Assignment - A2     4|403

- Title: Naive bayes classification
- Problem Statement: Download PIMA Indians Diabetes dataset. Use naive bayes algorithm for classification. Load the data from csv file & split it into training & test datasets. summarize the properties in the training dataset so that we can calculate probobilities & make predictions classify samples from a test dataset.
- Objective: Understand Naive bayes algorithm for classification.
- Outcome: predict whether the person has diabetes or not using Naive bayes classifica- based on parameters in dataset like blood pressure.
- S|W H|W requirements: 64 bit os (UNIX/ LINUX) python3 . Jypyler environment.

- Theory: Naive bayes classifiers are a family of simple probalistics classifiers
They are based on baye's theorem which describes the probability of a certain event occurring based on the prior knowledge of conditions that might be related to event.

bayes theorem is stated mathematically

$$P(A|B) = \frac{P(B/A)\ P(A)}{P(B)}$$

Scanned with CamScanner

where A,B are the events.
P(A|B) is a conditional prob. The likelihood
of event A occuring knowing that B is true
P(B|A) is a also conditional the likelihood
of B occuring knowing that A is true
P(A) & P(B) are marginal probabilities

Naive bayes is a technique for constructing
classifiers which applies the above theorem
with the strong assumption
These models assign class labels to
problem instances. represented as vectors of
features values

A family of algorithms based on one
common principle form the naive bayes
classifier the principle is that,
a particular feature is independent of the value
of any other feature. given the class
variable each feature contributes independently
to the probability of the positive outcomes

− About the data set :
The dataset is originally from the
National institute of Diabetes & Digestive
& Kidney Disease.

The objective of the dataset is to diagonostical
ly predict whether or not patient has
diabetes based on certain diognostic
measures included.
Several constraints were placed on the
selection of these instances from the
larger database in Particular all patients
here are at least 21 years old.

- conclusion: The Naive Bayes classifier was
successfully applied to the cleaned dataset
& the outcome was predicted with an
good accuracy.