

CLASSIFICATION

2020-04
Magsblust

Machine Learning

SUPERVISED

CLASSIFICATION REGRESSION

TARGET:

discrete
categorical

continuous
numeric, ordered

SUB-METHODOLOGIES:

time series
Anomaly detection
NLP

time series
anomaly detection

ALGORITHMS:

Logistic Regression
SVM
Decision Tree
Random Forest
KNN
Neural Network

OLS (linear regression)
GLM (tweedie regressor)
Polynomial Regression
Support Vector Regressor
Decision Tree Regressor
Neural Networks
LASSO, LARS

UNSUPERVISED

CLUSTERING DIMENSIONALITY REDUCTION

SUB-METHODOLOGIES:

Anomaly Detection
NLP

ALGORITHMS:

K-Means
DBSCAN
Hierarchical

PCA
t-SNE
factor Analysis

REINFORCEMENT

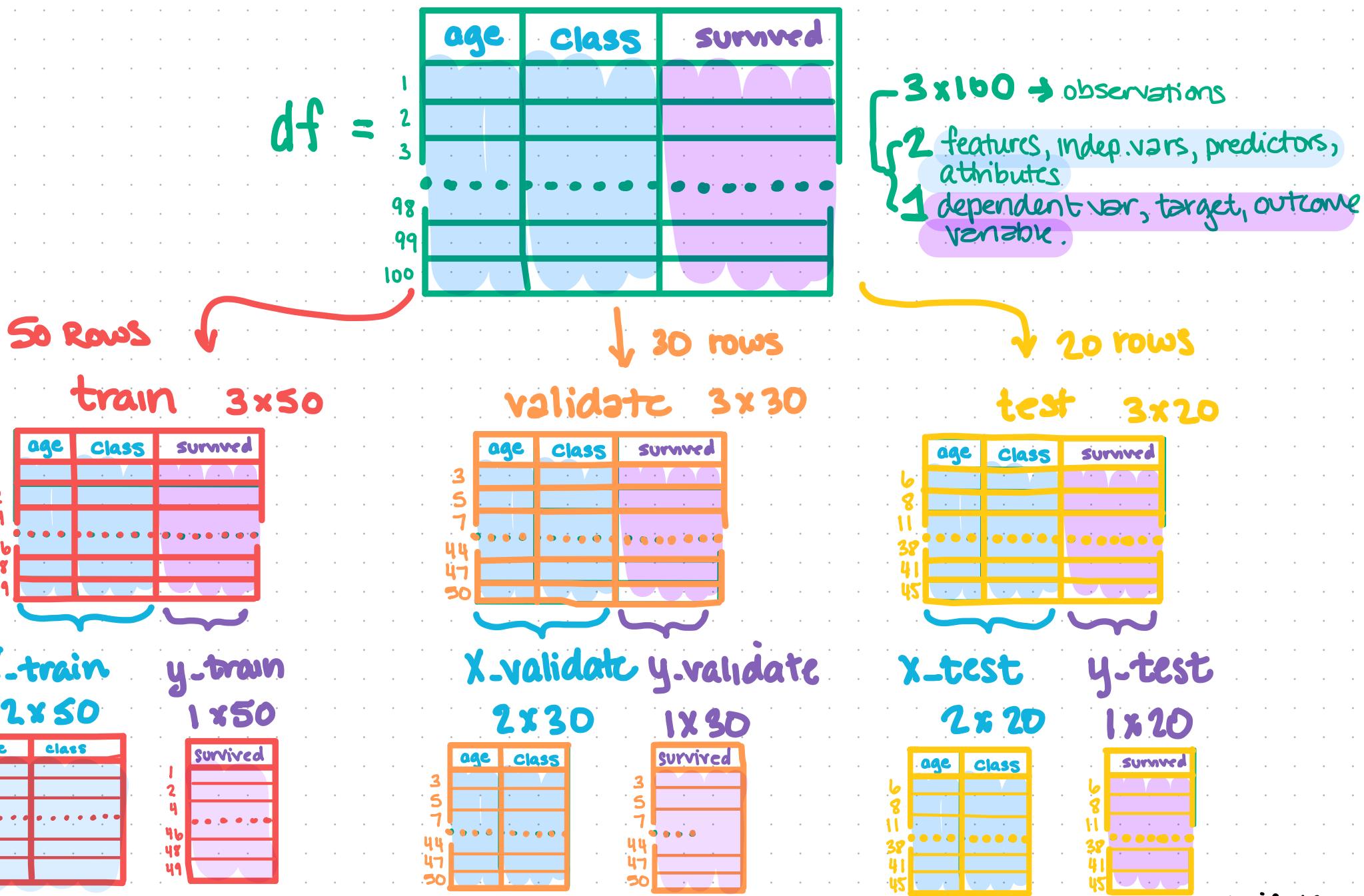
SUB-METHODOLOGIES

Anomaly Detection
NLP
Recommender Systems

ALGORITHMS:

Collaborative Filtering
Content Filtering
Neural Networks

Splitting & Sampling Dataframes



A SAMPLE FLOW

Acquire.py

```
myfile.candas as pd  
def acquire_titanic():  
    df = pd.read_csv('titanic.csv')  
    return df
```

Prepare.py

```
Import pandas as pd  
from sklearn.modelselection import train_test_split  
def clean_titanic():  
    <handle nulls>  
    <handle data errors>  
    <handle anomalies>  
    <new features>  
    <other things to prepare your data before splitting>  
    return df  
  
def split_titanic():  
    train_validate, test = train_test_split()  
    train, validate = train_test_split()  
    return train, validate, test  
  
def prep_titanic():  
    df = clean_titanic()  
    train, validate, test = split_titanic()
```

Explore.ipynb

```
Import pandas as pd  
Import acquire  
Import prepare  
Import matplotlib.pyplot as plt  
Import seaborn as sns
```

```
df = acquire.acquire_titanic()  
train, validate, test = prepare.prep_titanic()
```

EXPLORE

- 1. Univariate
- 2. Bivariate
- 3. Multivariate / Q&A

DRAW Conclusions
Select Features

MODEL

- 1. Fit on TRAIN
- 2. Evaluate on train, validate
- 3. Repeat with new Model (features, algorithm, hyperparams)
- 4. Select Best Model
- 5. Evaluate on test.

DELIVER

- 1. Create a final report with highlights, conclusions, recommendations
- 2. Create a slide presentation to deliver of analysis results, highlights, conclusions, recommendations.

YOU ARE
HERE!

EXPLORE

Univariate Exploration

- Exploration of 1 variable at a time

WHY?

① ID outliers — Anomalies
— Errors

② ID imbalanced data

/ target
\\ Features with
little - No Entropy

WHAT?

1. Document takeaways, findings
2. Document action needed, like
"Replace data error with —"
3. Document new Questions to ask the data
4. Return to prepare as needed.
5. Progress to Bivariate Exploration.

HOW?

Numeric DTypes

1. descriptive stats-

2. histogram



`df.describe()`
`Series.hist()`

3. boxplot



`plt.boxplot()`
`sns.boxplot()`

Categorical/Discrete Dtypes

1. Frequencies

2. Bar Plots

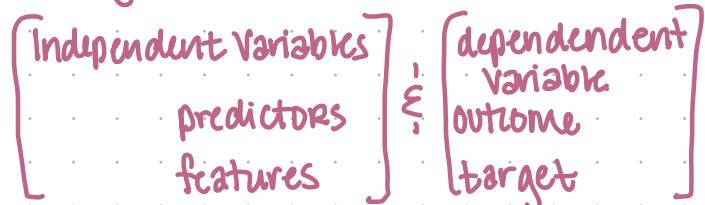
`series.value_counts()`
`sns.countplot()`

EXPLORE

Bivariate Exploration - Exploration of 2 variables at a time

WHY?

1. Identify Variable Interdependence
(You don't want that!)
2. Identify relationships between

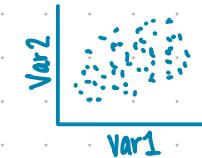


WHAT?

1. Document takeaways, findings
2. Document action needed, like
"Drop var1 due to Interdependence"
3. Document new Questions to ask the data
4. Return to prepare as needed.
5. Progress to Multivariate Exploration.

HOW?

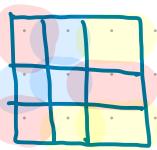
Scatterplot



Numeric x Numeric Dtypes

`sns.scatterplot()`

heatmap



`sns.heatmap()`

`sns.lineplot()`

`sns.lmplot()`

`sns.pairplot()`

Statistical tests → Pearson's R (linear corr.) & Spearman's (^{non-linear})

Discrete x Numeric Dtypes

Swarmplot



`sns.swarmplot()`

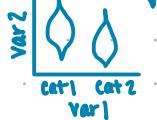
`sns.banplot()`

`sns.stripplot()`

`sns.violinplot()`

`sns.boxenplot()`

Violinplot



Statistical tests → t-test, ANOVA, Mann-Whitney U

Discrete x Discrete

Swarmplot



`sns.swarmplot()`

Crosstab

`pd.crosstab()`

Statistical test → χ^2 (chi-Square)

EXPLORE

Multivariate Exploration - Explore 3+ variables at a time

WHY?

Identify relationships between

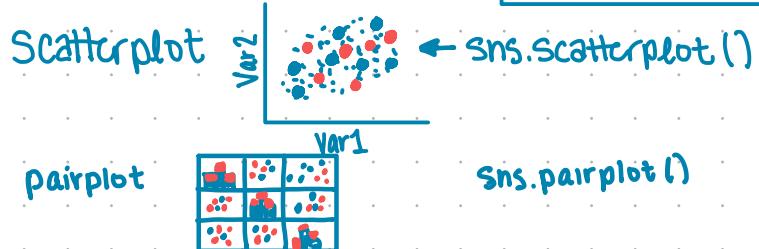


WHAT?

1. Document takeaways, findings, conclusions
2. Document initial recommendations
3. Document new Questions to ask the data
4. finalize features to move forward into modeling.

HOW?

Numeric x Numeric Dtypes



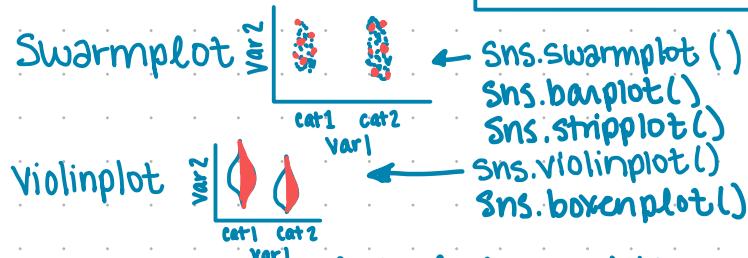
`sns.scatterplot()`

`sns.pairplot()`

Statistical tests → control for 1 group & run.

Pearson's R (linear corr.) & Spearman's (^{non-linear})

Discrete x Numeric Dtypes



`sns.swarmplot()`

`sns.banplot()`

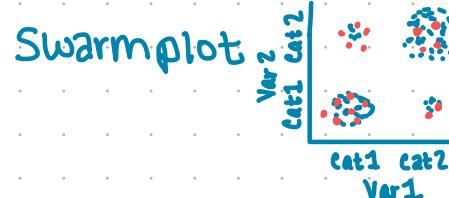
`sns.stripplot()`

`sns.violinplot()`

`sns.boxenplot()`

Control for 1 cat. variable, test the other.

Statistical tests → t-test, ANOVA, Mann-Whitney U



Discrete x Discrete

`sns.swarmplot()`

Statistical test → control for 1 categorical variable & test within the other 2

χ^2 (chi-Square)