

decide what to try next

what to do if getting large errors in prediction?

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features ($x_1^2, x_2^2, x_1^3, x_2^3, \dots$)
- Try decreasing/increasing λ

Machine Learning diagnostic

diagnostic: A test 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

Evaluate a hypothesis

Split training data into training set and test set to avoid overfitting

Training / testing procedure for linear regression

- Learn parameter θ from training set (70% of training data, e.g.)
- Compute test set error

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

(logistic regression is similar)

- Misclassification error.

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

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

Model selection and training/validation/test sets

Model selection:

choose right degree of polynomial:

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

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

$$\vdots \quad | \quad \vdots$$

$$d. h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_d x^d \rightarrow \mathbb{H}^{(d)}$$

d : degree of polynomial

choose d to minimize the error $J(\theta)$ in test set is NOT right, Since d might just be optimized for this specific test set.,

usually choose:

training data $\left\{ \begin{array}{l} 60\% \text{ Training set } (x^{(i)}, y^{(i)}) \\ 20\% \text{ cross validation set } (CV) \quad (x_{cv}^{(i)}, y_{cv}^{(i)}) \\ 20\% \text{ test set } (x_{test}^{(i)}, y_{test}^{(i)}) \end{array} \right.$

Errors: $J_{train}(\theta)$, $J_{cv}(\theta)$, $J_{test}(\theta)$

choose the degree of polynomial d to get minimum $J_{cv}(\theta)$
then use the θ from same d to calculate $J_{test}(\theta)$
which gives right test set error

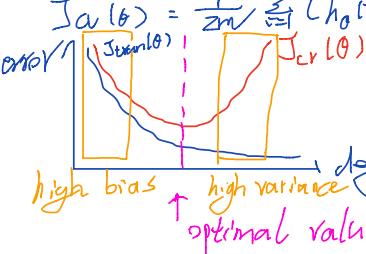
Bias vs. Variance

High bias: under fit

High variance: over fit

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

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



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

High variance (over fit): $J_{\text{train}}(\theta)$ will be low, $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$

Regularization and bias/variance

Linear regression with regularization

$$\begin{aligned} \text{Model: } h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3 + \theta_4 x_1^4 \\ J(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{m} \sum_{j=1}^m \theta_j^2 \end{aligned}$$

Large λ : $\lambda = 1000, \theta_1 \approx 0, \theta_2 \approx 0, \dots, \theta_4 \approx 0 \quad h_{\theta}(x) = \theta_0 \Rightarrow \text{high bias}$

Small λ : $\lambda = 0 \Rightarrow \text{high variance (over fit)}$

Intermediate $\lambda \Rightarrow \text{Just right}$

Choosing the regularization parameter λ

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

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

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

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

no regularization

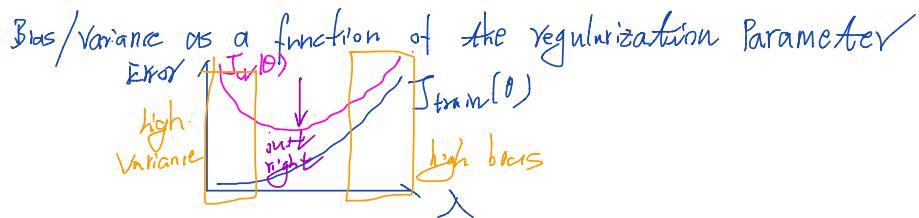
$$J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

J term involved

1. try $\lambda = 0 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(0)} \rightarrow J_{\text{cv}}(\theta^{(0)})$
2. try $\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$
3. try $\lambda = 0.02 \quad | \quad | \quad |$
4. try $\lambda = 0.04 \quad | \quad | \quad |$
5. try $\lambda = 0.08 \quad | \quad | \quad |$
- ⋮
12. try $\lambda = 10 \quad | \quad | \quad |$

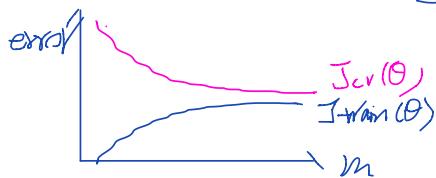
choose λ corresponding to smallest $J_{\text{cv}}(\theta)$, this is the proper value of λ

To see the error of this hypothesis, calculate $J_{\text{test}}(\theta)$

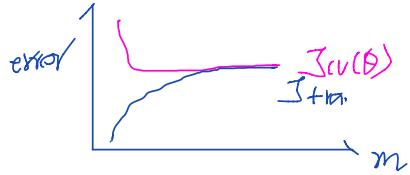


Learning Curves

Plot $J_{\text{train}}(\theta)$ and $J_{\text{cv}}(\theta)$ V.S. training set size m



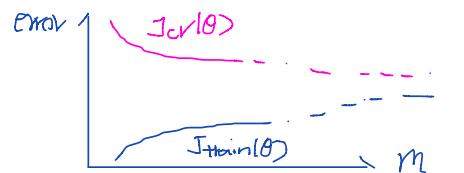
High bias



J_{cv} saturates pretty soon, $J_{train} \approx J_{cv}(\theta)$ by the end

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

High variance



* if a learning algorithm is suffering from high variance (overfitting)
, getting more training data is likely to help

Deciding what to try next - revisit:

debugging a learning algorithm

Suppose you implemented regularized linear regression, however when you apply this hypothesis on a new test set, you find it makes unacceptable large errors. What should you try next?

- Get more training examples
 - Try smaller sets of features
 - Try increasing λ
 - Try adding additional features
- $\left. \begin{array}{l} \\ \\ \end{array} \right\}$ fix high variance (overfitting) problems

- Try adding polynomial features
 - Try decreasing λ
- } fix high bias (underfitting) problems

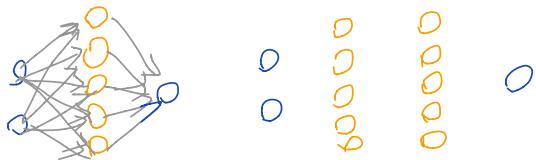
Neural Networks and overfitting

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

PREFERRED !!

Machine Learning System Design

Prioritizing what to work on: Spam classification example

Error Analysis

Recommended Approach

- Start with a simple algorithm that you can implement quick and dirty, 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 CV data set) where your algorithm made errors on, see if you spot any systematic trend in what type of examples it is making errors on.

Error metrics for skewed classes

Cancer classification example

Train logistic regression model $h(x)$, find that you got 1% error on test set (99% correct diagnostics)

BUT if in the training set only 0.5% of people have cancer, prediction $y=0$ (always) will give only 0.5% error, which seems better than "learned" predictions.



Skewed classes

Problem: 99.5% accuracy is better than 99% ??

Precision / Recall.

		<u>Actual class</u>	
		1	0
predicted class	1	True Positive	False Positive
	0	False negative	True negative

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

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

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

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

$$\frac{\text{True Pos}}{\# \text{ actual pos}} = \frac{\text{True Pos}}{\text{True pos} + \text{false neg}}$$

Trading off precision and recall

Normally:

Logistical regression: $0 \leq h_0(x) \leq 1$

predict 1 if $h_0(x) \geq 0.5$

predict 0 if $h_0(x) < 0.5$

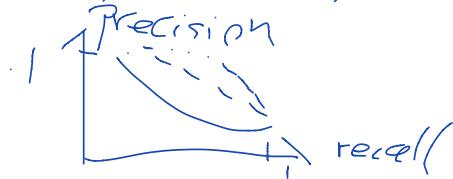
Suppose we want to predict $y=1$ only if very confident

→ higher precision, lower recall: $0.5 \rightarrow 0.7$

Suppose we want to avoid missing too many cases of cancer

→ higher recall, lower precision $0.5 \rightarrow 0.3$

Vary value of threshold \Rightarrow trade off precision vs recall



F₁ Score (on CV-set)

How to compare precision/^(P) recall numbers

Average: $\frac{P+R}{2}$ → not good

F₁ Score: $2 \frac{PR}{P+R}$ $P=0 \text{ or } R=0 \Rightarrow F_1=0$
 $P=1 \text{ & } R=1 \Rightarrow F_1=1$

Data for machine Learning

useful test: given the input x , can a human expert confidently predict y ? (x is features good enough?)

Large data rationale

- └ algorithm with many features \rightarrow low bias algorithms
- └ large training set \rightarrow unlikely to overfit
- └ $J_{\text{train}}(\theta)$ small
- └ $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

/ - - - - - /

Combine above two:

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

\Rightarrow have many features + Large training set