

→ Supervised Learning

→ 1 input column, 1 output column

→ Separate input & output column into different variables $\begin{matrix} x \\ y \end{matrix}$
from sklearn.model_selection import train_test_split
 $X_{train}, X_{test}, y_{train}, y_{test} = \text{train_test_split}(X, y, \text{test_size}=0.2, \text{random_state}=2)$

from sklearn.linear_model import LinearRegression
 $lr = \text{LinearRegression}()$

$lr.\text{fit}(X_{train}, y_{train})$

$lr.\text{predict}(X_{test}.\text{iloc}[:, \text{values}].\text{reshape}(-1,))$

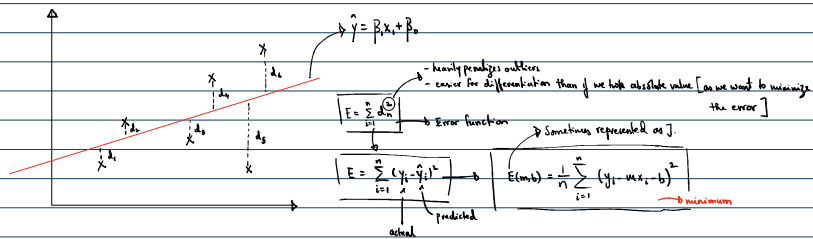
$m = lr.\text{coef}$
 $b = lr.\text{intercept}$
→ offset

weightage of input variable on output variable [how much the output column depends on the input column].
Dependence

→ line of best fit: a line which passes as close as possible to the data points

How to find m & b for line of best fit?
① closed form solution → Ordinary Least Squares (OLS)
② Non-closed form solution → Gradient descent

$$\textcircled{1} \quad b = \bar{y} - m\bar{x} \quad m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$



$$\begin{aligned} \frac{\partial E}{\partial b} &= \frac{\partial}{\partial b} \sum_{i=1}^n (y_i - mx_i - b)^2 = 0 \\ &= \sum_{i=1}^n 2(y_i - mx_i - b) = 0 \\ &= \sum_{i=1}^n (y_i - mx_i - b) = 0 \\ &= \sum_{i=1}^n (y_i - mx_i - b) = 0 \end{aligned}$$
$$\begin{aligned} \frac{\partial E}{\partial m} &= \frac{\partial}{\partial m} \sum_{i=1}^n (y_i - mx_i - b)^2 = 0 \\ &= \sum_{i=1}^n 2(y_i - mx_i - b)(-x_i) = 0 \\ &= \sum_{i=1}^n (y_i - mx_i - b)(-x_i) = 0 \end{aligned}$$
$$\begin{aligned} \frac{\partial E}{\partial b} &= \frac{\partial}{\partial b} \sum_{i=1}^n (y_i - mx_i - b)^2 = 0 \\ &= \sum_{i=1}^n 2(y_i - mx_i - b) = 0 \\ &= \sum_{i=1}^n (y_i - mx_i - b) = 0 \end{aligned}$$
$$\begin{aligned} \frac{\partial E}{\partial m} &= \frac{\partial}{\partial m} \sum_{i=1}^n (y_i - mx_i - b)^2 = 0 \\ &= \sum_{i=1}^n 2(y_i - mx_i - b)(-x_i) = 0 \\ &= \sum_{i=1}^n (y_i - mx_i - b)(-x_i) = 0 \end{aligned}$$

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \textcircled{2}$$

→ Regression Metrics
MAE ①
MSE ②
RMSE ③
 R^2 ④
Adjusted R^2 ⑤

①
 $MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$

Pro	Con
→ Same units	→ Modulus function graph isn't differentiable at 0
→ Robust outliers	

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

Pro	Con
→ loss function is differentiable	→ not same units (not easy interpret)
	→ Not robust to outliers

③
 $RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$

Pro	Con
→ loss function is differentiable	→ Not robust to outliers
→ Same units	

④ MAPE (Mean absolute percentage error)

$$y \rightarrow \hat{y} \rightarrow y - \hat{y} \rightarrow \frac{|y - \hat{y}|}{y}$$

→ Division by zero
→ Outliers

④ R^2 - How much better regression line is compared to the mean line
→ coefficient of determination / goodness of fit

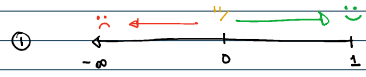
⑤ $R^2_{adj} = 1 - \frac{(1 - R^2)(n - 1)}{(n - 1 - k)}$, where $n = \# \text{ rows}$
 $k = \# \text{ of input cols}$

* irrelevant R^2_{adj} ↓
* relevant R^2_{adj} ↑
* when we have lots of cols
* there may exist cols which aren't as effective in predicting the output.

$$R^2 = 1 - \frac{SS_R}{SS_M} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

mean_absolute_error

$$R^2 = 1 - \frac{SS_R}{SS_M} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$



② — % of variance in the output column is being explained by the input columns.

solution

sklearn.metrics $\begin{cases} \text{mean_absolute_error} \\ \text{mean_squared_error} \\ r^2_score \end{cases} (y_test, y_pred)$

* penalizes original R^2 as you increase the # of features

Problem: if we add irrelevant columns, $R^2 \uparrow$ or — instead of \downarrow .

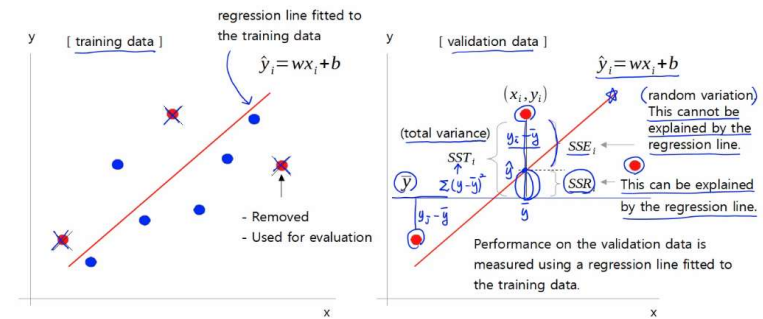
When SST is high (Data is very dispersed) there is a high chance the SSR will be high also but overall R^2 will not get penalized very much.

→ How much of the features are we able to exploit.

→ " " " " variance in the features are we able to utilize.

■ Performance evaluation metrics: MSE and R-squared

- The performance of the linear regression model can be measured by MSE or R^2 .
- Split the data set into training and validation data, and apply the regression line fitted to the training data to the validation data to evaluate the performance of the model.
- R-squared is a measure that determines the proportion of the variance of the dependent variable that can be explained by the independent variables. $R^2 = SSR / SST$ *



① mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$$

② R^2

(SS: Sum of Squares)

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (\text{Total Variance of } y)$$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (\text{Error, Residual})$$

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (\text{Regression})$$

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

($SST = SSE + SSR$)