

WEEK - 6

WEEK 6: Deciding what to try next

#1 Debugging a learning algorithm.

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

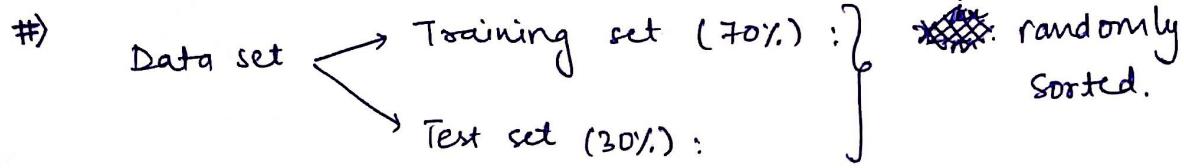
$$J = \frac{1}{2m} \left[\sum_{i=1}^m (h(\theta)x^i - y^i)^2 \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

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

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features.
- Try decreasing λ .
- Try increasing λ

#2 Evaluating a hypothesis.

Overfitting: fail to generalize to new examples not in training set.



So, in case of overfitting, $J(\theta)$ will be low and $J_{\text{test}}(\theta)$ will be high.

#) Training / testing procedure for LR / logistic regression :-

① Learn parameter θ from training data ($J(\theta)$) \rightarrow 70%.

② Compute test error: $J_{\text{test}}(\theta) = \frac{1}{2 * m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x^i) - y^i)^2$

$$\textcircled{3} \quad J_{\text{test}}(\theta) = -\frac{1}{m_{\text{t}}} \left[\sum_{i=1}^{m_{\text{t}}} y_{\text{test}}^i \log(h_{\theta}(x_{\text{test}}^i)) + (1 - y_{\text{test}}^i) \left(\log(\cancel{h_{\theta}(x_{\text{test}}^i)}) \right) \right]$$

↓

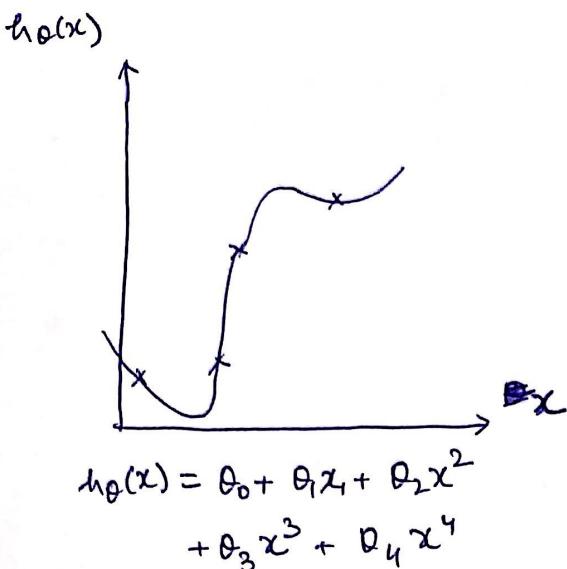
Another metric: Misclassification error :-

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

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

Model selection and training/Validation/test sets.

(parameters/λ)



*

Once parameters $\theta_0, \theta_1, \theta_2, \dots, \theta_4$ were fit to some set of data (training) set, the error of the parameters as measured on that data ($J(\theta)$) is likely to be lower than the actual generalization error.

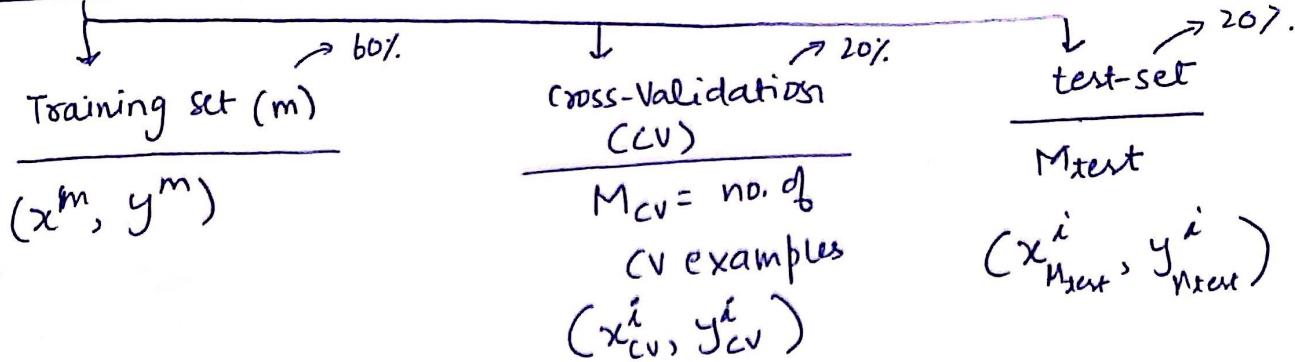
① Model Selection.

- degree of Polynomial ← $d=1$
1. $h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow J_{\text{test}}(\theta^1)$
 2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow J_{\text{test}}(\theta^2)$
 - ⋮
 10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{10} x^{10} \rightarrow J_{\text{test}}(\theta^{10})$

→ Measure the performance, $J_{\text{test}}(\theta^i)$

→ choose the lowest cost among the options.

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

Dataset

①

Training error:

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

Optimize the parameters in θ using the training for each polynomial degree

②

Cross-validation error:

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

Find the polynomial degree d with least error

③

Test error

$$J_{test}(\theta) = \frac{1}{2M_{test}}$$

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

Estimate J using test set.

Model Selection:

→ compute, $J_{cv}(\theta), J_{cv}(\theta^2) \dots J_{cv}(\theta^{10})$

→ select lowest order fact. for inst. $J_{cv}(\theta^4)$.

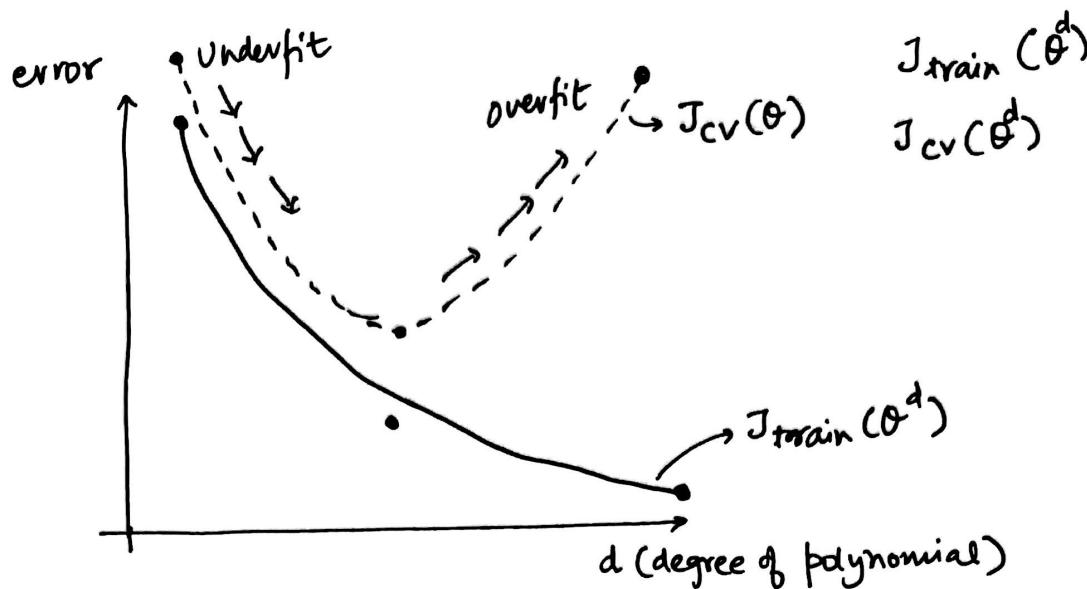
→ Estimate generalization error for test set

$$J_{test}(\theta^{(4)})$$

BIAS VS. VARIANCE

(Underfit)

(Overfit)



■ After having trained the data, for smaller d , $J_{\text{cv}}(\theta)$ may/may not be high. As we go on increasing d , we may find a suitable d for $J_{\text{cv}}(\theta)$ is low.
which

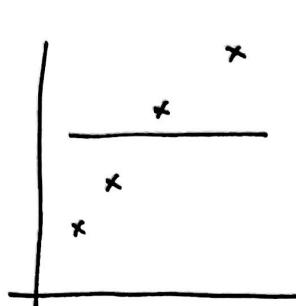
■ Bias: $J_{\text{train}}(\theta)$ will be high. } $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$
 $J_{\text{cv}}(\theta)$ will be high. }

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

Regularization & Bias/Variance (problem of Overfitting)

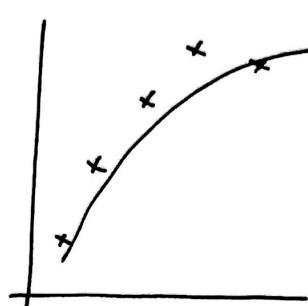
Model: $h_{\theta}(x) = \theta_0^T x$. $\theta \in \{\theta_0 \dots \theta_4\}$

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \lambda \theta_j^2 \right] \quad j \in \{1, \dots, n\}$$

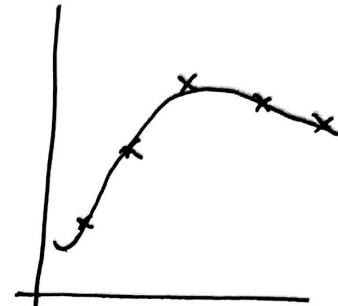


$\lambda = \text{very large}$

(High Bias)



$\lambda = \text{'just right'}$



$\lambda = 0$

(High Variance)

: How to choose λ :

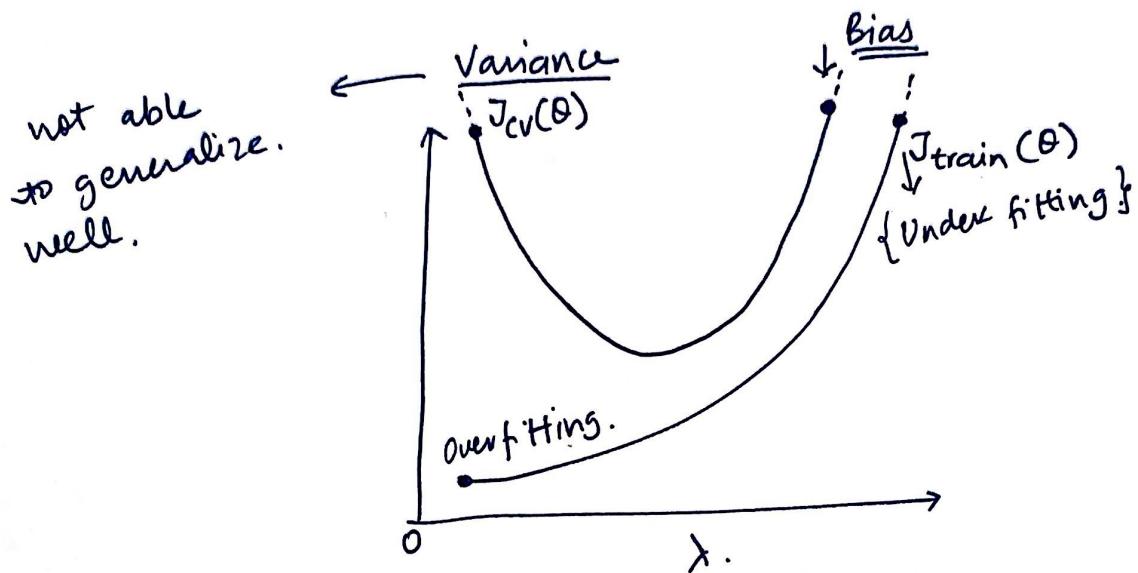
① $\left\{ \begin{array}{l} J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 \\ J_{\text{cv}}(\theta) \\ J_{\text{test}}(\theta) \end{array} \right\} \quad [\lambda = 0]$

② $\begin{array}{l} 1. \lambda = 0 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^1 \rightarrow J_{\text{cv}}(\theta^1) \\ 2. \lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^2 \rightarrow \dots \\ 3. \lambda = 0.02 \\ \vdots \\ 12. \lambda = 10 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{12} \rightarrow J_{\text{cv}}(\theta^{12}) \end{array} \quad \left. \begin{array}{l} \uparrow \\ \vdots \\ \uparrow \\ \rightarrow \text{lowest } J_{\text{cv}}(\theta^d) \\ \text{ex. } (\theta^5) \end{array} \right\}$

③ Pick θ^5 . Test error: $J_{\text{test}}(\theta^5)$

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

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



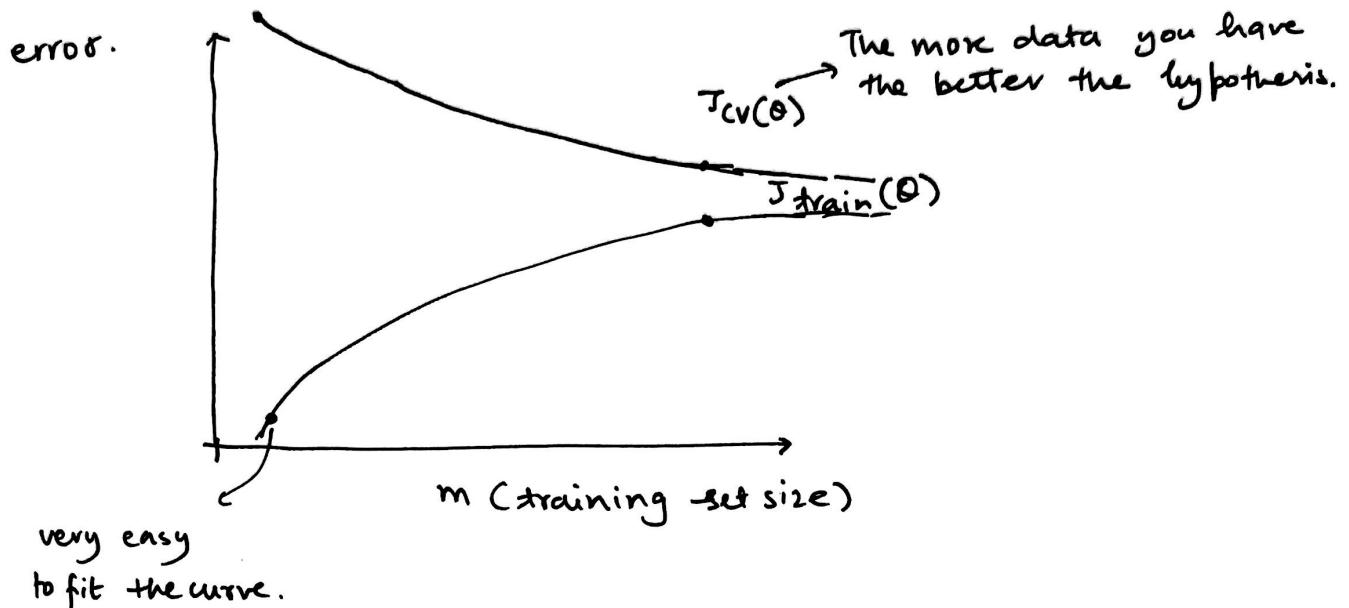
Steps to Model Selection:

1. Create a list of $\lambda \in \{0, 0.01, 0.02, \dots, 10.24\}$
2. Create a model ~~for~~ with different degrees
3. Iterate through the λ 's & for each λ compute θ^{learned}
4. Compute $J_{\text{cv}}(\theta^{\text{learned}})$ with $\lambda=0$ with diff degrees.
5. Select lowest $J_{\text{cv}}(\theta^{\text{learned}})$.
6. Using the best θ & λ , apply it on the $J_{\text{test}}(\theta)$ to see if it has good generalization of the problem.

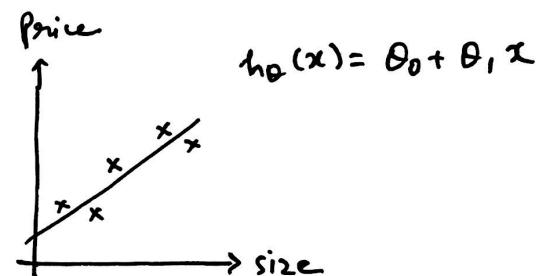
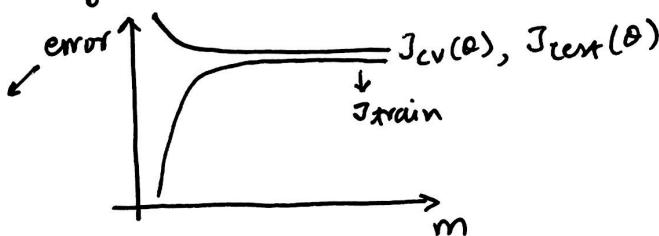
LEARNING CURVES

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

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

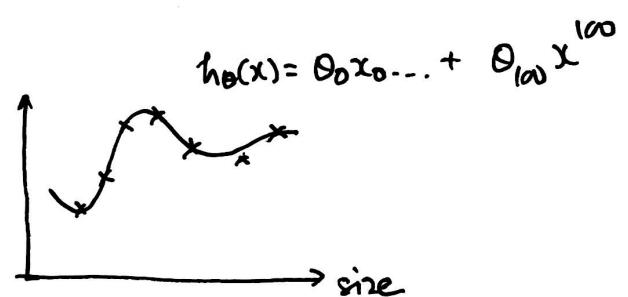
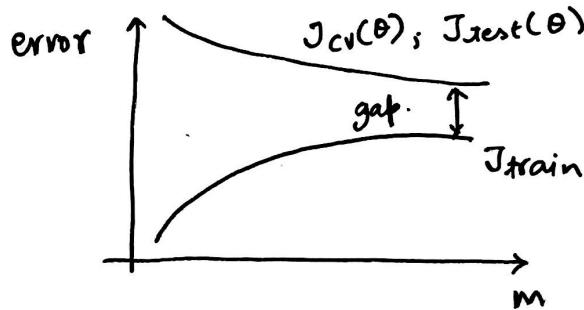


[#1] High Bias:



getting more data will not (by itself) help much.
training

[#2] High variance:



getting more ~~as~~ training data is likely to help.

Debugging a learning Algorithm

(75)

- ① Get more ~~training~~ training examples \rightarrow fixes high variance.
 - ② Try smaller sets of features \rightarrow fixes high variance.
 - ③ Try getting additional features. \rightarrow fixes high bias
 - ④ Try adding polynomial features \rightarrow fixes high bias.
 - ⑤ Try decreasing $\lambda \rightarrow$ fixes high bias.
 - ⑥ Try increasing $\lambda \rightarrow$ fixes high variance
-

EXAMPLE : BUILDING A SPAM CLASSIFIER :

what to prioritize

Non-spam (0)

Spam (1)

① Supervised learning:

$x = \text{features of email}$

$y \in \{0, 1\}$

{choose 100 words
indicative of spam/not spam}

Eg. deal, buy, discount, now, andrew, ...

$$x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \quad \begin{array}{l} \text{andrew} \\ \text{buy.} \\ \text{deal.} \\ \text{discount} \\ \text{now} \end{array} \quad x \in \mathbb{R}^{100}$$

Note: In practice, take most frequently occurring n words (10,000 to 50000) in training set, rather than manually pick 100 words

ERROR ANALYSIS

- ❖ Start with a simple algorithm that you can implement quickly. Implement it and test it on your CV data
- ❖ Plot learning curves to decide if more data, more features etc. are likely to help.
- ❖ Error Analysis:

Manually examine the examples that the algorithm made errors on (in your CV set). See if you spot any systematic trend in what type of examples it is making errors on. { if we develop new features by examining the test set then we may end up choosing features that work well w.r.t $\text{test}(\theta)$, and it no longer is a good est. }

Ex:

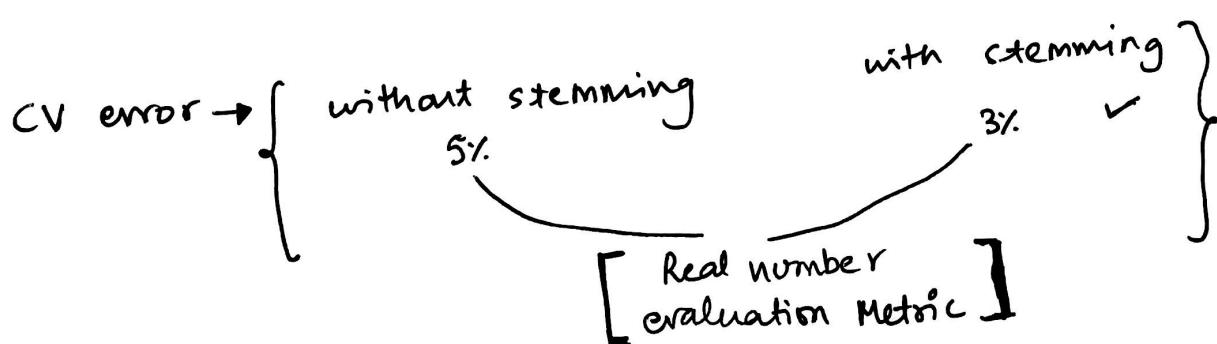
$$M_{CV} = 500$$

error = 100 emails.

what type of email it is?

what cues (features) you think would have helped the algorithm classify them correctly.

- ❖ 'stemming' Software : discount/discounted/discounting ✓
universe/university ✗



Handling Skewed data

Error metric for skewed classes.

① Precision/Recall:

$y=1$ in presence of rare class that we want to detect $\xrightarrow{\text{cancer}}$

Actual
Class.

Predicted
class

		1	0
1	True positive	False positive	
	False negative	True negative	

* Precision: of all patients where we predicted $y=1$, what fraction actually has cancer.

$$\left\{ \begin{array}{l} \frac{\text{True positives}}{\# \text{ predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{false pos.}} \end{array} \right\}$$

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

$$\left\{ \begin{array}{l} \frac{\text{True positives}}{\# \text{ Actual positive}} = \frac{\text{True Posit}}{\text{True post False neg.}} \end{array} \right\}$$

Trading off Precision & Recall :

→ logistic regression: $0 \leq h_\theta(x) \leq 1 = \begin{cases} 1 & \text{if } h_\theta(x) \geq 0.5 \\ 0 & \text{if } h_\theta(x) \leq 0.5 \end{cases}$

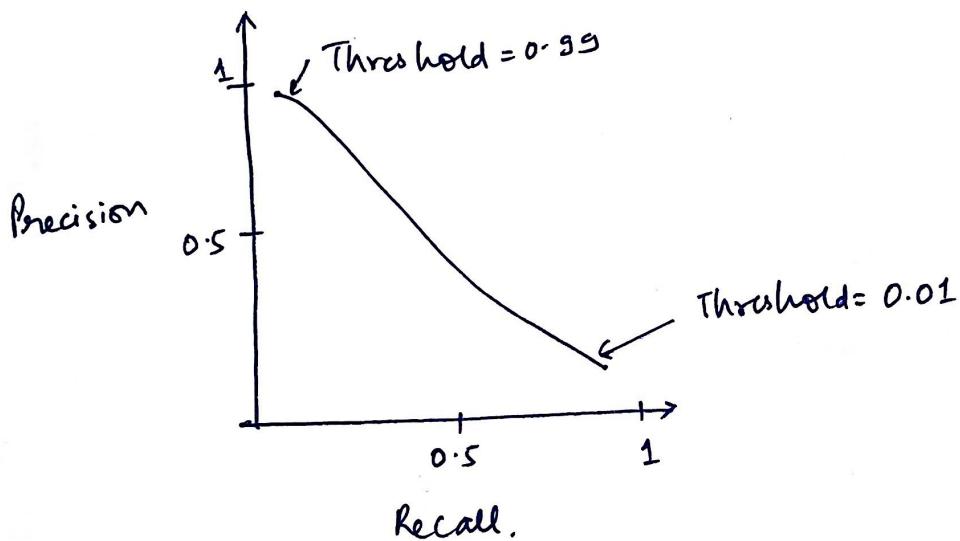
→ suppose we want to predict $y=1$ only if very confident

① if we change, $h_\theta(x) \geq 0.7 \Rightarrow$ Higher precision., lower recall.

② suppose we want to avoid missing too many cases of cancer (avoid false negative)

$$\left\{ \rightarrow \text{lower } h_\theta(x) < 0.3 \right\}$$

\Rightarrow higher recall, lower precision



(80)

F₁ Score : How to compare precision/recall numbers?

	(P)	(R)	Avg	F ₁ Score
Algorithm 1	0.5	0.9	0.75	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.51	0.0392

Predict y=1 all the time.

④ Avg. = $\frac{P+R}{2} \times$ (Not a good way)

⑤ F₁ Score = $\frac{2 * PR}{P+R}$