

py 03/03/21

Ch. 2 Con Notes

- (1) Big Picture
- (2) Frame the problem \rightarrow Pipelines
- (3) ~~Performance~~ & Notation & Performance Metrics

$x^{(i)}$ is the i th feature vector so
eg. If 4 ~~predictor~~ attributes in a given
feature labeled x_n

Then

$$(x^{(1)})^T = (x_{11}, x_{21}, x_{31}, x_{41}) = (1000, 423, 322,000, 4)$$

random #s

$X =$ matrix of features

$$\begin{pmatrix} (x^{(1)})^T = (x_{11}, x_{12}, \dots) \\ \vdots \\ (x^{(i)})^T = (x_{i1}, x_{i2}, \dots) \end{pmatrix}$$

$h \leftarrow$ hypothesis fcn (prediction fcn)
 $\hat{y}^{(i)}$ is a predicted value

- RMSE is sensitive to outliers (Euclidean ^{norm} distance, or ℓ_2 or $\| \cdot \|_2$ or $\| \cdot \|$)
- MAE uses Manhattan norm

$$MAE(X, h) = \frac{1}{m} \sum_{i=1}^m |h(x^{(i)}) - y^{(i)}|$$

(4) Check Assumptions

— what will data be used for what, are assumptions of models, etc.

py 39 **Note** Univariate regression means predict ONE $y^{(i)}$
multivariate regression means predict
a set(?) of $y^{(i)}$

pg 2

(5) Test set

(5) look @ data structure

- info() # Num nonnull & dtypes
- describe() # estimators
- value_counts() # Individual instances
- hist(bins=..., figsize=(,))

(6) Create a test set:

- don't look to avoid DATA SNOOPING
- avoid sampling biases by stratified samples

youtube

Types of CV (1) Random states

(1) leave one out: leave P out CV

exp1  Test Train

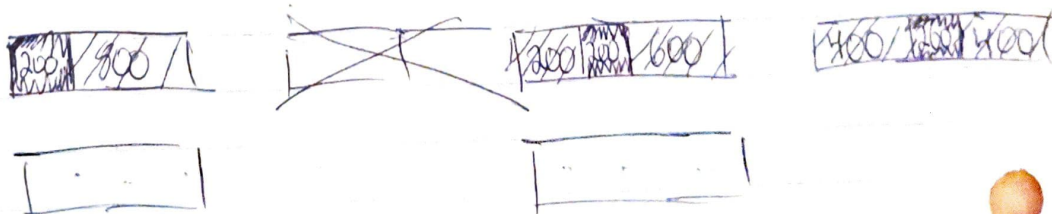
exp1 

notes: Requires many iterations to be representative
 low bias: you get better
 (1) only 1 model

(2) k-fold CV

k=5 (5 experiments)

+ folds are $\frac{m}{k}$ with $1000 = m + k=5$, each fold consists of k of 200 data units (?)



Final mean of ~~each~~ accuracy across k-folds or tell other estimators (min, max, mean)

② K-Fold - CV

Notes:

- If clustering of data (eg. Think classification) then you get unrepresentative results

③ Stratified CV



- At least some # of each class is contained ~~in~~ w/i training + testing
- Dependent on proportions in data set

★ ④ Time-series cross validation

T1

T2

T3

T4

T5

T6

T7

← Targets

Predict Stock eg:

~~D1 D2 D3 D4 D5~~

Stock

D1 | D2 | D3 | D4 | D5

Output

D6

D2 D3 D4 D5 D6

D7

pg 4

Encoding:
-- one is not ideal for categories
w/ large #'s of attr

e.g.: Professions \Rightarrow Cook, Police man, etc.
Country Code \Rightarrow "..."
species \Rightarrow animal 1, two, 3, ...

Feature Scaling

(1) Min-max scaling \leftarrow sklearn MinMaxScaler

(2) Standardization

\hookrightarrow Typical statistical procedure that accounts for $\sqrt{\sigma^2} = s_x$?

- does not bound to specific range
- less affected by outliers

★

~~WHERE TO APPLY TRANSFORMATIONS~~

Apply only on test ~~val~~ holdout validation set @ very eval. Otherwise it's fine to reuse this.