

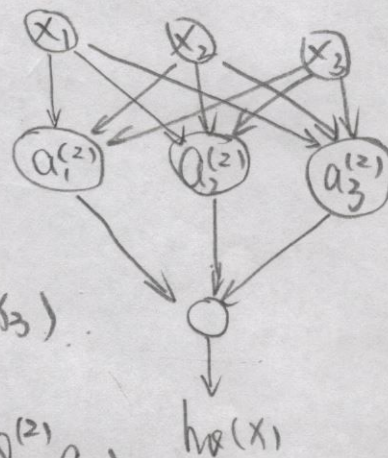
Neural Networks.

If network has S_j units in layer j , S_{j+1} units in layer $j+1$,

then $Q^{(j)}$ will be of dimension $S_{j+1} \times (S_j + 1)$

where $Q^{(j)}$ is the matrix of weights controlling function mapping from layer j to layer $j+1$.

$$\begin{cases} a_1^{(2)} = g(\theta_{10}^{(1)} \cdot x_0 + \theta_{11}^{(1)} \cdot x_1 + \theta_{12}^{(1)} \cdot x_2 + \theta_{13}^{(1)} \cdot x_3) \\ a_2^{(2)} = g(\theta_{20}^{(1)} \cdot x_0 + \theta_{21}^{(1)} \cdot x_1 + \theta_{22}^{(1)} \cdot x_2 + \theta_{23}^{(1)} \cdot x_3) \\ a_3^{(2)} = g(\theta_{30}^{(1)} \cdot x_0 + \theta_{31}^{(1)} \cdot x_1 + \theta_{32}^{(1)} \cdot x_2 + \theta_{33}^{(1)} \cdot x_3) \end{cases}$$



$$h_o(x) = g(\theta_{10}^{(2)} \cdot a_0^{(2)} + \theta_{11}^{(2)} \cdot a_1 + \theta_{12}^{(2)} \cdot a_2 + \theta_{13}^{(2)} \cdot a_3)$$

$$X^{(1)} = \begin{bmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{bmatrix}, \quad \theta^{(1)} = \begin{bmatrix} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \\ \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \\ \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{bmatrix}$$

$$a^{(2)} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = g(\theta^{(1)} \cdot X^{(1)}) \Rightarrow a = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix}, \quad \theta^{(2)} = \begin{bmatrix} \theta_{10}^{(2)} & \theta_{11}^{(2)} & \theta_{12}^{(2)} & \theta_{13}^{(2)} \end{bmatrix}$$

$$a_m^{(j+1)} = g(\theta_{m0}^{(j)} \cdot a_0^{(j)} + \theta_{m1}^{(j)} \cdot a_1^{(j)} + \dots)$$

$$h_o(x) = g(\theta^{(2)} \cdot a)$$

Applications:

- i). And. e.g: $h_o(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$, where $\theta_0 = -1.5, \theta_1 = \theta_2 = 1$.
- ii). Or. e.g: $h_o(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$, where $\theta_0 = -0.5, \theta_1 = \theta_2 = 1$.
- iii). Not. e.g: $h_o(x) = g(\theta_0 + \theta_1 x)$, where $\theta_0 = 0.5, \theta_1 = -1$.
- iv). XOR: $(\neg x_1 \wedge x_2) \vee (x_1 \wedge \neg x_2)$

Multiple Output Units:

$$\text{eg } y^{(i)} \in \left\{ \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \right\}$$

Cost Function of Neural Network:

$h_{\theta}(x) \in \mathbb{R}^K$, $(h_{\theta}(x))_i = i^{\text{th}}$ output, K output units.

m : the # of training examples, L : total # of layers in network.

S_L : # of units (excluding the bias unit) in layer L .

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K (y_k^{(i)} \cdot \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \cdot \log(1 - (h_{\theta}(x^{(i)}))_k)) \right] \\ + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2}_{\text{penalize all } \theta} \quad K \text{ is the \# of outputs units.}$$

$$\min_{\theta} J(\theta)$$

Need code to compute:

i). $J(\theta)$

ii). $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$.

Backpropagation algorithm:

1. Create a feed-forward network
2. Initialize all network weights to small random number
3. Until the termination condition is met, do
for $d \in D$ do

i). compute output o of every unit.

ii). for each output unit k , its error term is:

$$\delta_k = o_k(1 - o_k) \cdot (t_k - o_k)$$

% t_k is from y .

iii). for each hidden unit h ,

$$\delta_h = o_h(1 - o_h) \cdot \sum_{k \in \text{outputs}} w_{kh} \cdot \delta_k$$

iv). update all w : $w_{ji} = w_{ji} + \eta \cdot \delta_j \cdot x_{ji}$

Reshape function in matlab/octave.

Suppose: $\theta_1 \in \mathbb{R}^{10 \times 11}$, $\theta_2 \in \mathbb{R}^{10 \times 11}$, $\theta_3 \in \mathbb{R}^{1 \times 11}$

$$\text{thetaVec} = [\theta_1(:); \theta_2(:); \theta_3(:)]$$

To get $\theta_1, \theta_2, \theta_3$ back from thetaVec, we can Apply "reshape" function.

$$\text{e.g.: } \theta_1 = \text{reshape}(\underbrace{\text{thetaVec}(1:110)}_{\text{range}}, \underbrace{(10)}_{\text{row}}, \underbrace{(11)}_{\text{column}})$$

Gradient Checking:

$\theta \in \mathbb{R}^n$, the unrolled version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}, \dots$

$$\theta = [\theta_1, \theta_2, \dots, \theta_n]$$

then we have:

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \dots, \theta_n)}{2\epsilon}$$

\vdots

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_n - \epsilon)}{2\epsilon}$$

Apply this properties to check if the gradient get from backpropagation is correct.

- e.g.:
- Implement backpropagation to compute DVec (unrolled $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$).
 - Implement numerical gradient check to compute gradApprox.
 - Make sure they give similar values.
 - Turn off gradient checking, and using backpropagation for learning.

Random initialization to break symmetry.

Training a neural network.

- i). # of input units: dimension of features $x^{(i)}$
- ii). # of output units: number of classes.
- iii). Reasonable default: one or more hidden layer, and have same # of hidden units in every layer (usually the more the better).

Advice for Applying Machine Learning.

What you should try next if your hypothesis makes large errors on a new dataset?

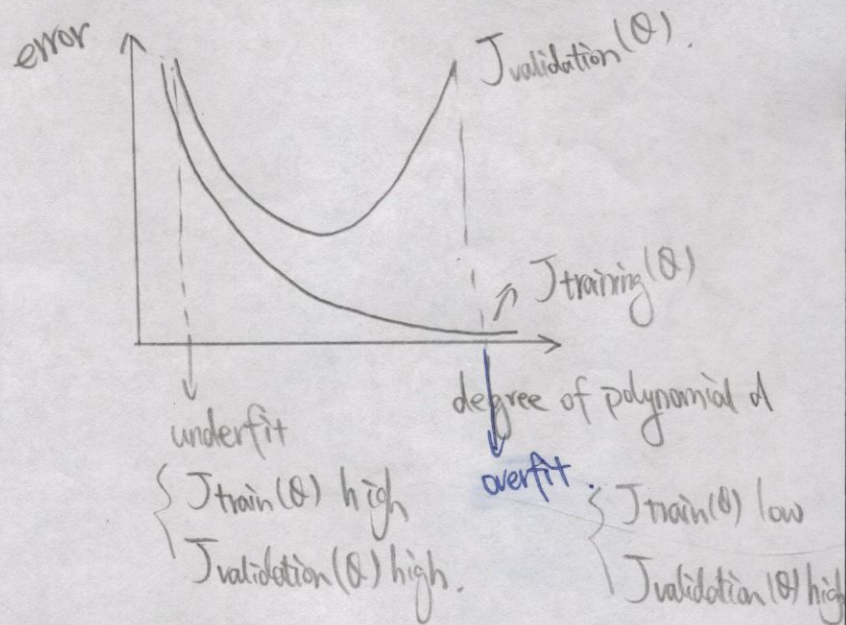
- i). Get more training examples.
 - ii). Try smaller sets of features.
 - iii). Try getting additional features.
 - iv). Try adding polynomial features.
 - v). Try decreasing λ or increasing λ .
- } fix overfitting (high variance)
- } fix high bias (underfitting).
- fix ~~underfitting~~ overfitting.

Evaluating your hypothesis: 60% Dataset for training,
20% Dataset for validation,
20% Dataset for testing.

Bias / Variance:

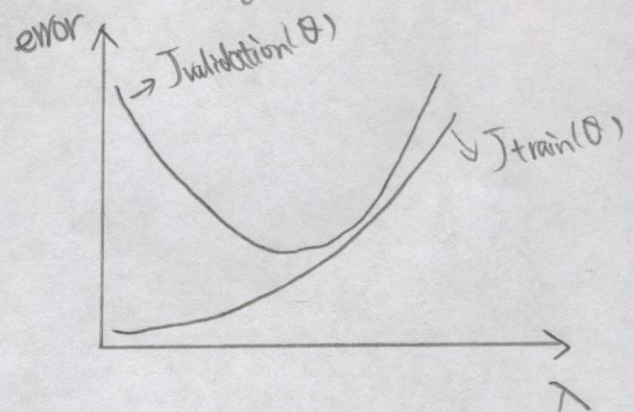
High bias \rightarrow underfit.

High variance \rightarrow overfit.



Linear Regression with Regularization:

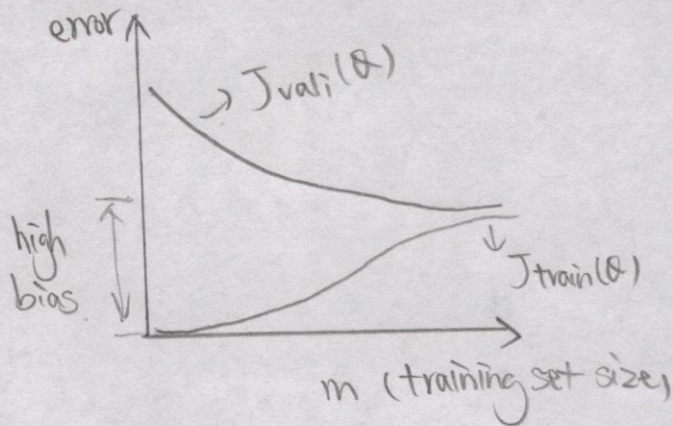
- i) Large $\lambda \rightarrow$ high bias, underfit.
- ii) Too small $\lambda \rightarrow$ high variance, overfit



Learning Curve.

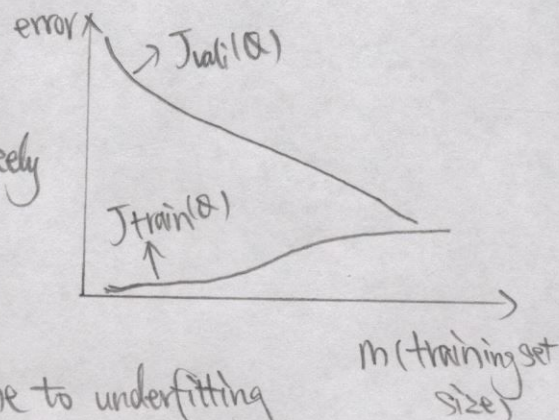
i). high bias.

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



ii). High variance

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



Small neural network: fewer parameters, more prone to underfitting.

Large neural network: more parameters, more prone to overfitting.

Recommended Approach:

1. Implement an algorithm and test it on your cross-validation data.
2. Plot learning curves to decide if underfit or overfit.
3. Error Analysis: manually check the examples your algorithm made error on to see what type of errors.

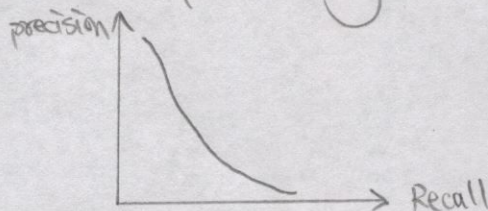
Error matrices for skewed classes.

Actual class			
	1	0	
predicted class	1	True pos	False Pos
	0	False neg	True neg

i). Precision: $\frac{\text{True pos}}{\# \text{pred Pos}} = \frac{\text{True Pos}}{\text{True pos} + \text{False Pos}}$

ii). Recall: $\frac{\text{True pos}}{\text{Actual Pos}} = \frac{\text{True Pos}}{\text{True Pos} + \text{False neg}}$

Trade off: predict 1 if $h_0(x) \geq \text{threshold}$.



How to compare precision/recall numbers?

\bar{F}_1 score : $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$, maximize \bar{F}_1 score.

Data for machine learning:

i). many parameters. \rightarrow low bias

ii). large training set. \rightarrow to avoid overfitting.

Combine i) + ii), always get low $J_{\text{test}}(\theta)$.