Dotta preprocessing using Scikit-learn

import sklean

standard format refers to data that has a mean and unit variance (ie. standard deviation = 1).
The process of converting data into this format is called data standarduation.

2 = 21 - U

2 - new data value n - old data value 11 - overall mean J- standard deviation

from sklearn. Meprocessing import scale col_standardinzed = scale (pizza - data)
scale also takes the axie parametal.

Range scaling

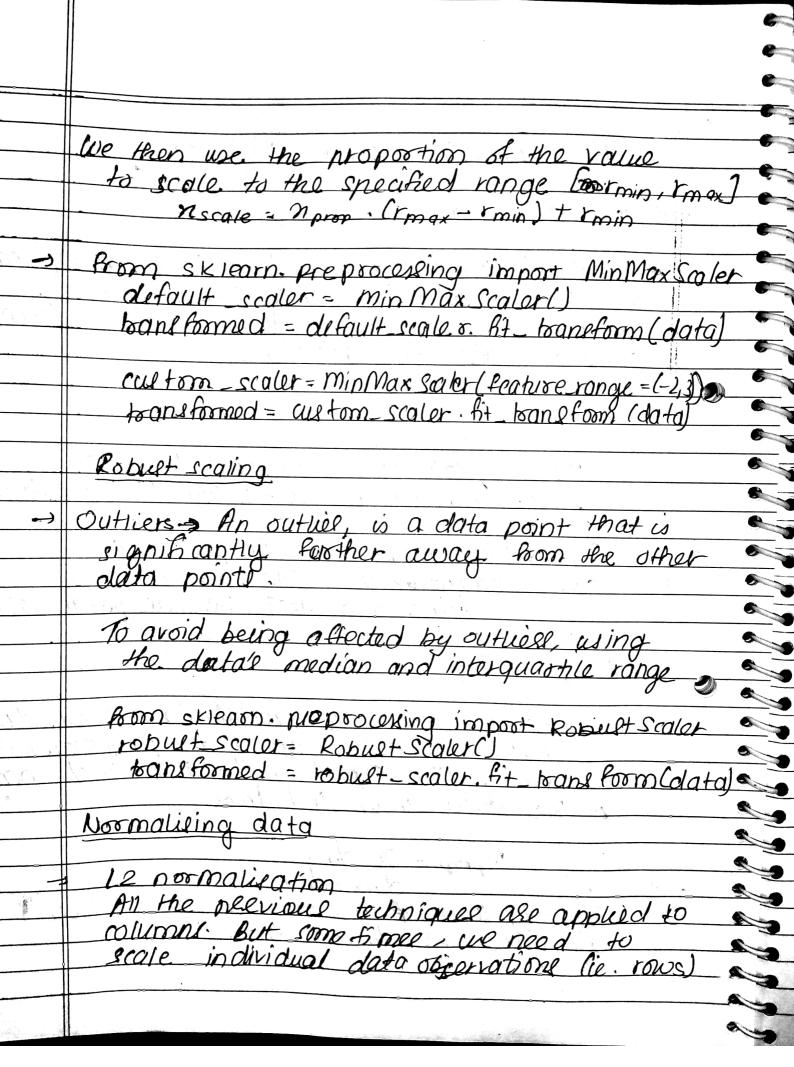
This is a two step process. For a given data value, n, we first compute the propostion of the value with respect to the min and max of the data dmin and dmax

Maron = n-dmin

dmox-dmin

This only worse is down not all data values a

are some is dmax & dmin



for instance, when clustering, we need to apply 12 parmalisation to each sow in order to calculate asine similarity scores. In general terme 12 nows of a sow is just the squall mot of the sum of squalled values of HD. 5000 $X = [n_1, n_2 \dots n_m]$ $X_{12} = \left[\frac{n_1}{l}, \frac{n_2}{l}, \frac{n_m}{l}\right], l = \left[\frac{m}{2}, \frac{n_2}{l}\right]$ L2 norm from sklearn. proprocessing import Normalizar normalizer= Normalizar() bansformed = normaliser. Fit transform (data) Data imputation. If only few values are missing, we can perform data imputation to substitute the missing data The four methods are · Using the room mean value · Using the median · Using the most frequent value defour method is using macan from sklean impute import simple imputer imp-mean = Simple Imputer () transformed = imp-mean fit transform (data)

veing the stoategy regroord, we can change the behavious stoate gy = g'mediah' 11 'most-frequent' for a constant value, Simple Imputer (5 trategy= 'constant', fill_value = -1) There are more advanced imputation methods such as K-nearest negligh bours and MICE but are rarely used in industriel since data of deaned. PCA - Poincip (com ponent analysis Used for reducing the features in the dataset which are sell important. 6 PCA extractl the principleal components of the dataset, which are uncorrelated set of latent 6 8 variables that encompass most of the information from the original dataset 6 we can use the mon n-components regulated to specify the number of principles components. default setting is to extract m-1 principal com ponente from sklearn decomparition import PCA • -PCO_Obj = PCA() pca_obj = PcA (n_componente = 3) 4 pc = pca obj. fit boardform (data). round(3) 1

| | II | |
|---|---|-------|
| | | |
| | | |
| | | |
| | | |
| | Lobelled data | |
| | | |
| | A big part of data science is to classify a data set into seperate categories. For example, we can classify a dataset of be turnous as either malignant or benign. | |
| | a data set into seperate caregories. For | |
| | example, we can downsify a doctaset of b | reast |
| | turnous as either mallament or benign. | |
| | . 0 | |
| | | |
| | | |
| _ | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | · | |
| | | |
| | | |
| | | |
| | | 1 |
| | | |
| | | |
| | | |
| | | |
| | | |