

1. Explain the linear regression algorithm in detail.

- A linear regression model attempts to explain the linear relationship between a dependent variable and independent variables (1 independent variable in case of Simple Linear Regression and more than 1 independent variables in case of Multiple Linear Regression.)
- Linear Regression algorithm for a given training data set, it will try to fit a best straight line in case of simple linear regression or a hyperplane for multiple regression using a metric called Residual Sum of Squares (RSS). The residuals are the difference between actual value and predicted values. Cost function is to minimize RSS using Ordinary Least Square
- We create a model using linear regression algorithm with different independent variables each time until we get a model which has significant variables which can explain target variable or dependent variable.
- $y_p = ax + b$ where b is the intercept and x is the independent variable and y_p is the predicted value for the dependent variable. This is the best fit line we got for our data.
- To understand how we got a and b parameters that make RSS become minimum we take partial derivative for each parameter and equate it with zero
- Please check the image for obtaining parameters using Ordinary Least Square method. Its nice to compute the values.

Ordinary Least Squares Method for simple Linear Regression

Let, x - data of independent variable
 \bar{x} - mean of x
 y - data of dependent variable
 \bar{y} - mean of y
 y_p - estimated (predicted) of y
 n - no. of observation

Let's find model parameters for model estimation $y_p = ax + b$ using Ordinary Least Square!

Let Residual Sum of Squares (RSS), $S = \sum (y - y_p)^2$

To get parameters that make RSS become minimum, take partial derivative for each parameter & equate it with zero.

$$S = \sum (y - ax - b)^2$$

$$\frac{\partial S}{\partial a} = 0 \quad (\text{we assume } b, x, y \text{ are const})$$

$$\frac{\partial (\sum (y - ax - b)^2)}{\partial a} = 2 \sum (y - ax - b)(-x) = 0$$

\Rightarrow we can cancel out const 2

$$\sum (-xy) + a \sum x^2 + b \sum x = 0$$

We know $\sum x = n\bar{x}$

$$b = \frac{\sum xy - a \sum x^2}{n\bar{x}} \rightarrow (1)$$

$$\frac{\partial S}{\partial b} = 0 \quad (\text{we assume } a, x, y \text{ const})$$

$$\frac{\partial (\sum (y - ax - b)^2)}{\partial b} = 2 \sum (y - ax - b)(-1) = 0$$

\Rightarrow we can cancel out 2 (constant)

$$-\sum y + a \sum x + b \sum 1 = 0$$

$\sum 1 = n, \sum x = n\bar{x}, \sum y = n\bar{y}$

$$-n\bar{y} + a n\bar{x} + nb = 0$$

$$= a\bar{x} + b = \bar{y} \rightarrow (2)$$

Insert (1) in (2)

$$a\bar{x} + \frac{\sum xy - a \sum x^2}{n\bar{x}} = \bar{y}$$

$$a \left(\bar{x} - \frac{\sum x^2}{n\bar{x}} \right) + \frac{\sum xy}{n\bar{x}} = \bar{y}$$

\rightarrow multiply $n\bar{x}$ on both sides

$$a (n\bar{x}^2 - \sum x^2) + \sum xy = \bar{y} \cdot n \rightarrow (3)$$

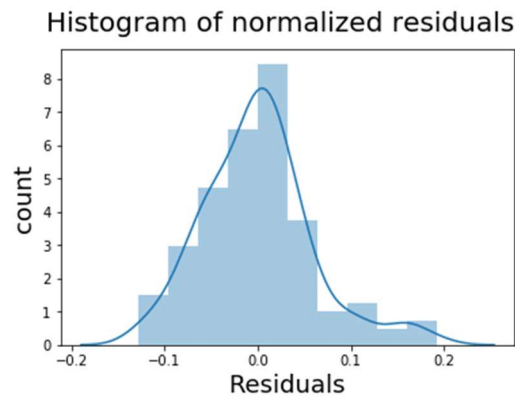
$$a = \frac{n \cdot \bar{x} \bar{y} - \sum xy}{(n\bar{x}^2 - \sum x^2)} \quad (\text{slope})$$

From (3) we get $b = \bar{y} - a\bar{x}$ (intercept)

2. What are the assumptions of linear regression regarding residuals?

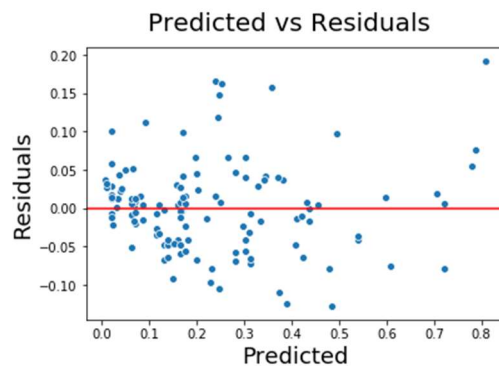
- a. **Residuals are normally distributed**, this is important to make inferences on the model that we have built, we need a notion of distribution of residuals. The residual/error terms generally follow a normal distribution with mean equal to zero. We can use histogram or q-q plot to plot the residuals.

For our car-predictions, the following histogram is generated which shows residuals are normalized.



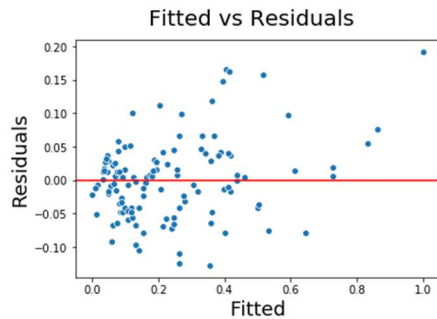
- b. **Residuals are independent of each other**, we should not see any visible patterns with the predicted and residuals.

For our car-predictions, the following graph is generated which shows residuals are independent.

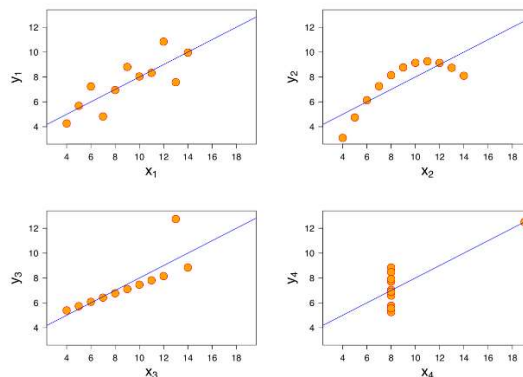


- c. **Residuals have constant variance (Homoscedasticity)**, the variance should not increase or decrease as error values changes. We will plot residuals vs actual y values.

For our car-predictions, the following graph is generated which shows most of the residuals are having constant variance.

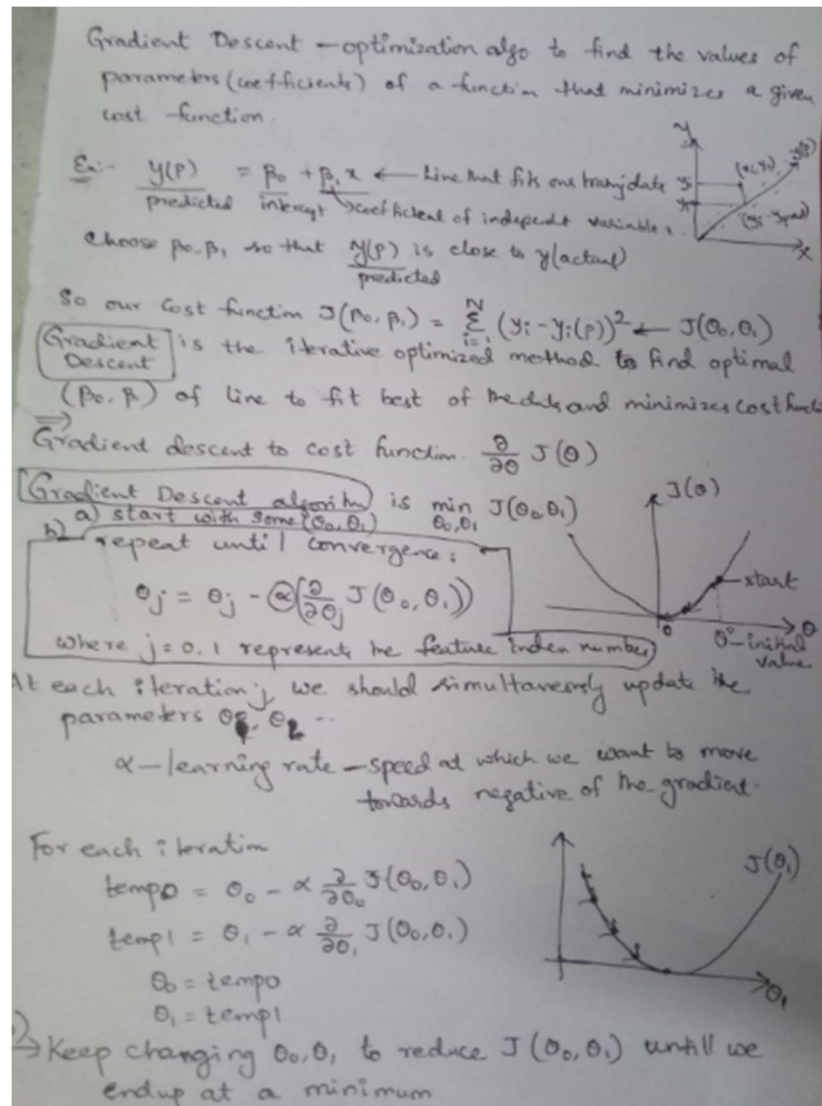


3. What is the coefficient of correlation and the coefficient of determination?
 - a. Coefficient of correlation is “R” value measures the strength and direction of linear relationship between two variables. The value of r is such that $-1 \leq r \leq +1$. The + and – signs are used for positive linear correlations and negative linear correlations, respectively.
 - b. Coefficient of determination is “ r^2 ” gives the proportion of the variance(fluctuation) of 1 variable that is predictable from the other variable. It is a measure that allows us to determine how certain one can be in making predictions from a certain model. It is the ratio of explained variation to total variation. The coefficient of determination is such that $0 < r^2 < 1$, and denotes the strength of the linear association between x and y.
4. Explain the Anscombe’s quartet in detail.
 - a. Anscombe’s Quartet was developed by statistician Francis Anscombe. It comprises four datasets, each containing eleven (x,y) pairs. The essential thing to note about these datasets is that they share the same descriptive statistics(ex: SUM, AVERAGE, STDDEV, MEAN, VARIANCE, CORRELATION COEFFICIENT,..).
 - b. But things change completely when they are graphed. Each graph tells a different story irrespective of their similar summary statistics. Below are the graphs for the 4 data sets.



- i. Dataset I ,appear to have clean and well-fitting linear models.
 - ii. Dataset II ,is not distributed normally.
 - iii. In Dataset III, the distribution is linear, but the calculated regression is thrown off by an outlier.
 - iv. Dataset IV ,shows that one outlier is enough to produce a high correlation coefficient.
- c. This quartet emphasizes the importance of visualization in Data Analysis. Looking at the data reveals a lot of the structure and a clear picture of the dataset.

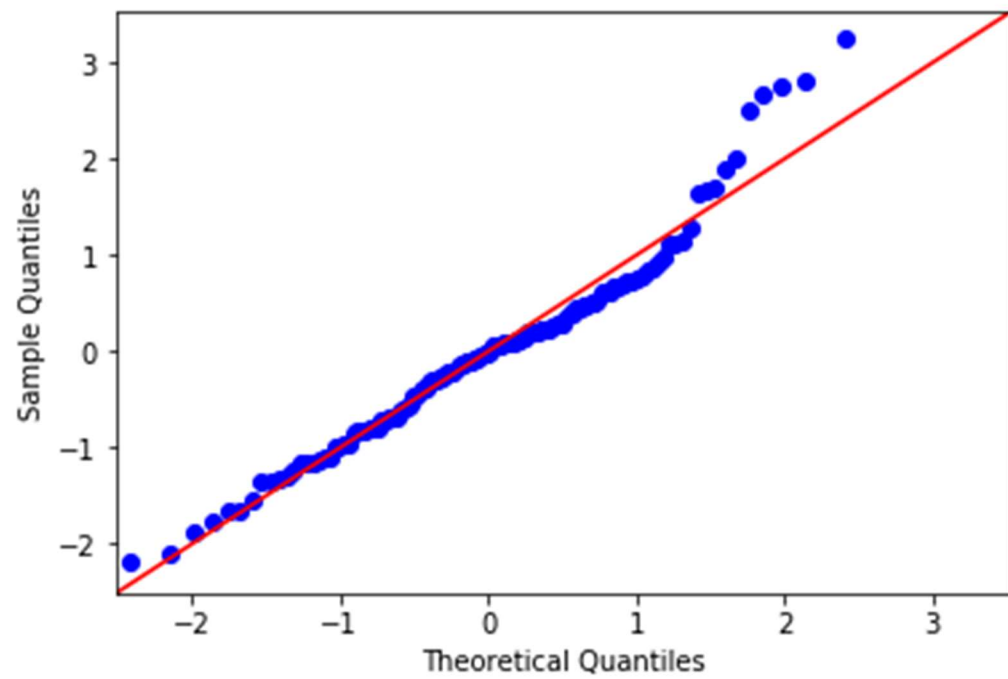
5. What is Pearson's R?
 - a. Pearson's correlation coefficient is also referred as Pearson's R. It is a measure of strength and direction of linear relationship between two variables. The value of r is such that $-1 \leq r \leq +1$. The + and – signs are used for positive linear correlations and negative linear correlations, respectively.
6. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?
 - a. Feature scaling or scaling is essential step of data processing which is applied to features to normalize when the range of raw data varies widely.
 - i. i.e, Scaling is important when we have features which vary in magnitudes, units and range and bring them to same level of magnitudes
 - ii. Scaling improve or speedup gradient descent algorithm.
 - iii. Scaling does not affect other parameters such as p-value and R-squared, it affects only the coefficients.
 - b. Standardized scaling brings all the data into standard normal distribution with mean 0 and standard deviation 1. Whereas Normalized or MinMax scaling, brings all the data in range of 0-1. The formulae used as given below:
 - i. Standardization : $X = (x - \text{mean}(x)) / \text{stddev}(X)$
 - ii. Normalization (MinMax Scaling): $X = (x - \min(x)) / (\max(x) - \min(x))$
7. You might have observed that sometimes the value of VIF is infinite. Why does this happen?
 - a. VIF, calculates how well one independent variable is explained by all other independent variables combined. $VIF(i) = 1 / (1 - R\text{-squared}(i))$.
 - b. So when $R\text{-squared}(i)$ becomes 1 then $VIF(i)$ becomes infinity. $R\text{-squared}(i)$ will be 1 when 100% variance of one variable is explained by another i.e, when two variables are exactly same. So we should look at variables with exact collinearity using correlation matrix or heatmap.
8. What is the Gauss-Markov theorem?
 - a. Gauss-Markov theorem states that ordinary least squares (OLS) estimator of the coefficients of a linear regression model is the best linear unbiased estimator (BLUE) if error terms satisfy below conditions.
 - i. Expected value of error term is zero for all observations
 - ii. Variance of the error term is constant
 - iii. Error term is independently distributed and not correlated.
9. Explain the gradient descent algorithm in detail.
 - a. Gradient descent mainly we start with some θ_0 , θ_1 for simple linear regression with intercept(θ_0) and one coefficient for independent variable(θ_1)
 - b. Keep changing θ_0 and θ_1 to reduce cost function $J(\theta_0, \theta_1)$ until we end up at a minimum



10. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression.

- The Q-Q plot is a graphical tool to assess if two data sets come from populations with a common distribution. Q-Q plot of the quantiles of X versus the quantiles/ppf of a distribution
- Importance of Q-Q plot in linear regression we can plot residuals to check if they are normally distributed instead of using histogram alone to check.
- Compare X against dist. The default is `scipy.stats.distributions.norm` (a standard normal).
 - `import statsmodels.api as sm`
 - `res = lm5.resid` #after training the model(lm5) we get residuals and plot q-q plot
 - `fig = sm.qqplot(res, stats.norm, fit=True, line='45')` # we can change `stats.norm` to `stats.t` or any distribution
 - `plt.show()`
- A perfectly normal distribution would exactly follow a line with slope = 1 and intercept = 0.

For our car-predictions, the following graph is generated which shows that residuals are normally distributed.



e.