

Machine Learning Regression

DR. BHARGAVI R

SCOPE

VIT CHENNAI

Regression - Fundamentals

Regression analysis enables us to -

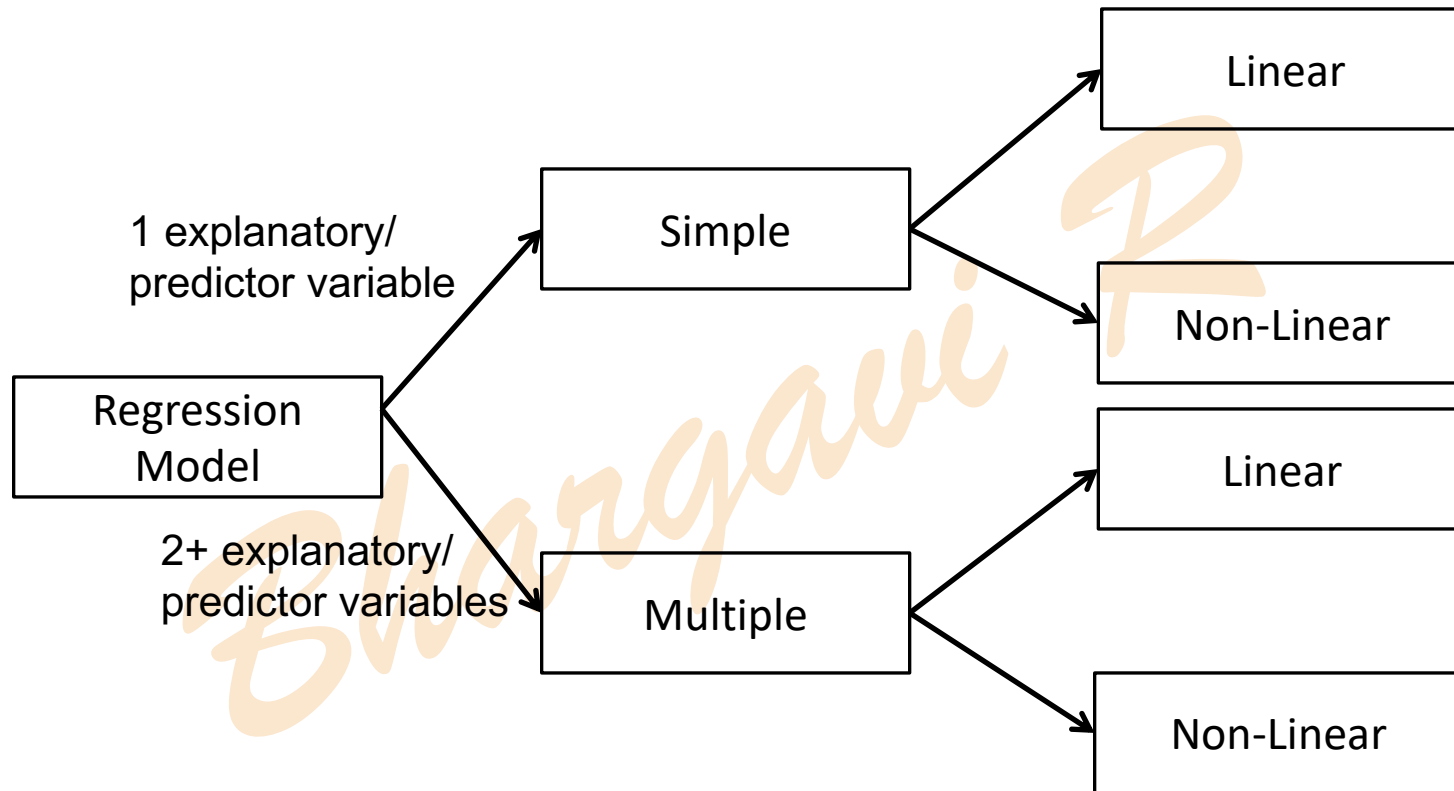
- Develop a mathematical function to predict the values of **Output/Dependent/Target/Response** variable, based on the value of **Input/Independent/ Predictor** variable(s).
- Quantify the effect of changes in the independent variable have on the dependent variable.
- Identify unusual observations.

Fundamentals (cont...)

Applications

- Sales prediction based on size of the shop
- Financial Forecasting: Predicting stock prices, interest rates, or economic indicators.
- Disease Progression: Predicting the progression of diseases based on patient data and medical history.
- Drug Response: Modeling the effectiveness of different drug dosages on patient outcomes.
- Quality Control: Analyzing factors that affect product quality in manufacturing processes.
- Signal Processing: Modeling and predicting signals in communication systems.

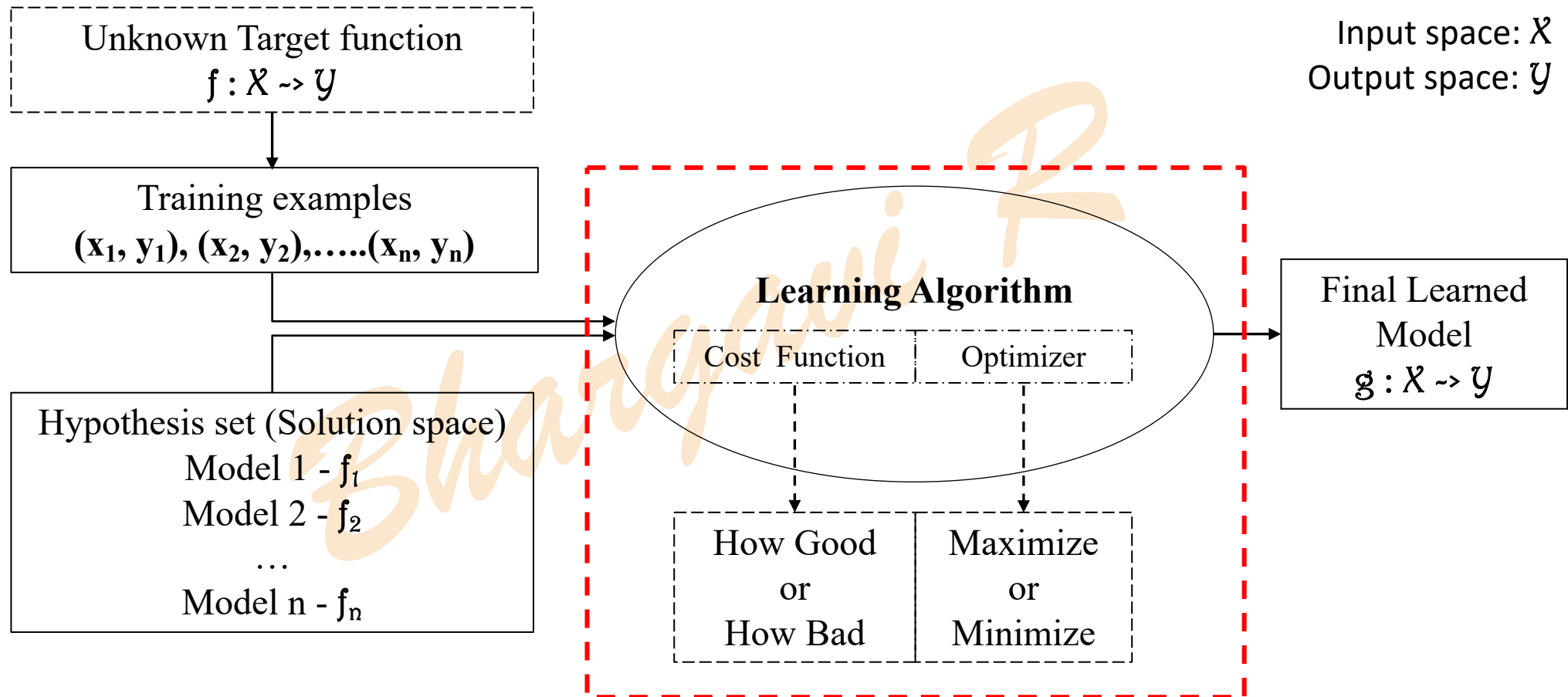
Regression Models - Types



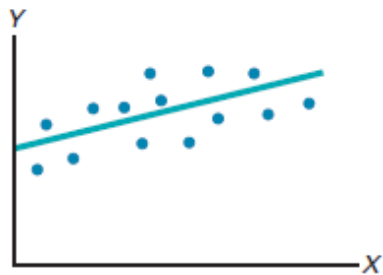
Simple Linear Regression

- Simple and straightforward.
- Predicts a *quantitative response* y based on *single predictor* variable x
- Assumes a linear relationship between x and y
- So **Hypothesis function** $h(x)$ is $y = h(x) = w_0 + w_1x$
- w_0 is the **intercept** and w_1 is the **slope**
- w_0 and w_1 are called the **linear regression coefficients** or **model coefficients** or **parameters** of the model
- Different values of w 's result in different hypotheses.

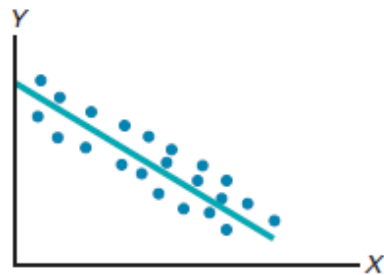
Machine Learning - Process



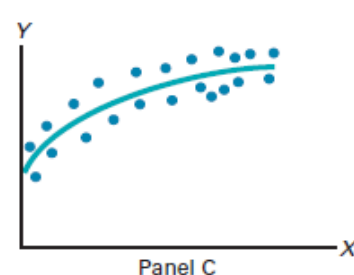
Relationship Between Variables



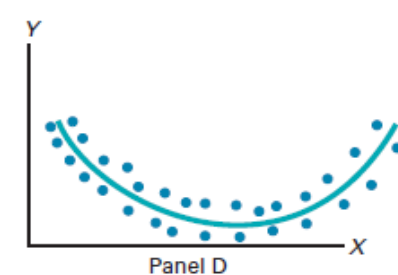
Panel A
Positive linear relationship



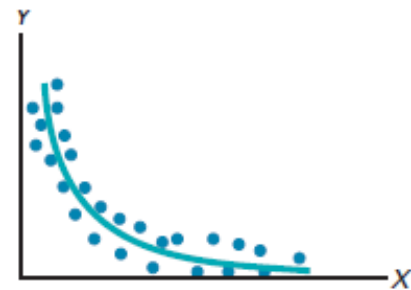
Panel B
Negative linear relationship



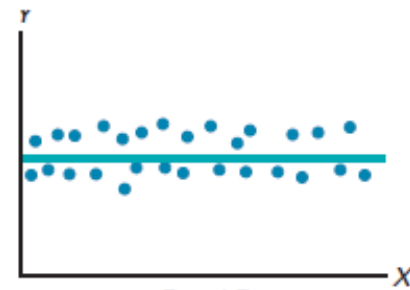
Panel C
Positive curvilinear relationship



Panel D
U-shaped curvilinear relationship



Panel E
Negative curvilinear relationship



Panel F
No relationship between X and Y

A Simple Use Case

“MyStyle” is a chain of apparel show rooms. Director of planning wants to forecast annual sales for all new stores for strategic planning, based on store size.

To examine the relationship between the store size in square feet and its annual sales, data is collected from already established stores.

Size of show room (Thousands of sq.ft)	Sales (millions of \$)	Size of show room (Thousands of sq.ft)	Sales (millions of \$)
1.7	3.7	1.1	2.7
1.6	3.9	3.2	5.5
2.8	7.7	1.5	2.9
5.6	9.5	5.2	10.7
1.3	3.4	4.6	7.6
2.2	5.6	5.8	11.8
1.3	3.7	3	4.1

Data

Showroom 1 (x_1 sq.ft, y_1 \$) (1.7 sq.ft , 3.7\$)

Showroom 2 (x_2 sq.ft, y_2 \$) (1.6 sq.ft , 3.9\$)

Showroom 3 (x_3 sq.ft, y_3 \$) (2.8 sq.ft , 7.7\$)

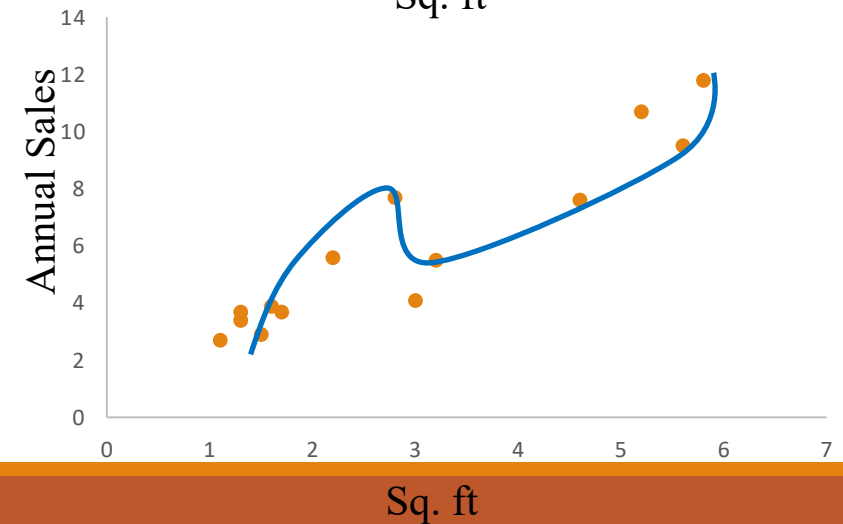
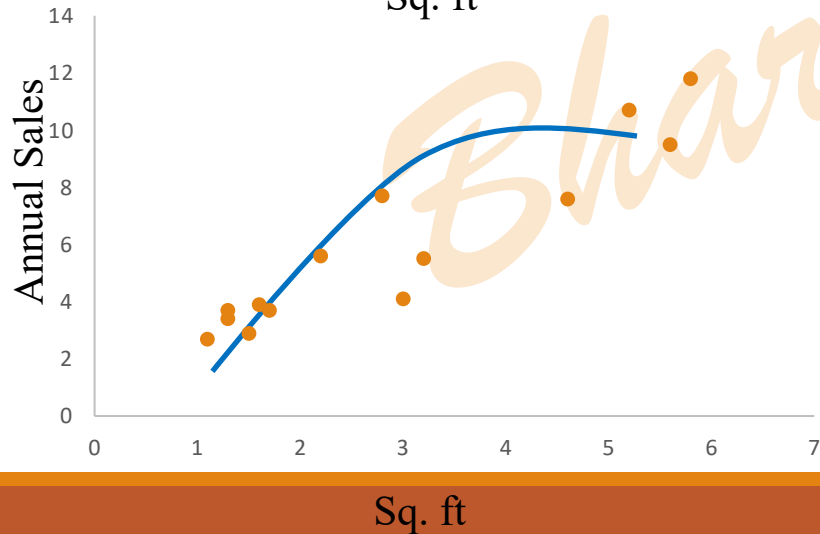
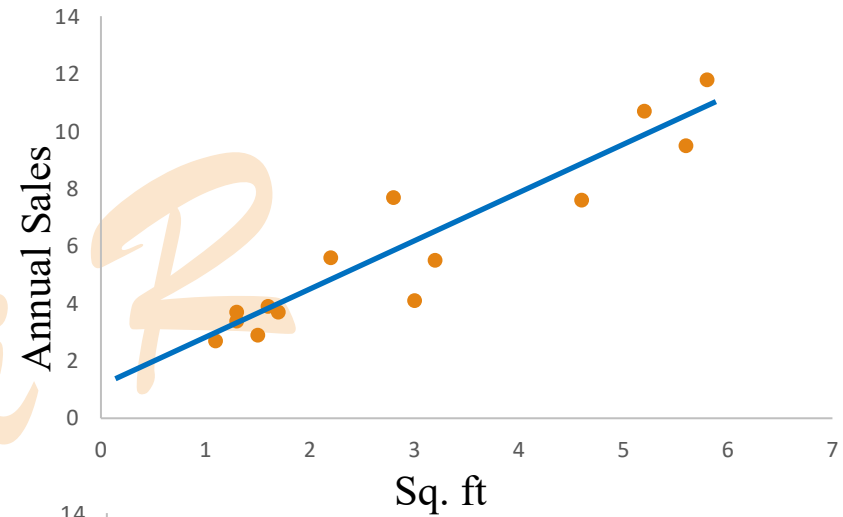
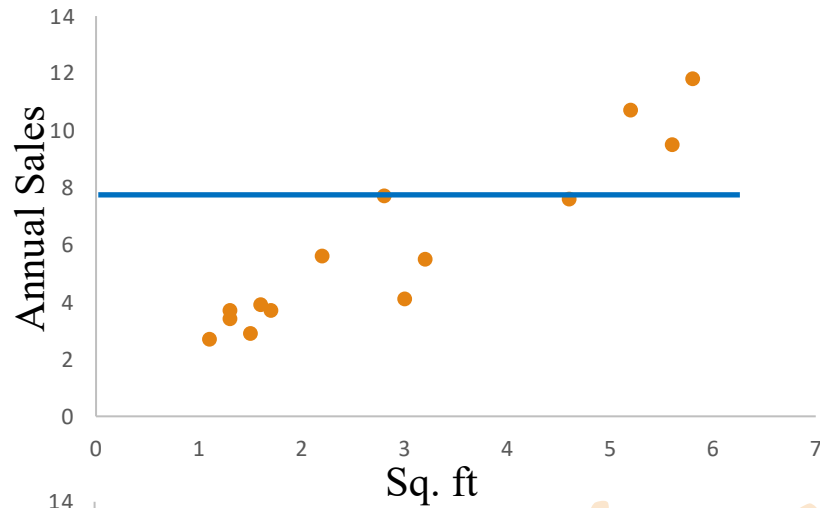
--

--

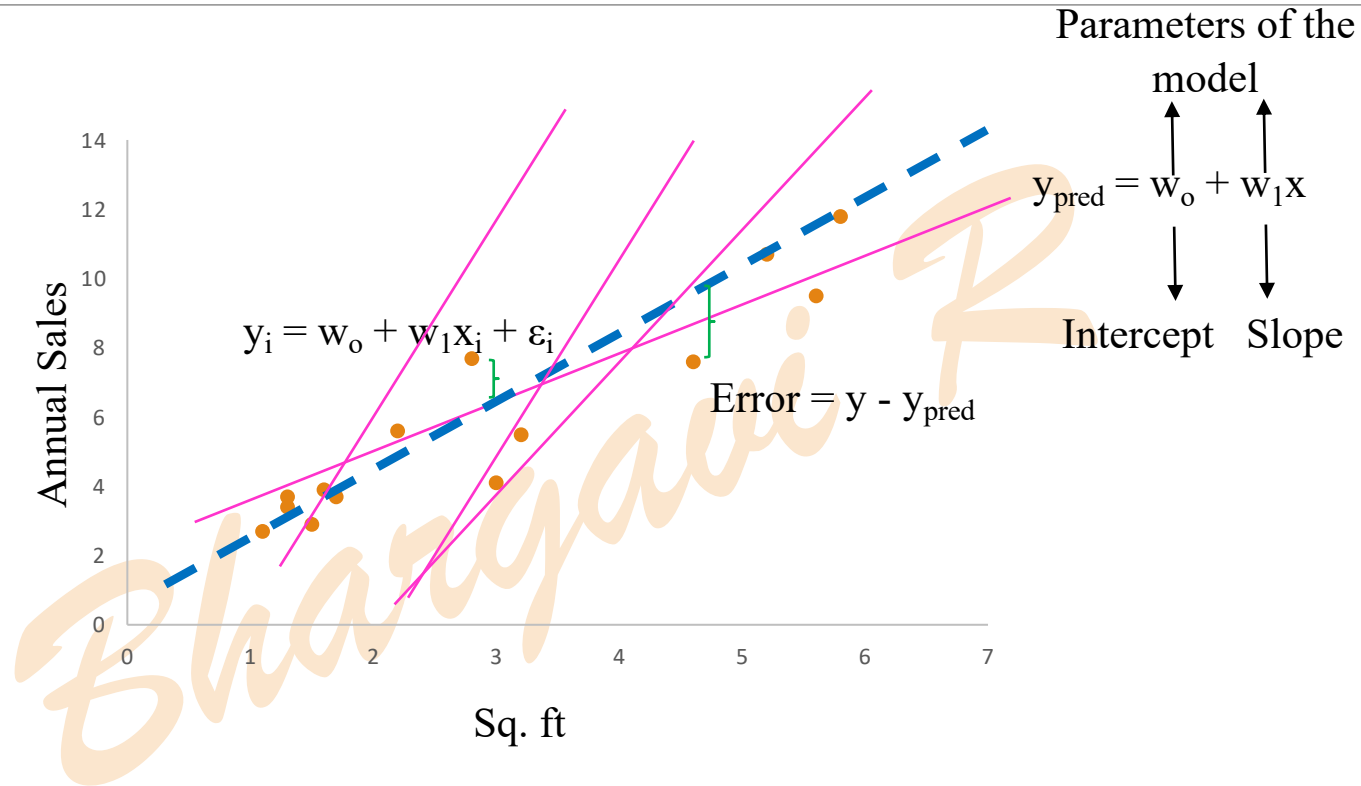
Showroom n (x_n sq.ft, y_n \$) (3 sq.ft , 4.1\$)


Input Output

Which Model



Estimate Best Function



Cost function

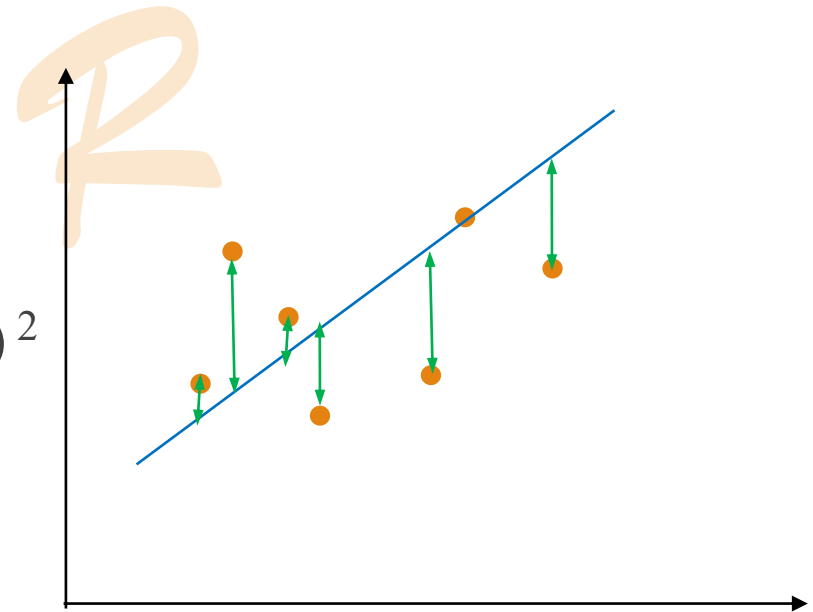
- RSS – Residual Sum of Squares

- $$\begin{aligned} \text{RSS} = & (\text{annual sales}_1 - (w_0 + w_1 \text{sq. ft}_1))^2 \\ & + (\text{annual sales}_2 - (w_0 + w_1 \text{sq. ft}_2))^2 \\ & + (\text{annual sales}_3 - (w_0 + w_1 \text{sq. ft}_3))^2 \\ & + \dots + (\text{annual sales}_n - (w_0 + w_1 \text{sq. ft}_n))^2 \end{aligned}$$

$$\text{RSS} = \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

$$\text{RSS} = \sum_{i=1}^n (y_i - h(x_i))^2$$

Cost function is $\sum_{i=1}^n (y_i - h(x_i))^2$



Objective Function

- Choose w_0 and w_1 such that y_{pred} is close to y for the training data.
- Minimize the Residual sum square (RSS)
- Minimize cost
- Objective function is $\min_{w_0 w_1} \sum_{i=1}^n (y_i - h(x_i))^2$

Or

$$\min_{w_0 w_1} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

Prediction

- Model parameters are obtained by minimizing the cost function.
- For any given test input x_t , value of y can be computed by substituting the values of x_t , \hat{w}_0 and \hat{w}_1 in the fitted regression line.
- let the estimated parameters be $\hat{w}_0 = 0.9645$, $\hat{w}_1 = 1.6699$. Then the annual sales of a show room of size 4000 sq,ft can be predicted as

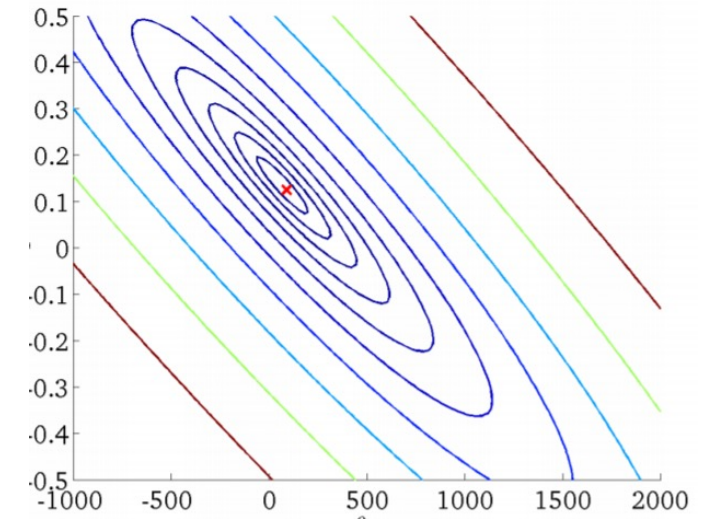
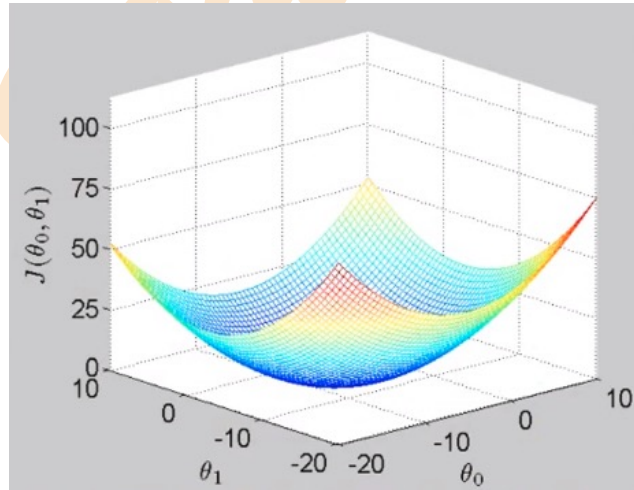
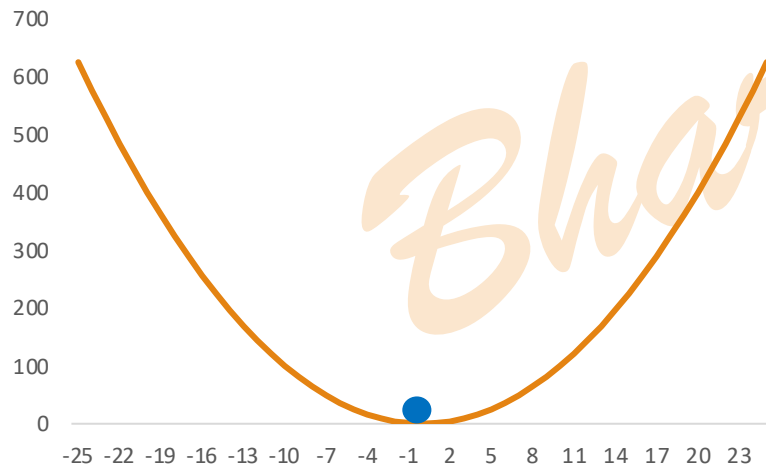
$$y_{\text{pred}} \text{ or } \hat{y} = 0.9645 + 1.6699 * 4 = 7.644 \text{ million \$}$$

Interpretation of Predicted coefficients

- Slope of the regression line i.e \hat{w}_1 can be interpreted as – with increase of 1 unit in x , the predicted value of y is estimated to increase by 1.6699 units
- Negative coefficients indicate the negative relation among the Predictor and Response variables.
- The intercept represents the predicted value of y when x equals 0

Optimizing the Cost Function

- RSS is a convex function
- Optimization problem = minimizing(RSS)
- Rate of change of RSS becomes/approaches zero at min(RSS)
- Rate of change of RSS is nothing but Gradient of RSS w.r.t parameters.



Optimizing the Cost Function (cont...)

- We can find solution (i.e w 's of line with minimum cost) in different approaches.
 1. Setting Gradient(RSS) = 0 and solving (closed form solution)
 2. Gradient descent (Iterative approach)

Bhargavi R

Find Gradient

$$RSS = \sum_{i=1}^n (y_i - (h(x_i)))^2 = \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

$$\text{Gradient}(RSS) = \nabla RSS(w_0, w_1) = \frac{\partial}{\partial \mathbf{w}} (RSS)$$

$$\frac{\partial}{\partial w_0} (\sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2) = \sum_{i=1}^n \frac{\partial}{\partial w_0} ((y_i - (w_0 + w_1 x_i))^2)$$

$$= -2 \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)]$$

$$\frac{\partial}{\partial w_1} (\sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2) = -2 \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)] x_i$$

$$\nabla RSS(w_0, w_1) = \begin{aligned} & -2 \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)] \\ & -2 \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)] x_i \end{aligned}$$

Approach 1

Setting Gradient of RSS to zero

$$\nabla RSS(w_0, w_1) = 0$$

$$-2 \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)] = 0 \quad (1)$$

$$-2 \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)] x_i = 0 \quad (2)$$

From (1) on solving for w_0 , we get $w_0 = \bar{y} - w_1 \bar{x}$

And Solving (2) after substituting the values of w_0 , we get, $w_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}) x_i}{\sum_{i=1}^n (x_i - \bar{x}) x_i}$

$$\text{Or equivalently } w_1 = \frac{\sum_{i=1}^n y_i x_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n}}{\sum_{i=1}^n x_i^2 - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n x_i}{n}}$$

Closed form Normal solution $\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$

Approach 2

A Simple Prediction Workout

- We know that $1 \text{ km} = 0.62137 * 1 \text{ Mile}$ (Ground truth)
- Can we learn this formula by looking at some data that we have in our hand?
- Let's see
- Assume that we have the following data (Training data)

Sl.No	Kms	Miles
1	100	62.137
2	0	0

- In practice we have many more observations

All We Know....

- Miles = f(kms)



- Data (Observations)
- There exists a linear relationship between input and output

i.e miles = $C * \text{kms}$ ----- (1)

Predict, Compare, Learn

- First assume that $C = 0.5$. Substituting in (1)

$$100 \text{ kms} = 50 \text{ Miles } (\neq 62.137)$$

$$62.137 - 50 = 12.137 \longrightarrow \text{Error}$$

- Let's now increase C to 0.6

$$100 \text{ kms} = 60 \text{ Miles } (\neq 62.137)$$

$$\text{Error} = 2.137 \text{ (Better than Previous guess)}$$

- Further increase C to 0.7

$$100 \text{ kms} = 70 \text{ Miles } (\neq 62.137)$$

$$\text{Error} = -7.863$$

Predict, Compare, Learn

- Let's now make C to 0.61

100 kms = 61 Miles (\neq 62.137)

Error = 1.137

----- and so on till we are okay with the error or till no further change in the error.

Bhargavi R

Gradient Descent (cont..)

While !(converge){

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla \mathbf{w}_t$$

$$t = t + 1$$

}

$\nabla \mathbf{w}$ is the partial derivative of the cost function w.r.t \mathbf{w} i.e $\nabla \mathbf{w} = \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}}$

η is a small constant called as Learning rate

Gradient Descent

Iterative approach

While not converged update

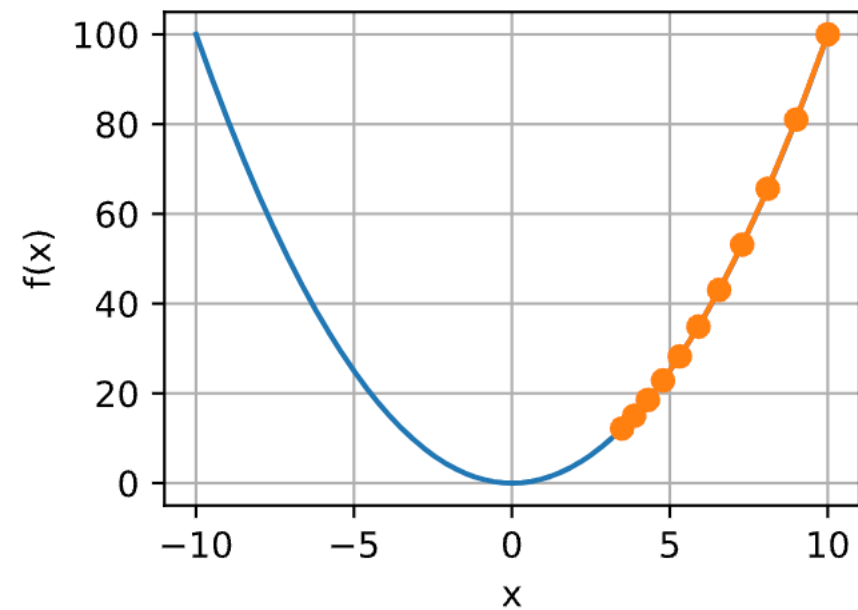
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla \text{RSS}(\mathbf{w})^t$$

$$\begin{bmatrix} w_0^{t+1} \\ w_1^{t+1} \end{bmatrix} = \begin{bmatrix} w_0^t - \eta(-2 \sum_{i=1}^n [y_i - (w_0^t + w_1^t x_i)]) \\ w_1^t - \eta(-2 \sum_{i=1}^n [y_i - (w_0^t + w_1^t x_i)]x_i) \end{bmatrix}$$

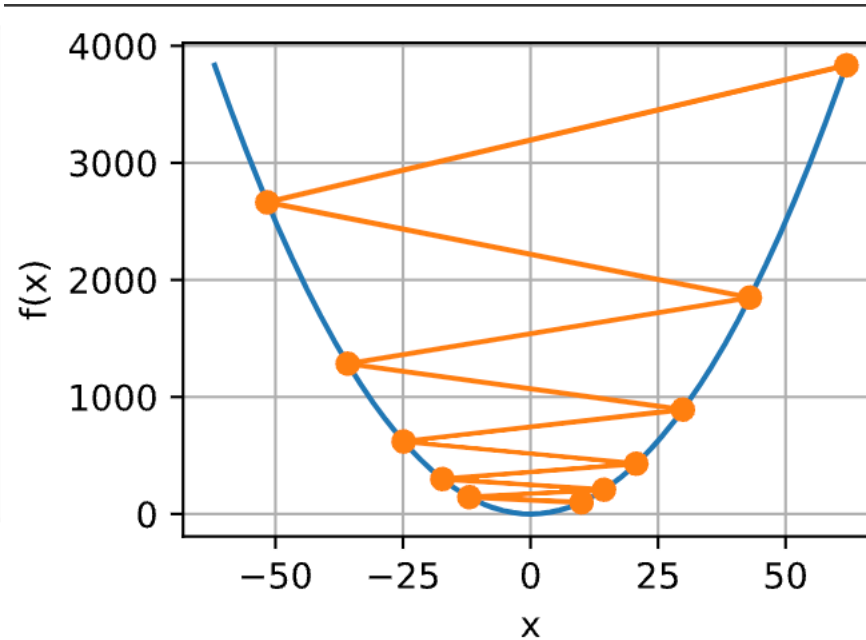
$$\begin{bmatrix} w_0^{t+1} \\ w_1^{t+1} \end{bmatrix} = \begin{bmatrix} w_0^t \\ w_1^t \end{bmatrix} + 2\eta \begin{bmatrix} \sum_{i=1}^n [y_i - (w_0^t + w_1^t x_i)] \\ \sum_{i=1}^n [y_i - (w_0^t + w_1^t x_i)]x_i \end{bmatrix}$$

Learning Rate Impact on Solution Convergence

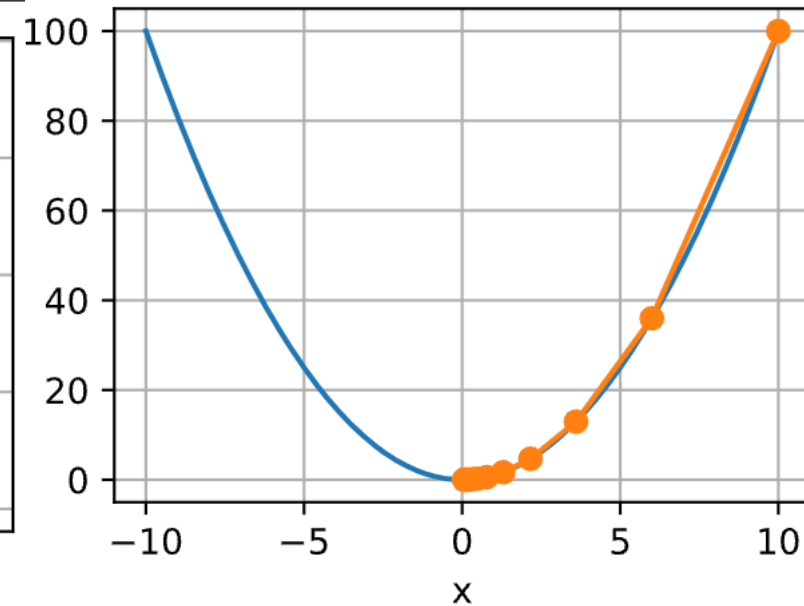
$\eta = 0.05$



$\eta = 1.1$



$\eta = 0.2$



Prediction Using More Than One Predictor

What should be the selling price of a used car based on the following features

- Color
- Age
- Model
- Mileage
- Distance travelled

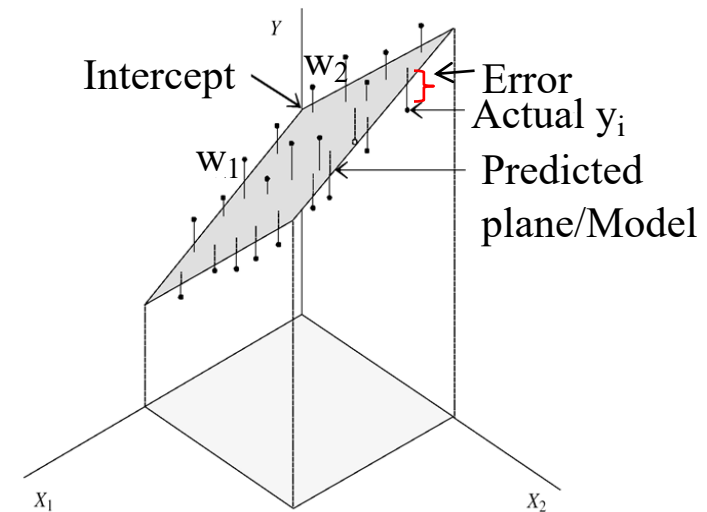
Bhargavi R

Multiple Linear Regression

Hypothesis function $h(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_mx_m$

A Multiple Regression model with n predictor variables fits a regression plane in $(n+1)$ dimensional space.

Ex: A multiple regression model with two explanatory variables fits a regression plane in 3-dimensional space



Data

Car1	x_{10}	x_{11}	x_{12}				x_{1m}	y_1
Car2	x_{20}	x_{21}	x_{22}				x_{2m}	y_2
Car3	x_{30}	x_{31}	x_{32}				x_{3m}	y_3
Car n	x_{n0}	x_{n1}	x_{n2}				x_{nm}	y_n

R

MLR Predicted Model and Interpretation

$$y = w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_mx_m \text{ where } x_0 = 1$$

- The intercept represents the predicted value of y when all x are set to 0.
- Slope for input feature x_i i.e w_i predicts the change in y per unit x_i holding all other input features constant.

Bhargavi R

Single Observation - Vector Notation

$$y_i = w_0 + w_1x_{i1} + w_2x_{i2} + \cdots + w_mx_{im} + \epsilon_i$$

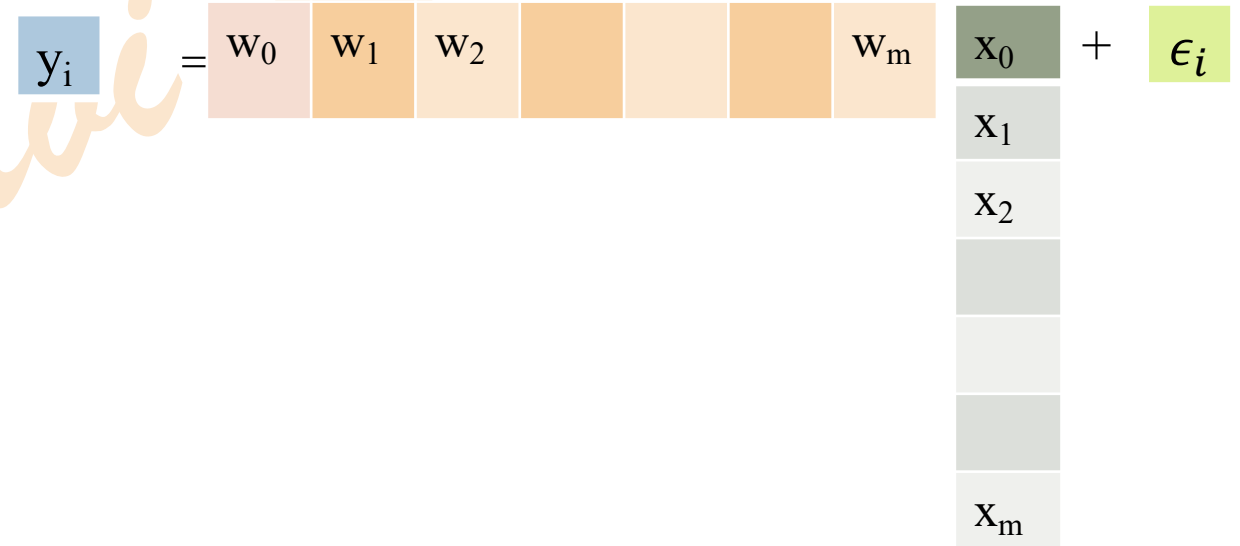
Introducing a dummy variable $x_{i0} = 1$

$$y_i = w_0x_{i0} + w_1x_{i1} + w_2x_{i2} + \cdots + w_mx_{im} + \epsilon_i$$

$$y_i = \sum_{j=0}^m w_j x_{ij} + \epsilon_i$$

$$y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i$$

Equivalently we can also write $y_i = \mathbf{x}_i^T \mathbf{w} + \epsilon_i$



All Observations - Matrix Notation

n observations with m attributes

The diagram illustrates the matrix notation for linear regression. It shows a vector y of n observations, a matrix X of n observations by $m+1$ attributes (including a bias attribute x_{i0}), a vector W of $m+1$ weights, and a vector ϵ of n error terms. The equation $y = XW + \epsilon$ is shown in both element-wise and matrix notation.

y_1	=	x_{10}	x_{11}	x_{12}			x_{1m}	w_0	+	ϵ_1
y_2		x_{20}	x_{21}	x_{22}			x_{2m}	w_1		ϵ_2
y_3		x_{30}	x_{31}	x_{32}			x_{3m}	w_2		ϵ_2
y_n		x_{n0}	x_{n1}	x_{n2}			x_{nm}	w_m		ϵ_n

$$y = XW + \epsilon$$

$$y = XW + \epsilon$$

MLR Cost Function - RSS

RSS for Simple Linear Regression is given by

$$RSS = \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

We will extend it to MLR

$$\begin{aligned} RSS(\mathbf{w}) &= \sum_{i=1}^n (y_i - (w_0 x_{i0} + w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im}))^2 \\ &= \sum_{i=1}^n (y_i - (\mathbf{w}^T \mathbf{x}_i))^2 \end{aligned}$$

$$RSS(\mathbf{w}) = (\mathbf{Y} - \mathbf{X}\mathbf{w})^T (\mathbf{Y} - \mathbf{X}\mathbf{w})$$

Gradient of RSS

$$\begin{aligned}\nabla RSS(\mathbf{w}) &= \nabla[(Y - X\mathbf{w})^T(Y - X\mathbf{w})] \\ &= -2X^T(Y - X\mathbf{w})\end{aligned}$$

Bhargavi R

Approach 1

Set Gradient = 0

$$\nabla RSS(\mathbf{w}) = -2X^T(Y - X\mathbf{w}) = 0$$

$$X^T(Y - X\mathbf{w}) = 0$$

$$X^T\mathbf{Y} - X^TX\mathbf{w} = 0$$

$$X^TX\mathbf{w} = X^T\mathbf{Y}$$

Multiplying both sides by $(X^TX)^{-1}$

$$(X^TX)^{-1}X^TX\mathbf{w} = (X^TX)^{-1}X^T\mathbf{Y}$$

But we know that $(X^TX)^{-1}X^TX = I$, the Identity matrix

$$\mathbf{w} = (X^TX)^{-1}X^T\mathbf{Y}$$

Closed form solution

Points to Remember – Closed form Solution

We can use closed form solution only if we can find the inverse of $X^T X$.

- $X^T X$ is a $m \times m$ matrix where m is the number of features.
- $X^T X$ is invertible only if the number of linearly independent instances $n > m$.
- Complexity of closed form solution is $O(n^3)$. Hence for huge data sets closed form solution is computationally very expensive

Approach 2 – Gradient Descent

While not converged {

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla \text{RSS}(\mathbf{w})^t$$

$$t = t+1$$

}

Substituting Gradient of RSS

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta(-2X^T(Y - X\mathbf{w}^t)) = \mathbf{w}^t + 2\eta X^T(Y - X\mathbf{w}^t)$$

While not converged {

$$\mathbf{w}^{t+1} = \mathbf{w}^t + 2\eta X^T(Y - \hat{Y}(\mathbf{w}^t))$$

$$t = t+1$$

}

Single Feature Update

Let's consider weight update for j^{th} feature

$$RSS(\mathbf{w}) = \sum_{i=1}^n (y_i - (w_0x_{i0} + w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im}))^2$$

Differentiating partially w.r.t w_j

$$\frac{\partial}{\partial w_j}(RSS) = \sum_{i=1}^n 2 (y_i - w_0x_{i0} - w_1x_{i1} - w_2x_{i2} - \dots - w_mx_{im})(-x_{ij})$$

$$= -2 \sum_{i=1}^n x_{ij} (y_i - \mathbf{x}_i^T \mathbf{w})$$

Single Feature Update (cont...)

Now Update of j^{th} feature weight

$$w_j^{t+1} = w_j^t - \eta \left(-2 \sum_{i=0}^n x_{ij} (y_i - \mathbf{x}_i^T \mathbf{w}^t) \right)$$

$$w_j^{t+1} = w_j^t + 2\eta \left(\sum_{i=1}^n x_{ij} (y_i - \hat{y}_i(\mathbf{w}^t)) \right)$$

Bhargavi R

Algorithm

Initialize $\mathbf{w}^{(1)} = 0$ (or random values)

While $\|\nabla RSS(\mathbf{w}^{(t)})\| > \varepsilon \{$

For $j = 0$ to m

$$\text{Partial}[j] = -2 \sum_{i=1}^n x_{ij} (y_i - \hat{y}_i(\mathbf{w}^t))$$

$$w_j^{t+1} = w_j^t - \eta * \text{Partial}[j]$$

$t = t+1$

}

Performance Metrics - Quality of Fit

Any regression model can be evaluated with the following metric

MAE – Mean Absolute Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

RMSE – Root Mean Squared

RSS – Residual Sum Squares = $\sum_{i=1}^n (y_i - \hat{y}_i)^2$

- RSS is Scale variant statistic

Metric - Quality of Fit (cont ...)

R² (Coefficient of determination) - represents the proportion of variance (of y) that has been explained by the independent variables in the model.

- R² is scale invariant statistic

$$R^2 = 1 - \frac{\text{Residual Sum Squares}}{\text{Total Sum Squares}} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

- Drawback - It either remains the same or increases with the addition of new independent variables.
- The value of R-squared does not decrease even when redundant variables are added.

Metric - Quality of Fit (cont ...)

Adjusted R-squared statistic – It takes into account the number of independent variables used for predicting the target variable.

- We can determine whether adding new variables to the model actually increases the model fit.

$$\text{Adjusted } R^2 = 1 - \left(\frac{(1-R^2)(n-1)}{n-m-1} \right)$$

n – Number of observations

m – Number of independent variables/features

- If R-squared does not increase significantly on the addition of a new independent variable, then the value of Adjusted R-squared will actually decrease.
- On the other hand, if on adding the new independent variable results in significant increase in R-squared value, then the Adjusted R-squared value will also increase.

Multicollinearity

- Multicollinearity refers to a situation in multiple linear regression where two or more independent variables are highly correlated.
- Multicollinearity causes several issues in the regression analysis like
 - Instability of Coefficient Estimates
 - Inflated Standard Errors of regression coefficients.
 - Interpretation Difficulties
- Correlation matrix gives collinearity between any two pairs of variables only it fails to capture multicollinearity.
- Variance Inflation Factor (VIF) can be used to identify multicollinearity and address the problem.

Variance Inflation Factor (VIF)

- VIF is a measure of multicollinearity among independent variables in a multiple regression model.
- VIF quantifies how much the variance of a regression coefficient is inflated due to multicollinearity.
- $VIF = \frac{1}{1 - R^2}$
- VIF starts at 1 and has no upper limit.
- $VIF = 1$, no correlation between the independent variable and the other variables.
- VIF exceeding 5 or 10 indicates high multicollinearity between this independent variable and the others.

Variance Inflation Factor (cont..)

- Fit the Initial Model: Fit a multiple linear regression model with all the potential predictors.
- Calculate VIF: Compute the VIF for each predictor.
- Identify High VIF: Identify predictors with high VIF values (commonly, a $VIF > 10$ is considered problematic).
- Remove High VIF Predictors: Remove the predictor with the highest VIF value.
- Refit the Model: Refit the model without the removed predictor.
- Recalculate VIF: Recalculate the VIF values for the remaining predictors.
- Repeat: Repeat the process until all VIF values are below the threshold (commonly 10).

Example

Consider a multiple regression problem with Three features x_1 , x_2 , x_3 and the target feature y .

Bhargavi R

To calculate the VIF for each predictor:

For X_1 :

1. Regress X_1 on X_2 and X_3 :

After performing this regression, let's assume we get $R^2 = 0.96$.

2. Calculate VIF for X_1 :

$$\text{VIF}_{X_1} = \frac{1}{1 - R^2} = \frac{1}{1 - 0.96} = \frac{1}{0.04} = 25$$

For X_2 :

1. Regress X_2 on X_1 and X_3 :

Suppose $R^2 = 0.97$.

2. Calculate VIF for X_2 :

$$\text{VIF}_{X_2} = \frac{1}{1 - R^2} = \frac{1}{1 - 0.97} = \frac{1}{0.03} = 33.33$$

For X_3 :

1. Regress X_3 on X_1 and X_2 :

Suppose $R^2 = 0.98$.

2. Calculate VIF for X_3 :

$$\text{VIF}_{X_3} = \frac{1}{1 - R^2} = \frac{1}{1 - 0.98} = \frac{1}{0.02} = 50$$

3. Address Multicollinearity by Removing Variables

From the VIF calculations:

- VIF for X_1 : 25
- VIF for X_2 : 33.33
- VIF for X_3 : 50

All VIF values are high, but X_3 has the highest VIF. We'll start by removing the predictor with the highest VIF.

Remove X_3 and Recalculate VIFs for Remaining Predictors:

1. Remove X_3 from the model:

The remaining predictors are X_1 and X_2 .

2. Recalculate VIF for X_1 and X_2 without X_3 :
-

- Regress X_1 on X_2 :

Suppose $R^2 = 0.95$.

$$\text{VIF}_{X_1} = \frac{1}{1 - R^2} = \frac{1}{1 - 0.95} = \frac{1}{0.05} = 20$$

- Regress X_2 on X_1 :

Suppose $R^2 = 0.94$.

$$\text{VIF}_{X_2} = \frac{1}{1 - R^2} = \frac{1}{1 - 0.94} = \frac{1}{0.06} = 16.67$$

4. Interpret Results

- New VIFs:
 - VIF for X_1 : 20 (still high but reduced)
 - VIF for X_2 : 16.67 (also high but reduced)

Even after removing X_3 , the VIF values are still high, which indicates some multicollinearity remains. If necessary, consider removing X_1 or X_2 next or using alternative methods like PCA or Ridge Regression for better model stability.

Bhargavi ←