

Lecture #19

Regression

over-view of previous lecture

• output would be in numerical value.

• target variable \Rightarrow numerical

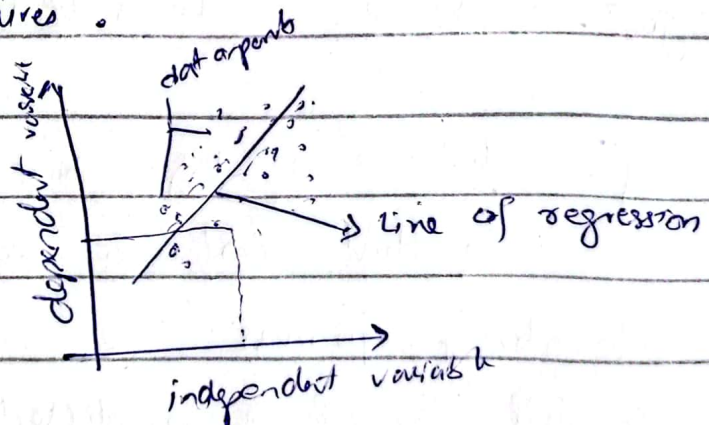
Algorithms

1. Linear Regression

\hookrightarrow supervised machine learning for predicting a continuous output based on one or more input features.

• a line touching/close to data points

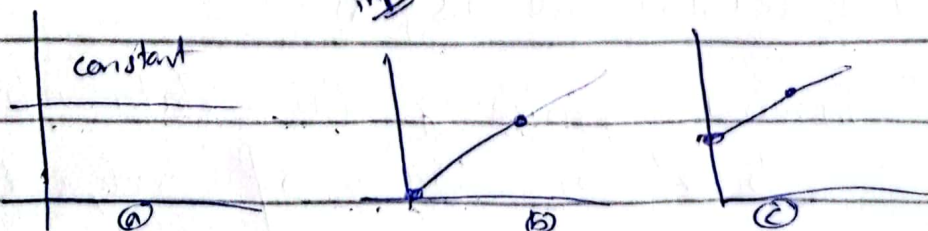
helps to predict actual data



$$\text{model: } f_{w,b}(x) = w \cdot x + b$$

x input

(w, b) parameters



(i) one feature \Rightarrow uni-variant regression

(ii) multiple features \Rightarrow multi-variant regression

$w \rightarrow$ weight
 $b \rightarrow$ bias

our goal \rightarrow finding those parameters (w, b) that
when value of x gives the output
close to our data-actual value.

\Rightarrow model training \Rightarrow bring the value of w and b
so that we can get the result.

Loss/cost function (Tells us difference b/w actual & predicted)
 \hookrightarrow distance function

input data \rightarrow plotted
at w randomly \rightarrow w and b \rightarrow example
Actual \rightarrow (Prediction Result) \rightarrow y and \hat{y}
 \rightarrow (cost function) \rightarrow y and \hat{y}

if distance is large then our value of w
so that cost is reduced.

iteratively process

until actual & predicted distance is minimum.

(2) Decision Tree Regressor

" involves partitioning the data into sub-sets
based on value of independent variable.

⇒ leaf node ⇒ average score/value

③ Random forest regressor

④ Gradient Boosting Regression

"build multiple decision trees sequentially,
each one correcting the error"

Implementation

dataset: USA housing

⇒ ensure data is correct, clean and according to
requirement

① Data Pre-processing

→ Visualizer (heat map)

- checking data
- checking shape (5000, 7)
rows → 5000, columns → 7
- view / check data
- info()
- handling null value
- drop / substitute
- describe() # 5-number summary

② ⇒ Setting model ready

(i) split the dataset \rightarrow input & target variable
 \bar{X} \bar{y}
grab everything except target variable we're trying to predict

(ii) using sklearn pre-processing (i.e standard scaling)

(iii) split the training data into $\left\{ \begin{array}{l} \text{train} \\ \text{test} \end{array} \right.$

\therefore (using train-test-split)

(iv) Model Training

$\left\{ \begin{array}{l} \text{import model} \\ \text{create instance} \\ \text{pass training data into model} \end{array} \right.$

(v) Model Testing

$\left\{ \begin{array}{l} \text{find predictions using test data} \\ \text{compare with actual \& predicted} \\ \text{(np.column_stack)} \end{array} \right.$

one way to evaluate regression model

Residual Analysis (note 3.6)

- in linear regression is a way to check how well the model fits the data
- involves looking at difference b/w actual data points & prediction from the model
- in good model \Rightarrow residual should be randomly scattered around zero on plot

(over estimation, under-estimation)

visualizer
→ distribution plot → symmetric → good regressor
→ scatter plot

(vi) Evaluator

— from metrics

- (i) mean square-error
- (ii) root mean square error

→ and some is → (vi) for all algorithms.