

a) Data Preprocess

b) EDA, Stats, Assumptions!

c) Split the dataset: 80 training – 20 testing or 70-30!

d) EQUATION:

$$y = mx + c$$

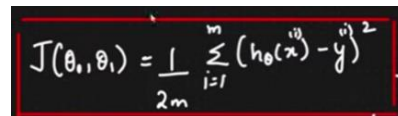
$$y = \beta_0 + \beta_1 x$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \text{i.e PREDICTED VALUE}$$

e) **BEST FIT LINE**: distance should be minimal between the predicted and real value!

i) will do it by applying the cost function

COST FUNCTION =


$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

ii) minimize the cost function to get the best fit line:

Method: GRADIENT DESCENT, will get the  $\theta_1$  values but which one is right among them?

Will get that by a term known as CONVERGENCE ALGORITHM

Later the best pt will get is known as GLOBAL MINIMA

f) **OVERFITTING** (Train: Good, Test: Poor; Low Bias, High Variance) or

**UNDERFITTING** (Train: Poor, Test: Good; High Bias, Low Variance)!

To solve OVERFITTING, will do **REGULARIZATION** (Ridge & Lasso)!

To solve UNDERFITTING, will do the following:

i) **Reduce Regularization**

ii) **Use a more complex model**: Polynomial Regression, Decision Trees, Random Forests, or Gradient Boosting Machines, AdaBoost, Neural Networks. These methods combine several weak learners to create a more powerful model that can capture complex relationships.

iii) **Increase the number of features or use feature engineering**:

Create new features, use feature transformations: Apply transformations like log transformations, polynomial features, or scaling to improve the model's ability to capture relationships.

iv) Train for more epochs (in case of neural networks) or increase the number of iterations in gradient-based algorithms like gradient descent. Sometimes the model hasn't had enough time to fully learn from the data, leading to underfitting.

v) **Add more data**

vi) **Remove noise or outliers**

vii) Try models that can handle nonlinear relationships between the features and the target, such as: Support Vector Machines (**SVM**) with nonlinear kernels.

viii) **Hyperparameter tuning**: Sometimes underfitting is due to suboptimal hyperparameters. Use techniques like Grid Search or Random Search to tune key parameters like learning rate, number of trees, maximum depth of trees, etc. in decision trees or ensemble models.

**g) Solve Overfitting:**

**1) REGULARIZATION:**

i) **Ridge Regression** (L2 regularization):

$\alpha$  value increases the coefficient of the magnitude decreases/shrinks/scale downs almost to 0

➔ Prevent overfitting!

ii) **Lasso Regression** (L1 regularization):

$\alpha$  value increases the coefficient of the magnitude decreases/shrinks/scale downs to exactly 0

➔ Prevent Overfitting & achieves Feature Selection!

**h) Solve Underfitting:**

**1) REDUCE REGULARIZATION:**

i) If we're using regularization techniques (e.g., Lasso or Ridge regression) and experiencing underfitting, it might be because our regularization strength is too high.

ii) Regularization shrinks the coefficients, and too much regularization can make the model too simple.

iii) Reduce the regularization strength by lowering the  $\lambda$  parameter. This allows the model to use more of the available features and coefficients, which can improve performance.

**i) CROSS – VALIDATION: K-Fold Cross-Validation!**

The dataset is divided into k subsets. The model is trained k times, each time using a different fold as the test set and the remaining k-1 folds as the training set.

**j) PERFORMANCE MATRIX:** helps to measure: How good our model is?

i)  $R^2$ : bigger the better

ii) Adjusted  $R^2$ : helps to remove independent variable which is least useful

iii) MAE

iv) RMSE

### **k) Common Problems & Solutions:**

#### **i) Multicollinearity:**

Occurs when independent variables are highly correlated, leading to instability in coefficient estimation.

**Solution:** Remove highly correlated features (**VIF**) or use **Principal Component Analysis (PCA)**.

#### **ii) Homoscedasticity**

Assumption that the variance of errors is constant across all levels of the independent variable(s).

**Solution:** If violated, try transforming the dependent variable (e.g., log transformation).

### **l) Model Evaluation on Test Data:**

Once the model has been trained, evaluate it using metrics such as **MSE**, **R<sup>2</sup>**, **RMSE**, etc. This helps assess how well the model is generalizing to unseen data.

**Keep checking all the below given assumption in the start, mid or in the end; whenever required!**

#### **m) Assumptions:**

**1) Linearity:** The relationship between predictors and the target is linear.

**2) Independence of Errors:** Residuals should not be correlated with each other.

**3) Homoscedasticity:** The variance of residuals is constant across all levels of the independent variable(s).

**4) Normality of Errors:** The residuals are normally distributed.

**5) No Multicollinearity:** Predictors should not be highly correlated with each other.

**6) No Significant Outliers/Influential Data Points:** Outliers or influential points should not unduly affect the model.

#### **n) How to Check and Validate Assumptions:**

**1) Linearity:** Visualize data and residuals (scatter plots).

**2) Independence:** Durbin-Watson test for autocorrelation.

**3) Homoscedasticity:** Plot residuals vs. fitted values.

**4) Normality:** Q-Q plot, Shapiro-Wilk test.

**5) Multicollinearity:** Correlation matrix, Variance Inflation Factor (VIF).

**6) Outliers:** Scatter plot, Cook's Distance, leverage statistics.

