

Debugging a Learning Algo

If we use a certain formula and get unacceptably large errors in prediction we can do these steps:

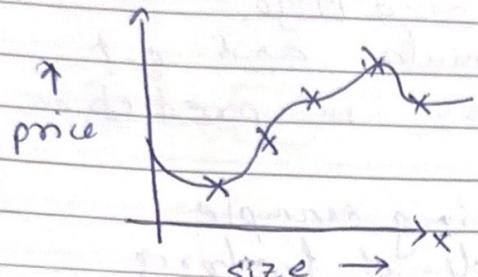
- Get more training example
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features
- Try increasing λ
- Try decreasing λ

But in reality, when applying these things on a DL model, some may turn out to be fruitful while others do not.

To help to choose / improve a DL, we run a diagnostic

Diagnostics take a lot of time to implement, but reduce the risk of spending time with no thing which do not help

Evaluating your Model



Model fits training data well but will fail to generalize to new examples not in training set.
(Overfitting)

$$x = x_1 \\ f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + \dots + w_n x^n + b.$$

here we are using only 1 feature and can see its not good.

x_1 = size in ft^2 x_2 = no. of bedrooms x_3 = no. of floors x_4 = age of home in years

We can evaluate our model by dividing our data set

Features	result
	$\left. \begin{array}{l} \text{training set} \\ \{x^1, y^1 \\ x^2, y^2 \} \end{array} \right\}$
	$\left. \begin{array}{l} \text{test set} \\ \{x^{m+\text{train}}, y^{m+\text{train}} \\ (x_{\text{test}}, y_{\text{test}}) \\ (x_{\text{test}}^{m+\text{test}}, y_{\text{test}}^{m+\text{test}}) \} \end{array} \right\}$

* Ex. Train/test procedure for linear model
Regression.

Fit Parameters

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

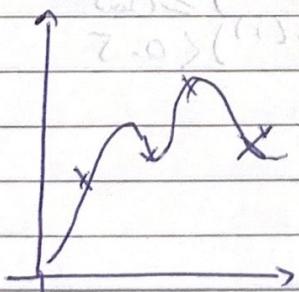
Compute test error:

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

Compute training error:

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

Both test & training error does not contain the regularization term.



$J_{\text{train}}(\vec{w}, b)$ will be low
 $J_{\text{test}}(\vec{w}, b)$ will be high

Train/Test procedure for classification problem

Fit parameters

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

Compute test/train error is same as fit parameters without the regularization term.

There is another that is commonly used for classification problems.

→ fraction of test set & fraction of training set the algorithm has misclassified

$$\hat{y} = \begin{cases} 1 & \text{if } f_{\vec{w}, b}(\vec{x}^{(i)}) \geq 0.5 \\ 0 & \text{if } f_{\vec{w}, b}(\vec{x}^{(i)}) < 0.5 \end{cases}$$

and count $\hat{y}(t_y)$

Model Selection

Once parameters \vec{w}, b are fit to the training set, the training error is most likely lower than the actual generalization error.

$J_{\text{test}}(\vec{w}, b)$ is a better estimate of how well the model will generalize to new data compared to $J_{\text{train}}(\vec{w}, b)$.

order

$$d=1 \quad 1. J_{\vec{w}, b} = w_1 x + b \rightarrow \vec{w}^{(1)}, b^{(1)} \rightarrow J_{\text{test}}(\vec{w}^{(1)}, b^{(1)})$$

$$d=2 \quad 2. J_{\vec{w}, b} = w_1 x + w_2 x^2 + b \rightarrow \vec{w}^{(2)}, b^{(2)} \rightarrow J_{\text{test}}(\vec{w}^{(2)}, b^{(2)})$$

$$d=3 \quad 3. J_{\vec{w}, b} = w_1 x + w_2 x^2 + w_3 x^3 + b$$

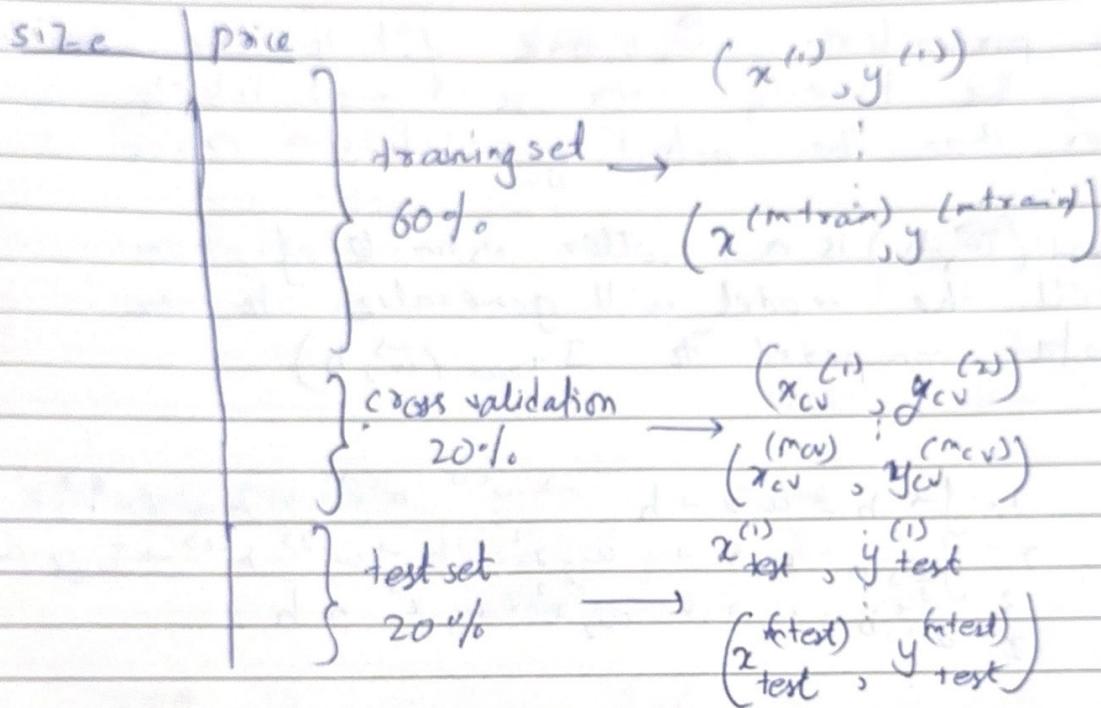
$$d=10 \quad 10. J_{\vec{w}, b} = w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_{10} x^{10} + b$$

We find the error of test set for each polynomial and take the lowest error model.

This method is flawed as

$J_{\text{test}} <$ generalization error and J_{test} is an optimistic estimate of generalization error

Model Selection - ML



Training error:

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

Cross validation error:

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

Test error:

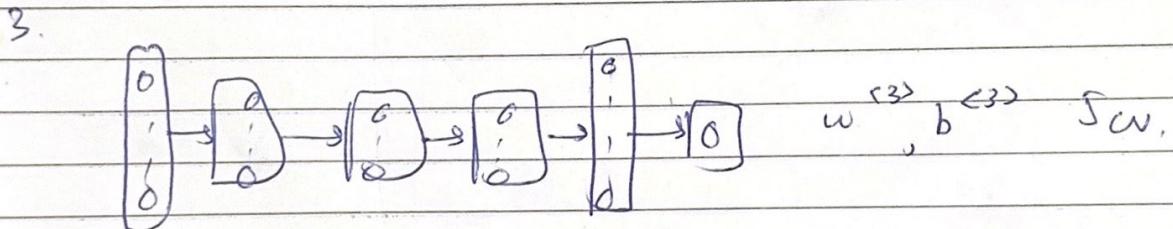
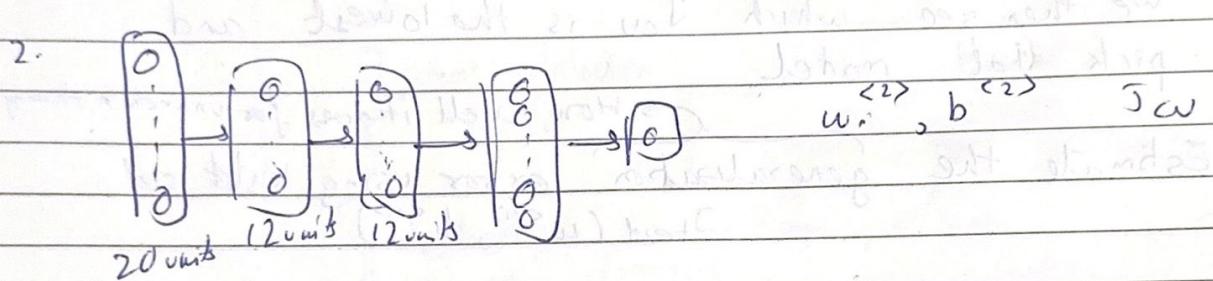
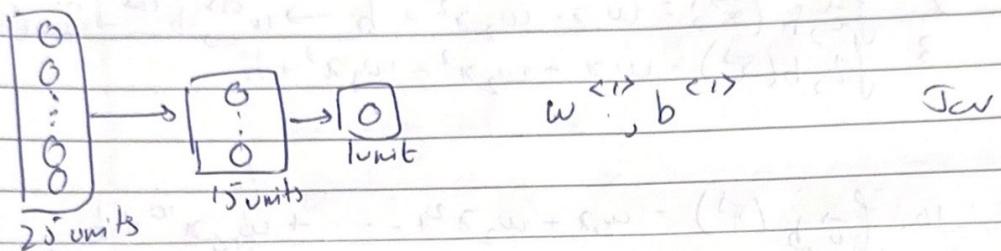
$$J_{\text{test}}(\vec{w}, b) = \frac{1}{2m_{\text{test}}} \left[\sum_{i=1}^{m_{\text{test}}} (f_{\vec{w}, b}(\vec{x}_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2 \right]$$

$$\begin{array}{ll}
 d=1 & 1. f_{\vec{w}, b}(\vec{x}) = w_1 x + b \rightarrow w^{<1>} b^{<1>} \rightarrow J_{CV}(w^{<1>} b^{<1>}) \\
 d=2 & 2. f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + b \rightarrow w^{<2>} b^{<2>} \rightarrow J_{CV}(w^{<2>} b^{<2>}) \\
 d=3 & 3. f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + w_3 x^3 + b \\
 & \vdots \\
 d=10 & 10. f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + \dots + w_{10} x^{10} + b
 \end{array}$$

We then see which J_{CV} is the lowest and pick that model

Estimate the generalization error using test set
 $J_{test}(w^{<4>} b^{<4>})$

Model Selection - NN

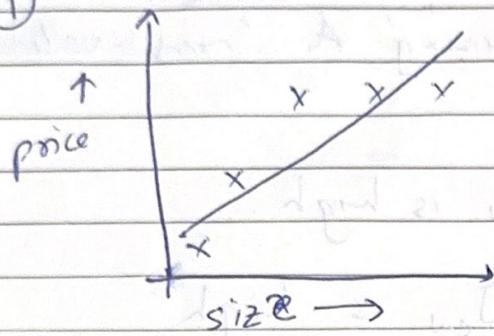


Pick model with lowest J_W .

Estimate generalization error using
 $J_{\text{test}}(w, b)$

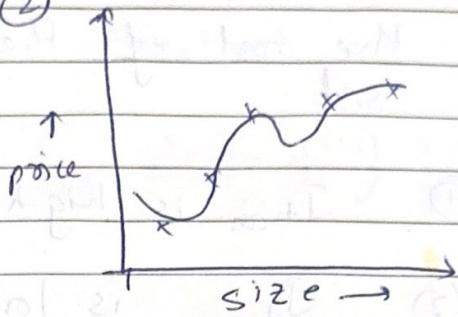
Bias & Variance

①



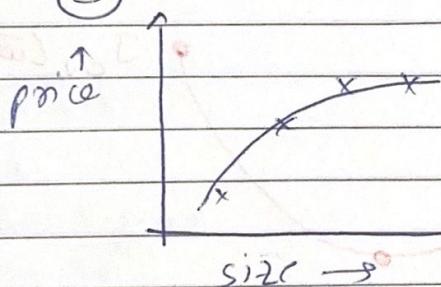
high bias
(underfit)

②



high variance
(overfit)

③



$$f(\vec{w}, b)(x) = w_1 x + w_2 x^2 + b$$

(just right)

This f is for 1 variable, so it is easy to plot.
we need another method for multiple features

→ add more weight (feature) or make depth

→ smooth level

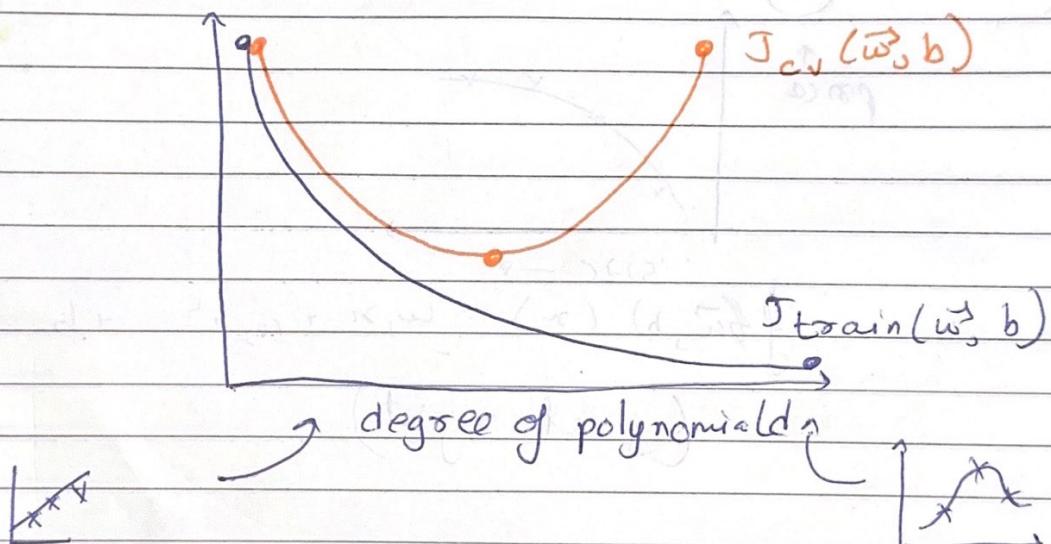
→ added a part of weight & dimension word with

To check the performance, we check the cost of the 'training' & 'cross-validation' set.

$d=1$ ① J_{train} is high, J_{cv} is high.

$d=4$ ② J_{train} is low, J_{cv} is high.

$d=2$ ③ J_{train} is low, J_{cv} is low.



\therefore High bias (underfit) $\rightarrow J_{\text{train}}$ will be high
 $(J_{\text{train}} \gg J_{\text{cv}})$

High variance (overfit) $\rightarrow J_{\text{train}}$ may be low

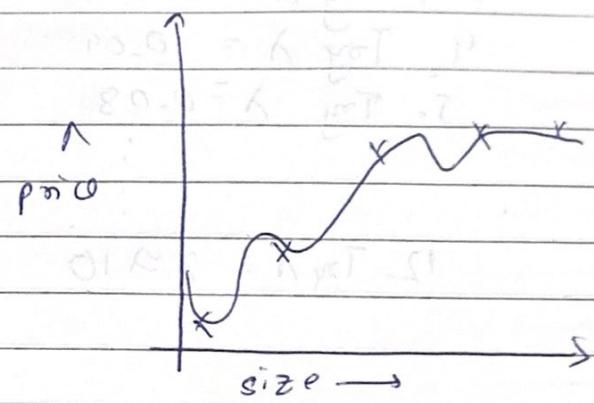
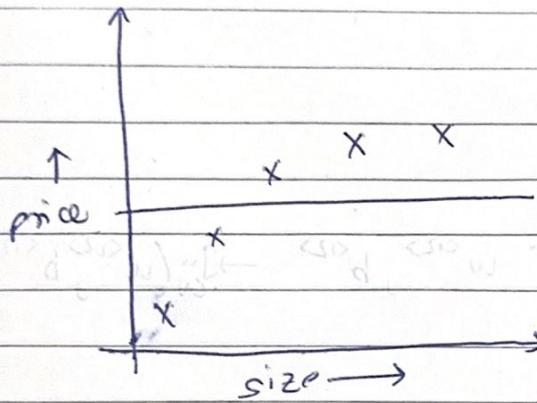
$$J_{\text{cv}} \gg J_{\text{train}}$$

High bias & variance $\rightarrow J_{\text{train}}$ is high & $J_{\text{cv}} \gg J_{\text{train}}$

Regularization & bias/variance.

Model: $f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

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



larger the value λ more
motivated will be the
regularization term to
keep \vec{w} small

small λ , less
general will be the
final graph, or
less motivation to
make it generalized

High bias (underfit)

$J_{\text{train}}(\vec{w}, b)$ is large
 $J_{\text{cv}}(\vec{w}, b)$ is high

High variance
(overfit)

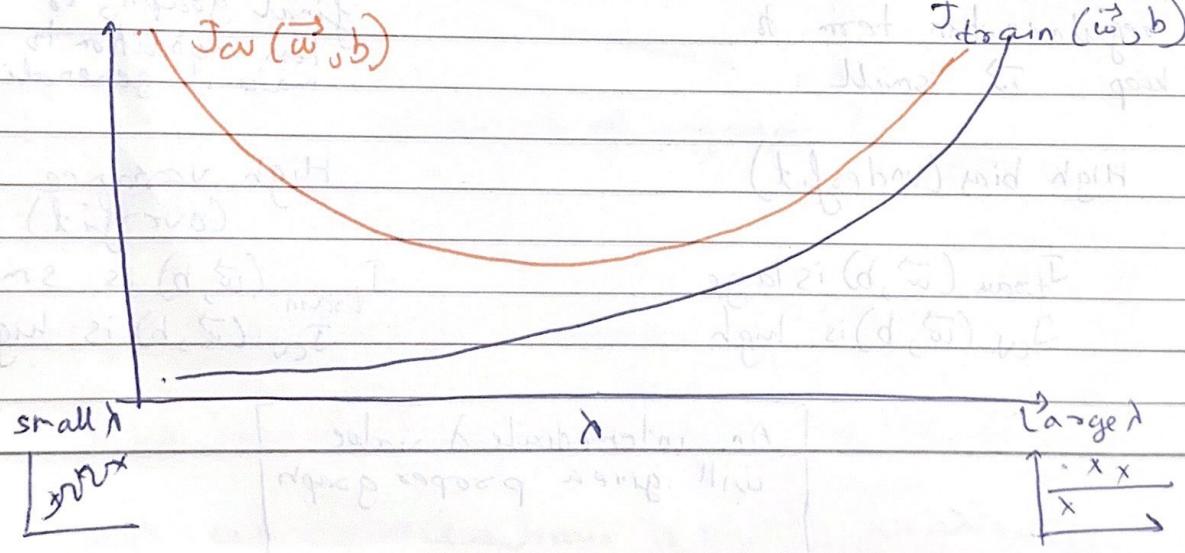
$J_{\text{train}}(\vec{w}, b)$ is small
 $J_{\text{cv}}(\vec{w}, b)$ is high

An intermediate λ value
will give a proper graph

To choose the proper λ value, we take various values for λ and the equation with the least J_{CV} will give the proper λ value.

- Ex. 1. Try $\lambda = 0 \rightarrow \min_{\vec{w}, b} J(\vec{w}, b) \rightarrow \vec{w}^{(1)}, b^{(1)} \rightarrow J_{CV}(\vec{w}^{(1)}, b^{(1)})$
 2. Try $\lambda = 0.01 \rightarrow \vec{w}^{(2)}, b^{(2)} \rightarrow J_{CV}(\vec{w}^{(2)}, b^{(2)})$
 3. Try $\lambda = 0.02 \rightarrow \vec{w}^{(3)}, b^{(3)} \rightarrow J_{CV}(\vec{w}^{(3)}, b^{(3)})$
 4. Try $\lambda = 0.04 \rightarrow \vec{w}^{(4)}, b^{(4)} \rightarrow J_{CV}(\vec{w}^{(4)}, b^{(4)})$
 5. Try $\lambda = 0.08 \rightarrow \vec{w}^{(5)}, b^{(5)} \rightarrow J_{CV}(\vec{w}^{(5)}, b^{(5)})$
 12. Try $\lambda = 0.2$ $\rightarrow \vec{w}^{(12)}, b^{(12)} \rightarrow J_{CV}(\vec{w}^{(12)}, b^{(12)})$

We take the best λ value and calculate generalization error using J_{test}



Summary: If baseline loss is equal/near to training cost, & high loss diff b/w training & cv then overfitting. If loss difference b/w baseline & test is high, then overfitting. If loss difference b/w baseline & test is low, then underfitting.

Whenever evaluating an algorithm, it is always good to measure it against another performance like:

- Human level performance
 - Competing algorithms performance
 - Guess based on experienced

Ex.

$$\textcircled{1} \quad \text{Human error} = 10.0\%$$

② Training error = 10.8%

③ Cross validation error = 14.8%

If we only see (2) $k \oplus$, we see underfitting, but checking (1), we see that it actually is an overfitting error.

Ex	Human error	J_{train}	J_{cv}	He	Te	Cve
		10.64%	10.64%	\uparrow	\downarrow	\downarrow
	Training errors (J_{train})	10.84%	10.84%	\uparrow	\downarrow	\downarrow
	Cross Validation Errors (J_{cv})	14.84%	14.84%	\downarrow	\downarrow	\downarrow

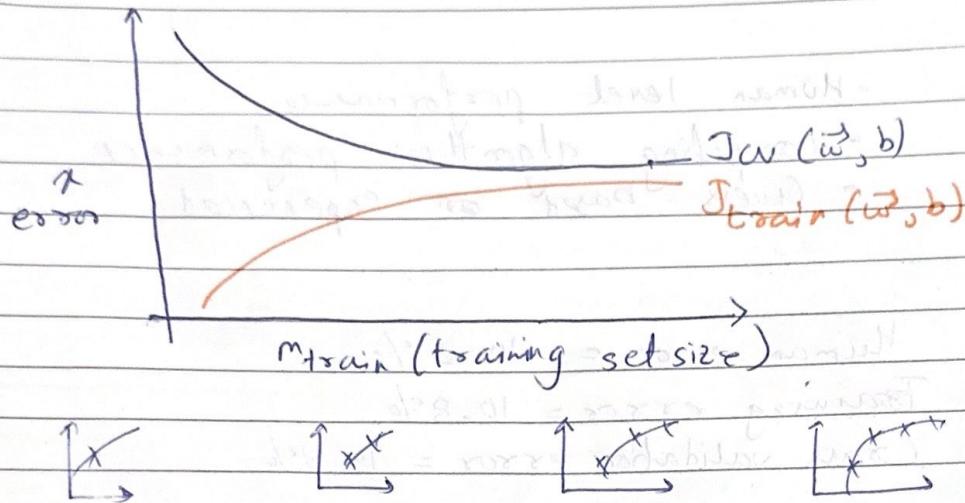
He \rightarrow 10.6 ↑
 Te \rightarrow 15.0 ↓ high variance
 Ce \rightarrow 19.7 ↓ high bias.

Learning Curves

Let's take $f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + b$.

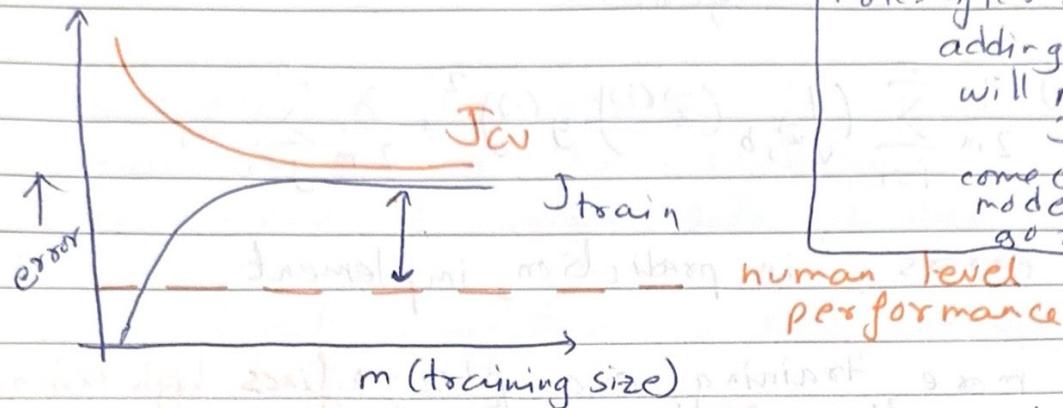
J_{train} = training error

and J_{cv} = cross validation error



If we have small no. of training examples, it is easy to keep the error low. However, if we increase the number of examples, it become more tricky to keep error low.

High Bias (Underfitting)

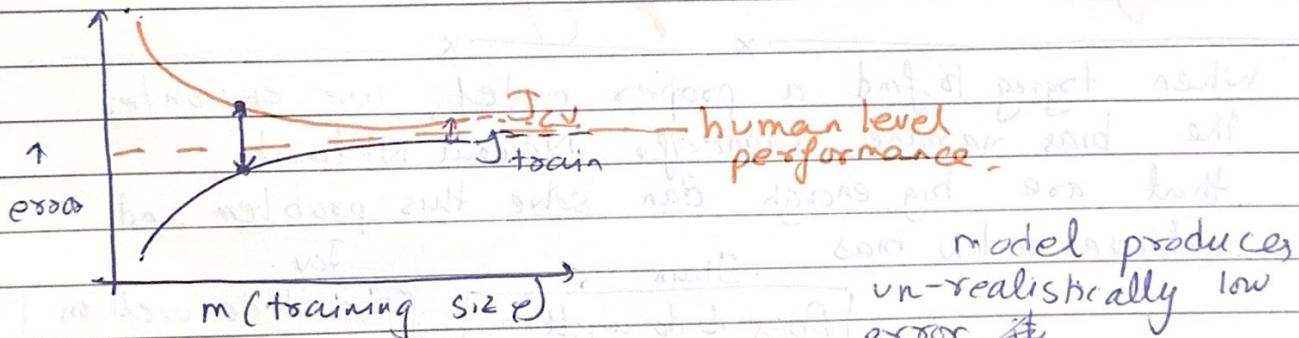


Note: After a point, adding more data will not make J_{CV} & J_{train} converge as the model cannot go to fit those examples.

Note:

If a learning algorithm suffers from high bias, getting more training data will not (by itself) help much.

High Variance (Overfitting)



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

Debugging a Learning Algorithm.

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

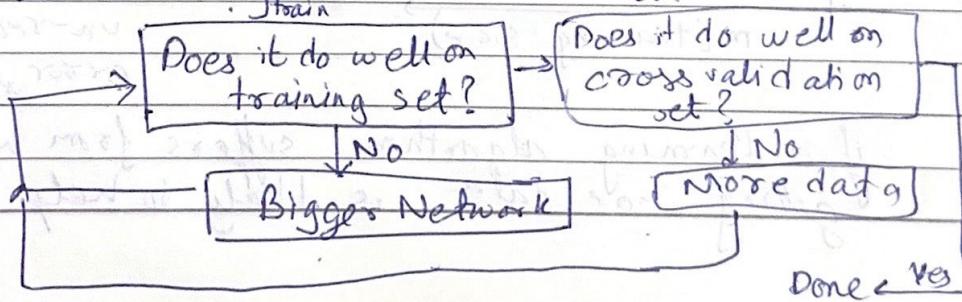
If it makes errors in prediction, implement

- HV • Get more training examples → fixes high variance
- HV • Try smaller feature set → fixes high variance
- HB • Try getting more features → fixes high bias
- HB • Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc)
- HB • Try decreasing λ → fixes high bias ↵
- HV • Try increasing λ → fixes high variance

HB → high bias = underfitting

HV → high variance = overfitting

When trying to find a proper model, we encounter the bias variance tradeoff. Neural Network that are big enough can solve this problem and achieve a low bias



Draw back is that bigger the network, more strain is put on the hardware
The number of unique data is also limited.

To prevent overfitting in NN, we use regularization.

In tensorflow, we have to add 'kernel regularizer' to the layer declaration.

layer1 = Dense (units=25, activation='selu',
kernel_regularizer=L2(0.0))