

PART 1 - Data Preprocessing.

workflow

Data Preprocessing

- import data
- clean data - * important step IRL.
- split into train & test sets.

Modelling

- Build Model
- Train Model
- Make Predictions.

Evaluation

- calculate performance matrix
- Make verdict.

Train-Test split.

usually 80/20 percent

↓
Training

→ Used trained model on test set that has known results to compare with models result.

Feature-Scaling

Remember (always applied to columns & ^{never on} ~~not~~ rows)

① Normalization
$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Most values b/w $[0, 1]$

where x is the value of column

x_{\min} is lowest & x_{\max} is highest

② Standardization
$$x' = \frac{x - \mu}{\sigma}$$
 where μ is mean/average

Most values except $[+3, -3]$ extremes or outliers.

σ is std. deviation

- with unscaled features, the lower valued/unit column become irrelevant & make small difference to large scale values.

Practical.

↪ allows to import a library or function or any module.

(1) import libraries

- (1) numpy - works with arrays for inputs.
- (2) matplotlib - for charts/graph/visualizations.
- (3) pandas - import dataset / create matrix of features & dependent variable vector.

◦ Libraries - ensemble of module containing ^{functions} ~~modules~~ & classes.

- import numpy as np
we can call with np ^{now} ~~now~~.
- import matplotlib.pyplot as plt.
↳ interested in this module.
- import pandas as pd.

(2) Import dataset

- (1) dataset = Pd.read_csv(filename.extension).
- create a dataframe with all values in data set.

(2) create 2 new entities

(1) matrix of features & (2) dependent variable vector

* features - the column with which we are going to predict
(independent variables) dependent variable.

* dependent variables - last column, that need to be predicted of dataset

full form - index location

feature of pandas, allows us to take index of columns/rows we want to extract.

* `x = dataset.iloc[:, :-1].values`

↓
Takes all rows (in python it means a range, without

LB or UB means all).

Here `[:, :-1]`

↓
Row
all

↓
column
(all except last).

so,

* `y = dataset.iloc[:, -1]`

③ Handling Missing Values.

Method I: Just delete the row of missing value.

- works good for large datasets, with low % of missing data.

Method II: Replace missing value by average of all data in column.

* Scikitlearn - has lot of tools & preprocessing tools for this task.

* we use a module called impute in scikitlearn & a class import the SimpleImputer from it.

`imputer = SimpleImputer(missing_values=np.nan,`

`strategy = mean)`

↓
value to replace
↑
to be replaced with

- Now, this object class has functions that will be used to find missing values & then impute it.

• `imputer.fit(x[:, 1:3])`

↑
expects all the column with numeric

• `imputer.transform(x[:, 1:3])` ← returns new updated matrix of features
↑
value
↑
upper bound is excluded in python (3-1) = 2 cols.

So to change original as it ~~go~~ have feature matrix

- $X[:, 1:3] = \text{imputer.transform}(X[:, 1:3])$

④ Encoding Categorical Data.

* It will be difficult to co-relate / compute these features with categorical string data, Thus we need to turn these categories to numbers.

Idea

Method I: Encode France to 0

Spain to 1

Germany to 2

- However model may interpret as order / numeric order & interpret as the order matters.

Better To-Do: ~~Turn~~ Turn to one hot encoding i.e.

turn the category into total columns with total no. of distinct category.

So for France Spain & Germany we get

for France 1 0 0

Spain 0 1 0

Germany 0 0 1

① Encode Independent Variable

- For this we use 2 classes

① ColumnTransformer class from composed module of skit

② OneHotEncoder class from preprocess module of skit

what kind of transformation
on which indexes

① $ct = \text{ColumnTransformer}(\text{transformers}, \text{remainder})$

to specify we are doing
encoding transformation

class name
to encode to

columns that won't
be applied for
transformations.

* $\text{transformers} = [(\text{'encoder'}, \text{'OneHotEncoder'})]$

$[0]] \rightarrow$ indexes to transform/encode

codename to say we wanna
keep columns as it is &
won't apply transformation.

Remainders z: 'passthrough'

↑
If we don't apply the remaining columns
won't be added.

• with this class, we can fit and transform at once
with function

~~data~~ $x = \text{ct.fit_transform}(x)$

↑
Machine learning use ~~np~~ numpy array so we
need to convert it to np array

We call $x = \text{np.array}(\text{ct.fit_transform}(x))$

② Encode Dependent Variable.

• we use LabelEncoder class from preprocessing from scikit

$\text{lbl} = \text{LabelEncoder}()$

↑
only one column with text
value, so its obvious what
needs to be performed

$y = \text{lbl.fit_transform}(y)$

↑
doesn't need to be numpy array since
its single dependent var. array.

⑤ Splitting Dataset into Train & Test set.

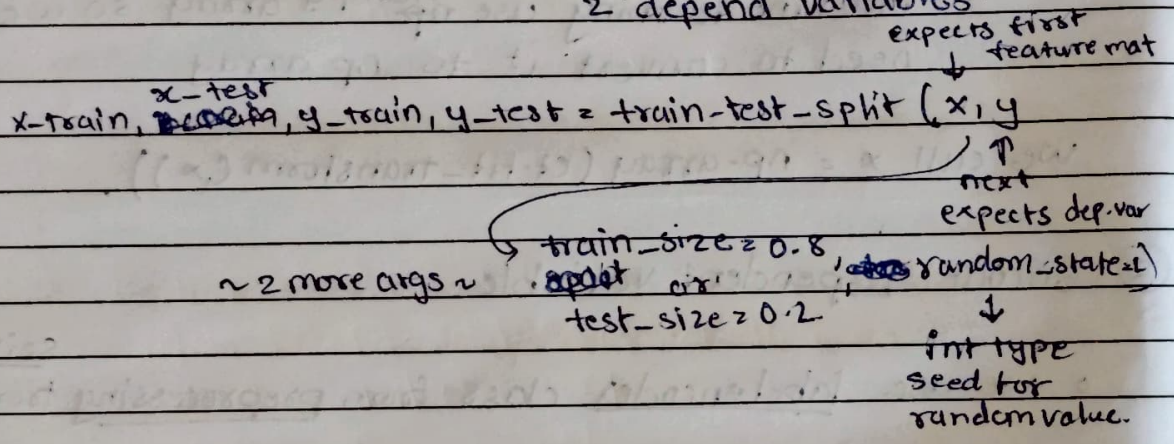
* We apply feature scaling after splitting data.

↓
because we only use train set for ML model while test set
represents real world data. & shouldn't be altered.
~~It just needs to be fed to model & get op.~~

- Also when applying feature scaling we use mean & standard deviation that might leak some info to train set which isn't supposed to happen. The ML model will get into on test set & may fit itself for it.

- We scale test later after splitting & scaling train set.
For this will use sklearn module that has.
model = Selection class with train-test-split fn.

• We split into 4. 2 feature matrices & 2 depend variables.



⑥ Feature Scaling (Not necessary for all models, works well all the time)

2 techniques (1) Standardization (with mean(μ) & deviation(σ)).
↳ (+3, -3).

(2) Normalization (with x_{min} & x_{max}).
↳ (0, 1)

works/recommended if we have normal distribution among most of our features.
(recommended for specific situations).

• For this step, we use same scikit module.
from sklearn.preprocessing import StandardScaler
↑
class.

SC = StandardScaler() ← automatically applies standardisation.

- Pummy variables \rightarrow variable columns after encoding categorical data.

- We don't need to apply standardisation to dummy variables since they are already b/w $(-3, 3)$ or $(0, 1)$

① First we fit our scales on training sets.

SD `x_train[:, 3:]` 2 `sc.fit_transform(x_train[:, 3:])`

↑
everything
after 2 with
3rd column

② Next we fit χ -test.

- we only apply transform method here since this data is like new data received IRL by model.

- The features of test set need to be scaled by the same scales used on training set.

- we do not calculate a new fit. we got that from fit-transform ^{earlier}

$x_test[:, 3:] = \text{sc.transform}(x_test[:, 3:])$

SUMMARY:

① All Libraries for preprocessing

- import pandas as pd [General]
- " numpy as np
- from sklearn.impute import SimpleImputer (for missing values)
- " sklearn.compose ColumnTransformer (for encoding categorical data)
- " sklearn.preprocessing OneHotEncoder (for encoding label/dependent data)
- " " " LabelEncoder (for splitting feature)
- " sklearn.model_selection train_test_split (for splitting feature)
- " sklearn.preprocessing StandardScaler (for scaling)

② To import dataset (and assign dep. & indep. variable)

```
dataset = pd.read_csv('./loadfilename')
```

```
x = dataset.iloc[:, :-1]
```

```
y = " "[:, 1]
```

③ Handle Missing data

```
imp = SimpleImputer(missing_values = np.nan, strategy = 'mean')
```

```
imp.fit(x[:, 1:3])
```

```
x[:, 1:3] = imp.transform(x[:, 1:3])
```

→ first fit to compute then transform to replace

④ Encoding Categorical Data

① Indep. variable

```
colTr = ColumnTransformer(transformers = [('encoding', OneHotEncoder(), [0])], remainder = 'passthrough')
```

```
x = np.array(colTr.fit_transform(x))
```

② Dep. variable

```
lbEnc = LabelEncoder()
```

```
y = lbEnc.fit_transform(y)
```

⑤ Split Data

```
x_tr, x_test, y_tr, y_test = train_test_split(x, y, test_size = 0.2,  
                                              random_state = 0)
```

⑥ Scaling (only when necessary)

```
sclr = StandardScaler()
```

```
x_train[:, 3:] = sclr.fit_transform(x_train[:, 3:])
```

```
y_train[:, 3:] = sclr.transform(x_test[:, 3:])
```