

My
Wing
know
our score.

Week 6

Date: _____

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices

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

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?

Note: There is a difference between someone that really knows how to powerfully & effectively apply that algorithm versus someone that's less familiar with the material

Solutions:

1) Get more training examples

↳ ~~may~~ may not work out at times

2) Try smaller set of features to avoid over fitting

Date: _____

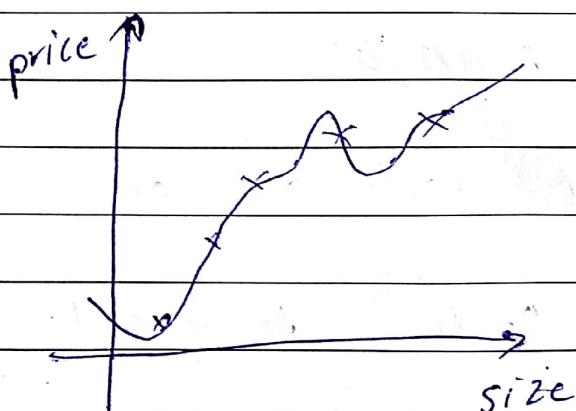
- 3) Try getting additional features
- 4) Try adding polynomial features
- 5) Try decreasing λ
- 6) Try increasing λ

Machine learning diagnostic

Diagnostic: a test that you can run to gain insight what is / isn't working with a learning algorithm, and gain guidance as to how best to improve its performance

↳ which ~~not~~ choose which avenue from above (out of 6)

Evaluating a hypothesis



Overfitting

Fails to generalize

to new examples

not in training set

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

Victory

Page No.

Date: _____

How to tell if hypothesis is overfitting with dataset with many features?

Size	Price	
2104	400	
1600	330	
2400	532	
1416	572	
3000	142	
1985	132	
1534	511	
1427	112	
1386	345	
1494	345	

} Training set (70%)

} Test set (30%)

Training set : $x^{(1)}, y^{(1)}$
 $x^{(2)}, y^{(2)}$

$x^{(m)}, y^{(m)}$

Test set : $x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}$
 $x_{\text{test}}^{(2)}, y_{\text{test}}^{(2)}$
 \vdots
 $x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})}$

m_{test} = no. of test examples

The records chosen for 70% & 30% must be random.

Testing Training / testing procedure for linear regression

- learn parameter θ from training data (minimizing training error $J(\theta)$)
- Compute test set error :

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

Training / testing procedure for logistic regression

- learn parameter θ from training data
- compute test set error

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

Date: _____

- Misclassification error (0/1 error misclassification error):

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

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

↳ ~~Fraction of~~ Fraction of
of mislabelled predictions
total predictions

Since $J_{\text{train}}(\theta)$ will always be low ~~be~~,
 $J_{\text{test}}(\theta)$ can be used to find out if there's
any overfitting: when $J_{\text{train}}(\theta)$ is low
 $\& J_{\text{test}}(\theta)$ is high

Model Selection & Training/Validation/Test sets

Once parameters $\theta_0, \theta_1, \dots, \theta_n$ are fit to some set of data (training set), the error of the parameters as measured on that data (the training error $T(\theta)$) is likely to be lower than the actual generalization error.

Consider the model selection problem

$$d=1 \quad 1) \quad h_{\theta}(x) = \theta_0 + \theta_1 x$$

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

⋮
⋮

$$d=10 \quad 10) \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$$

d = degree of

To choose a degree of polynomial & get estimate of how well hypothesis generalised to new examples:

~~Cheat~~ ~~page~~

first take models & minimise training error

it seems as if there is an additional parameter d to the

algorithm that is

being determined by the data set

Victory (continued on next page)

Date: _____

$$1) h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \theta^{(1)}$$

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

$$3) h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(10)}$$

theta vectors
after minimising
training error

$$\theta^{(1)} \rightarrow J_{\text{test}}(\theta^{(1)})$$

$$\theta^{(2)} \rightarrow J_{\text{test}}(\theta^{(2)})$$

Choose model with
lowest $J_{\text{test}}(\theta^{(i)})$

$$\theta^{(10)} \rightarrow J_{\text{test}}(\theta^{(10)})$$

Assume that $J_{\text{test}}(\theta^{(5)})$ has the lowest test set error then choose

$$\theta_0 + \theta_1 x + \dots + \theta_5 x^5$$

Suppose I want to ask how well does my model generalize?

↳ possible answer: report test set error

$$J_{\text{test}}(\theta^{(5)})$$

↳ not a fair estimate of hypothesis generalization

Date:

Problem: $T_{\text{test}}(\theta^{(5)})$ is likely to be an optimistic estimate of generalization error i.e. our extra parameter ($d = \text{degree of polynomial}$) is fit to data set

- ↳ we chose the value of d that gave the best possible performance on the test set
- ↳ hypothesis is likely to do better on test set than on new examples it hasn't seen before
- Parameter d is fit to the test set

To address this problem:

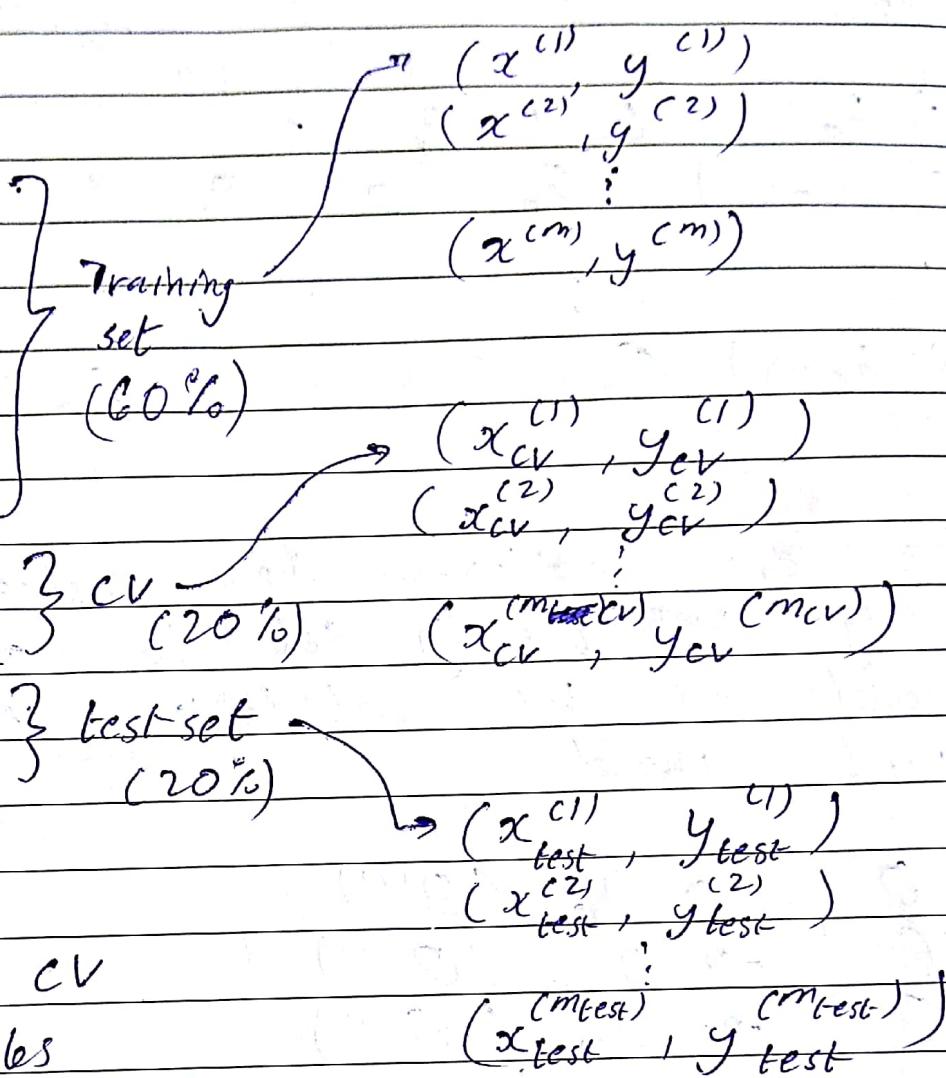
Given data set split into 3 pieces:

- 1) training set (60%)
- 2) cross validation set (CV) (20%)
- 3) test set (20%)

Date: _____

Data set:

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



m_{cv} = no. of CV examples

~~Train /~~

Thus the cost/error can be calculated for ~~the~~ training, cross validation & test data set: $J_{train}(0)$, $J_{cost}(0)$, $J_{test}(0)$

Model Selection

$$\min_{\theta} J(\theta) \rightarrow \theta^{(1)}$$

$$1) h_0(x) = \theta_0 + \theta_1 x$$

$$2) h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)}$$

$$10) h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)}$$

~~Use cross validation~~ Now instead of using ~~test~~ training set to select the model use the validation set to select model

$$\theta^{(1)} \rightarrow J_{CV}(\theta^{(1)})$$

$$\theta^{(2)} \rightarrow J_{CV}(\theta^{(2)})$$

$$\theta^{(3)} \rightarrow J_{CV}(\theta^{(3)})$$

$$\theta^{(10)} \rightarrow J_{CV}(\theta^{(10)})$$

Pick the ~~lowest~~ hypothesis with lowest cross validation error. Suppose for this example $h_0(x) = \theta_0 + \dots + \theta_4 x^4$ is chosen hence 4th order polynomial model is chosen

$$\hookrightarrow d=4$$

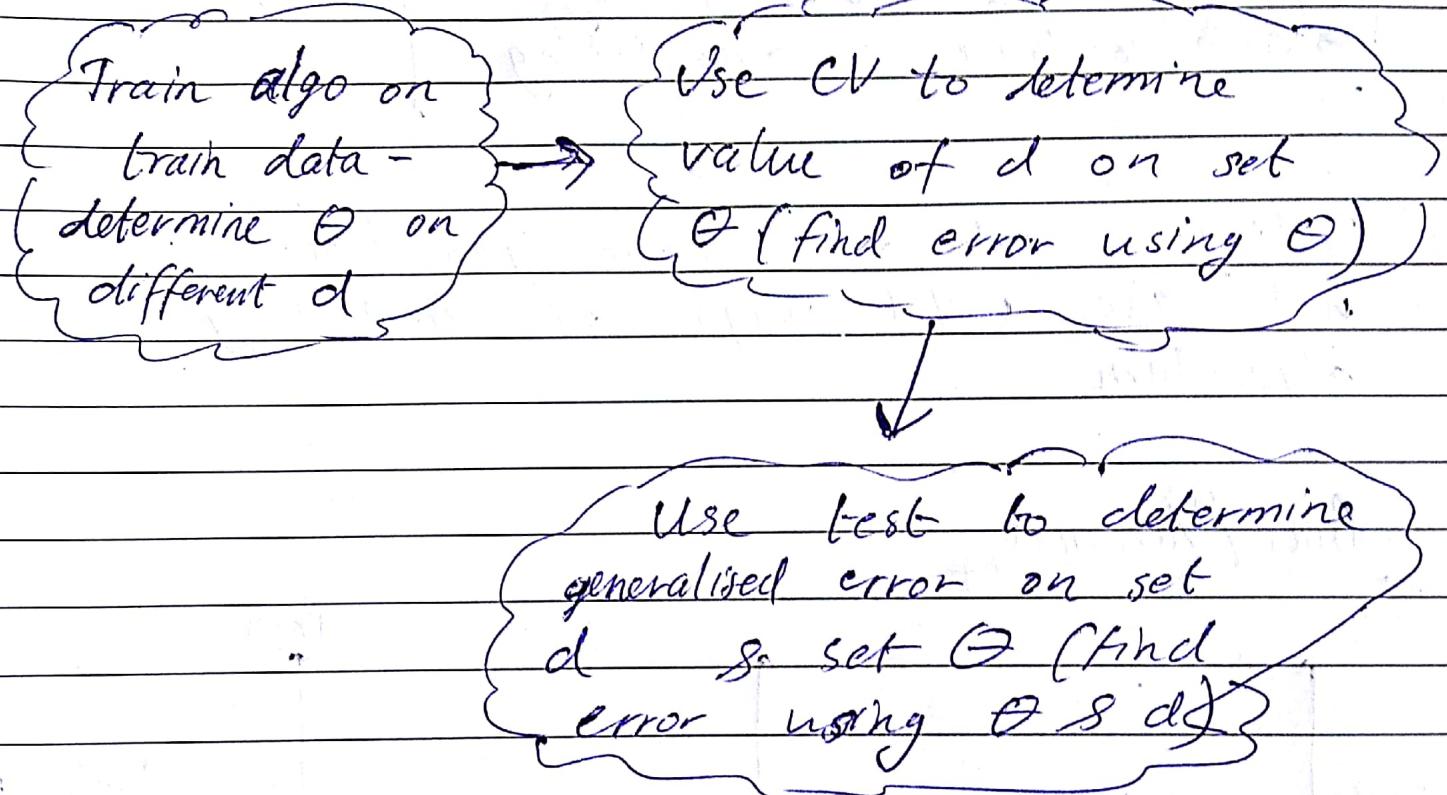
Date: _____

Use test set to estimate the generalization error of the model selected

Dataset

train	CV	test
-------	----	------

Note: CV can't be used to report generalised error as it was used to determine d



Using test set to get d vs then as a way of finding generalised error is wrong

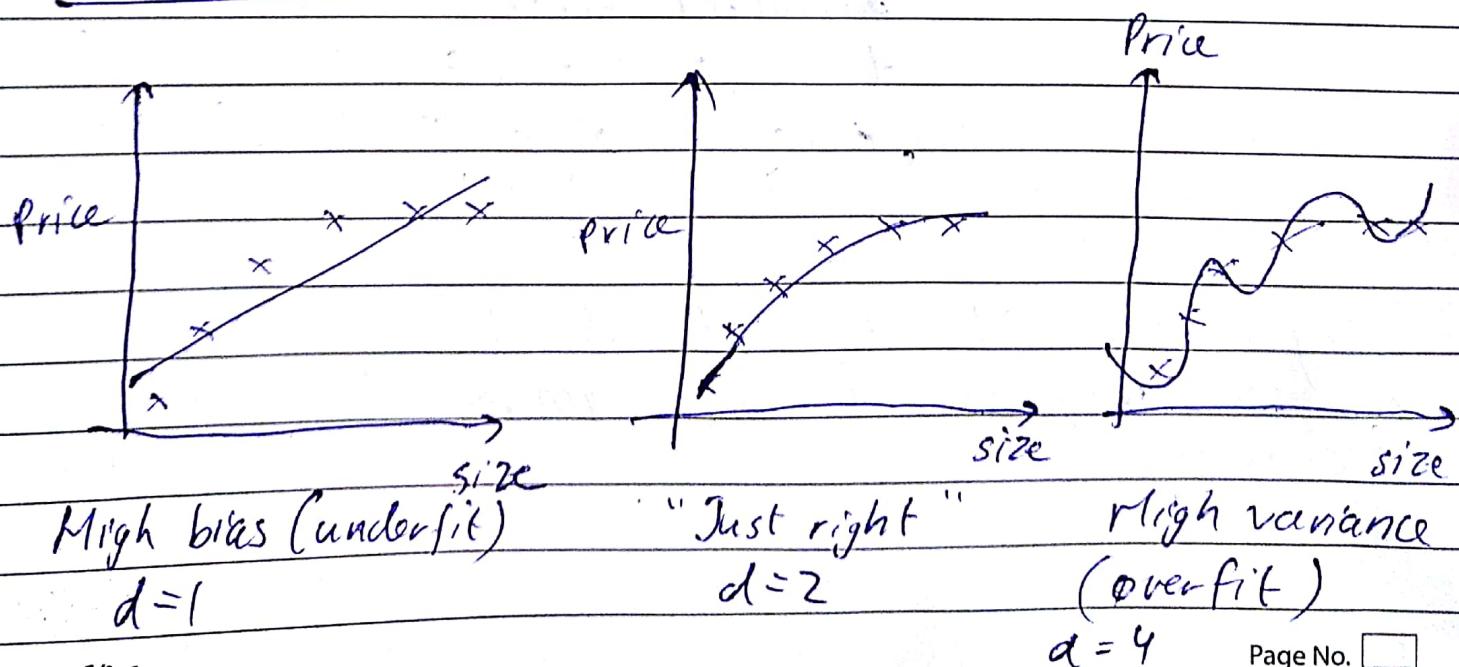
Diagnosing Bias Vs Variance

If a learning algorithm doesn't do well then almost all the time the model has a high bias problem or high variance problem \rightarrow underfitting or overfitting respectively.

High bias ~~is~~ = underfitting
high variance = overfitting

Finding out which one you have will give useful ways of improving algorithm

Bias / Variance



High bias (underfit)

$$d=1$$

"Just right"

$$d=2$$

High variance

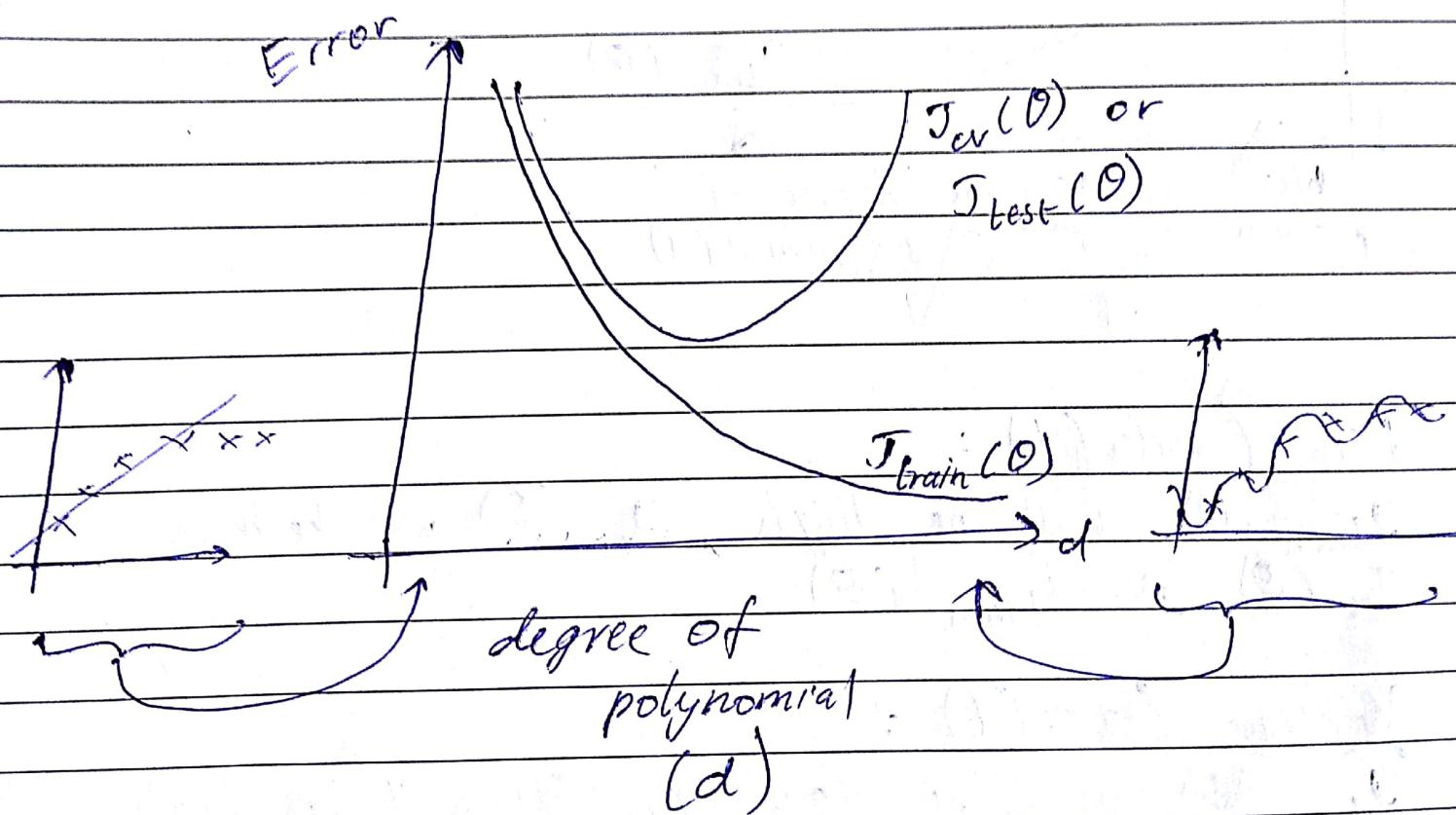
$$(over-fit) \quad d=4$$

Date:

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

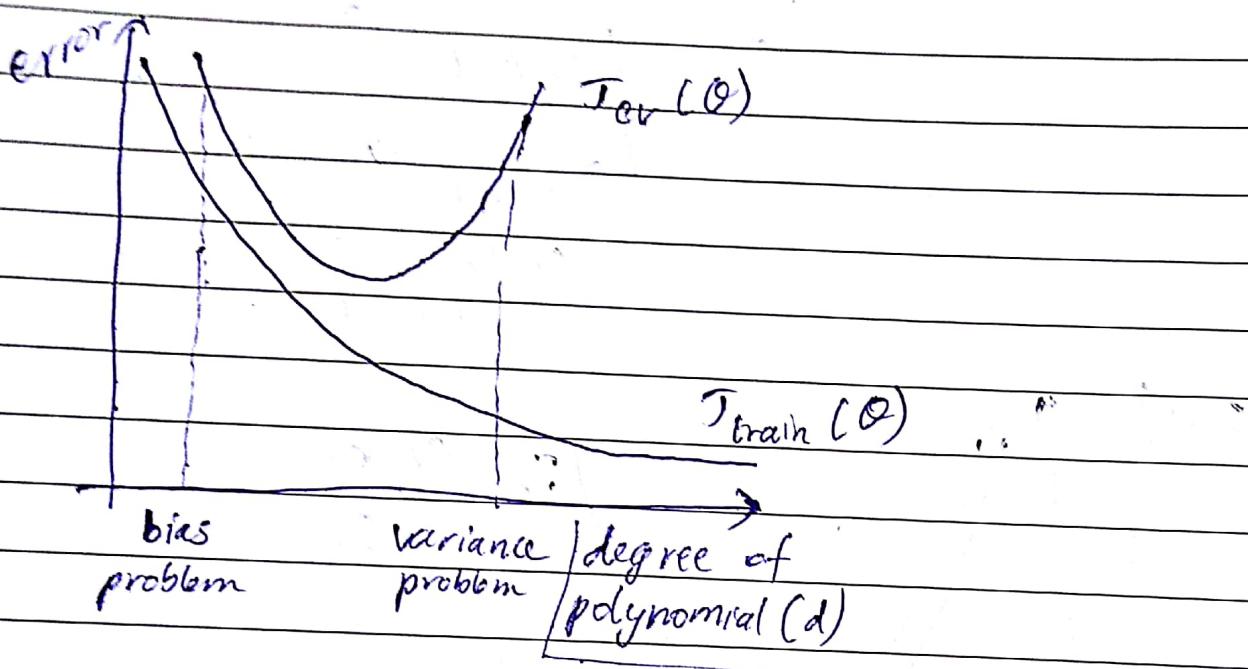
Cross validation error: $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} [h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)}]^2$
(or $J_{test}(\theta)$)

→ $J_{test}(\theta)$ will output similar results



Diagnosing bias vs Variance

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



Bias (underfit) :

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

Variance (overfit) :

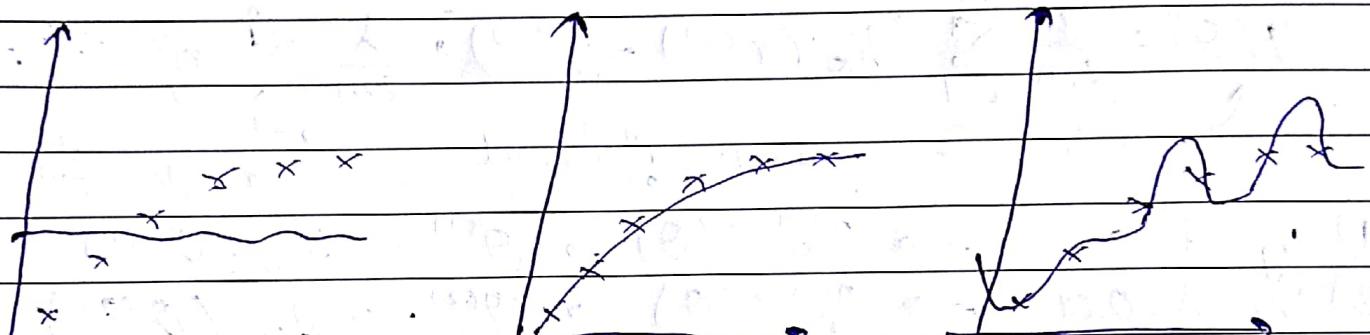
$J_{train}(\theta)$ will be low, $J_{cv}(\theta)$ will be high

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

Regularization & Bias/Variance

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

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



large λ
High bias (underfit)
 $\lambda = 1000, \theta_1 \approx 0, \theta_2 \approx 0\dots$
 $h_{\theta}(x) \approx \theta_0$

Intermediate
"Just right"

Small λ
High variance
(overfit)
 $\lambda \approx 0$

How is this intermediate λ chosen?

To choose the regularization parameter λ

~~use~~ the following :

- 1) $J_{\text{train}}(\theta)$ } all average sum
- 2) $J_{\text{cv}}(\theta)$ } of squared error
- 3) $J_{\text{test}}(\theta)$ } without regularization

Average sum of squared errors:

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

Choosing the regularization parameter

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

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

\curvearrowleft reg. applied on train set \curvearrowright reg. not applied

- 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 $\lambda = 0.02 \rightarrow \dots \rightarrow J_{cv}(\theta^{(3)})$
- 4) Try $\lambda = 0.04 \rightarrow \dots \rightarrow J_{cv}(\theta^{(4)})$
- 5) Try $\lambda = 0.08 \rightarrow \dots \rightarrow J_{cv}(\theta^{(5)})$
- ⋮
- 12) Try $\lambda = 10.24 \rightarrow J_{cv}(\theta^{(12)})$

Pick model with the lowest $J_{cv}(\theta)$

use that value of λ . Suppose ~~the~~ model 5 is picked

Use test set to find generalized error of the model chosen on previously unseen examples

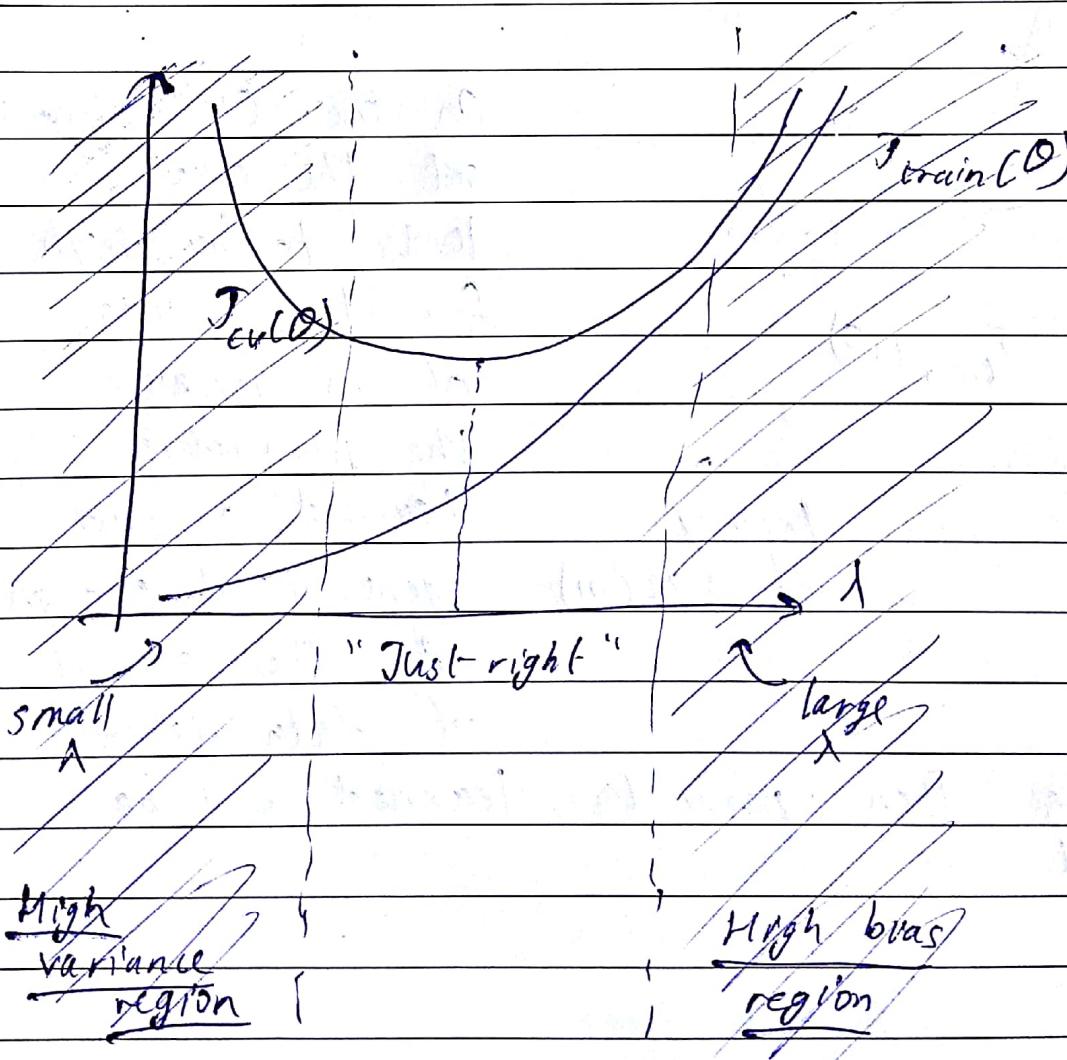
Date: _____

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

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

$$\textcircled{3} \quad J_{\text{cv}}(\theta) = \frac{1}{2m} \sum_{i=1}^{m_{\text{cv}}} (h_\theta(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

Find $\textcircled{1}$ first for different values of λ & train a specific $\theta^{(i)}$ then run $\textcircled{2}$ & $\textcircled{3}$ with the specific ~~$\theta^{(i)}$~~ $\theta^{(i)}$

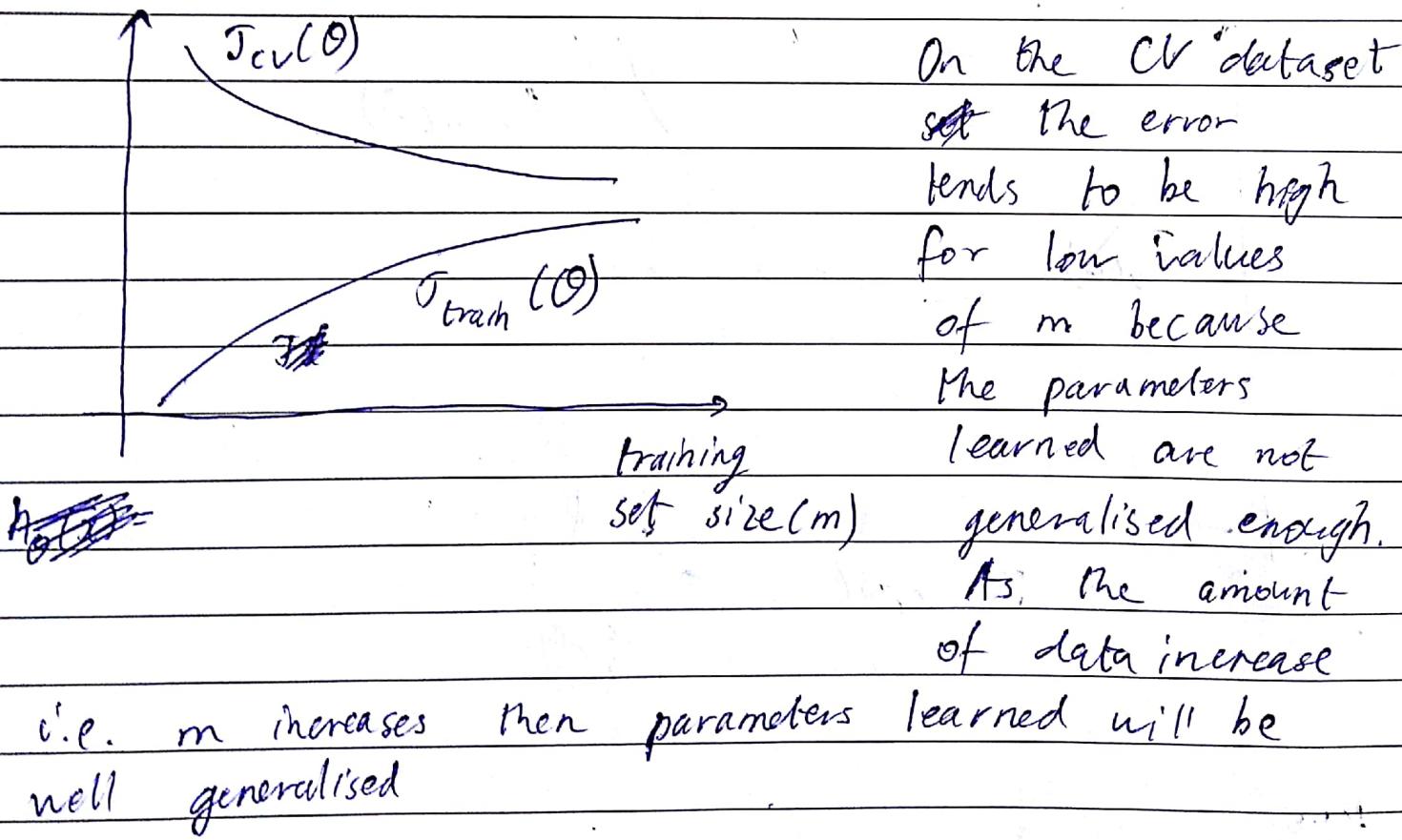


Learning Curves

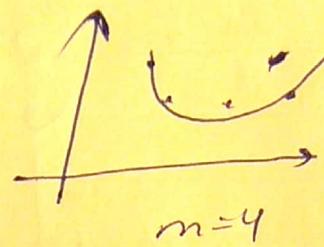
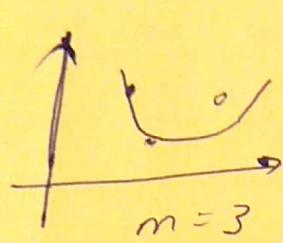
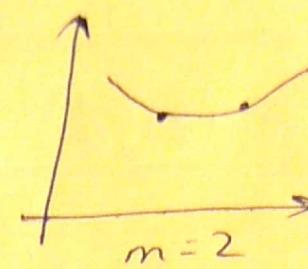
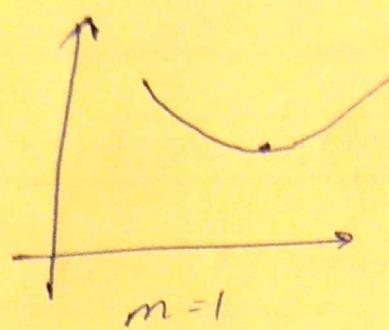
- ↳ used as sanity check that algorithm is working properly
- ↳ or want to improve performance of algorithm

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

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



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



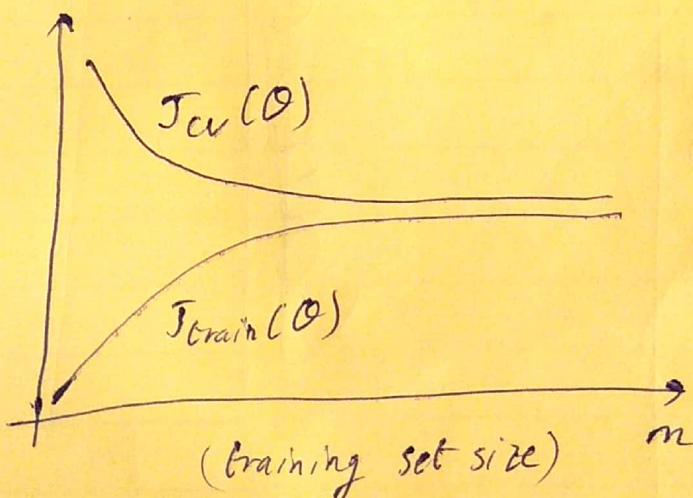
As the number of examples are low then fitting the data has ~~very~~ low cost

Hence for small values of m , $J_{\text{train}}(\theta)$ is low.

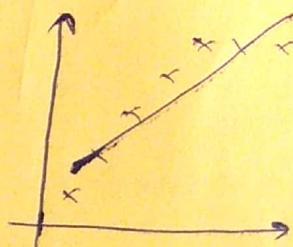
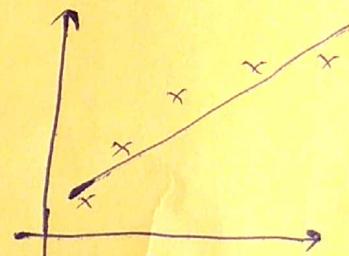
As ~~the~~ m increases

then fitting the dataset such that cost is low becomes difficult, thus $J_{\text{train}}(\theta)$ increases

High Bias



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



The shape of the fitted curve does not change much as more & more data points are added. Thus both J_{cv} & J_{train} plateau eventually.

Initially J_{cv} decreases as the parameters tend to learned tend to work better on generalised data but it plateaus eventually.

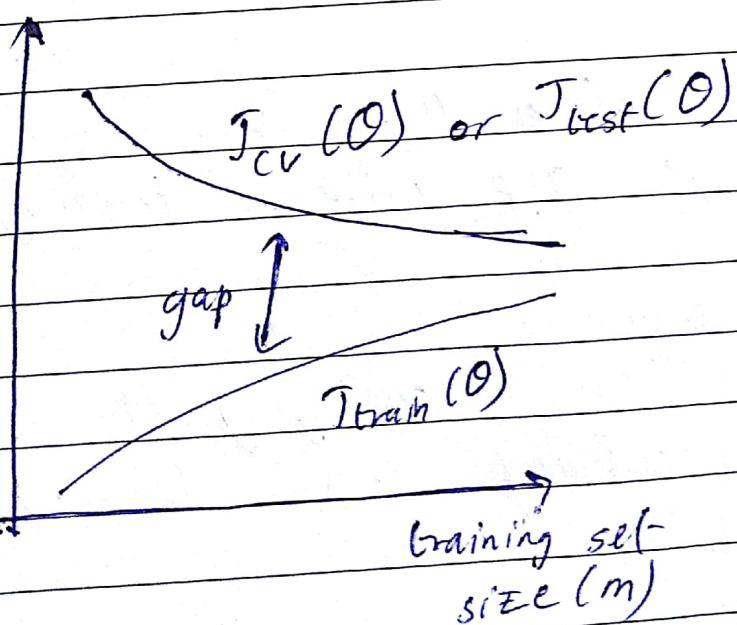
As before the training error will be small initially & increase as more examples are added. It plateaus as well.

~~This~~ $J_{train} \approx J_{cv}$ as they plateau.

If a learning algorithm has high bias then getting more training data will not help \rightarrow you will end up with similar values of J_{train} & J_{cv}

J_{cv} diagnostic
Both J_{cv} & J_{train} will be high in value]

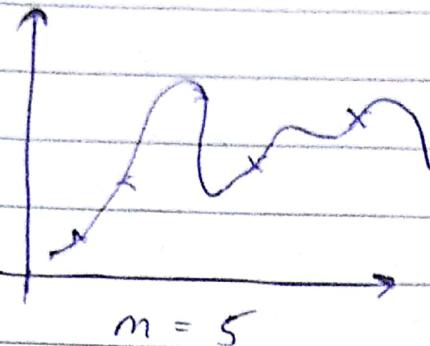
High Variance



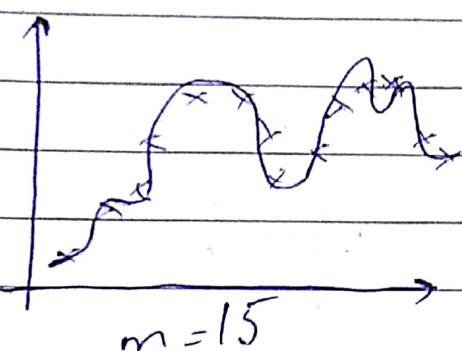
If a learning algorithm is suffering from high variance, getting training data is likely to help.

$$h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{100} x^{100}$$

Date:



When m is small then $h_0(x)$ (in this case) can fit it well. As the data set increases (m increases) then fitting the data set becomes more difficult, $J_{\text{train}}(\theta)$ increases then



On small values of m we also know that the model is being overfit hence the model will not generalize well

to unseen data therefore $J_{\text{cv}}(\theta)$ will be high. As more data is added to training, the less the degree of overfitting hence $J_{\text{cv}}(\theta)$ decreases.

Indicative diagnostic of a high variance problem is the large gap between training & cross validation.

Both J_{cv} & J_{train} are converging to each other hence adding more data helps.

The curves illustrated here are idealized, in reality you will encounter curves with more noise & messiness

Large training set size - $J_{\text{train}}(\theta)$ increases with training set size & $J_{\text{cv}}(\theta)$ continues to decrease without leveling off. Also, $J_{\text{train}}(\theta) < J_{\text{cv}}(\theta)$ but the difference between them remains significant.

Deciding What to do Next

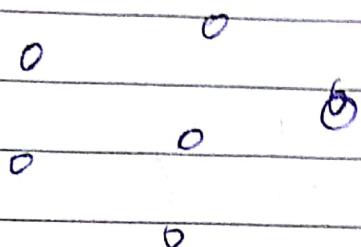
Debugging a Learning algorithm:

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 ~~an~~ unacceptably large errors in its prediction. What should you try next?

- Get more training examples \rightarrow fixes high variance
- Try smaller sets of features \rightarrow fixes high variance
- Try getting additional features \rightarrow fixes high bias
- Try adding polynomial features \rightarrow fixes high bias
- Try decreasing $\lambda \rightarrow$ fixes high bias
- Try increasing $\lambda \rightarrow$ fixes high variance

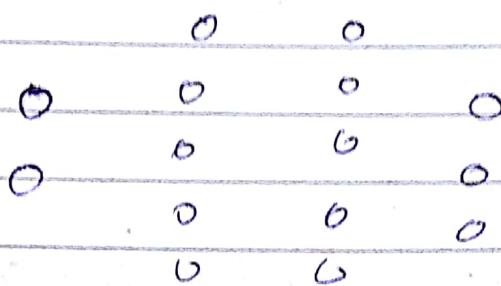
Neural Networks & overfitting

"Small" neural networks
 (few nodes per layer & few hidden layers \rightarrow fewer parameters) prone to underfitting



Computationally cheaper

"Large" neural networks
 (more parameters) prone to overfitting



Computationally more expensive

Use regularization (λ) to address overfitting

Usually it's best to use a "large" neural network with regularization as opposed to using "small" ~~neat~~ neural networks.

Date: _____

If you want to find how many hidden layers would work best then the best thing would be to split the data into train, CV & test. Check ~~what~~ what amount of layers deliver the best results on cross validation data set as it was used to determine the degree of linear regression (d)

Prioritizing what to work on

Suppose you want to build a spam ~~classification~~ classifier on emails.

Building a spam classifier

Supervised learning : x = feature of email, y = spam(1) or not spam(0)

Features x : Choose 100 words indicative of spam / not spam



Eg: deal, buy, discount, Andrew, etc.

Create a feature vector from the email depending on which words are present: 0 for not present, 1 for present (regardless of # of times a word is mentioned)

Date:

Note: In practice, take most frequently occurring n words (10'000 to 50'000) in training set, rather than manually picking 100 words.

How to spend time to make it have low error (spam detector)?

- Collect ~~test~~ lots of data
- develop sophisticated features based on email routing information (from email header)
- develop sophisticated features for message body e.g. should "discount" & "discounts" be treated as the same word? How about "deal" & "dealer"? features about pronunciation?
- develop sophisticated algorithm to detect misspellings ~~leg. misspelling~~

A research group will fixate on one of these problems & try to increase the accuracy of the model this way.

When faced with an ML problem there are usually many different ways to improve the algorithm

Error Analysis

↳ gives a more systematic way of choosing which & way to improve algorithm out of the many different options available as shown previously

Recommended approach

- Start with a simple algorithm that you can implement quickly. Implement it & test it on your cross-validation data. (may take up to 1 day - 24 hours - on implementation for this)
- Plot learning curves to decide if more data, more features etc are likely to help. In the absence of ~~the~~ the quick implementation & learning curves, it is very difficult to tell what to do next.
- Error analysis: Manually examine the examples (in cross validation sets) that your algorithm made errors on. See if you spot any systematic ~~trend~~ trend in what type of examples it is making errors on.

↳ This process will inspire you to design new features or tell of the current shortcomings of the system.

error analysis evaluated on CV so that when new features are made / decided we can check generalised error as new features may be CV data set dependent.

Date:

Suppose we have a spam classification ML problem:

$m_{cv} = 500$ 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
- (ii) What cues (features) you think would have helped the algorithm classify them correctly

→ spam from pharma companies, watch ~~replica~~ replica companies, phishing attempts & etc

Count ~~acc~~ according to each class

misclassified:

Pharma: 12

Replica: 4

phishing: 53

Other: 31

It would be the best use of your time to ensure that spam classifier correctly classifies phishing attempts, accuracy will increase greatly.

→ Create classes as such:

- deliberate misspellings: 5

- unusual email routing: 16

- unusual (spamming) punctuation: 32

Date:

It seems that most spammers use spamming punctuation when they write their emails hence it's best to create sophisticated features to address this misclassification.

Error analysis - process of manually examining the mistakes of the algorithm

↳ supports quick & dirty implementation of algorithm - we want to figure out what are the most difficult examples to classify so effort can be focused on these

The importance of numerical evaluation

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

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

↳ construct of NLP

"Porter stemmer" - tool that can be used to achieve the feature above.

Stemming software basically looks at the first few alphabets of the words & groups words together. This can lead to disadvantages:

- universe/university will be taken as some word when semantically they mean different things.

↳ Thus not easy to tell if stemming software should be used

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

↳ Need numerical evaluation (e.g. cross validation error) of algorithm's performance with & without stemming:

without stemming : 5% error } Thus using stemming
with stemming : 3% error } is a good idea

Gives you a very quick way of deciding to use stemming or not

The quick & dirty approach for creating algo allows you to use evidence to determine which direction to apply your time rather than using gut feeling

Error Metrics for Skewed Classes

~~Error metrics = numerical evaluation~~

numerical evaluation - error metrics

Cancer classification example

Train logistic regression model $h_\theta(x)$
($y=1$ if cancer, $y=0$ otherwise)

You find that you got 1% error on test set
(99% correct diagnoses)

However only 0.50% patients have cancer

Suppose you use a non-learning function:

function $y = \text{predict_cancer}(x)$ *

$y = 0$; % ignore x !

return

This function now has 0.50% error,
better than the learning algo applied!

Skewed classes - ratio of positive
to negative examples is at extremes

of +ve examples is much smaller
than -ve in this example

Date: _____

If skewed classes are present in the data set it becomes hard to use simple classification accuracy:

Version 1: 99.2% (0.08% error)

Version 2: 99.5% (0.05% error)

Technically according to simple classification Version 2 is an improvement over version 1. But Version 2 can just simply be setting $y=0$ for all x . as previously shown.

↳ difficult to say if there is any actual improvement with the learning classifier

Other ~~error~~ metrics are used for skewed classes

Date: _____

Precision / Recall

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

~~Actual class~~ ^{Actual}

~~Predicted class~~

		1	0
		True positive	false positive
Predicted class	1		
	0	false negative	True negative

false positive - prediction is 1 but it is wrong

false negative - prediction is 0 but it is wrong

Precision (of all patients where we predicted $y=1$, what fraction actually has cancer?):

$$\frac{\text{True positives}}{\# \text{ predicted } \text{positive}} = \frac{\text{True +ve}}{\text{True +ve} + \text{false +ve}}$$

Recall (of all patients that actually have cancer, what fraction did we correctly detect as having cancer?):

$$\frac{\text{True +ve}}{\# \text{ actual } \text{positive}} = \frac{\text{True +ve}}{\text{True +ve} + \text{false -ve}}$$

The higher the value of both, the better.

If our algo predicts 0 all the time ($y=0$) then recall will be 0

↳ there will be no true +ves

If precision & recall are high then we know that even if we have skewed classes, the learning algorithm is doing well.

Trading off Precision & Recall

We want to control the trade-off between precision & recall

We continue with cancer example as before.

Here logistic regression : $0 \leq h_\theta(x) \leq 1$

Predict 1 if $h_\theta(x) \geq 0.7$

Predict 0 if $h_\theta(x) < 0.7$

→ we predict $y=1$ (cancer) only if we are very confident

→ higher precision now because there are more true tress as compared to false tress → predictions are only made when we are more confident

→ lower recall now because we're predicting $y=1$ on a small number of people who actually have cancer

Suppose we want to avoid missing too many causes of cancer (avoid false negatives)

↳ when in doubt we want to predict that they have cancer so that at least they look further into it

Set the probability ~~threshold~~ threshold lower:

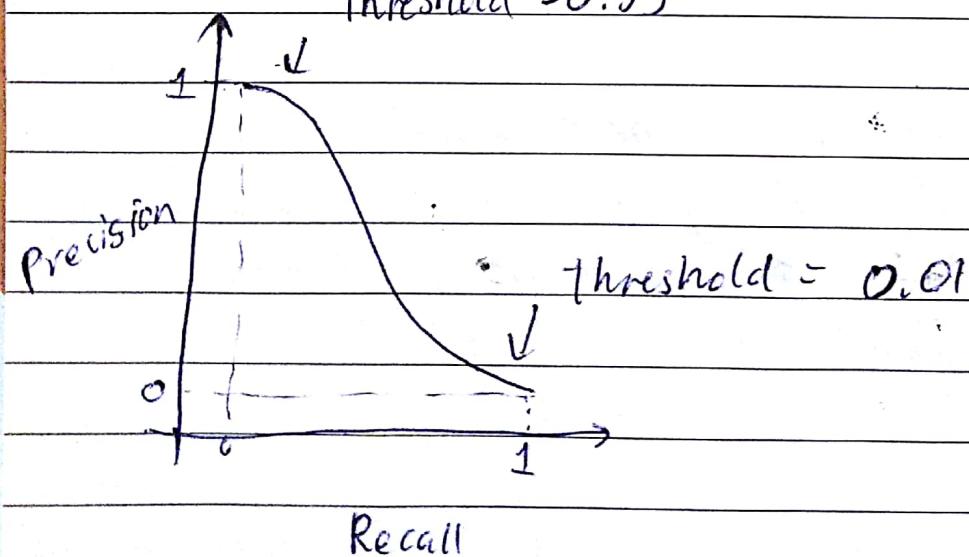
Predict 1 if $h_\theta(x) \geq 0.3$

Predict 0 if $h_\theta(x) < 0.3$

↳ higher recall & lower precision

More generally: Predict 1 if $h_\theta(x) \geq \text{threshold}$

Threshold = 0.99



F₁ score (F score)

How to compare precision/recall numbers?
Which one is the best?

By use precision/recall we lost. The simplicity of a single real number metric \rightarrow it is difficult to decide which algorithm with which threshold is better.

\hookrightarrow slows down the process of decision making \hookrightarrow when evaluating different changes

	Precision	Recall	Avg	F ₁ score
Algorithm 1	0.5	0.4	0.45	0.444
Algo 2	0.7	0.1	0.4	0.175
Algo 3	0.02	1.0	0.51	0.0392

It can ~~not~~ be clearly said that algo 3 is the worst as Recall can be 1 if we predict $y=1$ all the time. Thus there will be no such instances where we predicted 0 hence there will be no false -ves. Hence Recall = 1.0

\hookrightarrow the 2 extremes of a high threshold or a low threshold tend to be bad classifiers.

Date: _____

Now to make a single value metric from P & R we can try a normal average & F_1 score:

$$\text{Average} = \frac{P+R}{2}$$

$$F_1 \text{ score} = 2 \frac{RP}{P+R}$$

With the average, algo 3 is considered the best which is obviously wrong & hence average is a bad measure. However the F_1 score correctly determines Algo 3 the worst.

→ F score is widely used in the ML community

~~F~~ score extremes:

- $P=0, R=0$ then F score = 0
- $P=1, R=1$ then F score = 1

There are many different ways of combining precision & recall, F score is just one way

$$F \text{ score} \in [0, 1]$$

if enough data is given to an inferior algorithm,
it can beat a superior algorithm ~~not~~ ↗

Date:

Try out various thresholds & on
cross validation set & pick whichever
gives the greatest f score.

Data for ML

There have been many studies that
show that many different learning
algorithms tend to give similar ranges
of performance ~~but performance is given a~~
~~driven by training data~~

"It's not who has the best algorithm
that wins. It's who has the most
data"

This statement is true under certain
conditions

Large data rationale

Assume feature $x \in \mathbb{R}^{n+1}$ has sufficient
information to predict y accurately

Example: For breakfast I ate two eggs
choose either two, too or too into
blank space

Counter example : predict housing price from only size (feet²) & no other features

Useful test : given the input x , can a human expert confidently predict y ?

↳ so a human expert will be able to fill in the black space in the example (a normal english speaker) while a human house expert will not be able to give an estimate price for the house with just the area.

Use a ~~large~~ learning algorithm with many parameters

↳ thus becomes low bias algorithm hence $J_{\text{train}}(\theta)$ will be small

Use a ~~very~~ very large training set (unlikely to overfit) hence it has low variance thus:

$$J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$$

hence $J_{\text{test}}(\theta)$ will be small as well

Date: _____

Conti

Condition: can a human expert look at
the features or & confidently predict
 y