==Performance metrics used in Regression==: To calculate the accuracy

of a Regression model we use performance metrics.

==There are two techniques:==
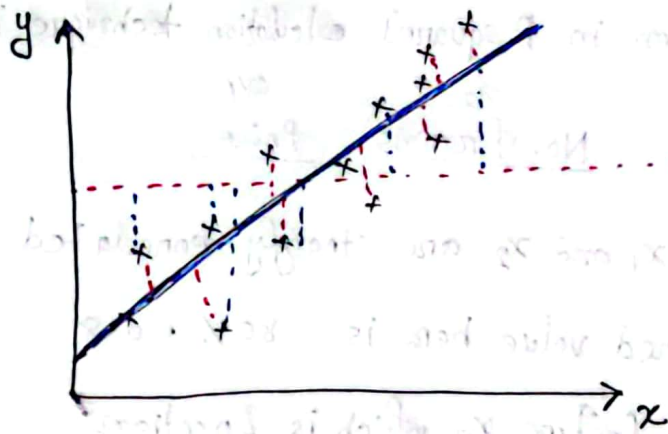
   ① R squared  ② Adjusted R squared

==R squared Formula:==

                       {Error}

$$\text{Rsquared} = 1 - \frac{SS_{Res}}{SS_{Total}}$$

               SS Res = Sum of square residual ↑

               SS Total = Sum of square total

==SS Res== = The total error $\sum(y_i - \hat{y_i})^2$   $\left[\sum(\text{actual point} - \text{predicted point})^2\right]$



==SS Total== → We take a line (straight) from the average value of $y$.

Now the summation of $(y_{avg} - actualpoints)^2$ would be SS Total

$$\sum(y_{avg} - \infty\ y_i)^2$$

So, we can write the R squared Formula like below →

$$R \text{ squared} = 1 - \sum_{j=1}^{n} \frac{(y_i - h_\theta(x))^2}{(y_i - \bar{y_i})^2}$$

$h_\theta(x) = \hat{y_i}$ (predicted point)

$\bar{y_i} = y \text{ mean}$

$$= 1 - \sum_{j=1}^{n} \frac{(y_i - \hat{y_i})^2}{(y_i - \bar{y_i})^2}$$

==R squared ranges between 0 to 1.==

==(2) Adjusted R squared:==

One of the problem in R squared calculation technique is. Souppose

| $x_1$ | $x_2$ | O/P |
|---|---|---|
| House size | No. of rooms | Price |

Both the feature $x_1$ and $x_2$ are strongly correlated with output (Price)

Suppose, R squared value here is = 80% ≈ 0.8

Let's add another feature $x_3$ which is Location

| $x_1$ | $x_2$ | $x_3$ | O/P |
|---|---|---|---|
| House Price | No. of rooms | Location | Price |

Here 3 of features are strongly correlated with output. For that R squared will increase even more. Suppose that become → 90%

Now, add another feature, which is gender ($x_4$)

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | O/P |
|-------|-------|-------|-------|-----|
| House size | No. of rooms | Location | Gender | Price |

In the case of gender, this feature is not highly correlated or important for output prediction. Although for this feature also R squared a bit.

Let's say it Increased 1% and became 91%.

So, no matter what feature we are adding, R squared is increasing

That is not right for the model accuracy. For gender feature R squared should not increase. Adjusted R squared solve this specific problem.

==Formula of Adjusted R squared:==

$$1 - \frac{(1-R^2)(N-1)}{N-P-1}$$

N = Num of datapoints

$R^2$ = R squared

P = Num of independent features.

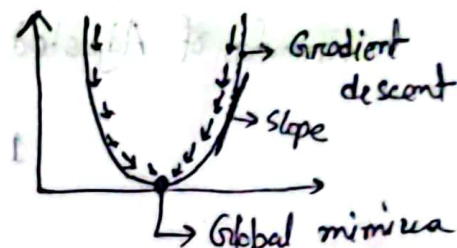We should perform both R squared and adjusted R squared to know that not every feature is important.

① Mean Squared Error (MSE)

② Mean Absolute Error (MAE)

③ Root Mean Squared Error ( RMSE)

Previously we already discussed about Mean squared Error whose Formula

was → $J(\theta_0, \theta_1) = \frac{1}{n} \sum_{j=1}^{n} \left( y_i - h_\theta(x)_j \right)^2$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \hat{y}_i \right)^2 \rightarrow \text{quadratic Equation}$$
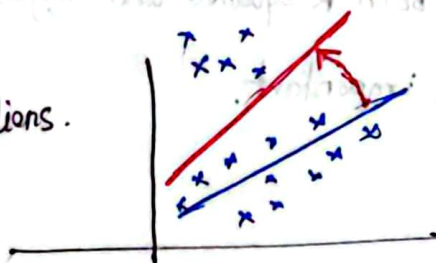
**Advantage of MSE:**

① Equation is differentiable.

② It has only one local/Global minima

**Disadvantage of MSE:**

① Not Robust to outliers.

If dataset has outliers, Best fit line moves away from where it
should be to ⟨⟩ the side of outliers a bit. (can't handle situation)
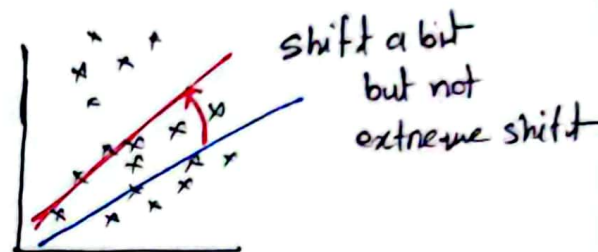
② It is not in the same unit.

As we are squaring the errors, then the unit will also be squared.

② Mean Absolute Error (MAE):

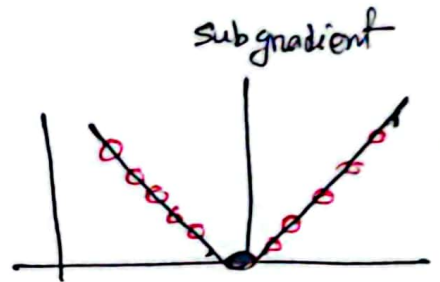Formula: $\frac{1}{n} \sum\limits_{i=1}^{n} |y_i - \hat{y}_i|$

Advantage:

① Robust to outliers

② It will be in the same unit

shift a bit
but not
extreme shift

Disadvantage:

① Convergence usually take more times.

subgradient

③ Root Mean squared Error (RMSE):

Formula, $RMSE = \sqrt{MSE}$

$$= \sqrt{\frac{1}{n} \sum\limits_{i=1}^{n} (y_i - h_\theta(x)_i)^2}$$

Advantages:

① Same Unit

② Differentiable

③ 1 Global minima

Disadvantage:

① Not Robust to outliers

When you have outliers → USE MAE

When you don't have outliers → USE MSE, RMSE