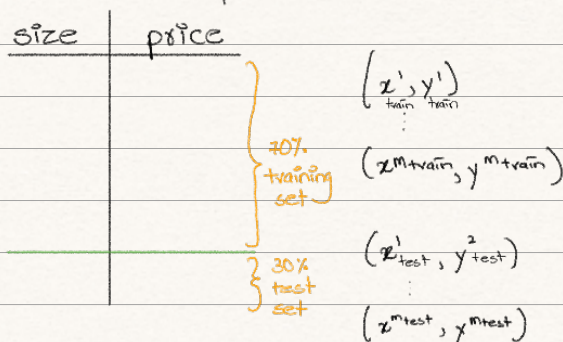


Subject: Advice For applying ML

//

Evaluate a model:

- Split data into parts



Linear Regression:

Minimize Cost:

$$J(\vec{w}, b) = \left[\frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m_{train}} \sum_{j=1}^n w_j^2 \right]$$

Compute test error:

$$J_{test}(\vec{w}, b) = \frac{1}{2m_{test}} \left[\sum_{i=1}^{m_{test}} (f_{\vec{w}, b}(\vec{x}_{test}^{(i)}) - y_{test}^{(i)})^2 \right]$$

Compute training error:

$$J_{train}(\vec{w}, b) = \frac{1}{2m_{train}} \left[\sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}_{train}^{(i)}) - y_{train}^{(i)})^2 \right]$$

Model Selection:

$d=1$ 1) $f_{\vec{w}, b} = w_1 x + b \longrightarrow J_{test}(w^{<1>}, b^{<1>})$

$d=2$ 2) $f_{\vec{w}, b} = w_1 x + w_2 x^2 + b \longrightarrow J_{test}(w^{<2>}, b^{<2>})$

\vdots

$d=10$ 10) $f_{\vec{w}, b} = w_1 x + w_2 x^2 + \dots + w_{10} x^{10} + b \longrightarrow J_{test}(w^{<10>}, b^{<10>})$

→ since, we are also using α as a parameter, then this α might overfit the test set

Subject:

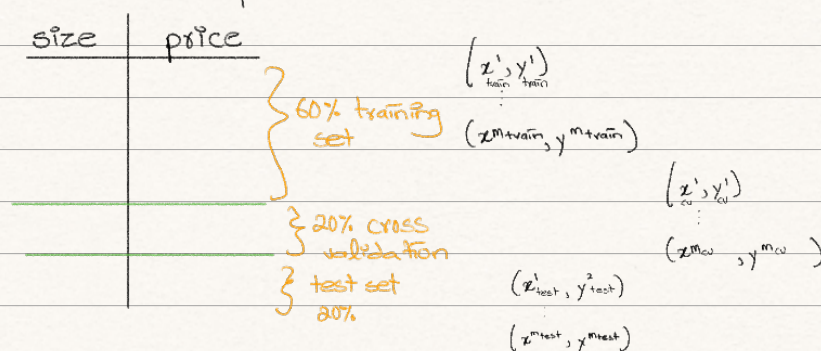
- Choose one with lowest J_{test}

→ This is NOT good as the model will overfit the test set and will NOT generalize truly

Training / cross validation / test set

→ validation / development / dev

- Split data into parts



Compute test error:

$$J_{\text{test}}(\vec{w}, b) = \frac{1}{2m_{\text{test}}} \left[\sum_{j=1}^{m_{\text{test}}} (f_{\vec{w}, b}(\vec{x}_{\text{test}}^{(j)}) - y_{\text{test}}^{(j)})^2 \right]$$

Compute training error:

$$J_{\text{train}}(\vec{w}, b) = \frac{1}{2m_{\text{train}}} \left[\sum_{j=1}^{m_{\text{train}}} (f_{\vec{w}, b}(\vec{x}_{\text{train}}^{(j)}) - y_{\text{train}}^{(j)})^2 \right]$$

Compute CV error:

$$J_{\text{cv}}(\vec{w}, b) = \frac{1}{2m_{\text{cv}}} \left[\sum_{j=1}^{m_{\text{cv}}} (f_{\vec{w}, b}(\vec{x}_{\text{cv}}^{(j)}) - y_{\text{cv}}^{(j)})^2 \right]$$

Subject: / /

Model Selection:

1) $f_{w,b} = w_1x + b \longrightarrow J_{cv}(w^{(1)}, b^{(1)})$

2) $f_{w,b} = w_1x + w_2x^2 + b \longrightarrow J_{cv}(w^{(2)}, b^{(2)})$

⋮

10) $f_{w,b} = w_1x + w_2x^2 + \dots + w_{10}x^{10} + b \longrightarrow J_{cv}(w^{(10)}, b^{(10)})$

- Choose one with lowest J_{cv} .
- Estimate generalization error using J_{test}

Diagnosing Variance and bias:

High bias (underfit)

"Just right"

High variance (overfit)

J_{train} is high

J_{train} is low

J_{train} is low

J_{cv} is high

J_{cv} is low

J_{cv} is high

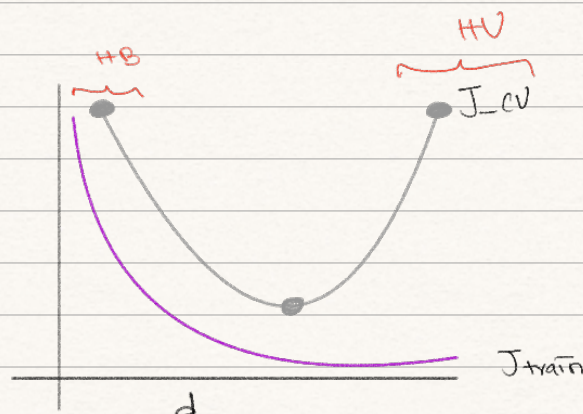
$J_{train} \approx J_{cv}$

$J_{cv} \gg J_{train}$

High variance & bias

J_{train} is high

$J_{cv} \gg J_{train}$



Subject: / /

Regularization and bias/Variance:

Model: $f_{\vec{w},b}(x) = w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$

$$J(\vec{w},b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w},b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m w_j^2$$

Large λ

small λ

Intermediate λ

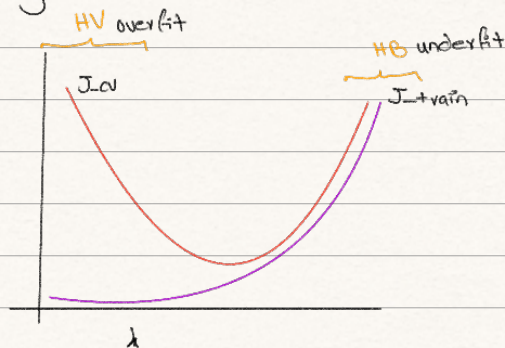
high bias (under fit)

High variance (over fit)

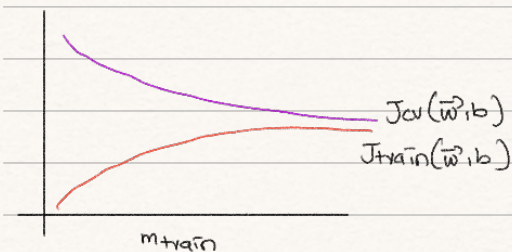
1. Try $\lambda=0 \rightarrow \min_{\vec{w},b} J(\vec{w},b) \rightarrow w^{(1)}, b^{(1)} \rightarrow J_w(w^{(1)}, b^{(1)})$

\vdots

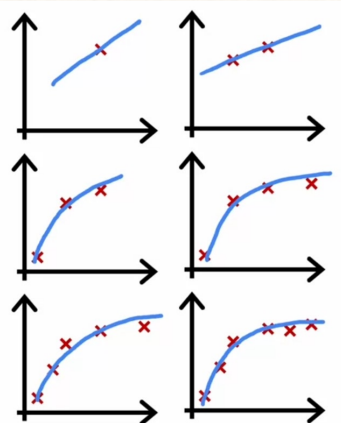
2. Try $\lambda \approx 10 \rightarrow \min_{\vec{w},b} J(\vec{w},b) \rightarrow w^{(12)}, b^{(12)} \rightarrow J_w(w^{(12)}, b^{(12)})$



Learning Curves:

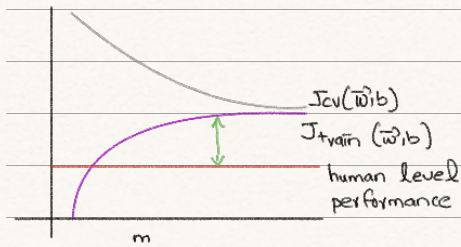


Why J_{train} gets bigger?

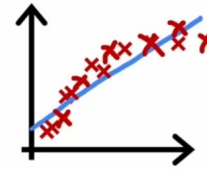
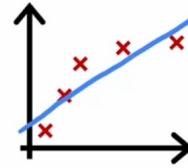


Subject:

High bias:

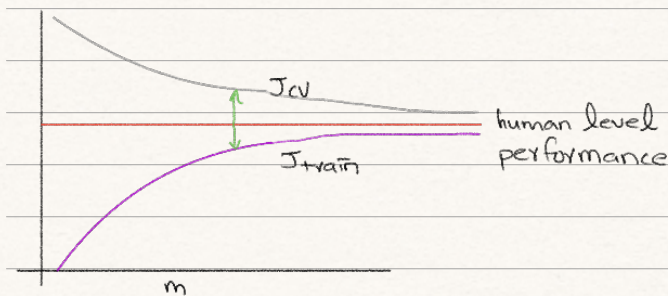


Why curves flatten?



Regardless of amount of data, a linear model is too simple & won't change much

High Variance:



• More data likely to help

What to do next?

High Bias:

- Try using more features
- Try adding polynomial features
- Decrease λ

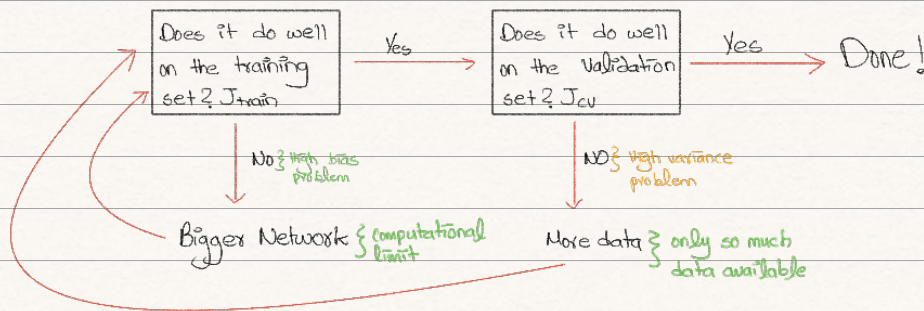
High Variance:

- Get more data
- Try using less features
- Increase λ

Subject: / /

Bias & variance in neural network:

Large networks trained on small/moderate data are low bias machines



• A larger Neural Network which is regularized correctly will perform just as good or better than smaller. So, it never hurts to go bigger!

Regularized:

layer_1 = Dense (25, activation = "relu", kernel_regularizer = L2 (0.01))

shows down the algo ☹️

Machine learning development process:

Error analysis:

$m_{cv} = 500$

Algo misclassifies 100 of them

Manually examine the 100 examples and group them based on common traits.

if misclassified amount is very large, try random sampling

Data Augmentation:

• Modifying existing training example to create a

Subject: / /

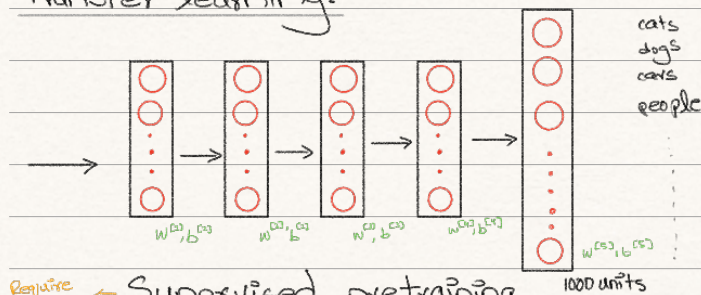
new training examples.

- Distortion introduced should be represented of the type of noise/distortions in the test set. $\left. \begin{array}{l} \text{Distortions in test set} \\ \text{must be meaningful} \end{array} \right\}$

Data synthesis:

- Create new data on your own.

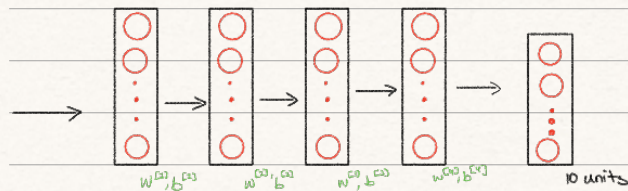
Transfer learning:



Require large data
1 million examples

Supervised pretraining

↳ easily available on internet



fine tuning

You want to train model to recognize digits. But you don't have data.

Train model on some other data like recognize dogs, cats etc. Now, copy the neural network &

use all the previous parameters but train last layer.

Option 1: only train output layer's parameters.

Option 2: Train all parameters. Initialize with above nn parameters.

Why does this work?

layer 1: detect edges

layer 2: detect corners

layer 3: detect curves/basic shapes

} learn generic features

- Use the same input type.

Subject: / /

Skewed datasets:

- Ratio of + & - is NOT 50/50

Precision/Recall:

$y=1 \rightarrow$ presence of a rare class

Confusion matrix:

		Actual class	
		1	0
Predicted class	1	15	5
	0	10	70

Annotations:
- True positive: 15
- False positive: 5
- False negative: 10
- True negative: 70

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

↳ Of those we predicted as $y=1$, how many actually has the disease?

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

↳ Of those that have the disease, how many did we correctly predict?

Trading off Precision & Recall:

Raising threshold increases precision but lower recall and vice versa

F1 score:

How to compare precision/recall numbers?

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

Subject:	P	R	Avg	F1 Score	/ /
Algo 1	0.5	0.4	0.45	0.444	→ Better of the three
Algo 2	0.7	0.1	0.4	0.175	
Algo 3	0.02	1.0	0.501	0.0392	

$$F1 \text{ score} = \frac{1}{\frac{1}{P} + \frac{1}{R}} = 2 \frac{PR}{P+R}$$