

1 Bias \rightarrow It is the amount that a model prediction differs from the target value, compared to the training dataset.

(A high level bias can lead to underfitting)
(Bias always with training data)

A linear algorithm often has a high bias, which make them learn fast.

2 Variance \rightarrow

Low Bias ML Algorithms \rightarrow Decision Tree, K-NN, SVM

High Bias ML Algorithms \rightarrow Linear Regression, Logistic Regression

2 Variance error \rightarrow It is the amount that the estimate of the target function will change if different training data was used.

(Variance can lead to overfit)

(Variance always with testing data)

A model with high Variance may reflect random noise in the training dataset instead of target function. A model with low variance mean sample data is closer to prediction.

Low-Variance ML algo - Linear Regression, Logistic Regression

High-Variance ML algo - Decision Tree, KNN, SVM

Note → The goal of any Supervised ML algorithm is to achieve low bias & low variance.

Note - increasing the bias will decrease the variance, increasing the variance will decrease the bias.

3 Total error → It is the sum of the bias error and variance error.

The total goal is to balance bias & variance, so the model does not underfit & overfit the data.

4 When Variance Occur

If the ML model performs well with the training database, but does not perform well with the test database, then Variance Occur.

Absolutely measure of variation from the mean of the data.

Variance and Standard Deviation

Variance = $\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$

Standard Deviation = $\sqrt{\text{Variance}}$

It measures how much data is spread out from the mean.

$$\text{Variance} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

Overfitting & Underfitting

4 Underfitting — When ML is underfitting, it means it isn't learning much from the training data. They didn't have enough information on the target variable or (it only performs ~~bad~~ on training data as well ~~as~~ performs poor on testing data.

— Reasons of Underfitting —

- High bias & low variance.
- The size of training data set used is not enough.
- The model is too simple so user by performs.
- Training data is not clear & also contains noise in it.

— Techniques to reduce underfitting —

- Increases model complexity by using hyperparameter.
- Remove noise from the data.
- Increase the duration of training to get better result.

5 Overfitting — Model is said to be overfitted when the model does not account prediction on testing data.

When model get trained with so much data, it start learning from the noise of inaccurate data entries in our data set.

Ex - Overfitting might cause our ~~the~~ model to predict that every person coming to our site will purchase.

- Reason of Overfitting —

- High Variance & Low bias

- a) The model is too complex

- Small size of ~~Ex~~ training data

Note → Overfitting ~~give~~ gives good ~~accuracy~~ accuracy in training dataset but bad accuracy in testing data set whereas

Underfittness give bad accuracy in testing as well as training dataset

-> Static model

>> import statsmodels.formula.api as smf

>> lm = smf.ols(formula='sales ~ TV + radio', data=df)
• fit()

>> lm.summary()

Residual = actual - predicted

$$r = y - \hat{y}$$

R-Squared & Adjusted R-Squared for Regression

Residual Sum Square = RSS = $\sum (y_i - \hat{y}_i)^2$ = Loss
° The lower RSS, the better is the model
Prediction,

Total Sum Square = TSS = $\sum (y_i - \bar{y})^2$

$$R^2 = 1 - \frac{RSS}{TSS}$$

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

n = no. of data points

p = no. of independent variable

∴ Adjusted R^2 says that it will give more reliable result as compare to R^2 .

~~Value~~ For Linear Regression or the Regression algorithm it is better to check Adjusted R^2 instead of R^2 as in Adjusted R^2 it take n (no. of data point) as well as p (no. independent variable) to control "Overfitting".

Note : If we add useless data or column in our dataset then also R^2 is going to increase by small amount which is not necessary so we use or see Adjusted R^2 accuracy score as it divided $\frac{1}{n-p}$, no. of data point (n) & the no. of independent variable (p)

6 VIF \Rightarrow Variance inflation factor to the multi-collinearity among the feature of the dataset.

$$\text{VIF} = \text{Variance inflation factor} = \frac{1}{1 - R^2}$$

(If $\text{VIF} > 5$ (greater than 5) then the features are multi-collinear)