

① Deciding what to try next - - -

* Debugging a learning algorithm.

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

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

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? low accuracy.

• Get more training examples.

• Try smaller set of features

• Try getting additional features.

• Try adding polynomial features ($x_0^1, x_1^2, x_2^3, \dots$)

• Try changing λ . (Increase & decrease)

Machine learning diagnostic.

Diagnostic:

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

Diagnoses can take time to implement, but doing so can be a very good use of your time.

Q Which of following statements about diagnostic are true?

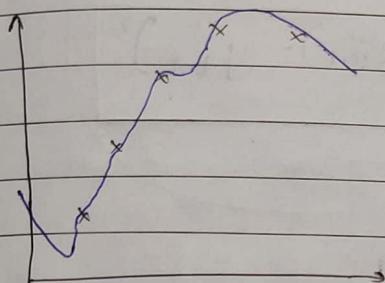
- ✓ It's hard to tell what will work to improve a learning algorithm, so the best approach is to go with gut feeling and just see what works.
- ✓ Diagnostics can give guidance as to what might be more fruitful things to try to improve a learning algorithm.
- ✓ Diagnostics can be time consuming to implement and try, but they can still be a very good use of your time.
- ✓ A diagnostic can sometimes rule out certain courses of action (changes to learning algo) as being unlikely to improve its performance significantly.

Evaluating a hypothesis -

Getting low value of training error \rightarrow good thing
 \rightarrow underfitting

Not generalized for new examples \rightarrow
(Training error $J(\theta) \rightarrow$ low)

(Test error $J_{\text{test}}(\theta) \rightarrow$ high)

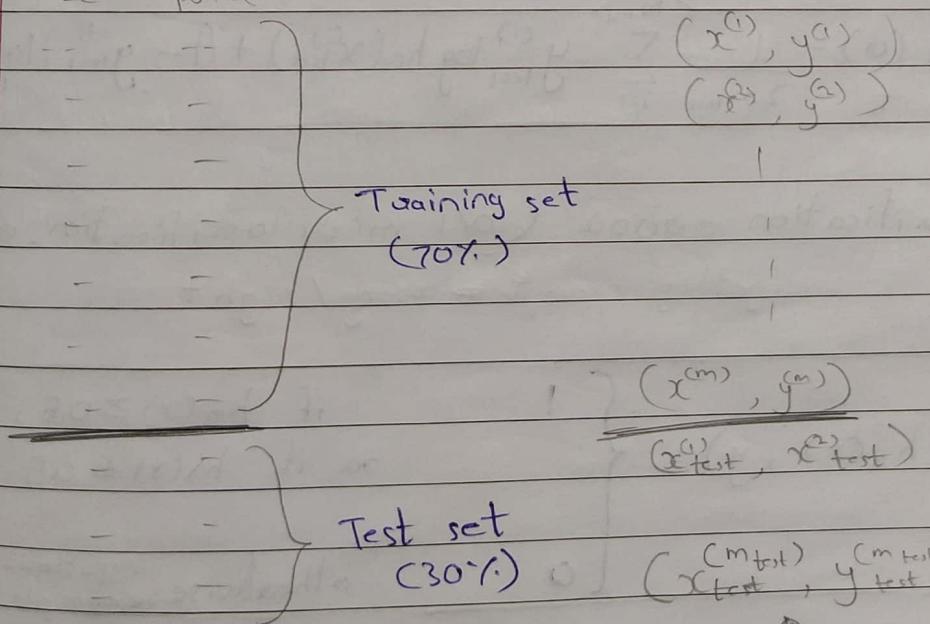


$h_{\text{can}}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$, Here we can plot & get to know about overfitting.

But when we have many features
We cannot plot how hypothesis looks.

Dataset:

size price



Just remember to randomly choose train, test.

representation.

* Training / testing procedure for linear regression.

↳ Learn parameters θ from training data
Minimizing training error $J(\theta)$ → for 70% of data

2) Compute test set error

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

plug θ here, from 1)

* Training / testing procedure for logistic regression.

↳ Learn parameters θ from training data

2) Compute test set error

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

Misclassification error (0/1 misclassification error)

getting example → wrong/right

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

Hypothesis classified examples → y correct

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

count of how many examples classified wrongly

① Model Selection and Train/Validation/Test Sets

* Model Selection process

- ↗ Degree of polynomial?
- What features to include?
- λ^2

Once your parameter is what fit to some set of data, maybe the training set, (may be something else), then the hypothesis error of your hypothesis as measured on that same data set, (such as training error). That's unlikely a good estimate of your actual generalization error.

i.e. How well the hypothesis will generalize to new examples.

The training error^{($J(\theta)$)} is likely to be lower than the actual generalization error.

★ Model Selection.

Minimizing training error gives some parameter vector θ

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

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

$$2) h(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

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

$$3) h(x) = \theta_0 + \theta_1 x + \theta_2 x^3$$

$$\theta^{(3)} \rightarrow J_{\text{test}}(\theta^{(3)})$$

$$10) h(x) = \theta_0 + \dots + \theta_{10} x^{10} \quad \theta^{(10)} \rightarrow J_{\text{test}}(\theta^{(10)})$$

$d = \text{degree of polynomial} - \text{one more parameter}$

Here we choosed model with lowest test error ($d=5$)

i.e. choose $\theta_0 + \dots + \theta_5 x^5$

How well does the model generalize?

Report test set error $J_{\text{test}}(\theta^{(5)})$.

This will not be a fair estimate of how well my hypothesis generalizes.

because, we chose fit value of d ($d=5$) according to test set → that gave us the best performance on test set

So, Hypothesis of $d=5$, is likely to do better on this test set than it would on new examples that it hasn't seen before

Not Predictive

[How well]

[How well]

generalizes to

New examples

fit a_0, a_1, \dots on \rightarrow training set \rightarrow Hypothesis a_0, a_1, \dots will go well on training set, but not well on new examples.fit d on \rightarrow test set \rightarrow Hypothesis

Examples

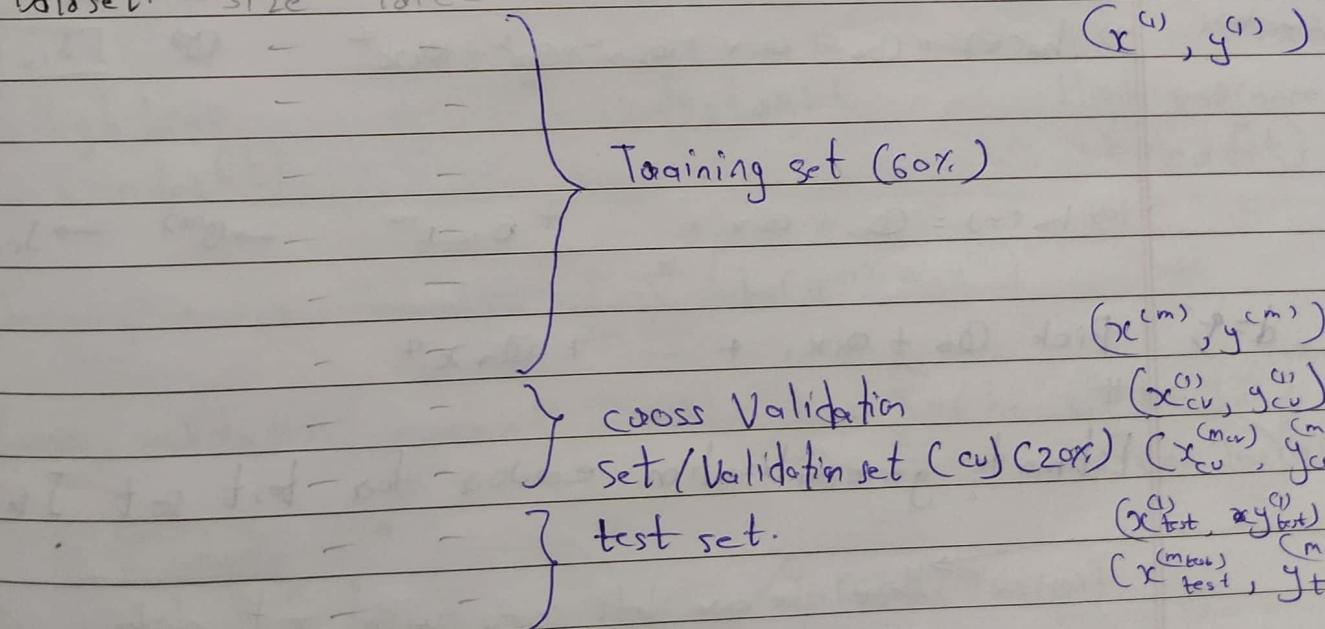
we haven't

seen before.

Performance of hypothesis on that test set may not be fair estimate of how well the hypothesis is likely to do on

★ Evaluating your hypothesis.

Dataset: size Price



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

Cross Validation error: $J_{cv}(\alpha) = \frac{1}{2m_{cv}} \sum_{t=1}^{m_{cv}} (h_{\alpha}(x_{cv}^{(t)}) - y_{cv}^{(t)})^2$

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

* Model Selection

~~We instead of test set to select 'd'~~,
 we will use cv set to select 'd'.
 i.e. to select model.

1) $h_{\alpha}(x) = \alpha_0 + \alpha_1 x \rightarrow \min_{\alpha} J(\alpha) \rightarrow \alpha^{(1)} \rightarrow J_{cv}(\alpha^{(1)})$

2) $h_{\alpha}(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 \rightarrow \alpha^{(2)}$

3) $h_{\alpha}(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \alpha_3 x^3 \rightarrow \alpha^{(3)}$

4) $h_{\alpha}(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \alpha_3 x^3 + \alpha_4 x^4 \rightarrow \alpha^{(4)} \quad | J_{cv}(\alpha^{(4)})|_{\text{lower}}$

10) $h_{\alpha}(x) = \alpha_0 + \alpha_1 x + \dots + \alpha_{10} x^{10} \rightarrow \alpha^{(10)} \rightarrow J_{cv}(\alpha^{(10)})$

$d=4$, Pick $\alpha_0 + \alpha_1 x + \dots + \alpha_4 x^4$

Estimate generalization error for test set $J_{\text{test}}(\alpha^{(4)})$

→ Here we are left over with test data.

We will use test data to measure generalization errors of model.

We might generally expect $J_{cv}(\alpha) < J_{test}(\alpha)$?.

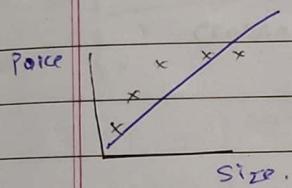
→ An extra parameter d has been fit to cross validation test set, not according to test set.

-- ① Diagnosis Bias vs. Variance. - - - - -

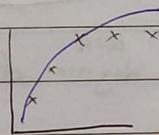
Most of time algo is in underfitting / overfitting problem

Bias or Variance or a bit of both?

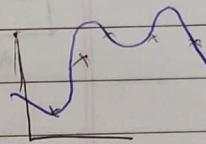
* Bias / Variance.



$$\alpha_0 + \alpha_1 x$$



$$\alpha_0 + \alpha_1 x + \alpha_2 x^2$$



$$\alpha_0 + \alpha_1 x + \alpha_2 x^2 + \alpha_3 x^3 + \alpha_4 x^4$$

High bias
(underfit)

Just right
 $d=2$

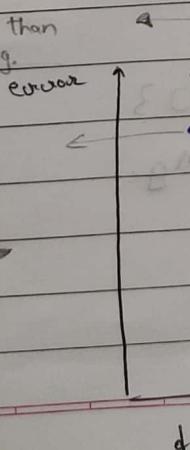
High variance.
(overfit)

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

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

This will generate more error than for training.

This will be high as, high bias



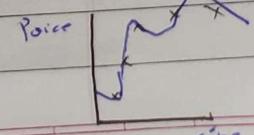
$d=1$

$J_{train}(\alpha)$

$J_{cv}(\alpha) \approx J_{test}(\alpha)$

$J_{cv}(\alpha) \approx J_{test}(\alpha)$

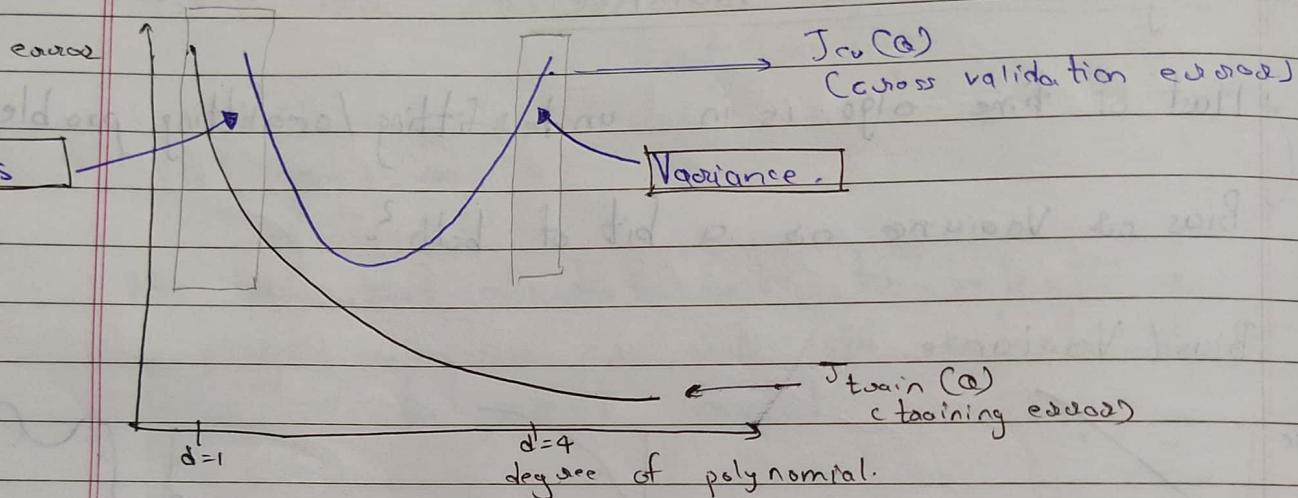
This increases because Train data overfits testing data
- underfits.



$d=4$

* Diagnosing Bias vs Variance.

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



Bias (Underfit)

$J_{train}(\alpha)$ will be high

$$J_{cv}(\alpha) \approx J_{train}(\alpha)$$

Variance (Overfit)

$J_{train}(\alpha)$ will be low

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

All $J_{train}(\alpha) = 0.10$, $J_{cv}(\alpha) = 0.3$
 High Variance (overfitting).

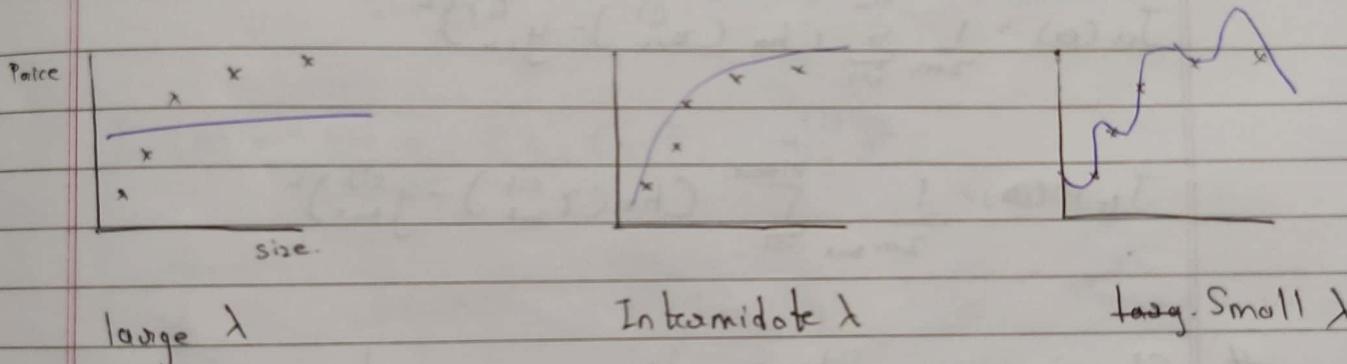
- - @ Regularization and Bias / Variance.

→ How regularization affects Bias & Variance.

* Linear regression with regularization.

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

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



large λ

High bias
(underfit)

Intermediate λ

Just right

too small λ

High variance.
(overfit)

$$\lambda = 1000, \theta_0 \approx 0, \theta_1 \approx 0, \dots$$

$$h_{\theta}(x) \approx \theta_0$$

$$\lambda = 0 (I \rightarrow 0)$$

So, no regularization term

* Choosing a regularization parameter. λ .

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

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

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad \text{--- (without regularization)}$$

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

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

Choosing regularization parameter.

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

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

we can go even less

1) Try $\lambda=0 \rightarrow \min J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$

2) Try $\lambda=0.01 \rightarrow \min J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{\text{cv}}(\theta^{(2)})$

3) Try $\lambda=0.02$

4) Try $\lambda=0.04$

5) Try $\lambda=0.08$

$\rightarrow J_{\text{cv}}(\theta^{(5)})$ lowest $J_{\text{cv}}(\theta)$

1) Try $\lambda=10.24 (\theta)$

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

x_2

~~Training data \rightarrow~~ choosing x | choosing α

Page No: 88
Date: / /

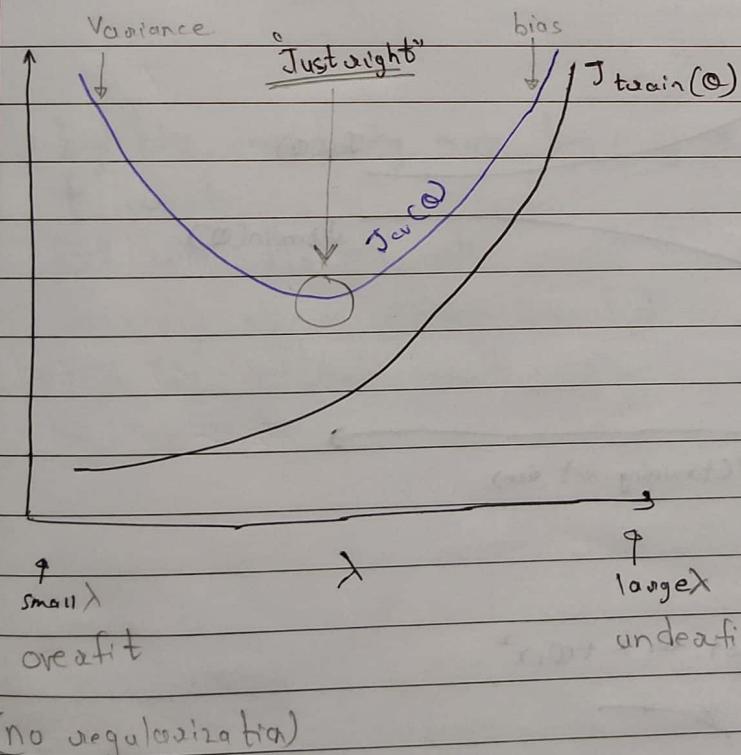
Pick (say) $\alpha^{(i)}$. Test error: $J_{\text{test}}(\alpha^{(i)})$

* Bias/Variance as a function of the regularization parameter λ .

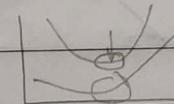
$$J(\alpha) = \frac{1}{2m} \sum_{i=1}^m (h_\alpha(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \alpha_j^2$$

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

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



Here, we don't have considered regularization term in J_{train} & J_{cv} . I think it's by mistake. We should have considered regularization.



And Also, for slight, J_{train} will also be lowest.

learning curve will tell about bias/Variance or both.

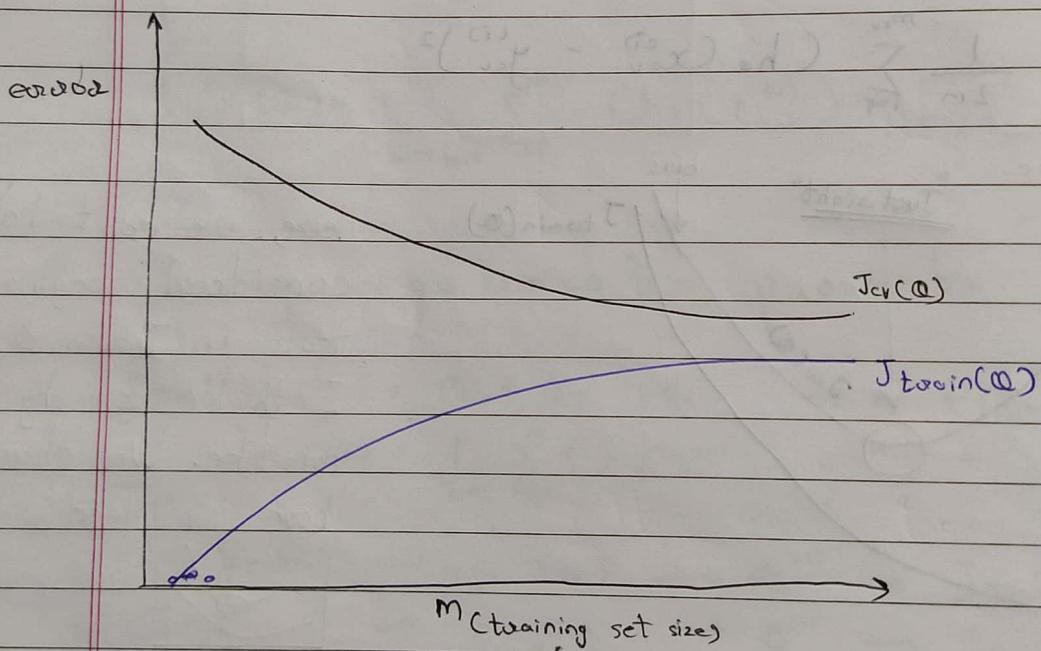
Learning Curves

To diagnose \rightarrow Algo suffering from bias or variance problem

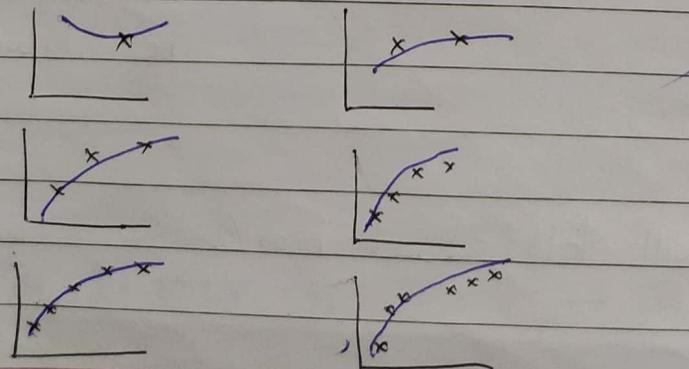
learning curves

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

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



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

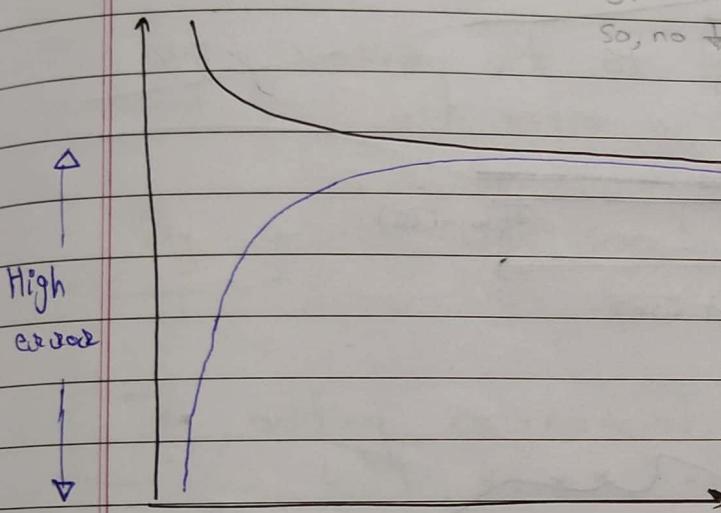


So, $m = 1/2/3$, { training $J_{\text{train}}(\theta) = 0$

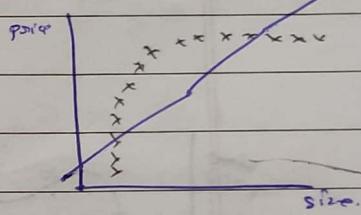
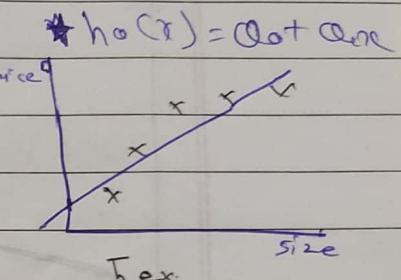
$J_{\text{train}}(\theta) \approx 0 \ (\rightarrow C > 0)$ Regularization

Without regularization

★ High bias



increasing m will give nearby same hypothesis, so same line,
so, no \downarrow in $J_{\text{cv}}(\theta)$.



but same line

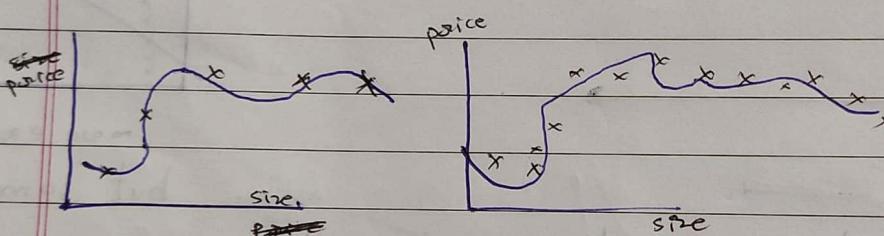
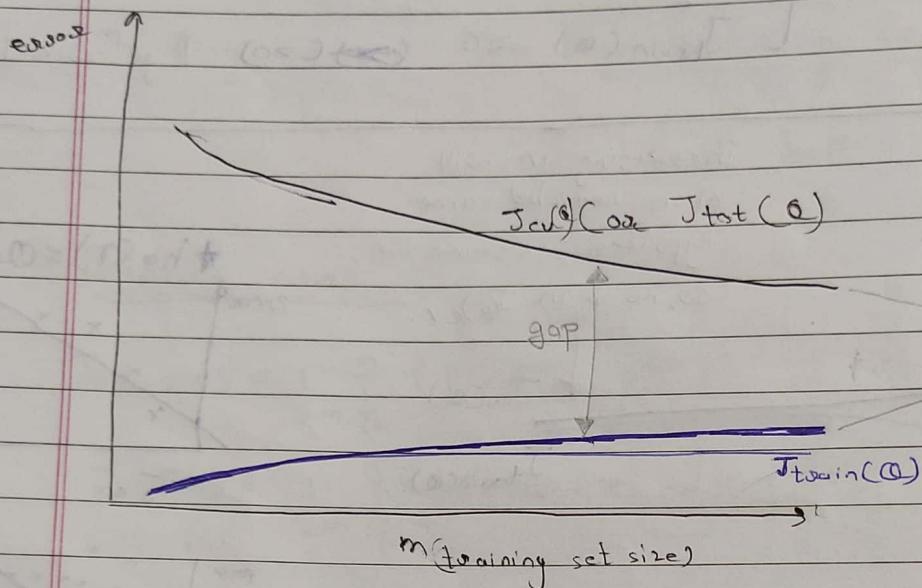
In high bias, considering more training ex. will \leftarrow (poorly close to previous)

not work.

$J_{\text{cv}}(\theta)$ will flatten out.

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

★ High Variance.

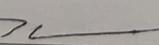


$$h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

(and small λ)

In high variance, getting more training data is likely to help.

→ $J_{\text{train}}(\theta) \uparrow$, $J_{\text{cv}}(\theta) \downarrow$ → This is what actually matters.

- Q1 Getting more ~~data~~ training data will help learning algo's performance
- ✓ • Algo is suffering from high bias. ✓ ?
 - ✓ • Algo is  Variance. ✓
 - ✓ • $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$ ✓
 - ✓ • $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$ ✓ ?

Deciding What to do Next Revisited --

If model is making a very high error than expected

Get more training examples

→ fixes high variance, not fixes high bias.

Try smaller sets of features

→ fixes high variance, not fixes high bias

Try getting additional features

→ fixes high bias (current hypothesis may be too simple)

Try adding polynomial features ($x_1^2, x_2^2, x_1x_2, \text{etc}$)

→ fixes high bias, not fixes high variance.

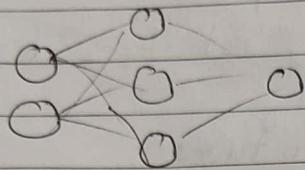
decreasing $\lambda \rightarrow$ fixes high bias

Increasing $\lambda \rightarrow$ fixes high variance

★ Neural networks and overfitting.

"small" neural network

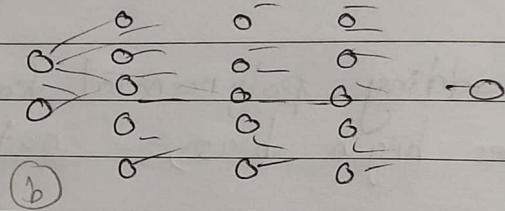
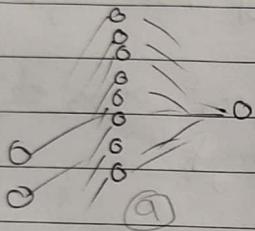
fewer parameters; more prone to underfitting



(computationally cheap)

"large" neural network.

more computation parameters; more prone to overfitting



(computationally more expensive)

The best way is using large neural network with regularization (disadvantage \rightarrow Time complexity)

a or b

Use train, test, cv,

Try training with one or 2/3 hidden layers.
See $J_{cv}(\theta)$

Q1 We fit a neural network, $J_{cv}(\theta) \gg J_{train}(\theta)$
Is increasing no of hidden layers likely to help?

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

→ Better to skip this video.

④ Prioritizing What to Work on

* Building a spam classifier.

Supervised learning $x =$ features of small

$y = \text{spam}(1)$ or not $\text{spam}(0)$.

Features x: choose 100 words indicative of

e.g. deal, buy, discount spam / not spam,
andrew, now ↗

$$x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 1 \\ 1 \\ 0 \end{bmatrix} \text{ and new } x \in \mathbb{R}^{100}$$

$$x_j = \begin{cases} 1 & \text{if word } j \text{ appears in } \underline{\text{email}} \\ 0 & \text{otherwise.} \end{cases}$$

^a Deal of the Day! Buy now!

In practice, take most significant (frequently) occurring words (10,000 to 50,000) in training set rather than manually pick 100 words & use them as features.

• Collect lots of data

"Honey Pot" Proj.

features → email routing

- Is discount \hookrightarrow Discounts same. , Punctuation's.

• Detecting malpractice, medicine, watches

~~61~~ visit ~~cooperative~~ P.T.O.

* Recommended approach.

- Take simple algo; implement & use CV. set.
- Plot curves (more data/ features)
- Error analysis Manually examine e.g. examples.
e.g.

$$M_{cv} = 500$$

wrongly classified CV = 100

examine 100 \rightarrow Type

\Rightarrow which features will help

In NLP \rightarrow Porter stemmer is used

Q Which of the following do you agree? Check all that apply.

For some learning applications, it is possible to imagine coming up with many different features. But it can be hard to guess in advance which feature will be most useful.

For spoken classification, algo to detect and correct deliberate misspellings will make a significant improvement in accuracy.

Because spam classifⁿ uses very high dimensional feature vectors, ($n=50,000 \dots$), a significant effort to collect a massive training set will always be a good idea

There are often many possible ideas for how to develop a high accuracy learning system; "gut feeling is not recommended". To choose among alternatives.

① Error Metrics for skewed example.

Train logistic regression model $h_0(x)$.
 $(y=1 \rightarrow \text{cancer} / y=0 \stackrel{\text{else}}{\rightarrow} \text{otherwise})$

Found that we got 100% error on test set.

→ (99.7% correct diagnoses)

but from 100

only 0.5% had actually cancer → 0.5% penalty

If we have said, no one has cancer, → 0.5% penalty.

Skewed Class

Accuracy is not a really good measure.

99.2% accuracy (0.8% error)

99.5% accuracy (0.5% error)

→ we cannot tell which one gives more quality to classifier.

Because predicting 0 all times is not a good classifier.

★ Precision (Recall)

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

$y=1$	Positive
$y=0$	Negative

		Actual Class		Predicted Class		Positivity
		1	0	True Positive	False Positive	
Predicted Class	1	True Positive	False Positive	True Positive	False Positive	Positivity
	0	False negative	True negative	False negative	True negative	

FN → must be as low as possible

Telling cancer patient that you don't have cancer
False negative must be ↓,
False pos, not must, but should be ↓.

★ Precision

Of all patients where we predicted $y=1$, what fraction actually have cancer.

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

★ Recall

of all patients where that actually have cancer, what fraction did we correctly detect as having cancer.

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

- If we make a classifier which predicts $y=0$ all the time.
 \rightarrow ~~recall = 0~~ as these will no true positive.
As these will no y predicted as 1.

(ii) Precision = ?.

$$\begin{array}{|c|c|} \hline & 80 & 20 \\ \hline 80 & | & 820 \\ \hline \end{array} = \frac{80}{80+20} = 0.8$$

(iii) Recall = ?.

$$\begin{array}{|c|c|} \hline & 20 & 20 \\ \hline 80 & | & 820 \\ \hline \end{array} = \frac{80}{80+80} = 0.5$$

• Classifier with high precision / high recall is actually a good classifier -- (cannot decide by accuracy)

• While defining precision/recall,
we assign $y=1$ to class that we are trying to detect

An algorithm cannot cheat (predicts $y=1$ all the time)
and get high precision / recall

High precision / High recall \rightarrow good classifier

① Trading off Precision & Recall - - -

for some applications we will want to somehow control trade-off b/w precision & recall.
 ↗ balance. maybe

★ Trading off precision & Recall.

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

Predict $y=1$ if $h_\theta(x) \geq 0.5$ or 0.9

Predict $y=0$ if $h_\theta(x) < 0.5$ or 0.9

Suppose we want to predict $y=1$ (cancer)
 ↗ only if very confident.

High precision, lower recall

We are classifying $y=1$,
 only if we are more
 confident.

So, more no. of
 predictions will be
 correct.

We are predicting $y=1$, on
 smaller no of patients

(False positive ↓, maybe) (False negative ↓, maybe)

2) Suppose we want to avoid missing too many cases of cancer (avoid false negative).
false negative must be as low as possible.

In doubt we want to predict $y=1$. (low confident)

predict 1 if $h(x) \geq 0.3$

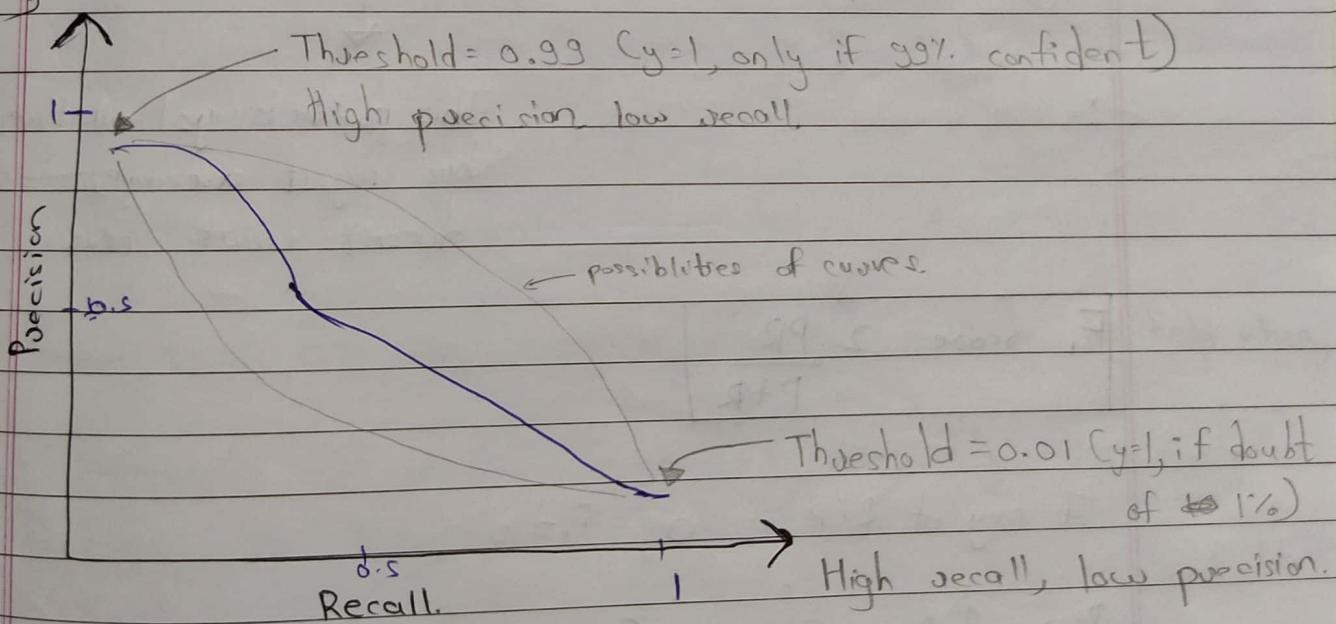
predict 0 if $h(x) \leq 0.3$

High recall

Low precision

We have to choose according to situation 1) / 2)

(A general trade off b/w precision vs Recall)



Note: Try diff thresholds & pick which maximizes F_1 .

Sadguru
Page No: 101
Date: / /



F_1 score (F score)

How to compare precision/recall numbers.

	Precision (P)	Recall (R)	Avg.	F_1 Score
Algorithm 1	0.5	0.4	0.45	0.45
Algo 2	0.7	0.1	0.4	0.175
Algo 3	0.02	1.0	0.51	0.035

Precision/Recall → 2 nos

F_1 score → 1 no

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

This predicts $y=1$ all time, still has highest averages → average doesn't work here

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

If any of R/P is 0, $F_1 = 0$

$$P=0 \& R=0 \Rightarrow F_1 \text{ score}=0$$

$$P=1 \& R=1 \Rightarrow F_1 \text{ score}=1$$

We measure R & P of Cross Validation set
which maximizes $\frac{R+P}{2}$

Data for Machine learning -

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

A inferior algo can beat superior algo, if more data is given, (more data to both algo, still inferior can win)

* Large data advantage.

→ To test this question, can human predict y given x .
Assume x has all features to predict y correctly.

- Use learning algorithm with many parameters
(e.g. logistic/linear, nn)

low bias algo. → $J_{\text{train}}(\theta)$ will be small.

- Use a very large training set
(unlikely to overfit)

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

→ $J_{\text{test}}(\theta)$ will be small.

③ Having large data can improve algo, but unlikely to help when,

- ✓ feature x do not predict y accurately. (linear algo?)

- ✓ Algo with large no of features. (low bias)

- ✓ feature x cannot predict y accurately (nn, big #hidden layers)

- ✓ Not using regularization.