

Normal equation based Approach

⇒ Using linear algebra. All this problem of selecting 'x', how many iteration and other

Advantages of NE over Linear Regression

- Single step
- No 'x' selection

Disadvantages

- when n becomes big more Computational Power required.

Derivation

$$h(x) = \theta^T x \quad \text{in vectorize format. } x^T x$$

lets put x in different format

$$x^i = \begin{bmatrix} 1 \\ x_1^i \\ \vdots \\ x_n^i \end{bmatrix} \quad (n+1) \times 1$$

$$X = \begin{bmatrix} x_0^1 & x_1^1 & x_2^1 & \dots \\ x_0^2 & x_1^2 & \dots & \dots \\ \vdots & \vdots & \vdots & \vdots \\ x_0^m & x_1^m & \dots & \dots \end{bmatrix} \equiv \begin{bmatrix} (x^1)^T \\ (x^2)^T \\ \vdots \\ x^m \end{bmatrix}$$

$$\dim(X) = m \times (n+1)$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$\dim(y) = m \times 1$$

$$y' = \theta_0 x_0' + \theta_1 x_1' + \dots$$

$$y' = \theta^T x' = x'^T \theta$$

$$y^1 = (x^1)^T \theta$$

$$y^2 = (x^2)^T \theta$$

$$y^m = (x^m)^T \theta$$

$$\Rightarrow y = X\theta = Xw$$

$$w = \theta$$

$$\dim(\theta) = (n+1) \times 1$$

$$h_\theta(x) = \theta^T x = x^T w$$

$$\text{error} = e = xw - y \quad \begin{matrix} \text{Linear} \\ \text{Error in each training} \end{matrix} \quad \text{eg}$$

$$E = e^T e \quad \text{(Cumulative error)}$$

Also known as MMSE (

Minimized Mean Squared error.

$$E = (xw - y)^T (xw - y) \equiv (xw^T - y^T) (xw - y)$$

$$(1 \times 1)^T \times (1 \times 1) = (1 \times 1) = \text{Scalar value.}$$

$$(A+B)^T = A^T + B^T \quad \text{Property of matrices.}$$

$$E = (xw)^T (xw) - (xw)^T (y) - y^T (xw) + y^T y$$

$$(xw)^T (y) = ((xw) y)^T$$

$$y^T (xw) = ((xw)^T y)^T \quad \text{(Property of Transpose)}$$

Now look at dim of matrices.

$$y^T(xw) = ((xw)^T y)$$

$$\dim(xw) = m \times 1$$

$$\dim y = m \times 1$$

$$\dim((xw)^T y) = (1 \times m) \times (m \times 1) = 1 \times 1$$

~~[x]~~ Transpose of scalar is equal to original

$$x^T = x \quad (\text{where } x \text{ is scalar})$$

$$E = (xw)^T(xw) - 2(xw)^T y + y^T y$$

$$E = w^T x^T x w - 2(xw)^T y + y^T y$$

$$E = w^T x^T x w - 2w^T x^T y + y^T y$$

$$\frac{\partial w^T}{\partial w} = 1$$

$$\frac{\partial 2w^T x^T y}{\partial w} = 2 \begin{bmatrix} x_{10}^1 & x_{11}^1 & x_{12}^1 & \dots \\ x_{20}^2 & x_{21}^2 & x_{22}^2 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \end{bmatrix}$$

$$= \begin{pmatrix} x_{10}^1 y_0 + x_{11}^1 y_1 + x_{12}^1 y_2 + \dots \\ x_{20}^2 y_0 + x_{21}^2 y_1 + x_{22}^2 y_2 + \dots \\ \vdots \end{pmatrix}$$

$$= x_{10}^1 y_0$$

$$\frac{\partial}{\partial w} \left[(x_0' w_0 + x_1' w_1 + \dots) y_1 + \dots + (x_n' w_0 + \dots) y_m \right]$$

$$= \frac{\partial}{\partial w}$$

$$P = 2 \sum_{i=1}^m y_i (x_0' w_0 + x_1' w_1 + \dots)$$

$$P = 2 \sum_{i=1}^m y_i \sum_{c=0}^n x_c^i w_c$$

$$\frac{\partial P}{\partial w_0} = 2 (x_0' y_1 + x_0^2 y_2 + \dots)$$

$$\frac{\partial P}{\partial w_1} = 2 (x_1' y_1 + x_1^2 y_2 + \dots)$$

$$\frac{\partial P}{\partial w_n} = 2 (x_n' y_1 + x_n^2 y_2 + \dots)$$

$$\frac{\partial 2 w^T x^T y}{\partial w} = 2 x^T y$$

$$P_2 = w^T x^T x w$$

$$= (w_0 \ w_1 \ w_2 \ \dots \ w_n) \begin{pmatrix} x_0^1 & x_0^2 & x_0^3 & \dots \\ x_1^1 & x_1^2 & x_1^3 & \dots \\ x_2^1 & x_2^2 & x_2^3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

$$\begin{pmatrix} x_0^1 & x_1^1 & x_2^1 & \dots \\ x_0^2 & x_1^2 & x_2^2 & \dots \\ x_0^3 & x_1^3 & x_2^3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}$$

$$= \begin{pmatrix} w_0 & w_1 & w_2 & \dots & w_n \end{pmatrix} \begin{pmatrix} (x_0^1)^2 w_0 + \dots & (x_1^1)^2 w_1 \\ \vdots & \vdots \\ (x_0^n)^2 w_0 & \dots & (x_n^n)^2 w_n \end{pmatrix}$$

$$P2 = \cancel{w_0} ((x'_0)^2 w_0 + \dots (x'_n)^2 w_n) + \dots (w_n (x'_n)^2 w_0 \dots (x'_n)^2 \cancel{w_n})$$

$$\frac{\partial P2}{\partial w_0} = 2 \cancel{w_0} (x'_1)^2 + 2 \cancel{w_0} (x'_{n2})^2 \dots 2 \text{ on } x'_n$$

$$\frac{\partial P2}{\partial w} = 2 X^T X w$$

$$\frac{\partial F}{\partial w} = 2 X^T X w - 2 X^T y$$

Now for minimum error derivative should be zero at that point

$$\frac{\partial F}{\partial w} = 0$$

$$\frac{\partial y^T y}{\partial w} = 0$$

As there is no 'w' in it

$$2 X^T X w = 2 X^T y$$

$$w = (X^T X)^{-1} X^T y$$

As ~~the~~ Dimension increases ^{increases} time taken ~~to~~ or no. of computation required to process to find inverse is increased so for high Dimension (Generally > 2048) it takes high amount of time so we prefer Regression Model there.

When normal equation is preferred

~~when dimension < 1024 to~~

$\dim(w) < 1024$

NE is preferred

Clustering

Unsupervised Learning

In supervised learning we are given a label associated with each data point, which helps in classification of data.

In unsupervised we don't have any labels associated with it and tasked to find some structure

One such approach could be to divide the given data into various distinct cluster based on some properties

- Use of to find pattern in stars
- Dividing Customer base into segment on some criteria

Extra

Can we labeled data for clustering?

The thing is that we can use whatever we want, to use, what matters is results. Results for labelled data is not as good as that of some algorithm for labelled data such as regression. They perform better.

Why?

Because Clustering tries to find relation b/w data point where as in labelled data we want relation b/w data point and its label. So they perform poorly as they fail to do so.

for eg. if some person visits some website e-commerce label when he buy something

Classifier tries to predict whether he will buy this time or not (or what he will buy) whereas Clustering may try to cluster it type of user with its parameter based on ~~time~~ data points

K-Means Clustering

This is an iterative algorithm which perform algorithm which two step in each

first: randomly takes k - data point as center

Second Then

for every iteration, colour the points based on its nearest centroid

Then move centroid to center of all colored points

then repeat until not a single ~~centroid~~ point change its color.

Can k-means only segment well separated data?

Even where it is not seen to segment data. k-means clustering separate our data points. This use heavily on market segmentation

k-means clustering always give non overlapping k-segments

In Linear Regression we have cost function to reduce. In Clustering we take similar approach to minimize the avg. distance of all the cluster centroids with the datapoint in that cluster

Formally let

$c^{(i)}$ = index of cluster to which data point $x^{(i)}$ is assigned

μ_k = cluster centroid

$\mu_{c^{(i)}}$ = centroid of cluster to which $x^{(i)}$ is assigned

$$J = \frac{1}{n} \sum_{i=1}^n \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

And our objective is to minimize this cost function

This cost function is also called distortion.