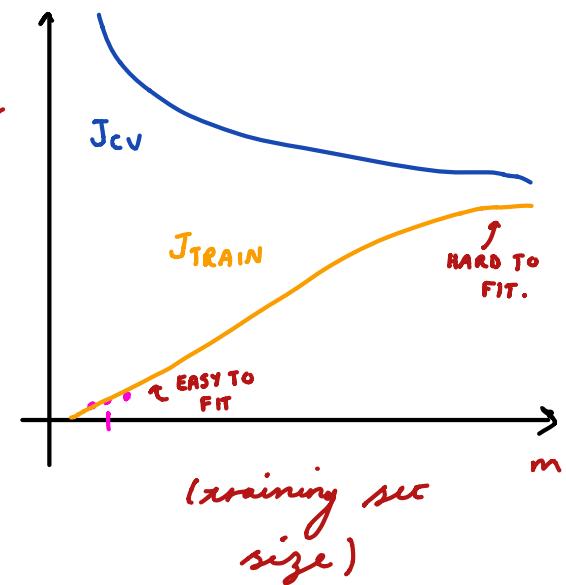
Learning Curves

PLOT:

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

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

- Artificially reduce the training set size and plot above EQNs as $J(m)$.

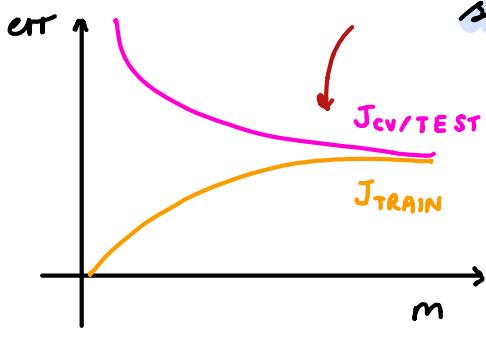


- * $J_{\text{TRAIN}} \approx 0$ w/ REG for low m
- * As m gets larger, it becomes harder to fit a quadratic function to the data
- * Whereas for $J_{\text{CV}}(\theta)$, this shows the opposite trend ::

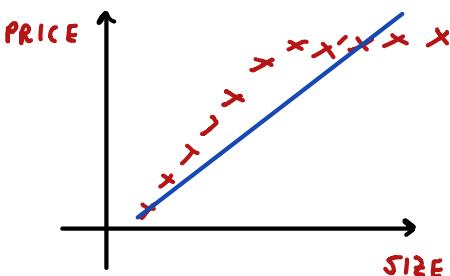
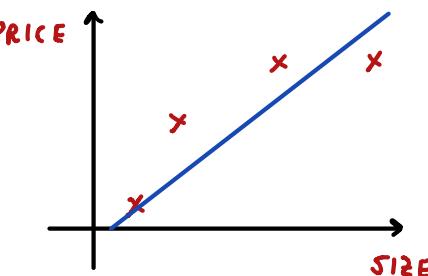
MORE TRAINING EX. \Rightarrow BETTER GENERALIZATION

HIGH BIAS

eventually reaches saturation



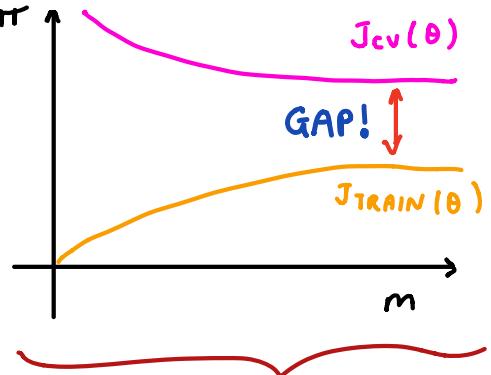
$$h_{\theta} = \theta_0 + \theta_1 x$$



- Even though m has increased,

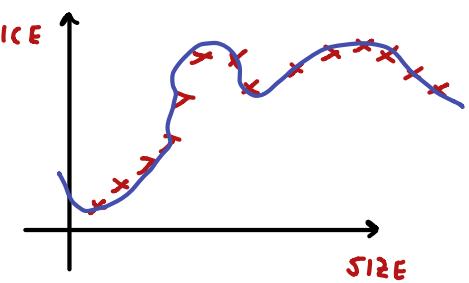
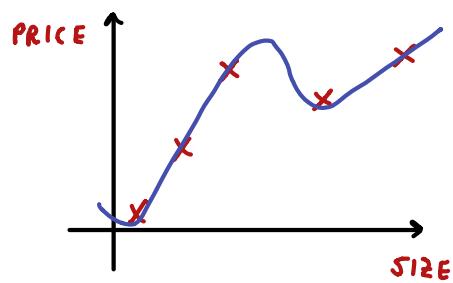
- High errors in general by insisting on a linear hypothesis we get a similar $h_{\theta}(x)$. $\} J_{\text{CV}} \approx J_{\text{TRAIN}}$
- Getting more TRAINING DATA WONT HELP FOR HIGH BIAS

HIGH VARIANCE



(small λ)

$$h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$



- Since we are overfitting on the training set, $J_{cv}(\theta)$ will remain large:

$$\left. \begin{array}{l} \\ \end{array} \right\} J_{cv} \gg J_{train}$$

- As m increases, $J_{cv}(\theta)$ will KEEP DECREASING

∴ **IF OVERFITTING AND J_{cv} IS HIGH,
GET MORE DATA**

SUMMARY :

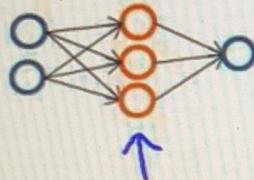
Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features ($x_1^2, x_2^2, x_1 x_2$, etc) → fixes high bias.
- Try decreasing λ → fixes high bias
- Try increasing λ → fixes high variance

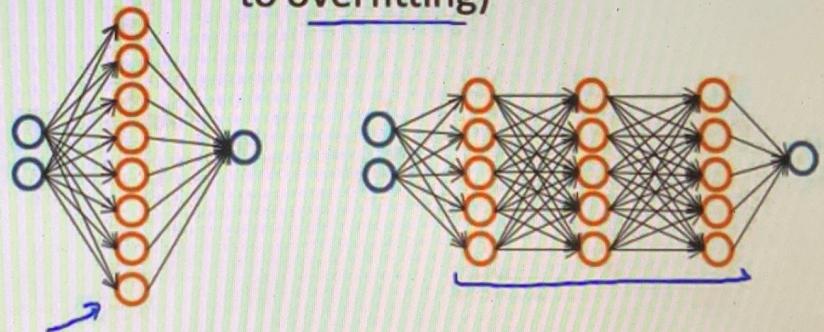
Neural networks and overfitting

→ “Small” neural network
(fewer parameters; more
prone to underfitting)



Computationally cheaper

→ “Large” neural network
(more parameters; more prone
to overfitting)



Computationally more expensive.

Use regularization (λ) to address overfitting.



Machine Learning, Design

EXAMPLE: Spam Classifier

<p>From: cheapsales@buystufffromme.com To: ang@cs.stanford.edu Subject: Buy now!</p> <p>Deal of the week! Buy now! Rolex <u>w4tchs</u> - \$100 <u>Med1cine</u> (any kind) - \$50 Also low cost <u>M0rgages</u> available.</p> <p style="text-align: center; color: red;">Spam (1)</p>	<p>From: Alfred Ng To: ang@cs.stanford.edu Subject: Christmas dates?</p> <p>Hey Andrew, Was talking to Mom about plans for Xmas. When do you get off work. Meet Dec 22? Alf</p> <p style="text-align: center; color: blue;">Non-spam (0)</p>
--	--

SUPERVISED LEARNING

- Let there be a labelled training set with spam $\rightarrow 1$ and non-spam $\rightarrow 0$.
- X = features of e-mail
 - ↳ e.g. deal, buy, discoune, Calum
 - ↳ Take list of ~ 100 words and check whether the word appears in the email
 - ↳ $1 = \text{PRESENT}, 0 = \text{NOT-PRESENT}$

$$X = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{100}; x_j = \begin{cases} 1 & \text{if word } j \text{ present} \\ 0 & \text{otherwise} \end{cases}$$

FEATURE
 VECTOR
 $\in \mathbb{R}^{100}$

NOTE: In practice, take most frequently used n words ($10,000$ to $50,000$) in training set, rather than manually pick 100 words

How to spend your time to minimize error?

- ① Collect lots of data e.g. "honeypot" project
- ② Develop sophisticated features based on c-mail
ROUTING information
- ③ Develop sophisticated features for message body
- ④ Develop algorithm to spot mis-spellings

Error Analysis

Recommended Approach:

- ① Start with a simple algorithm that you can implement quickly
- ② Plot learning curves to decide if more data, more features etc. will help.
- ③ ERROR ANALYSIS: Manually examine the examples (in CV set) that your algorithm made errors on. See if SYSTEMATIC TREND can be spotted on errors.

$m_{CV} = 500$ examples in cross validation set

Algorithm misclassifies 100 emails.

Manually examine the 100 errors, and categorize them based on:

- (i) What type of email it is pharma, replica, steal passwords, ...
- (ii) What cues (features) you think would have helped the algorithm classify them correctly.

WORTH FOCUSING ON

Pharma: 12

Replica/fake: 4

Steal passwords: 53

Other: 31

→ Deliberate misspellings: 5
(m0rgage, med1cine, etc.)

→ Unusual email routing: 16

→ Unusual (spamming) punctuation: 32

- Always have a numerical evaluation of your algorithm e.g.

Should discount/discounts/discounted/discounting be treated as the same word?

Can use "stemming" software (E.g. "Porter stemmer")
universe/university.

Error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try it and see if it works.

Need numerical evaluation (e.g., cross validation error) of algorithm's performance with and without stemming.

Without stemming: 5% error With stemming: 3% error

Skewed Classes

Ex :

* Train logistic regression model $h_0(x)$ { $y=1$ if cancer,
 $y=0$ otherwise)

* Find 1% error on test set

* ONLY 0.50% HAVE CANCER

↳ SKEWED CLASSES

Far greater no. of patients without cancer than with it.

• In the skewed classes regime it is not always clear whether an improvement in accuracy was due to a beneficial change to the algorithm, or by mathematical luck due to the skewed class

ALTERNATIVE: Precision / Recall

* $y=1$ in presence of rare class that we want to detect

ACTUAL CLASS

		1	0
1	TRUE POSITIVE	FALSE POSITIVE	
	FALSE NEGATIVE	TRUE NEGATIVE	
0			

PRECISION

- Of all patients we predicted $y=1$, what fraction actually have cancer?

$$\text{P} = \frac{\text{True Positives}}{\text{True pos} + \text{False pos.}}$$

RECALL

- Of all patients that actually have cancer, how many did we predict correctly?

$$\text{R} = \frac{\text{True Positives}}{\text{True pos} + \text{False neg.}}$$

$$\text{Precision} = \frac{\text{True positives}}{\# \text{predicted as positive}} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

$$\text{Recall} = \frac{\text{True positives}}{\# \text{actual positives}} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

- Having an algorithm with both high recall and high precision confrom we have a classifier that will generalize well to new data.

—> FIXES SKEWED CLASSES

P
R
E
D
I
C
T
E
D