

Week 6

• Evaluating a Learning Algorithm:

a. Deciding What to Try Next:

Debugging a learning algorithm:

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

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

- Get more training examples(sometimes it doesn't work)
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features(x_1^2, x_2^2, x_1x_2 , etc)
- Try decreasing λ
- Try increasing λ

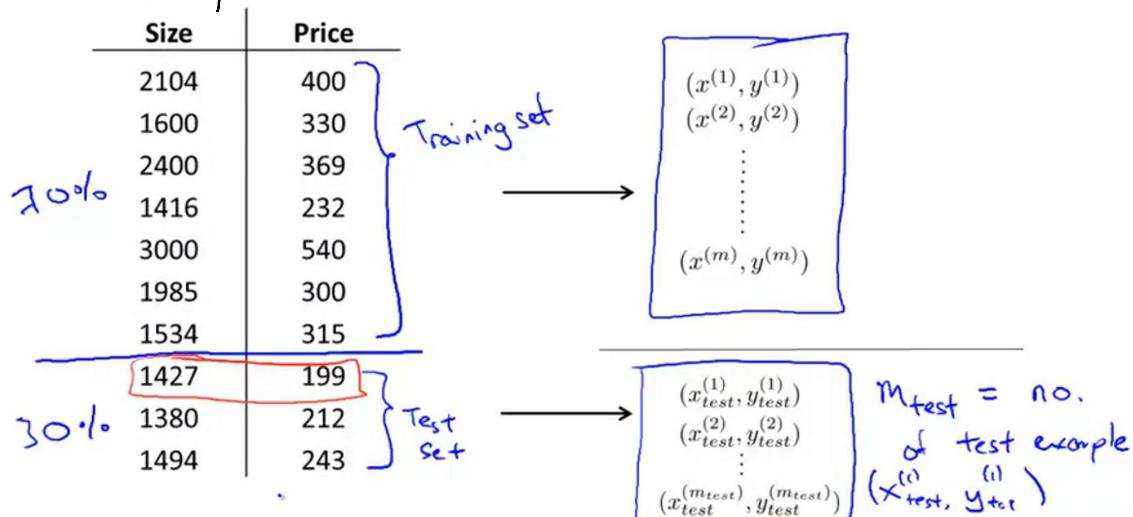
Machine learning diagnostic:

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.

b. Evaluating a Hypothesis:

How to tell a hypothesis is overfitting?

Dataset: Split datasets



Training/testing procedure for linear regression

- ⇒ - Learn parameter θ from training data (minimizing training error $J(\theta)$) 70%

- Compute test set error:

$$\xrightarrow{\text{linear}} J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x^{(i)}_{\text{test}}) - y^{(i)}_{\text{test}})^2$$

$$\xrightarrow{\text{logistic}} J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y^{(i)}_{\text{test}} \log h_{\theta}(x^{(i)}_{\text{test}}) + (1-y^{(i)}_{\text{test}}) \log (1-h_{\theta}(x^{(i)}_{\text{test}})).$$

Sometimes using misclassification error.

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

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x^{(i)}_{\text{test}}), y^{(i)}_{\text{test}}).$$

c. Model Selection and Train/Validation/Test Sets:

one more parameter d denote degree of polynomial.

Model selection

$\rightarrow d = \text{degree of polynomial}$

$$d=1 \quad 1. \quad h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \Theta^{(1)} \rightarrow J_{\text{test}}(\Theta^{(1)})$$

$$d=2 \quad 2. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \Theta^{(2)} \rightarrow J_{\text{test}}(\Theta^{(2)})$$

$$d=3 \quad 3. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \Theta^{(3)} \rightarrow J_{\text{test}}(\Theta^{(3)})$$

⋮

$$d=10 \quad 10. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \Theta^{(10)} \rightarrow J_{\text{test}}(\Theta^{(10)})$$

Choose $\boxed{\theta_0 + \dots + \theta_5 x^5}$ ↪

↗

How well does the model generalize? Report test set

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

⑤ We may choose an overly optimistic estimate

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

We choose the best J_{test} on test set,

it by test set.

Evaluating your hypothesis

Dataset:

Size	Price	
2104	400	60% Training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	20% Cross validation set (CV)
1427	199	
1380	212	20% Test set
1494	243	

split into 3 pieces.

$m_{cv} = \text{number of CV examples}$

$(x_{cv}^{(1)}, y_{cv}^{(1)})$, $(x_{cv}^{(2)}, y_{cv}^{(2)})$, ..., $(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

$(x_{test}^{(1)}, y_{test}^{(1)})$, $(x_{test}^{(2)}, y_{test}^{(2)})$, ..., $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

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

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

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

When we face a model selection problem:

Model selection

- $\hat{\theta}_1: h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)}$
- $\hat{\theta}_2: h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)}$
- $\hat{\theta}_3: h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)}$
- \vdots
- $\hat{\theta}_{10}: h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)}$

find $\min J_{\text{cv}}(\theta^{(i)})$.

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

$$J_{\text{cv}}(\theta^{(2)})$$

$$\vdots$$

$$J_{\text{cv}}(\theta^{(4)})$$

$$J_{\text{cv}}(\theta^{(6)})$$

$$d=4$$

Pick $\theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4$ ←

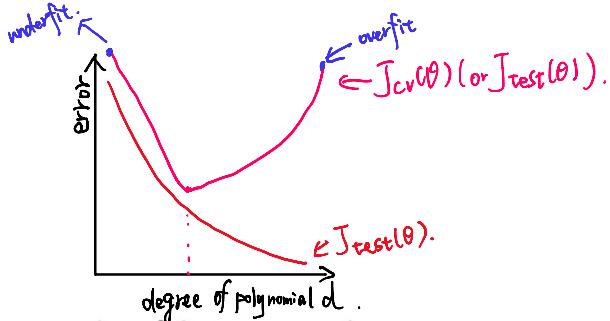
Estimate generalization error for test set $J_{\text{test}}(\theta^{(4)})$ ←

• Bias vs. Variance:

a. Diagnosing Bias vs. Variance:

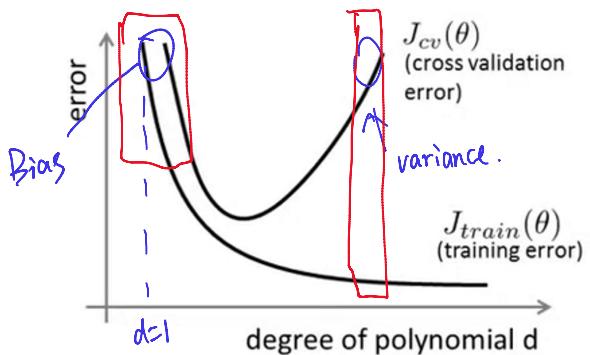
High bias or high variance → underfitting or overfitting

It's important to figure out which of these two problems is bias or variance or a bit of both. It would give a strong indicate for whether the useful and promising ways to try to improve your algorithm.



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 problem or a variance problem?



Bias (underfit) :
 $J_{train}(\theta)$ will be high.
 $J_{cv}(\theta) \approx J_{train}(\theta)$

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

b. Regularization and Bias/Variance:

Choose the regularization parameter λ

first :

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

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

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

} no regularization.

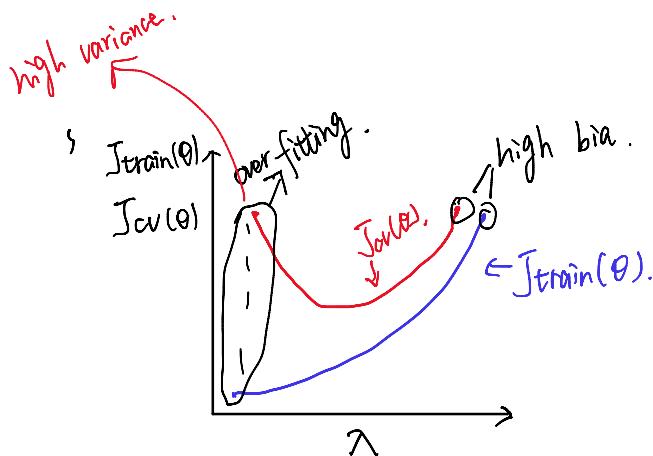
Second : try different λ value.

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

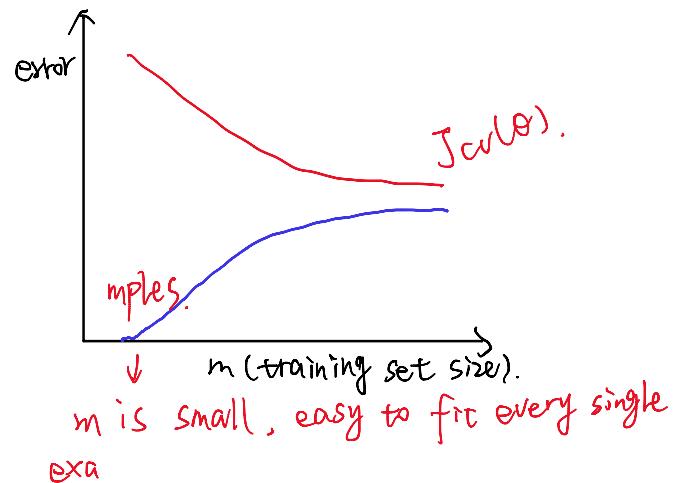
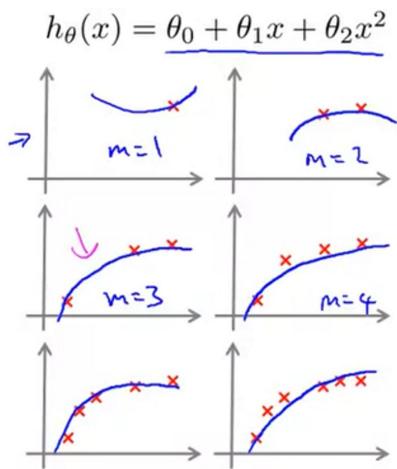
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

- | | | |
|---|---|--|
| <ol style="list-style-type: none"> 1. Try $\lambda = 0$ 2. Try $\lambda = 0.01$ 3. Try $\lambda = 0.02$ 4. Try $\lambda = 0.04$ 5. Try $\lambda = 0.08$ | $\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
$\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
$\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
\vdots
$\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(s)}$ | <p>find $\min J_{cv}(\theta)$.</p> <p>in this case, we suppose $s=m$.</p> <p>$J_{cv}(\theta^{(s)})$</p> |
|---|---|--|

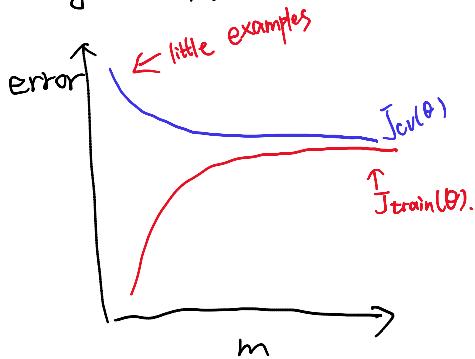
5. Try $\lambda = \underline{0.08}$ $\rightarrow \theta^{(5)}$ $J_w(\theta^{(5)})$
 :
 12. Try $\lambda = \underline{10}$ $\rightarrow \theta^{(12)}$ $J_w(\theta^{(12)})$
 $\uparrow \underline{10.24}$ Pick (say) $\theta^{(5)}$. Test error: $J_{\text{test}}(\theta^{(5)})$



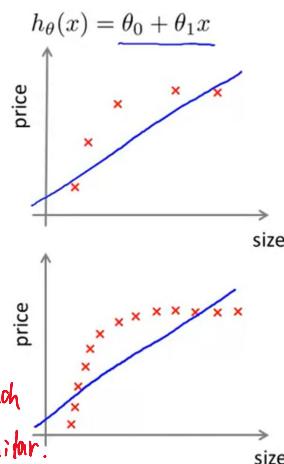
c. Learning Curves:
 a very useful thing to plot, check algorithm working correctly or improve the performance of the algorithm.



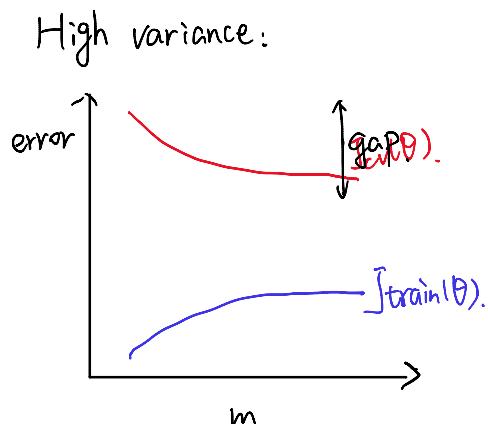
High bias:



because you have few parameters and so much data, when m is large, train and CV will be similar.

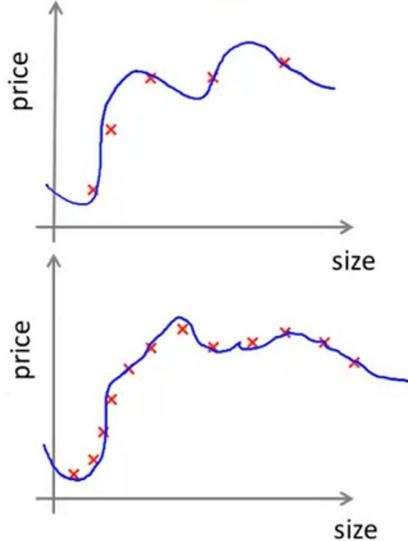


if a learning algorithm is suffering from high bias, getting more training data will not help much.



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

(and small λ)



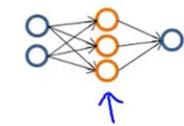
If a learning algorithm is suffering from high variance, getting more training data is likely to help.

b. Deciding What to Do Next Revisited:

- **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 λ :** Fixes high bias
- **Increasing λ :** Fixes high variance

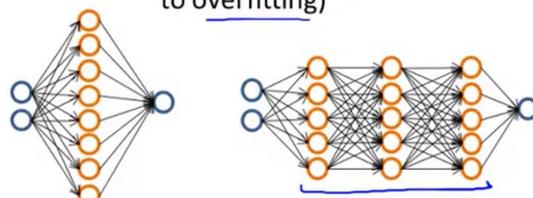
Neural networks and overfitting

→ “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.

using a larger neural network by using regularization to address overfitting that's often more effective than using a smaller neural network.

• Building a Spam Classifier:

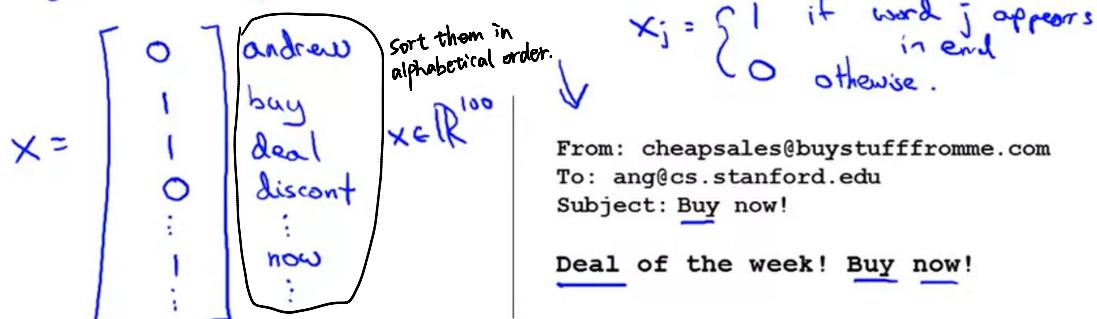
- a. Prioritizing What to Work On:
How do we want to represent x ?

Building a spam classifier

Supervised learning. $x = \text{features of email}$. $y = \text{spam (1) or not spam (0)}$.

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

E.g. deal, buy, discount, andrew, now, ...



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

How to spend your time to make it have low error?

- Collect lots of data
 - E.g. "honeypot" project.
- Develop sophisticated features based on email routing information (from email header).
- Develop sophisticated features for message body, e.g. should "discount" and "discounts" be treated as the same word? How about "deal" and "Dealer"? Features about punctuation?
- Develop sophisticated algorithm to detect misspellings (e.g. m0rtgage, med1cine, w4tches.)

b. Error Analysis:

Recommended approach

- - Start with a simple algorithm that you can implement quickly. Implement it and test it on your cross-validation data.
- - Plot learning curves to decide if more data, more features, etc. are likely to help.
- - Error analysis: Manually examine the examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

it will tell you current things or shortcomings of the system and give you inspiration to come up with improvements.

We should let evidence guide our decisions on where to spend our time rather than use gut feeling, which is often wrong.

Example:

$m_{CV} = 500$ 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.

Pharma: 12

may not worth to spend time.

Replica/fake: 4

→ Deliberate misspellings: 5

Steal passwords: 53 *★*

(m0rgage, med1cine, etc.)

Other: 31

→ Unusual email routing: 16

★ Unusual (spamming) punctuation: 32

work on this problem.

The importance of numerical evaluation

through this evaluation, we'll know our method is good or not.

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

Error Analysis on Cross Validation Set.

• Handling Skewed Data:

a. Error Metrics for Skewed Classed:

Cancer classification example

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

Find that you got 1% error on test set.

(99% correct diagnoses)

Only 0.50% of patients have cancer.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
    return
```

0.5% err.

In this case, means even though we just make $y=0$, it only has 0.5% error. This circumstance sometimes occurs when the ratio of

positive to negative examples is very close to one of two extremes.(above case: the number of positive is much much smaller than negative.)——skewed class.

Precision/Recall

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

Actual class		
		1
		0
Predicted class	1	True positive False positive
0	False negative	True negative

$$y = 0 \\ \text{recall} = 0$$

Precision

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

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

Recall

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

$$\frac{\text{True positives}}{\#\text{actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

b. Trading Off Precision and Recall:

Trading off precision and recall

→ Logistic regression: $0 \leq h_\theta(x) \leq 1$

Predict 1 if $h_\theta(x) \geq 0.5$ ~~0.7~~ ~~0.9~~ 0.3

Predict 0 if $h_\theta(x) < 0.5$ ~~0.7~~ ~~0.9~~ 0.3

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

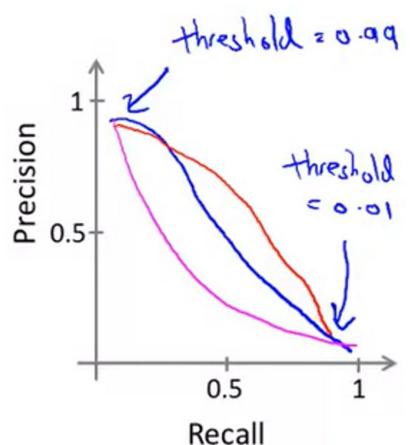
→ Higher precision, lower recall.

→ Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

→ Higher recall, lower precision.

$$\rightarrow \text{precision} = \frac{\text{true positives}}{\text{no. of predicted positive}}$$

$$\rightarrow \text{recall} = \frac{\text{true positives}}{\text{no. of actual positive}}$$



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

How to compare precision/recall numbers?

	Precision (P)	Recall (R)	Average	F. Score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.51	0.0392

Average ~~$\frac{P+R}{2}$~~

$$F. \text{Score}: 2 \frac{PR}{P+R}$$

predict $y=1$ all the time.

$P=0$ or $R=0 \rightarrow F_1=0$ worst

$P=1$ and $R=1 \rightarrow F_1=1$. best.

• Using Large Data Sets:

a. Data for Machine Learning:

Designing a high accuracy learning system

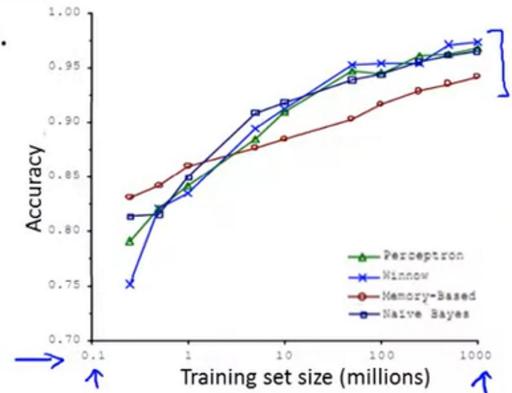
E.g. Classify between confusable words.

{to, two, too} {then, than}

→ For breakfast I ate two eggs.

Algorithms

- - Perceptron (Logistic regression)
- - Winnow
- - Memory-based
- - Naïve Bayes



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

Large data rationale

→ Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features; neural network with many hidden units). low bias algorithms. ←

→ $J_{train}(\theta)$ will be small.

Use a very large training set (unlikely to overfit)

low variance ←

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

→ $J_{test}(\theta)$ will be small

Ask yourself can a human experts look at the features x and **confidently predict the value of y** ? Because that's sort of a certification that y can be predicted accurately from the features x . secondly can we actually **get a large training set**, and train the learning algorithm **with a lot of parameters** in the training set.