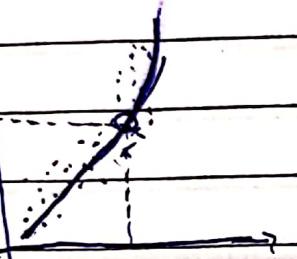


are main make em' valid. While I'm all well about a wise words with just weaving up words on my workshift. X-ray, my x-radiation holes, x-height letters the xylophones tones. Yellow back & out are yesterday's lawn yard sell our yarn. zigzag zombies, zoom into the zenith. zero in ZEBRAS!

8

Linear Regression



$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n \quad | \quad x_0 = 1 \text{ & } x_i = \text{feature i-th column}$$

$$\text{cost}(h_\theta(x)) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \quad | \quad x^{(i)} \& y^{(i)} \rightarrow i^{\text{th}} \text{ row/training example}$$

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad | \quad \theta_j = \text{Feature j-parameter}$$

VECTORIZED IMPLEMENTATION

Given matrixes ~~of size~~ ~~size~~ ~~matrix~~ ~~matrix~~

$m \rightarrow$ no. of examples

$n \rightarrow$ no. of features

- $[X]^{m \times (n+1)}$
- $[y]^{m \times 1}$
- $[\theta]^{(n+1) \times 1}$

Hypothesis = $X^* \theta$; % Gives a $m \times 1$ vector.

Cost($h_\theta(X)$) = sum((hypothesis - y).^2) / (2*m); % 1x1 var.

Gradient = $(1/m) * (X^* (\text{hypothesis} - y))$; % Returns a $(n+1) \times 1$ vector

GRADIENT DESCENT

$\theta = \theta - ((\alpha) * \text{gradient})$

$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

LOGISTIC REGRESSION:

Sigmoid function: $g(z) = \frac{1}{1+e^{-z}}$



Hypothesis:

$$h_\theta(x) = g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n) \text{ or } g(\theta^\top x)$$

Decision Boundary:

$$x = 0$$

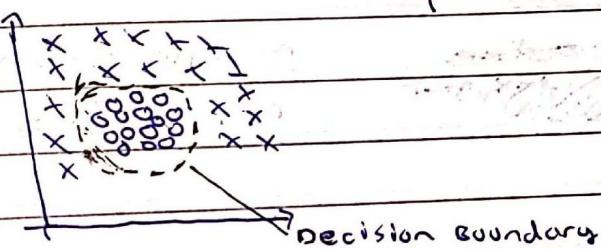
$$x = 1$$

$$g(x) < 0.5 \text{ or } g(x) > 0.5$$

$$g(x) \geq 0.5 \text{ or } g(x) > 0.5$$

$$\text{when, } \theta^\top x < 0 \text{ or } \theta^\top x \leq 0$$

$$\text{When, } \theta^\top x \geq 0 \text{ or } \theta^\top x > 0$$



Cost function: ~~cost~~ cost($\theta^\top x$)

$$\text{cost}(h_\theta(x)) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))]$$

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

$$[x]^{m \times (n+1)}, [y]^{m \times 1}, [\theta]^{(n+1) \times 1}$$

$n \rightarrow$ no. of features

VECTORIZED IMPLEMENTATION

$\geq = x^* \theta$; % Returns a $m \times 1$ vector

$\text{sigmoid} = 1. / (1 + \exp(-z))$; % Returns our hypothesis. ($m \times 1$) vector. 0.5 sigmoid

$\text{hypothesis} = \text{sigmoid};$

$\text{cost} = (-1/m) * ((y^* \log(\text{hypothesis})) + (1-y)^* (1 - \log(\text{hypothesis})))$

$\text{gradient} = (1/m) * (x^* \theta - y)^* \text{hypothesis};$

REGULARIZATION

$$\text{cost}'(h(x)) = \text{cost}(h(x)) + \left(\frac{\lambda}{2m}\right) \sum_{j=1}^n \theta_j^2$$

* note: Not minimizing / regularizing θ_0 since it has no effect on the fn. IMO would negatively affect it.

LINEAR REGRESSION:

$$\text{GRADIENT} = \frac{1}{m} \left[\sum_{i=1}^m ((h(x^{(i)}) - y^{(i)})^2 x_j^{(i)}) + \lambda \theta_j \right] \text{For } j \neq 0$$

IMPLEMENTATION: $[x]^{m \times (n+1)}, [y]^{m \times 1}, [\theta]^{(n+1) \times 1}$

$\text{gradient}(1) = (1/m) * (x(:, 1)' * (\text{hypothesis} - y));$

$\text{gradient}(2: \text{end}) = (1/m) * ((x(:, 2:\text{end})' * (\text{hypothesis} - y)) + \lambda \theta_j);$

LOGISTIC REGRESSION:

$$\text{GRADIENT} = \frac{1}{m} \left[\sum_{i=1}^m ((h(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}) + \lambda \theta_j \right]$$

IMPLEMENTATION:

$\text{gradient}(1) = (1/m) * (x(:, 1)' * (\text{hypothesis} - y));$

$\text{gradient}(2: \text{end}) = (1/m) * ((x(:, 2:\text{end})' * (\text{hypothesis} - y)) + \lambda \theta_j);$

$[x]^{m \times (n+1)}, [y]^{m \times 1}, [\theta]^{(n+1) \times 1}$

LINEAR REGRESSION VIA NORMAL EQUATIONS!

$$[\theta] = (X^T X)^{-1} X^T y$$

↓
pseudo inverse.

$$\rightarrow \text{theta} = \text{pinv}(X^T X) * X^T y;$$

REGULARIZATION:

$$[\theta] = (X^T X + \lambda I)^{-1} X^T y.$$

I is an identity matrix with element $(1,1) = 0.1$ $[I]^{(n+2) \times (n+2)}$

$$L = \begin{bmatrix} 0 & & \\ 1 & 1 & & \\ 0 & 1 & 1 & \ddots & \\ \vdots & \vdots & \ddots & \ddots & \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix}$$

The regularized $[\theta]$ is always invertible.

$$L = \text{ones}(n+2, n+2); L(1,1) = 0;$$

$$\rightarrow \text{theta} = \text{pinv}(X^T X + \lambda L) * X^T y;$$

GRADIENT DESCENT

VS

NORMAL EQUATIONS

- Need to choose alpha

+ No need to choose alpha

- Needs many iterations

+ No need to iterate

+ $O(kn^2)$

- $O(n^3)$, need to calculate $\text{inv}(X^T X)$

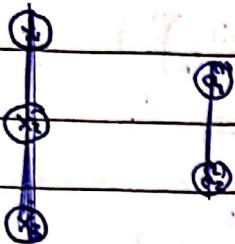
+ Works well when n is large

- Slow if n is very large

+ $n > 10,000$

+ $n < 10,000$

NEURAL NETWORKS



$L=1$
Layer 1
Input layer

$L=2$
Layer 2

$L=3$
Layer 3
Output layer

$$(a_0^{(1)}) \quad (a_1^{(1)}) \quad (a_2^{(1)})$$

$$(a_0^{(2)}) \quad (a_1^{(2)})$$

$$(a_1^{(2)}) \quad (a_2^{(2)})$$

$$(a_1^{(2)}) \quad (a_2^{(2)})$$

$$(a_1^{(2)}) \quad (a_2^{(2)})$$

FORWARD PROPAGATION

$a \rightarrow$ Is the activation function.

Ex: Sigmoid, ReLU, tanh, etc.

Let L denote the layer number; i denote the node number & j denote individual components.

Then $a_i^{(l)} = a \left(\sum_{j=0}^{n-1} \theta_{ij}^{(l-1)} a_j^{(l-1)} \right)$

$n \rightarrow$ no. of nodes in the previous layer.
 $a(x) \rightarrow$ activation function.

VECTORIZED IMPLEMENTATION:

$$[a^l] = \begin{bmatrix} a_0^{(l)} \\ \vdots \\ a_n^{(l)} \end{bmatrix} = \begin{bmatrix} a^{(l)}(z_0) \\ \vdots \\ a^{(l)}(z_n) \end{bmatrix} = a([z])$$

IF $s_l \rightarrow$ no. of nodes in layer l ,

$$[z]^{s_l \times 1} = [\theta]^{s_l \times s_{l-1}} \times [x^{(l-1)}]^{s_{l-1} \times 1}$$

COST FUNCTION

$$J(\theta) = \frac{1}{m} \left(\sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log(h_\theta(x^{(i)})_k) + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k) \right] \right) + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j=1}^{s_l} (\theta_{ij}^{(l)})^2$$

$m \rightarrow$ No. of training examples
 $K \rightarrow$ No. of classes
 $L \rightarrow$ No. of layers
 $s_l \rightarrow$ No. of nodes in layer l .

$y^{(i)}$ refers to a vector of dimensions $K \times 1$.

For a classification problem where the third class is the correct answer for a given $x^{(i)}$ (where $K = 5$),

$y^{(i)}$ would look like this.

$$y^{(i)} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Here $y_k^{(i)}$ refers to the k -th element of vector y .

$$\text{Ex: } y_2^{(i)} = 0, y_3^{(i)} = 1$$

Similarly,

$h_\theta(x^{(i)})$ refers to the $K \times 1$ vector which is the collection of outputs of the final (output) layer of our neural network.

$$\text{Ex: } h_\theta(x^{(i)}) = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.06 \\ 0.21 \\ 8.31 \end{bmatrix}, \text{ Here } h_\theta(x^{(i)})_k \text{ refers to the } k\text{-th element of vector } h_\theta(x^{(i)}).$$

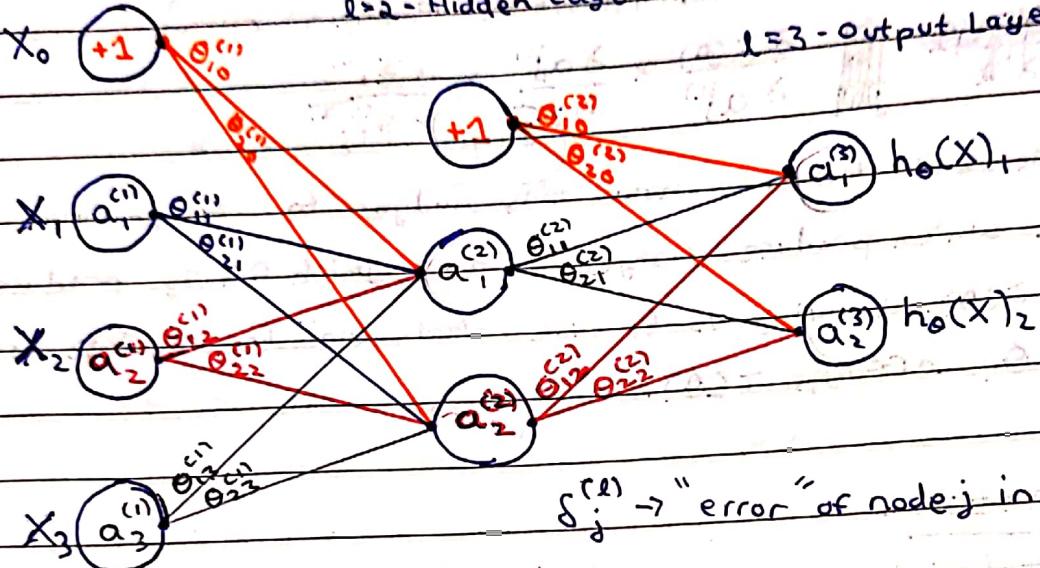
$$\text{Ex: } h_\theta(x^{(i)})_2 = 0.2, h_\theta(x^{(i)})_3 = 0.06.$$

Hence, the first part of our cost function calculates the sum of the error in the predictions by our neural network, i.e., error in each $h_\theta(x^{(i)})_k$ for its corresponding $y_k^{(i)}$.

$l=1$ - Input Layer (3 units)

$l=2$ - Hidden Layer (2 units)

$l=3$ - Output Layer (2 units)



For $l=1$ (output layer) - [Vectorized implementation]

$$\delta^{(1)} = a^{(1)} - y$$

For $l \neq 1$, $l > 1$

$$\Theta^{(2)} = \begin{bmatrix} x_0 & x_1 & x_2 \\ a_1^{(3)} & \theta_{10} & \theta_{11} & \theta_{12} \\ a_2^{(3)} & \theta_{20} & \theta_{21} & \theta_{22} \end{bmatrix}$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(2+1)} * g'(z^{(2)})$$

Intuition via example. On the right

we have matrices $\Theta^{(2)}$, $(\Theta^{(2)})^T$ & $\delta^{(2)}$

displayed with their respective values

as per the above neural network.

Let us now understand the term,

$$(\Theta^{(2)})^T \times \delta^{(2+1)} | [\Theta^{(2)}]^T \times \delta^3 = 1$$

$$\delta^{(3)} = \begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \end{bmatrix} \rightarrow \delta_1^{(3)} = h_\Theta(x)_1 - y_1^{(i)}$$

$$\delta^{(3)} = \begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \end{bmatrix} \rightarrow \delta_2^{(3)} = h_\Theta(x)_2 - y_2^{(i)}$$

$$\delta^{(2)} = \begin{bmatrix} \theta_{10} \delta_1^{(3)} + \theta_{20} \delta_2^{(3)} \\ \theta_{11} \delta_1^{(3)} + \theta_{21} \delta_2^{(3)} \\ \theta_{12} \delta_1^{(3)} + \theta_{22} \delta_2^{(3)} \end{bmatrix}$$

$\delta^{(2)}$ simply calculates the weighted sum of error due to θ in $\delta^{(3)}$ for every node in

The $*g'(z)$ term comes from the calculus derivative ~~term~~ of the cost function. [Try to derive? Application of chain rule]

$$\text{Ex: } \frac{\partial J(\theta)}{\partial \theta_{ij}^{(l-1)}} \rightarrow \cancel{\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l-1)}}} \quad \cancel{\frac{\partial J(\theta)}{\partial a_j^{(l)}}} \times \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \times \frac{\partial z_j^{(l)}}{\partial \theta_{ij}^{(l-1)}}$$

useful to think about them as how tiny nudges to θ produce changes in z which produce changes in a , which finally produce changes in $J(\theta)$ allowing us to use chain rule to calculate how a nudge to θ , affects $J(\theta)$.

BACK PROPAGATION:

1. Iterate over training examples $(x^{(i)}, y^{(i)})$ from $i=1$ to m .

2.005

1. Set $\Delta_{ij}^{(l)} = 0$ for all $i, j \in l$.

2. Iterate over training examples $(x^{(i)}, y^{(i)})$ from $i=1$ to m .

for $i=1:m$ [Vectorized Implementation]

(i) Set $a^{(1)} = x^{(i)}$ & add bias term.

(ii) Perform forward propagation to compute $a^{(l)}$ for $l=2, 3 \rightarrow L$

(iii) Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

(iv) Compute $\delta^{(L-1)}, \delta^{(L-2)} \rightarrow \delta^{(2)}$ using

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * g'(z^{(2)})$$

(v) Sum all these $\delta_{ij}^{(l)}$ values per example iteration to $\Delta^{(l)}$

$$\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} S^{(l+1)} (a^{(l)})^T$$

end for

3. Calculate gradient terms from Δ .

$$\text{If } i=0, D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)}$$

$$\text{If } i \neq 0, D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}$$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = \frac{1}{m} D_{ij}^{(l)}$$

EVALUATING A TRAINED ALGORITHM

- A common practice is to evaluate the trained algorithm by splitting the dataset into a training dataset & testing dataset.

General practices: Training dataset - 70%, Testing dataset - 30%.
(or)

Training dataset - 80%, Testing dataset - 20%.

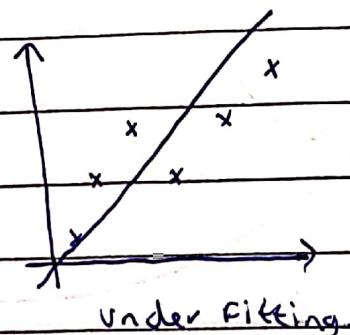
The algorithm and even the person's knowledge of the algorithm nor the creator of the algorithm should ever observe / infer any information from the test set. It should be a closely guarded secret.

In no way should it influence the algorithm.

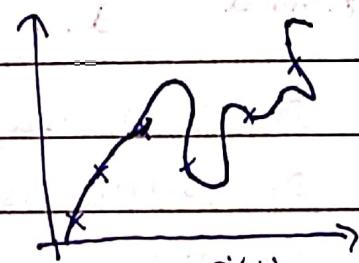
BIAS vs VARIANCE

High Bias \rightarrow Underfitting

High Variance \rightarrow Overfitting



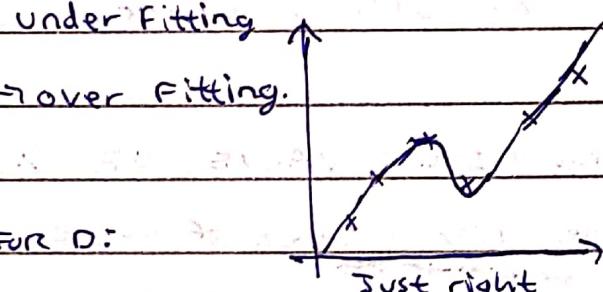
Choosing the right order of the polynomial terms in the hypothesis can help correct under & over fitting.



If d is the order of the polynomial,

Generally: If d is very less, \rightarrow underfitting

d is very large \rightarrow overfitting.



① CHOOSING THE RIGHT VALUE FOR D:

We can split our dataset as follows:

60% \rightarrow Training data

20% \rightarrow Cross validation data

20% \rightarrow Testing data

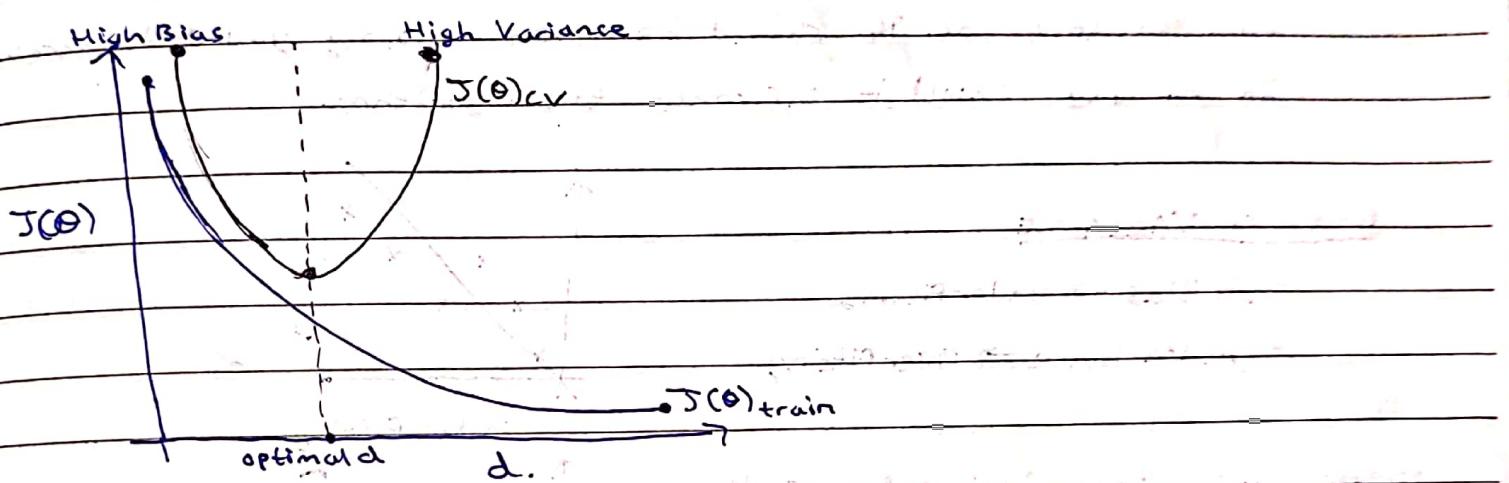
We can iterate over a range of values for λ . Ex: [1 to 10]

$J(\theta)_{\text{cross-validation}}$ is our cost function of choice.

Squared error cost function or logistic error cost function etc.

From the cross validation dataset we can select an optimal value for λ .

The best value for λ is the $J(\theta)_{\text{cross-validation}}$ which gives the least value. for $(x^{(i)}, y^{(i)})$ from our cross-validation (CV) dataset.



If $J(\theta)_{\text{training}}$ is large \rightarrow High bias (underfitting)

& $J(\theta)_{\text{cv}}$ is large

If $J(\theta)_{\text{training}}$ is very small \rightarrow High Variance (overfitting)

& $J(\theta)_{\text{cv}}$ is large

② CHOOSING THE RIGHT VALUE OF λ

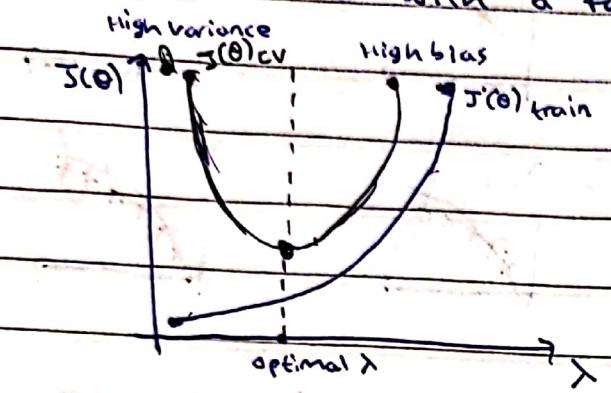
We can iterate over a range of values for λ . [Ex: $0.1 \rightarrow 0.2, 0.4, 0.8, \dots, 10$]

$J(\theta)_{\text{train}}$ is the cost function of our choice over our training examples.

and the regularization terms. $J(\theta)_{\text{cv}}$ is just the cost function over our CV examples.

without the regularization term. The $[\theta]$ parameters our cost function calculates its value upon ~~is~~ the theta calculated by minimizing

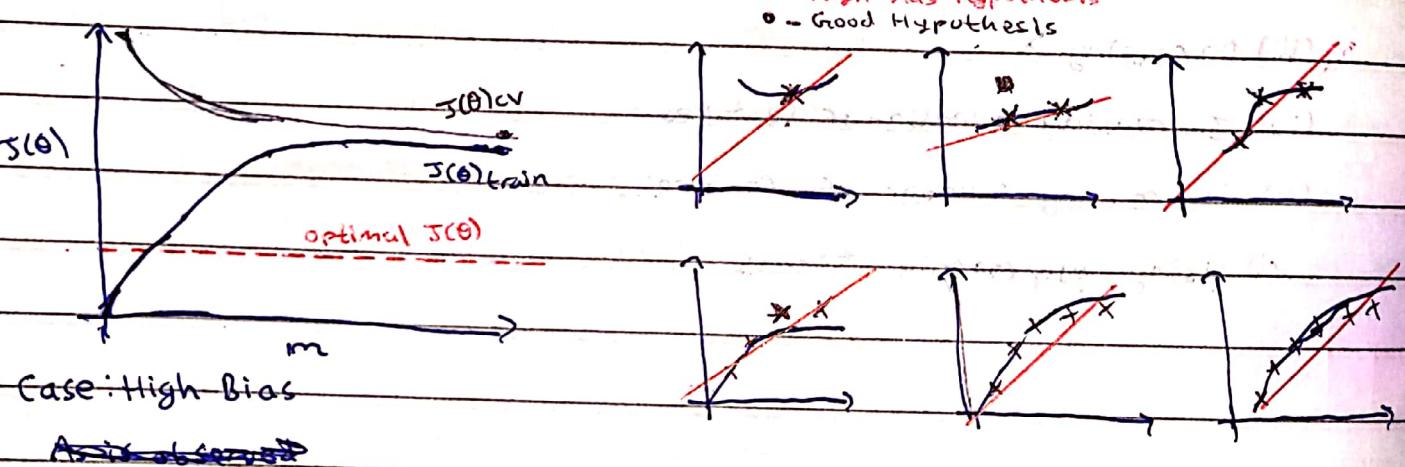
$J(\theta)$ which is the sum of the cost function of our choice and the regularization term with a factor of ' λ '.



Following the same rules as when determining d , we can diagnose the problem as high bias or high variance and by iterating over different values for λ we can pick our optimal value for our regularization parameter ' λ '.

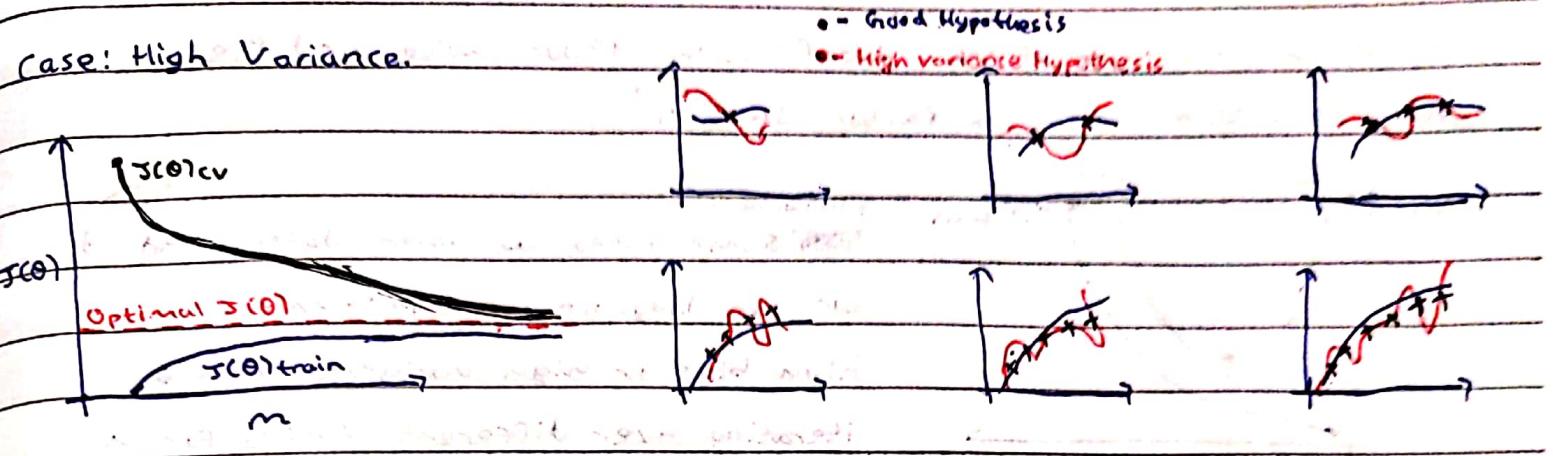
③ LEARNING CURVES (Iterating over m)

Here $J(\theta)_{\text{train}}$ is the cost function on examples $(x^{(i)}, y^{(i)})$ from $i=1 \dots m$. θ is the parameter vector obtained by training on examples $(x^{(i)}, y^{(i)})$ from $i=1 \dots m$. $J(\theta)_{\text{cv}}$ is the cost function on examples $(x^{(i)}, y^{(i)})$ from all the examples in the cv dataset.



From the above curve, we can infer that increasing the number of training examples alone cannot correct high bias problems.

case: High Variance



From the above observed curves, we can infer that increasing training examples does indeed help lower $J(\theta)_{cv}$. The 2 curves do converge eventually.

OPTIMIZING THE LEARNING ALGORITHM

Potential methods & when to use them:

- Problem it fixes

(i) Increasing training examples - High variance

(ii) Increasing λ - High Variance

(iii) Decreasing λ - High Bias

(iv) Increasing number of features - High Bias

(v) Decreasing number of features - High Variance

(vi) Using polynomial terms - High Bias

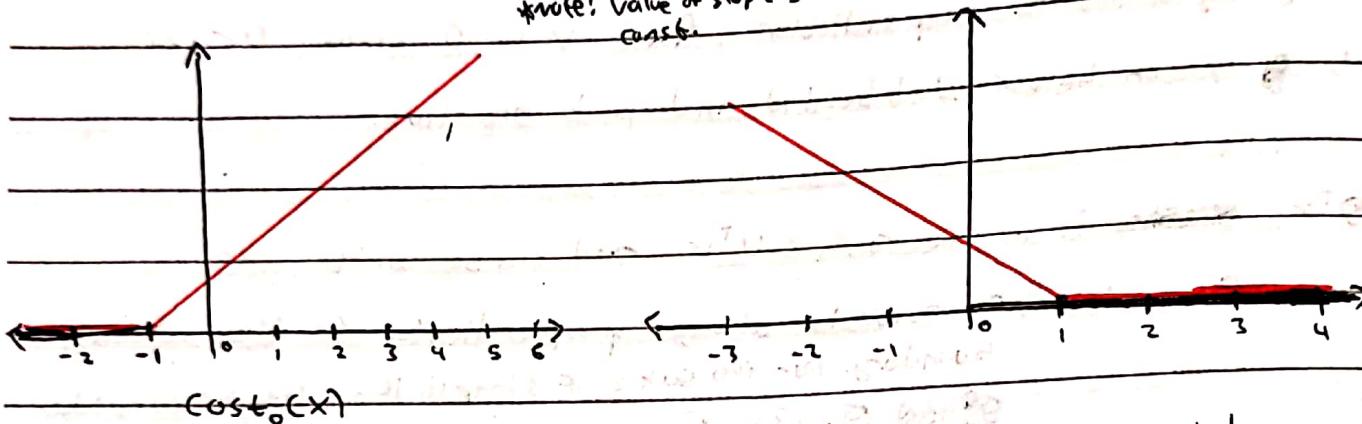
SVMs (SUPPORT VECTOR MACHINES)

Also called Large Margin Classifiers.

COST FUNCTION:

$$J(\theta) = C \sum_{i=1}^m [y^{(i)} \text{cost}_0(x^{(i)}) + (1-y^{(i)}) (\text{cost}_0(x^{(i)}))] + \sum_{j=1}^n \theta_j^2$$

*note: value of slope is const.



The slope of the ~~green~~ line does not matter much.

$$\text{h}_0(\theta^T x) \begin{cases} 0 & \text{when } h_0(x) < 0 \\ 1 & \text{when } h_0(x) \geq 0 \end{cases}$$

This cost function encourages our SVM to not 'just' predict the answer but instead predict it ~~more~~ with more certainty.

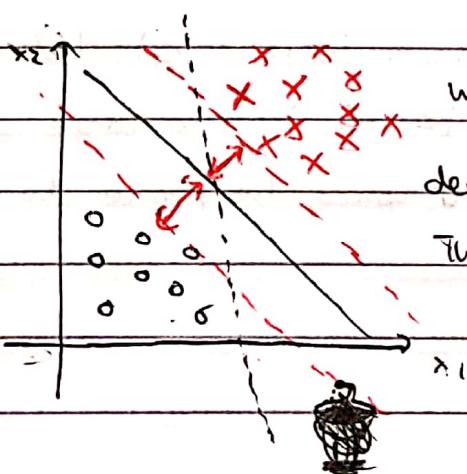
If we make C very large it will try to minimize the first term to 0.

$$\therefore \theta^T x \geq 1 \quad \text{or} \quad \theta^T x \leq -1 \quad \text{Then the minimization objective becomes}$$

$$\text{h}(x)=1 \quad \text{h}(x)=0 \quad \rightarrow C \times 0 + \sum_{i=1}^n \theta_i^2 = \sum_{i=1}^n \theta_i^2$$

While both decision boundaries work, the SVM chooses the decision boundary marked by the ~~green~~ non-dotted line.

This is because it chooses the decision boundary with the maximum margin between it and the datapoints.



$$U = \begin{bmatrix} U_1 \\ U_2 \end{bmatrix} \quad V = \begin{bmatrix} V_1 \\ V_2 \end{bmatrix}$$

$$U^T V = U_1 V_1 + U_2 V_2 = p \|U\| \|V\|$$

? can be proved.

p is the length of the perpendicular from V to U when $\theta < 180^\circ$.

$\theta > 180^\circ$, $-V$ can be extended behind and p is negative.

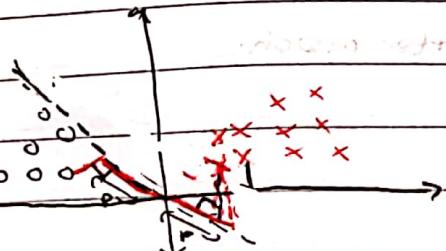
$\Theta^T x$ is of the form $U^T V$ and can be written as

$p \|U\|$. Note: Θ -vector is always perpendicular to the decision boundary. For the sake of simplicity, let us consider

$$\Theta_0 = 0$$

Minimization objective:

$$\min \sum_{i=1}^n \Theta_i^2 \text{ when } \Theta^T x \leq -1 \text{ or } \Theta^T x \geq 1$$



$$\therefore \Theta^T x = p \cdot \|\Theta\|$$

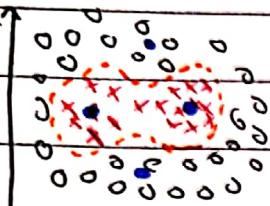
and if we are trying to minimize $\|\Theta\|^2$ then we must maximize p

$$\sum \Theta_i^2 = (\|\Theta\|)^2$$

p is the "distance between the examples

and the decision boundary. Therefore, for large values of C it will choose the decision boundary which provides a large margin between itself and the training examples. However, if C is too large then it may be hypersensitive to outliers and hence be prone to overfitting.

KERNELS: They provide a new way for us to create complicated decision boundaries without using very high order polynomial features.



→ Decision boundary we wish to draw.

Let us choose these 4 landmarks.

We can now compute a new set of features

f_1, f_2, f_3, f_4 which are a measure of the

similarity [distance] between $x^{(i)}$ & l_1, l_2, l_3, l_4 .

We can measure similarity through a 'kernel' function or 'similarity function'.

$$f_j^{(i)} = \text{similarity}(x^{(i)}, l_j) = e^{\left(-\frac{\|x^{(i)} - l_j\|^2}{200}\right)}$$

If we use no kernels & just use $\Theta^T X$, it is also called a Linear kernel.

Other kernels include the polynomial kernel

standard form : $f_j = (x^{(i)} \cdot l_j + c)^n$

& other kernels which are specific to different use cases.

SVM training libraries have well optimized methods to train using

using features ' f ' to train on instead of (computed using..

a gaussian kernel) instead of X , we can compute complicated decision boundaries.

While the concept of kernels can be used with other learning algorithms,

kernels are specially optimized to be used with SVMs.

BIAS VS VARIANCE

Large C - High Variance	Small C - Low Variance
---------------------------	--------------------------

Large σ - High bias	Small c - Low bias
----------------------------	----------------------

Note: Gaussian kernels can be slow to compute when using a large number of features since the no. of features generated = m .

While using SVM's landmarks are placed at all positions occupied by all training examples.

$$x^{(i)} = f^{(i)} \text{ or } X = \emptyset$$

~~so landmarks will~~

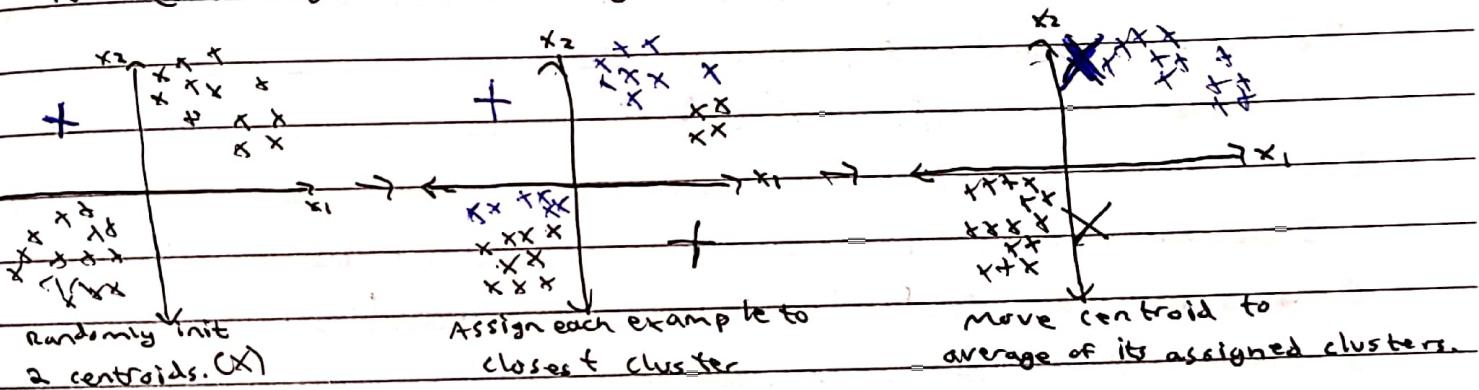
UNSUPERVISED LEARNING

The primary difference between supervised and unsupervised learning is that the training examples / data provided to an unsupervised learning algorithm is unlabelled.

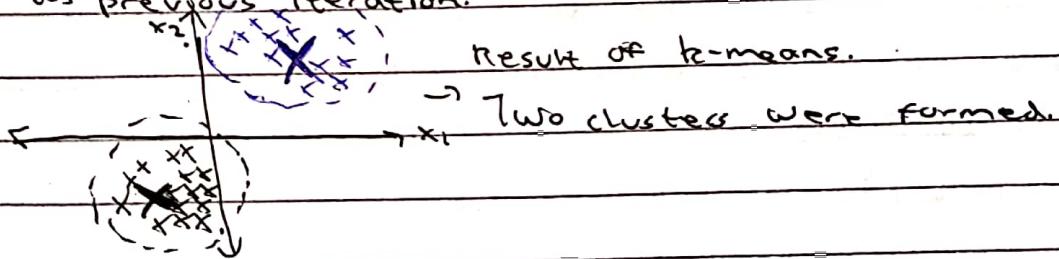
Training data $\Rightarrow \{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$, i.e., there is no $y^{(i)}$

Unsupervised learning algorithms try to find structure in unlabelled data.

K-Means Algorithm - Widely used to find clusters in unlabelled data.



→ Repeat algorithm until next position of each cluster is the same as previous iteration.



K-Means Algorithm

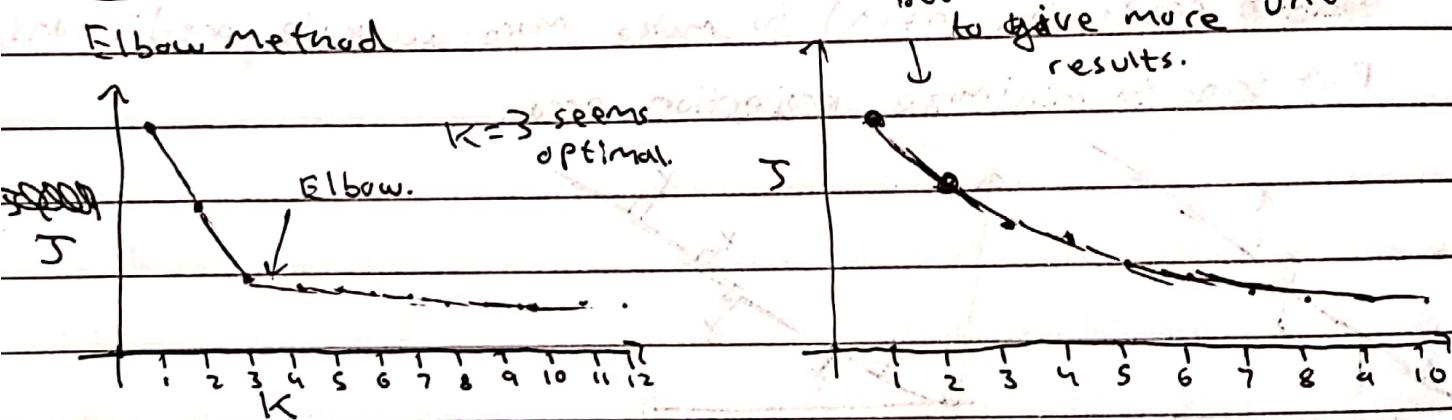
- Randomly initiate K-cluster centroids at randomly chosen $x^{(i)}$ (M_k)
- Assign $c^{(i)}$ to cluster k per example $x^{(i)}$ which minimizes the $\|x^{(i)} - M_k\|^2$ squared distance between them, i.e., the cluster closest to $x^{(i)}$
- Move M_k to the average of all the examples labeled to cluster k

OPTIMIZATION OBJECTIVE

$$\text{Cost function: } J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, M_1, M_2, \dots, M_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - M_k\|^2$$

It is possible for k-means to land in local optima for different random initializations of clusters M_1, \dots, M_K .
So it is advisable to run k-means multiple times with different random initializations for the clusters.

CHOOSING K



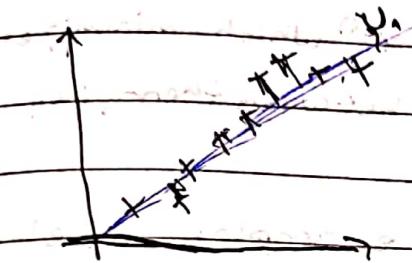
Hence the elbow method might not always work and it is better to handpick K depending on the situation / purpose for classification.

DEBUGGING: Cost function J should always decrease with increase in K. If not then possible reasons are local optima & faulty code / implementation.

PRINCIPAL COMPONENT ANALYSIS (PCA)

It is used to reduce the dimensions of some given data.

Ex: Given 2D data



PCA helps find the vector $\mathbf{w} \in \mathbb{R}^2$ such that the projection error, i.e., the perpendicular distance from the data point to the surface provided

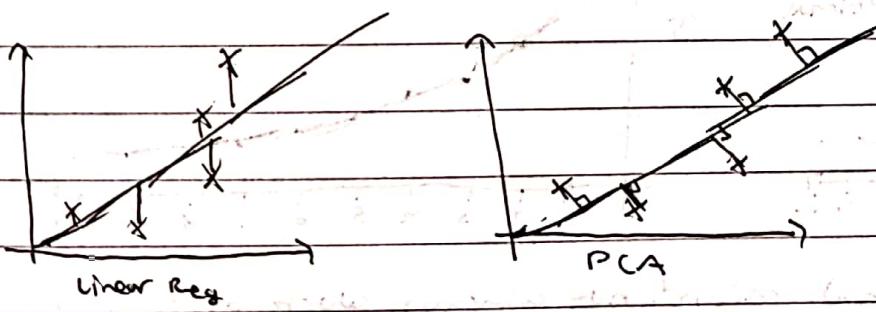


projections of data x on U , which provide

1D vector of features $\mathbf{z} \in \mathbb{R}^d$

In general, PCA finds the direction / surface (upto 3D) onto which to project the data X such that projection error between the data points & the vectors are minimum.

*Note: While linear regression tries to reduce the error between label y and ~~excess~~ prediction $h_\theta(x)$ to make more accurate predictions, PCA tries to minimize projection error.



THE PCA ALGORITHM

To reduce ~~the~~ data of n-dimensions to K-dimensions.

Given $X^{(i)}$ is a row vector of all the features, i.e., $X \in \mathbb{R}^n$

- Preprocess data by applying mean normalization. Also feature scaling if necessary
- Compute covariance matrix (Σ)

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (X^{(i)}) (X^{(i)})^\top$$

Vectorized implementation: If $X = \mathbb{1}^{m \times n}$

Then covariance matrix

$$\Sigma = \left(\frac{1}{m}\right) * X' * X$$

$$\begin{bmatrix} \dots & X^{(1)\top} & \dots \\ \dots & X^{(2)\top} & \dots \\ \dots & X^{(m)\top} & \dots \end{bmatrix}$$

- Compute eigen vectors of matrix Σ

- Using Singular Value Decomposition

$$[U, S, V] = \text{svd}(\text{covariance matrix});$$

- Using eigenvector function [only applies to square matrices]

$$[U, S, V] = \text{eig}(\text{covariance matrix});$$

Matrix U is a $n \times n$ matrix which contains the vectors

U_1, U_2, \dots, U_n as its columns.

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ U_1 & U_2 & \dots & U_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

SVD can also be used directly on the data matrix X , and matrix V would contain our eigen vectors. We apply SVD on the Σ matrix since $n \times n$ is generally more computably reasonable than $n \times n$ and potentially for the right S matrix. For a covariance matrix, U & V are the same.

Matrix U :

$$U = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ U_1 & U_2 & U_3 & \dots & U_n \\ 1 & 1 & 1 & \dots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix}$$

we select the first K vectors.

$$U_{\text{reduce}} = U(:, 1:K)$$

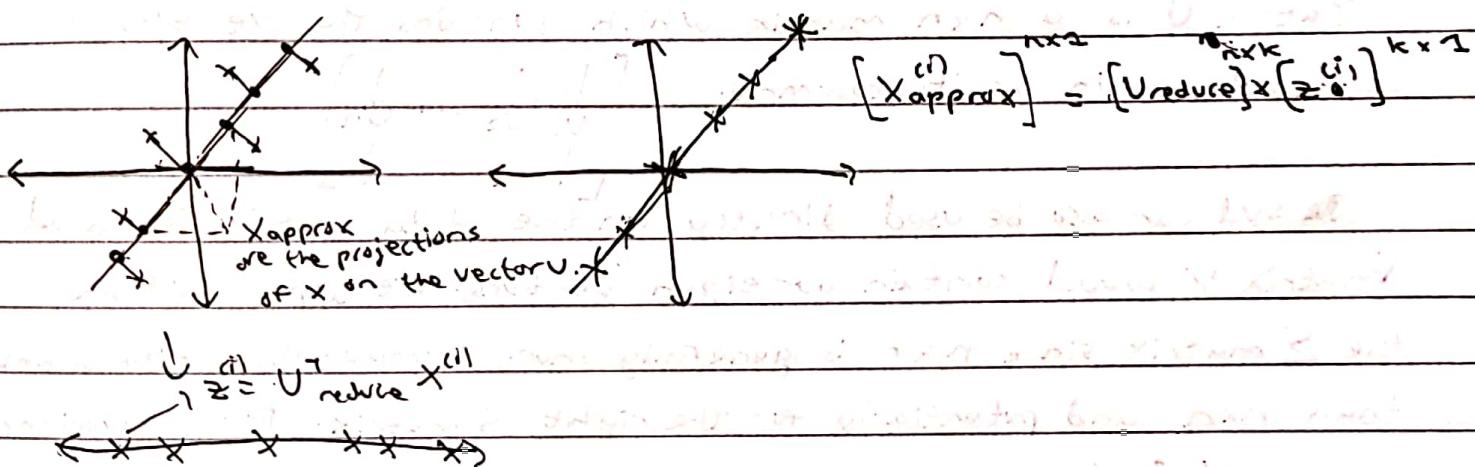
$$x^{(i)} \in \mathbb{R}^n \quad z^{(i)} \in \mathbb{R}^K$$

$$z^{(i)} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ U_1 & U_2 & \dots & U_K \\ 1 & 1 & \dots & 1 \end{bmatrix}^T \begin{bmatrix} x^{(i)} \\ x^{(i)} \\ x^{(i)} \end{bmatrix} = \begin{bmatrix} -U_1 \\ -U_2 \\ \vdots \\ -U_K \end{bmatrix} \times \begin{bmatrix} x^{(i)} \\ x^{(i)} \\ x^{(i)} \end{bmatrix} = \begin{bmatrix} z^{(i)} \end{bmatrix}$$

$$\text{Ex: } \begin{bmatrix} U_{11} & U_{12} & U_{13} \\ U_{21} & U_{22} & U_{23} \end{bmatrix} \times \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \end{bmatrix} = \begin{bmatrix} U_{11} \cdot x_1^{(i)} + U_{12} \cdot x_2^{(i)} + U_{13} \cdot x_3^{(i)} \\ U_{21} \cdot x_1^{(i)} + U_{22} \cdot x_2^{(i)} + U_{23} \cdot x_3^{(i)} \end{bmatrix}$$

It maps the raw data point $x^{(i)}$ onto the eigenvector.

RECONSTRUCTION FROM REDUCED DIMENSION COMPRESSION



CHOOSING A GOOD VALUE OF K FOR DIMENSIONALITY REDUCTION

$$\text{Average squared projection error} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \hat{x}^{(i)}\|_2^2$$

$$\text{Total variation in data} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|_2^2$$

We choose K to be the smallest natural number which satisfies the condition \Rightarrow Average Squared Projection Error ≤ 0.01

Total Variation in Data

$$0.01 \rightarrow 1\%, i.e., \text{we can say that } 99\% \text{ of variance was retained}$$
$$\frac{\sum_{i=1}^m \|x^{(i)} - \hat{x}^{(i)}\|_2^2}{\sum_{i=1}^m \|x^{(i)}\|_2^2} \leq 0.01.$$

Acceptable values for variance retained are values upward of 85%.

While applying svd,

$$[U, S, V] = \text{Svd}(\text{co-variance matrix } \Sigma)$$

$$S = \begin{bmatrix} \sigma_{11} & & & \\ & \sigma_{22} & & \\ & & \ddots & \\ & & & \sigma_{kk} \\ & & & & \sigma_{nn} \end{bmatrix}_{n \times n} \rightarrow \text{it is a diagonal matrix.}$$

For a given K, we can calculate variance retained by.

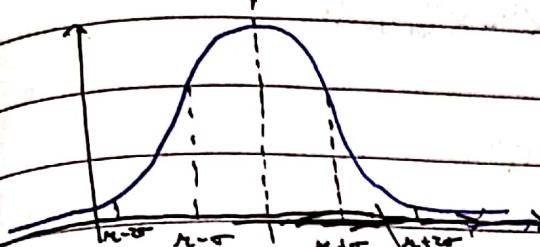
85% variance retained

$$1 - \frac{\sum_{i=1}^k \sigma_{ii}}{\sum_{i=1}^n \sigma_{ii}} \geq 0.85$$

ANOMALY DETECTION ALGORITHM

GAUSSIAN DISTRIBUTION

μ = Mean



σ = Standard deviation

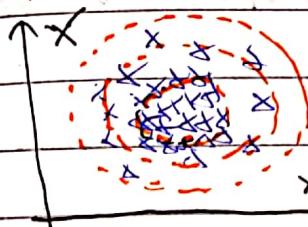
σ^2 = Variance

$$\text{Eqn} \rightarrow p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

parametrised by

Area under graph is always equal to one.

If we have features that can be mapped as a gaussian distribution we can use them to create an anomaly detection algorithm. Use `hist(x)` command in octave to see the histogram plot of our data. We can use transformation such as $\log(x+c)$ or x^c or $x^{\frac{1}{2}}$ to make our data distribution more similar to a normal gaussian distribution. We can also create new features which involve an operation between 2 existing features.

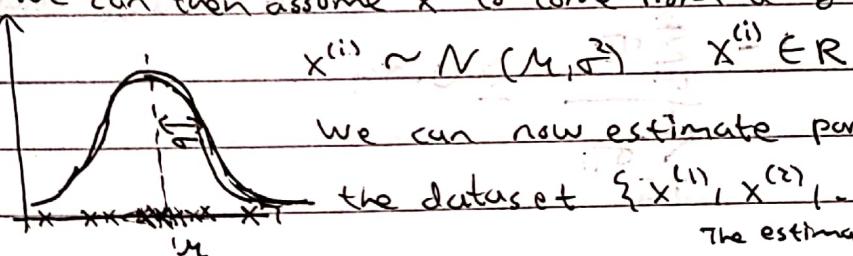


X - Anomalous Data.

If $h_\theta(x) < \epsilon$ ∈ threshold constant, then it is classified as anomalous.

If our training data resembles a gaussian distribution, we can assume the data source to be distributed normally.

We can then assume x to come from a gaussian distribution.



$x^{(i)} \sim N(\mu, \sigma^2) \quad x^{(i)} \in \mathbb{R}$

We can now estimate parameters μ & σ using

the dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

The estimated μ

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \bar{x})^2$$

Now when $x^{(i)} \in \mathbb{R}^n$ Then we find parameters M_j & σ_j^2 for each feature of example $x^{(i)}$

Then $x_1 \sim N(\mu_1, \sigma_1^2)$, $x_2 \sim N(\mu_2, \sigma_2^2)$, ..., $x_n \sim N(\mu_n, \sigma_n^2)$

$$\therefore p(x) = p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

ALGORITHM:

- Choose features x_j which you think are indicative of anomalous examples. It helps if these features also resemble a gaussian distribution. Use functions to transform them to resemble gaussian distributions if necessary & possible.

- Calculate parameters μ_1, \dots, μ_n & $\sigma_1^2, \dots, \sigma_n^2$

$$\text{where } \mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \& \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

- Given new example x , calculate $p(x)$ using estimated parameters.

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

- Classify as anomaly if $p(x) < \epsilon$ | $y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases}$

Given a training set of ~~minority~~ anomalous & majority normal data

60% normal examples \rightarrow training set

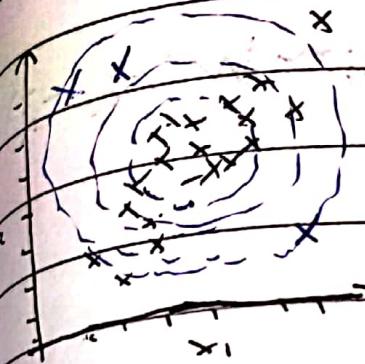
20% normal examples + 50% anomalous examples \rightarrow CV-set

20% normal examples + 50% anomalous examples \rightarrow test set.

Use metrics such as precision, recall, F-score, etc to measure accuracy.

Can use CV-set to choose value of ϵ

MULTIVARIATE GAUSSIAN DISTRIBUTION



Given two features which vary linearly, the features in blue ^{look like} ~~are~~ anomalies. But the normal model cannot classify them as anomalies.

We need to create a feature $x_3 = \frac{x_2}{x_1}$ or something similar.