

# Lesson #06

## Decision Trees



- Introduction to Decision Tree
- Converting categorical variables
- Splitting Data
- Decision Trees as flows of data
- Entropy & Gini
- Information gain
- Applying Decision Trees
- Overfitting problem
- Case study: classification problem

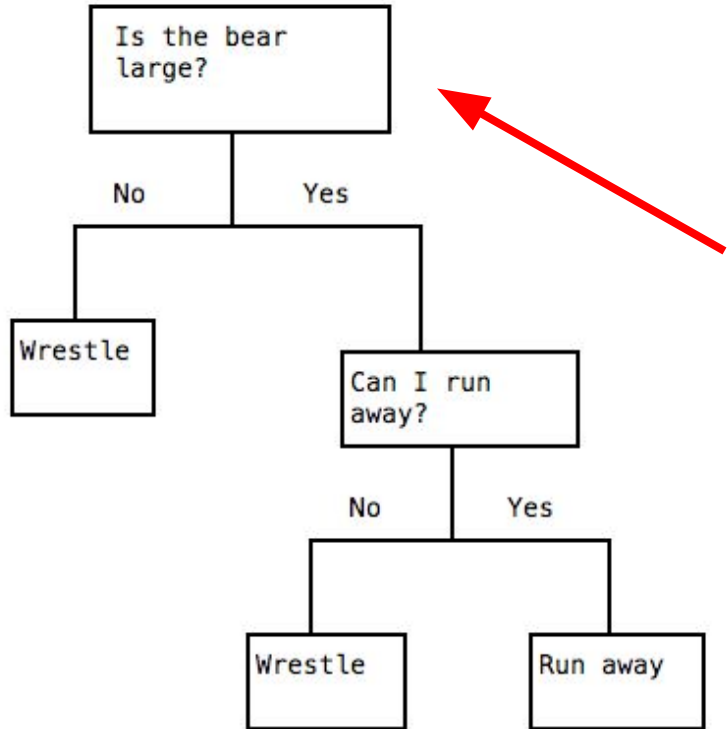
How can an  
algorithm be  
represented as a  
tree?





# Decision Tree (classification)

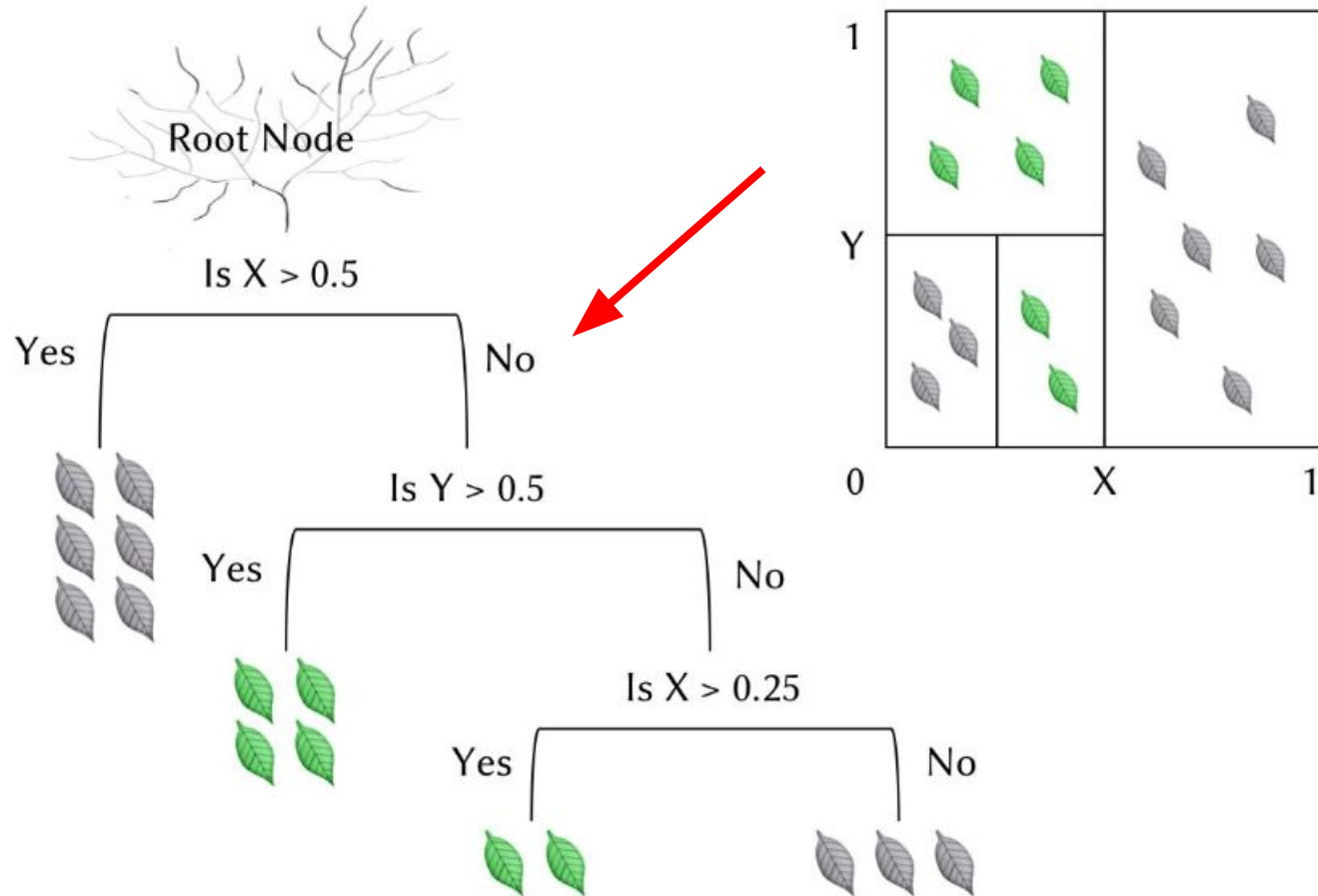
Should I wrestle this bear?



Bear name	Size	Escape possible?	Action
Yogi	Small	No	Wrestle
Winnie	Small	Yes	Wrestle
Baloo	Large	Yes	Run away
Gentle Ben	Large	No	Wrestle

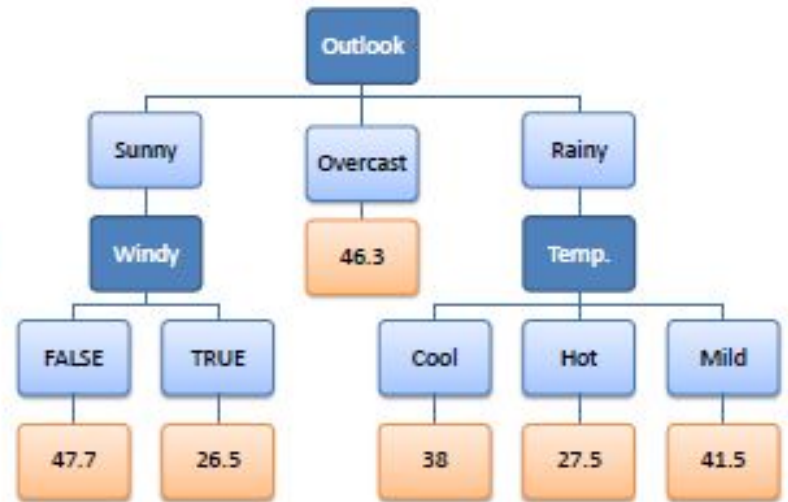
# Decision Tree (classification)

5



# Decision Tree (regression)

Predictors				Target
Outlook	Temp.	Humidity	Windy	Hours Played
Rainy	Hot	High	False	26
Rainy	Hot	High	True	30
Overcast	Hot	High	False	48
Sunny	Mild	High	False	46
Sunny	Cool	Normal	False	62
Sunny	Cool	Normal	True	23
Overcast	Cool	Normal	True	43
Rainy	Mild	High	False	36
Rainy	Cool	Normal	False	38
Sunny	Mild	Normal	False	48
Rainy	Mild	Normal	True	48
Overcast	Mild	High	True	62
Overcast	Hot	Normal	False	44
Sunny	Mild	High	True	30

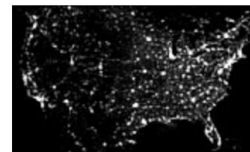




**CASE STUDY**

**Download:** [Data Folder](#), [Data Set Description](#)

**Abstract:** Predict whether income exceeds \$50K/yr based on census data. Also known as "Census Income" dataset.



Data Set Characteristics:	Multivariate	Number of Instances:	48842	Area:	Social
Attribute Characteristics:	Categorical, Integer	Number of Attributes:	14	Date Donated	1996-05-01
Associated Tasks:	Classification	Missing Values?	Yes	Number of Web Hits:	1564397

**Source:**

Donor:

Ronny Kohavi and Barry Becker  
Data Mining and Visualization  
Silicon Graphics.  
e-mail: ronnyk '@' live.com for questions.

### Data Set Information:

Extraction was done by Barry Becker from the 1994 Census database. A set of reasonably clean records was extracted using the following conditions: ((AAGE>16) && (AGI>100) && (AFNLWGT>1)&& (HRSWK>0))

Prediction task is to determine whether a person makes over 50K a year.

**high\_income**

>50K

$\leq 50K$

$\leq 50K$

>50K

$\leq 50K$



## Dataset info

Number of variables	15
Number of observations	32561
Total Missing (%)	0.0%
Total size in memory	3.7 MiB
Average record size in memory	120.0 B

Variables types

Numeric	6
Categorical	9
Boolean	0
Date	0
Text (Unique)	0
Rejected	0
Unsupported	0

## Warnings

- [capital\\_gain](#) has 29849 / 91.7% zeros Zeros
- [capital\\_loss](#) has 31042 / 95.3% zeros Zeros
- Dataset has 24 duplicate rows Warning

Be careful

# Converting Categorical Variables

---

```
1 income.workclass.head()
```

```
0          State-gov
1    Self-emp-not-inc
2          Private
3          Private
4          Private
```

```
1 # Convert a single column from text categories to numbers
2 col = pd.Categorical(income["workclass"])
3 income["workclass"] = col.codes
4 income.workclass.head(5)
```

```
0    7
1    6
2    4
3    4
4    4
```

```
1 col.categories[7]
```

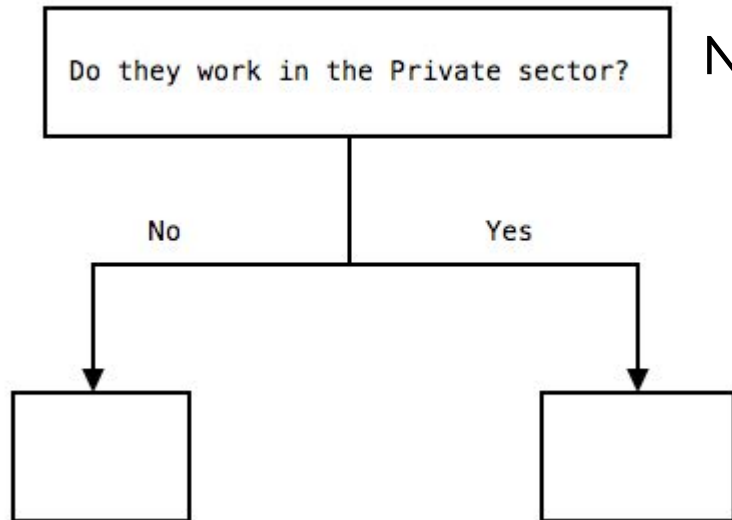
```
' State-gov'
```

```
1 col.categories
```

```
Index(['?', ' Federal-gov', ' Local-gov', ' Never-worked', ' Private',
       ' Self-emp-inc', ' Self-emp-not-inc', ' State-gov', ' Without-pay'],
      dtype='object')
```

# Splitting Data

What income do people make?

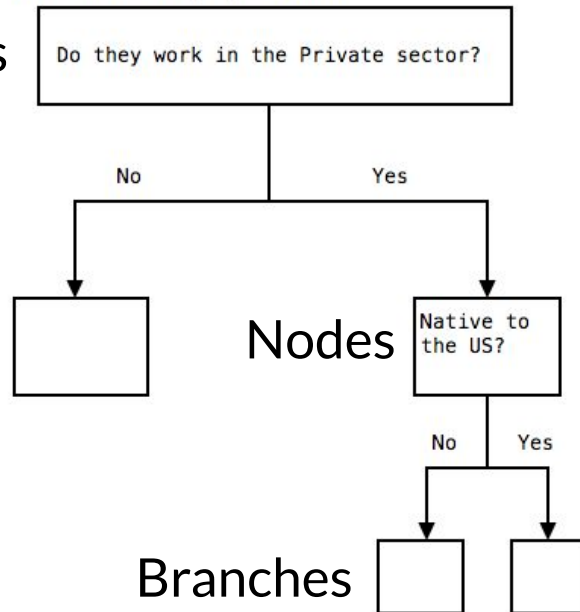


Nodes

Branches

What income do people make?

Nodes

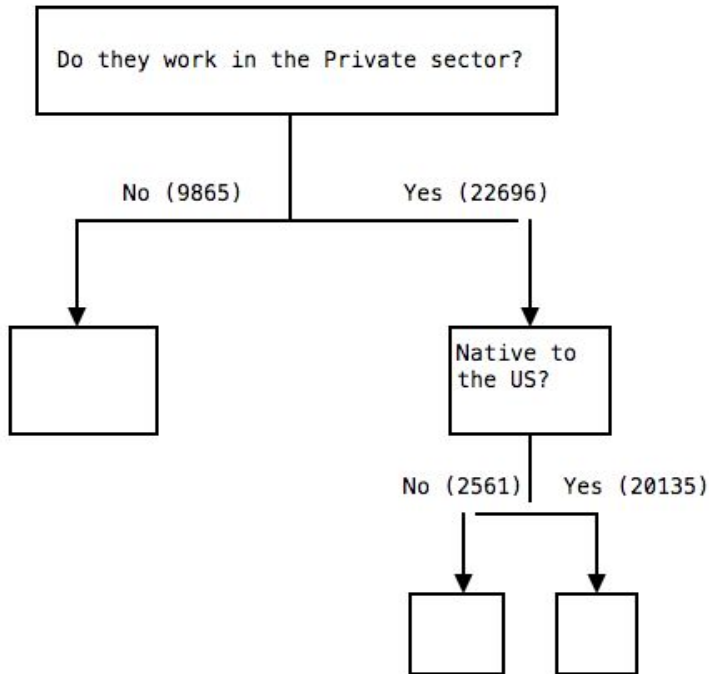


Nodes

Branches

# Decision Tree as Flows of Data

What income do people make?



We'll need to continue splitting nodes until we get to a point where all of the rows in a node have the same value for **high\_income**.

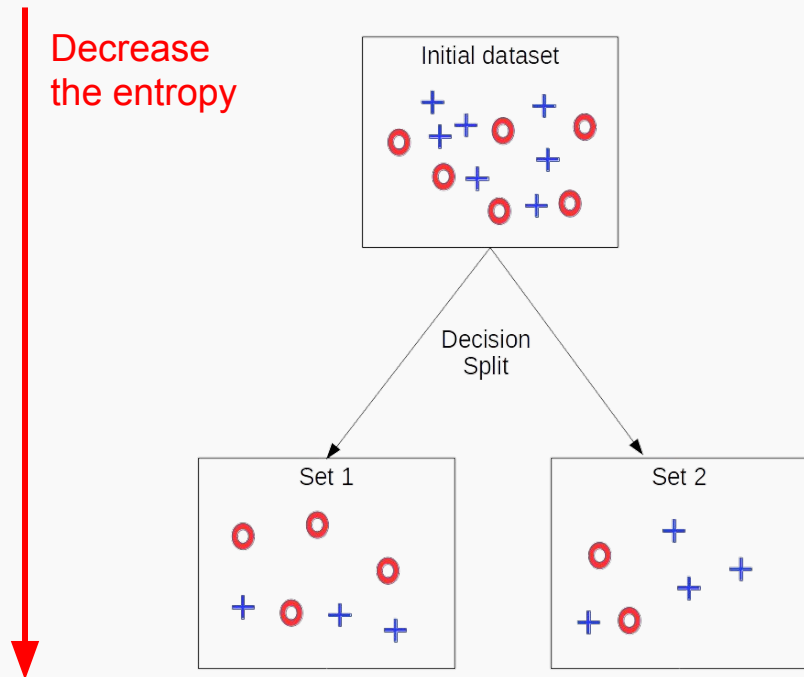




# Entropy

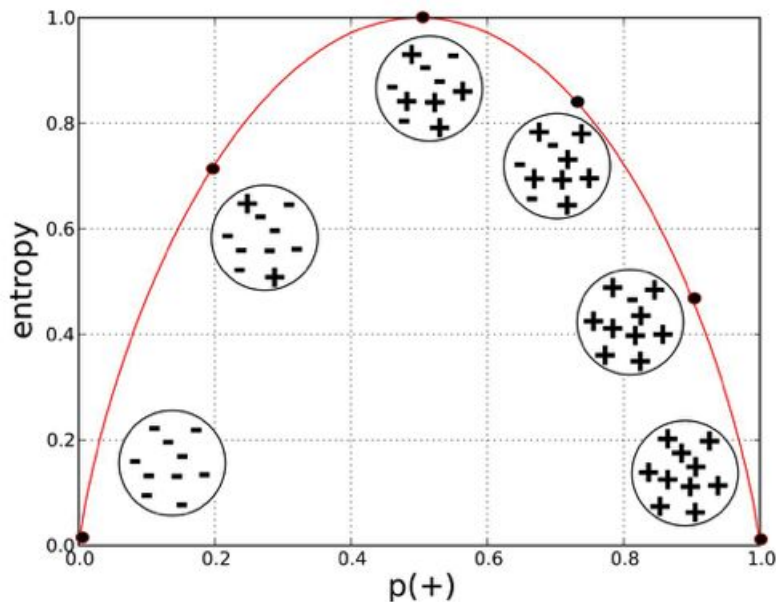
Entropy is an indicator of how messy your data is.

# Why Entropy in Decision Trees?



- The goal is to tidy the data.
- You try to separate your data and group the samples together in the classes they belong to.
- You maximize the purity of the groups as much as possible each time you create a new node of the tree
- Of course, at the end of the tree, you want to have a clear answer.

# Mathematical definition of entropy



- Suppose a set of  $N$  items, these items fall into two categories:
  - $+$ :  $n$  items
  - $-$ :  $m$  items
  - $S = n+m$
- $p = n/S, q = m/S$
- $p + q = 1$
- $E = -p \log(p) - q \log(q)$

# Generalization

---

Feature X

$$E(x) = - \sum_{i=1}^c P(x_i) \log_b P(x_i)$$

$P(x_i)$  is the fraction of  
examples in a given class  $i$

```
c = pd.AA.unique()
```



# Entropy using the frequency table of one attribute

---

high\_income

1

1

0

0

1

$$E(x) = - \sum_{i=1}^c P(x_i) \log_b P(x_i)$$

$$E(\text{high\_income}) = - \left( \frac{2}{5} \times \log_2 \frac{2}{5} \right) - \left( \frac{3}{5} \times \log_2 \frac{3}{5} \right)$$

$$E(\text{high\_income}) = 0.53 + 0.44$$

$$E(\text{high\_income}) = 0.97$$

# Entropy using the frequency table of two attributes

age	high_income	split_age
25	1	0
50	1	0
30	0	0
50	0	0
80	1	1

split\_age is based on median of age (suppose equal to 50)

$$E(T, X) = \sum_{c \in X} \frac{|X_c|}{|X|} E(T|X_c)$$

$$c_0 = \frac{4}{5} \times E([1, 1, 0, 0])$$

$$c_1 = \frac{1}{5} \times E([1])$$

$$E(T, X) = c_0 + c_1$$

$$= 0.17$$

# Information Gain

---

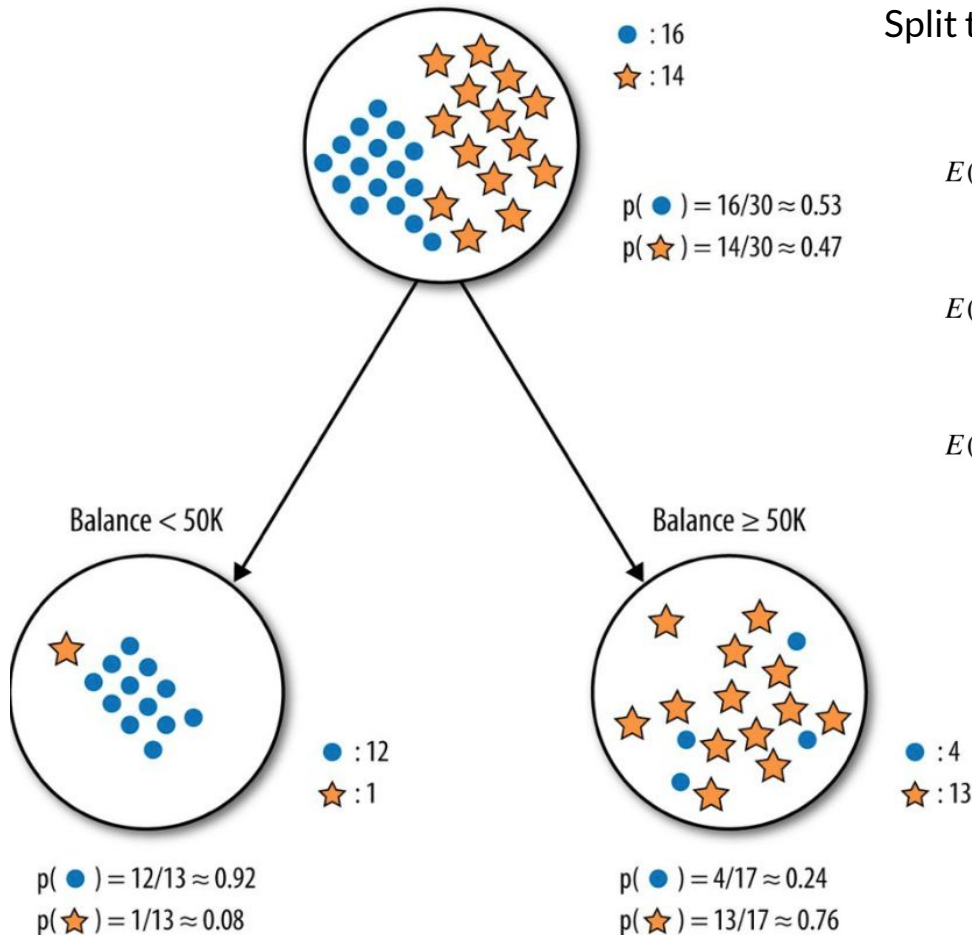
The information gain is based on the **decrease in entropy after a dataset is split** on an attribute.

$$IG(T,X) = E(T) - E(T|X)$$

Information Gain from X on T

Constructing a decision tree is all about finding attribute that returns the **highest information gain** (i.e., the most homogeneous branches).

Entire population (30 instances)



Split the Decision Tree from feature "Balance"

$$E(\text{Parent}) = -\frac{16}{30} \log_2 \left( \frac{16}{30} \right) - \frac{14}{30} \log_2 \left( \frac{14}{30} \right) \approx 0.99$$

$$E(\text{Balance} < 50K) = -\frac{12}{13} \log_2 \left( \frac{12}{13} \right) - \frac{1}{13} \log_2 \left( \frac{1}{13} \right) \approx 0.39$$

$$E(\text{Balance} \geq 50K) = -\frac{4}{17} \log_2 \left( \frac{4}{17} \right) - \frac{13}{17} \log_2 \left( \frac{13}{17} \right) \approx 0.79$$

$$E(\text{Balance}) = \frac{13}{30} \times 0.39 + \frac{17}{30} \times 0.79$$

$$= 0.62$$

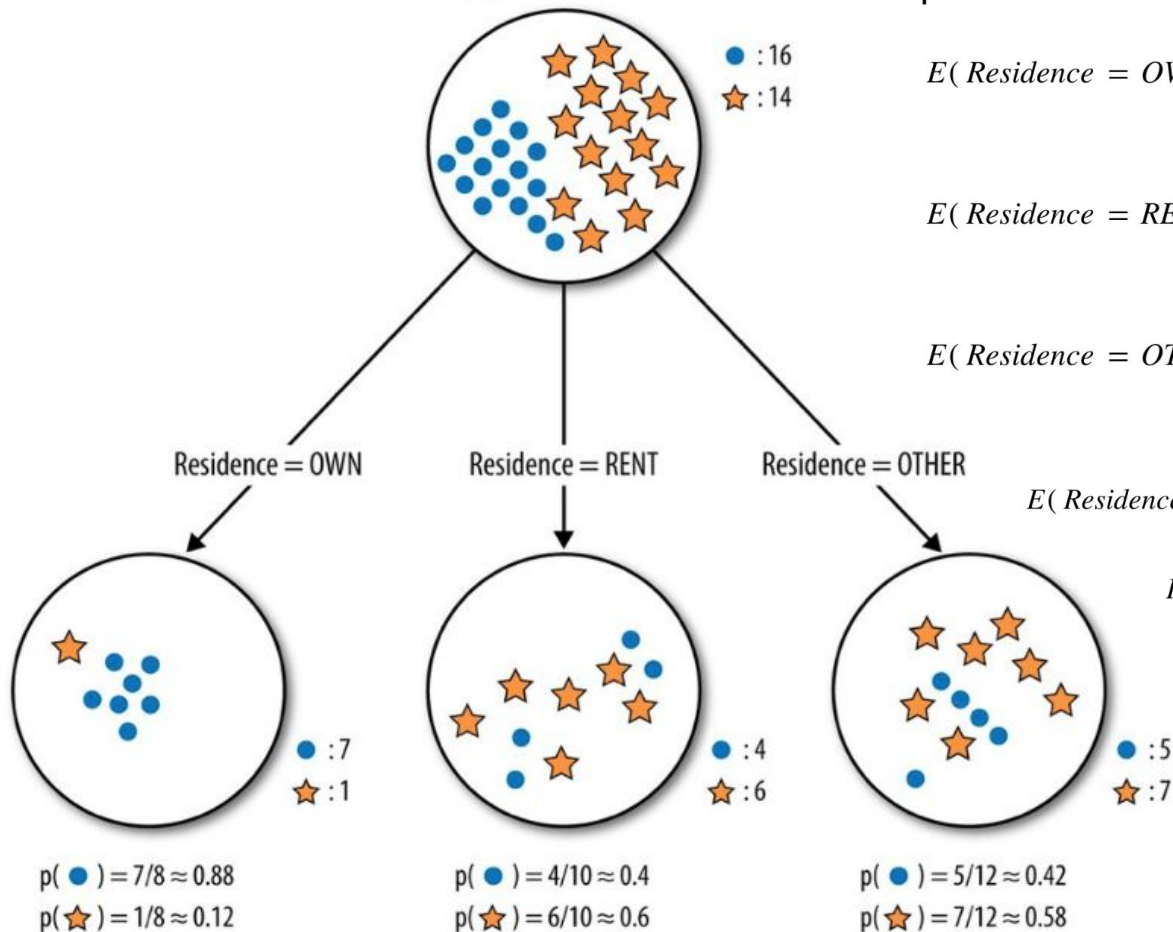
$$IG(\text{Parent}, \text{Balance}) = E(\text{Parent}) - E(\text{Balance})$$

$$= 0.99 - 0.62$$

$$= 0.37$$



Entire population (30 instances)



Split the Decision Tree from feature "Residence"

$$E(\text{Residence} = \text{OWN}) = -\frac{7}{8} \log_2 \left( \frac{7}{8} \right) - \frac{1}{8} \log_2 \left( \frac{1}{8} \right) \approx 0.54$$

$$E(\text{Residence} = \text{RENT}) = -\frac{4}{10} \log_2 \left( \frac{4}{10} \right) - \frac{6}{10} \log_2 \left( \frac{6}{10} \right) \approx 0.97$$

$$E(\text{Residence} = \text{OTHER}) = -\frac{5}{12} \log_2 \left( \frac{5}{12} \right) - \frac{7}{12} \log_2 \left( \frac{7}{12} \right) \approx 0.98$$

$$E(\text{Residence}) = \frac{8}{30} \times 0.54 + \frac{10}{30} \times 0.97 + \frac{12}{30} \times 0.98 = 0.86$$

$$\begin{aligned} IG(\text{Parent}, \text{Residence}) &= E(\text{Parent}) - E(\text{Residence}) \\ &= 0.99 - 0.86 \\ &= 0.13 \end{aligned}$$

$$Gini(x) = 1 - \sum_{i=1}^c P(x_i)^2$$

$$Entropy(x) = - \sum_{i=1}^c P(x_i) \log_b P(x_i)$$

Gini index or Entropy is the criterion for calculating **Information Gain**. Both of them are measures of impurity of a node.

# Evaluating Classifiers



Predicted Values

Actual Values

1

0



$$\text{precision} = \frac{TP}{(TP + FP)}$$

$$\text{precision} = \frac{TN}{(TN + FN)}$$

## Evaluating Binary Classifiers

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FN + FP}$$

$$\text{TPR} = \frac{\# \text{true positives}}{\# \text{true positives} + \# \text{false negatives}}$$

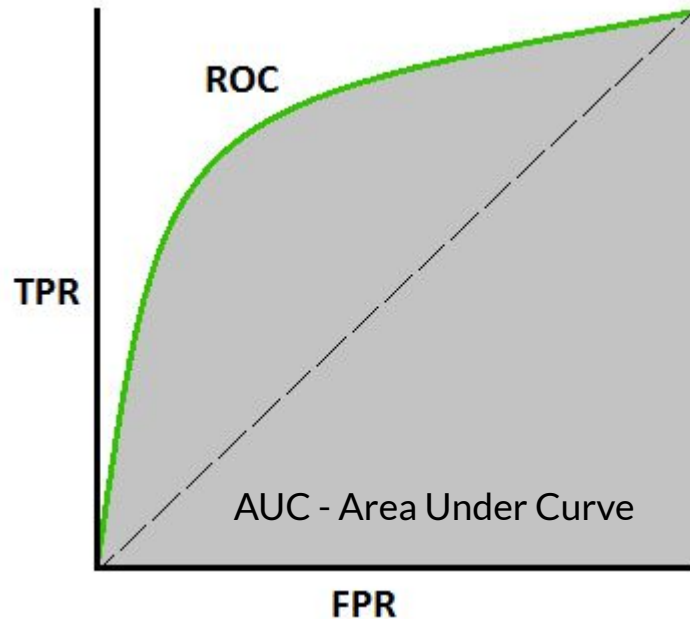
Recall

$$\text{TNR} = \frac{\# \text{true negatives}}{\# \text{true negatives} + \# \text{false positives}}$$

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$



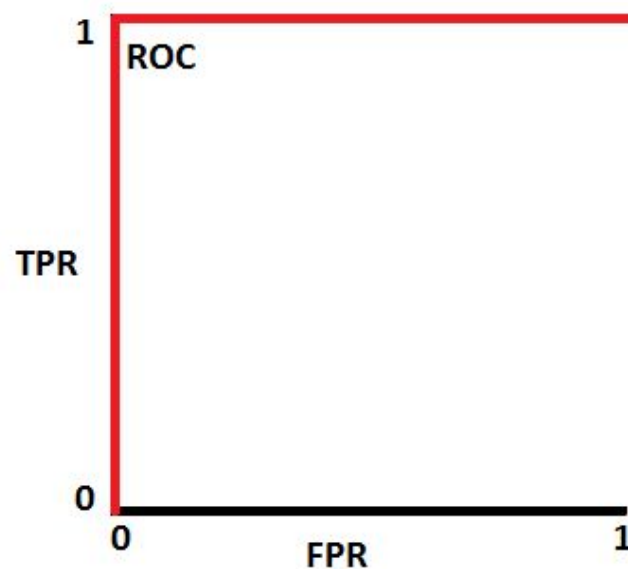
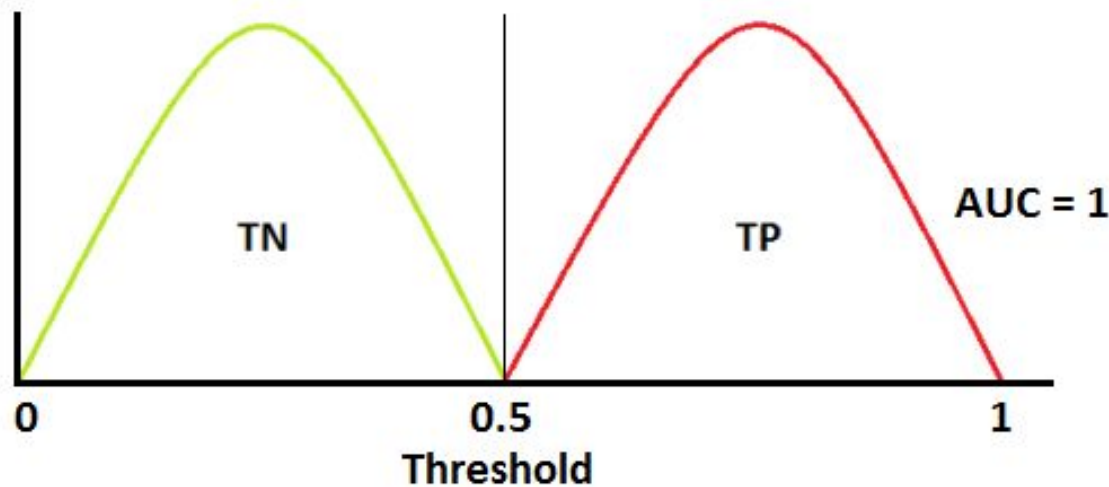
# Receiver Operating Characteristic (ROC)



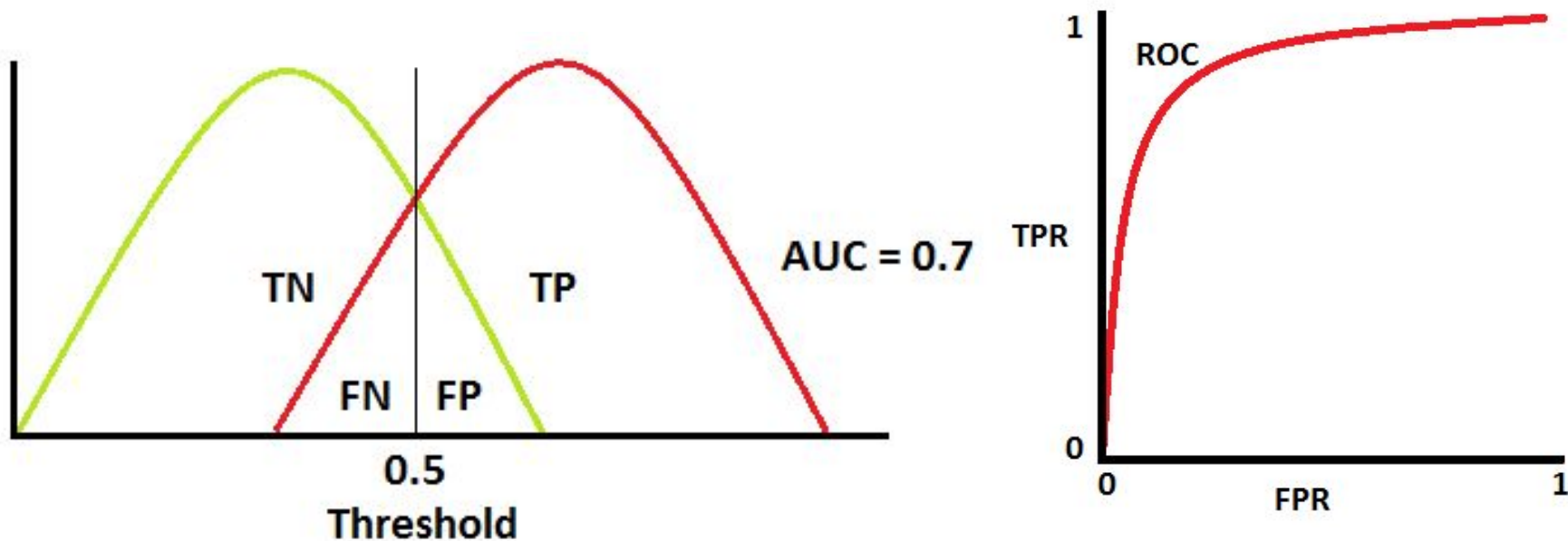
TPR  
(measure the impact of  
True Positive)

$FPR = 1 - TNR$   
(measure the impact of  
False Positive)

# How to speculate the performance of model?

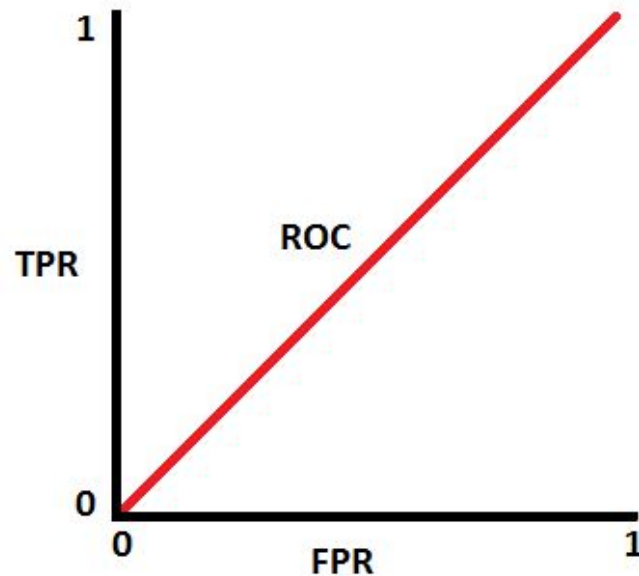
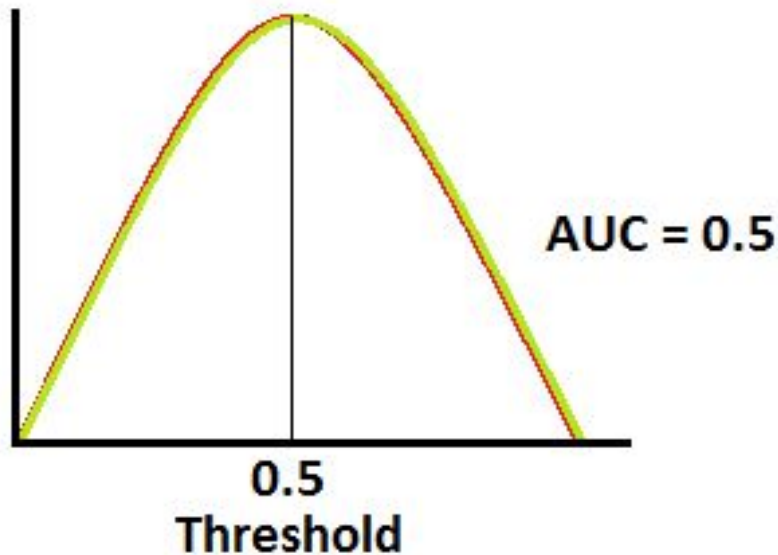


# How to speculate the performance of model?

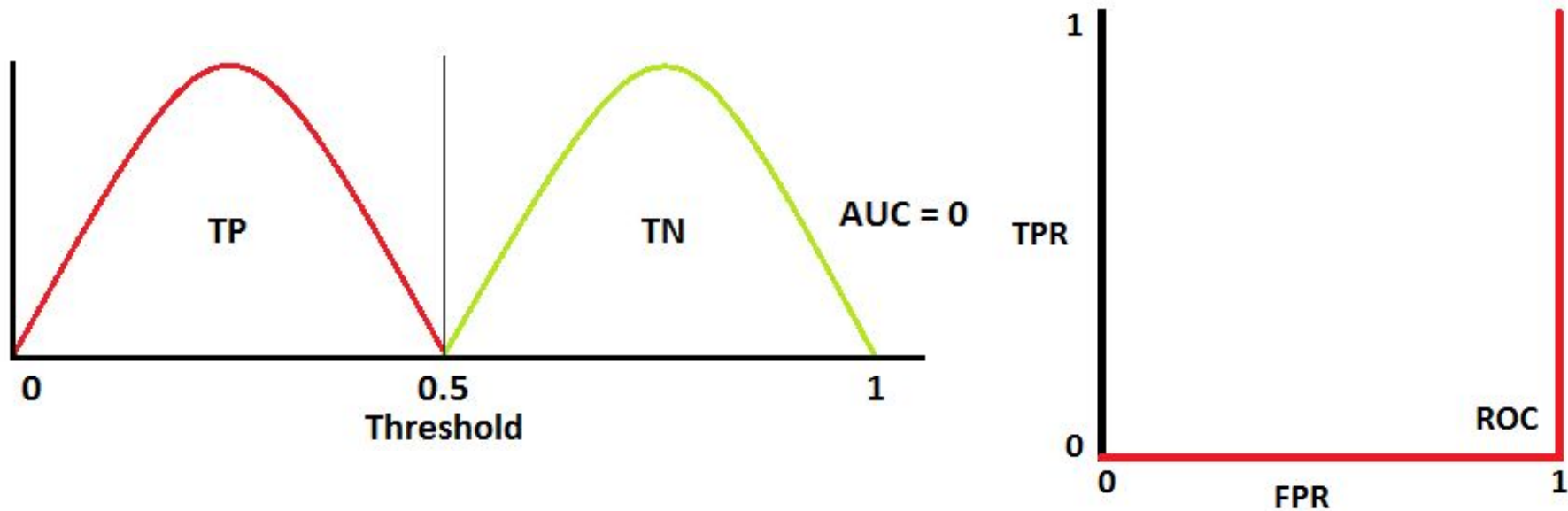


# How to speculate the performance of model?

---



# How to speculate the performance of model?





# Applying Decision Tree

---

## `sklearn.tree`.**DecisionTreeClassifier**

```
class sklearn.tree. DecisionTreeClassifier (criterion='gini', splitter='best', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None,  
max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, class_weight=None, presort=False)
```

[\[source\]](#)

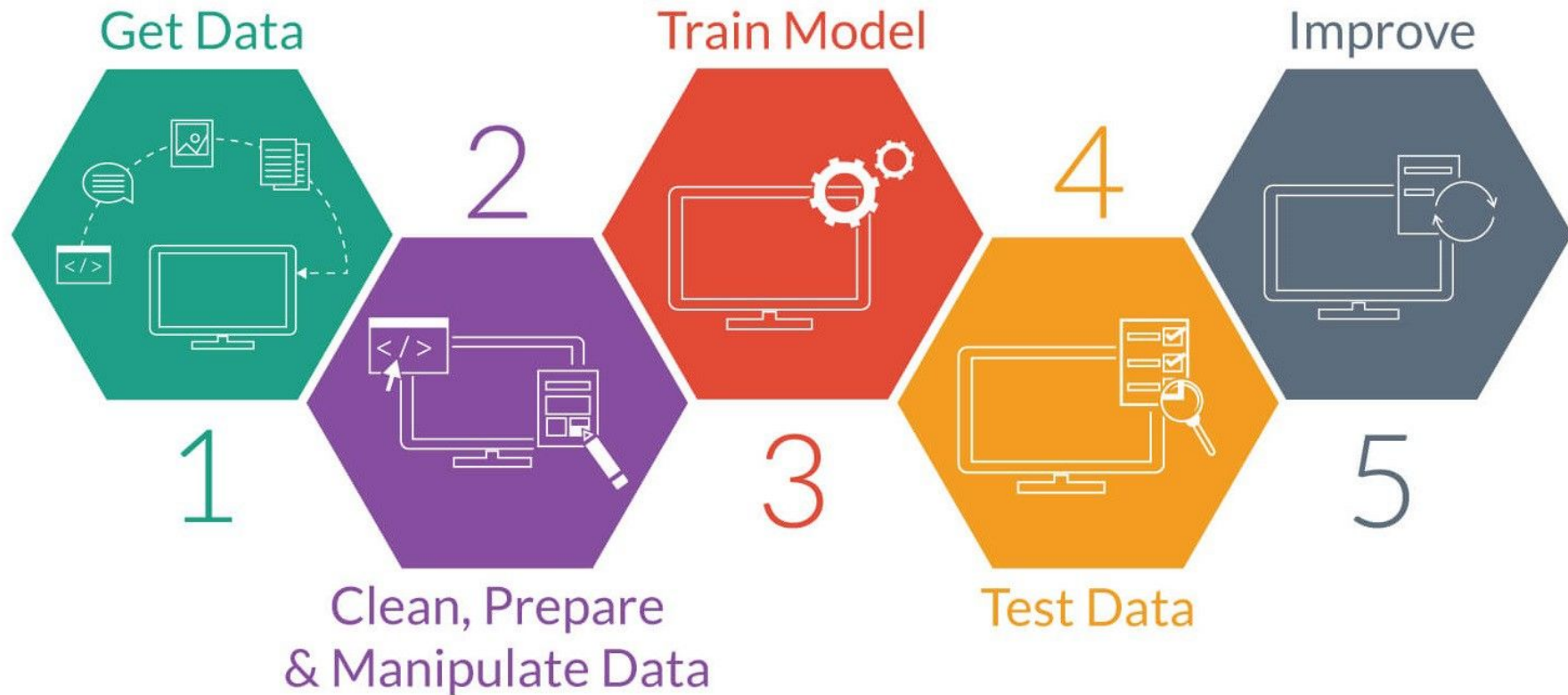
## `sklearn.tree`.**DecisionTreeRegressor**

```
class sklearn.tree. DecisionTreeRegressor (criterion='mse', splitter='best', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None,  
max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, presort=False) ¶
```

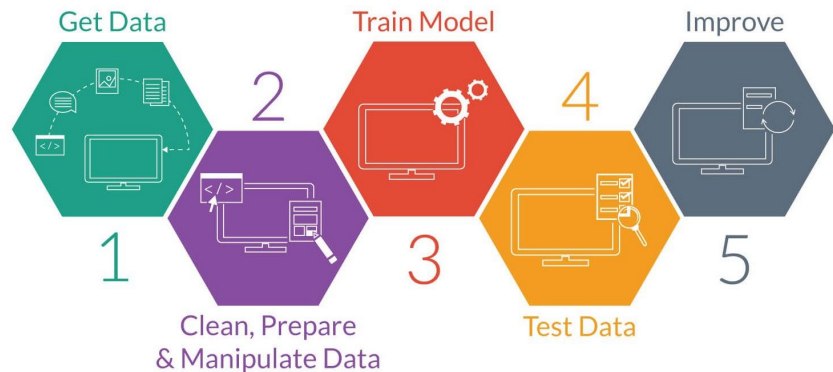
[\[source\]](#)



# A general ML workflow



# A general ML workflow

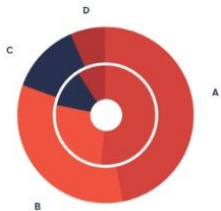


1. Load Libraries
2. Get data, including EDA
3. Clean, prepare and manipulate data (feature engineering)
4. Modeling (train and test)
5. Algorithm Tuning
6. Finalizing the Model

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn.model_selection import train_test_split
6 from sklearn.tree import DecisionTreeClassifier
7 from sklearn.neighbors import KNeighborsClassifier
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.model_selection import KFold
10 from sklearn.model_selection import cross_val_score
11 from sklearn.pipeline import Pipeline
12 from sklearn.model_selection import GridSearchCV
13 from sklearn.metrics import accuracy_score
14 from sklearn.metrics import classification_report
15 from sklearn.metrics import confusion_matrix
16 import pydotplus
17 from IPython.display import Image
18 from sklearn import tree
```

# Load Libraries

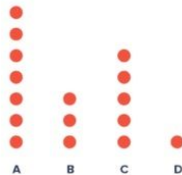
Multi-level Donut Chart



Angular Gauge



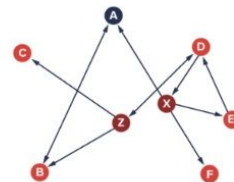
Dot Plot



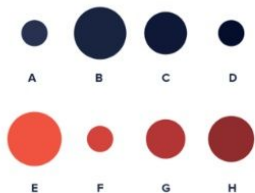
Pie Chart



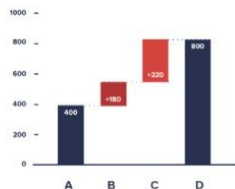
Sociogram



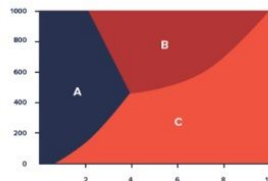
Proportional Area Chart (Circle)



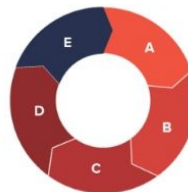
Waterfall Chart



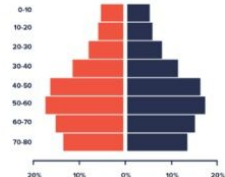
Phase Diagram



Cycle Diagram



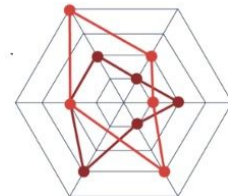
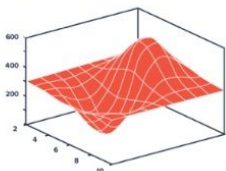
Population Pyramid



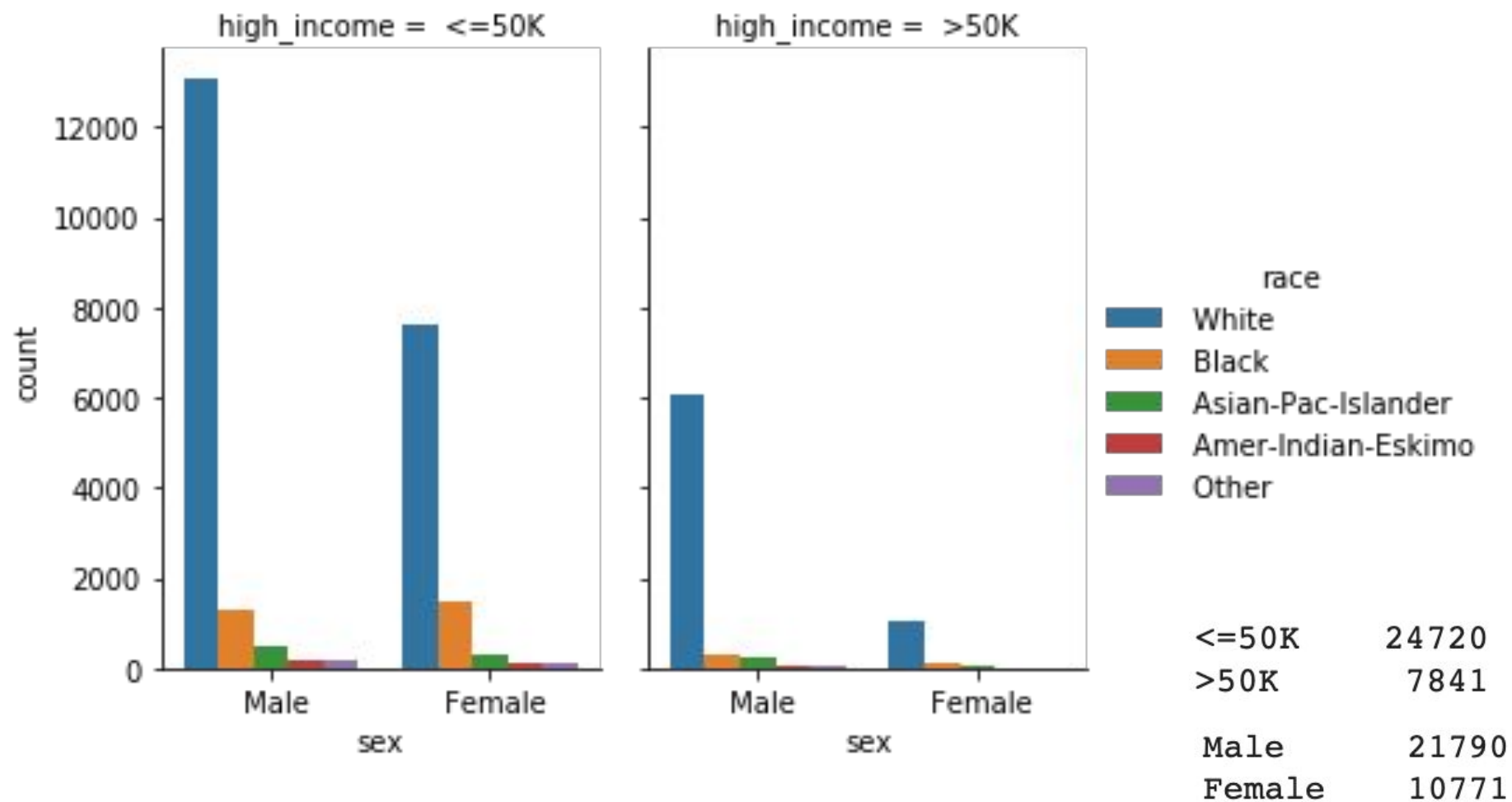
Boxplot



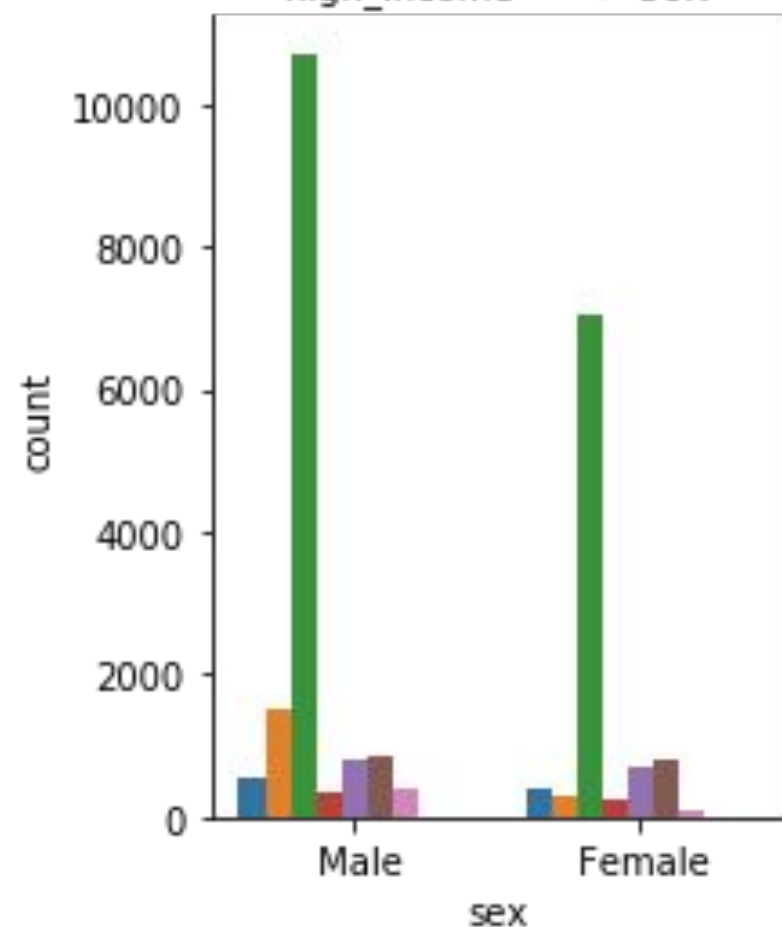
Graph



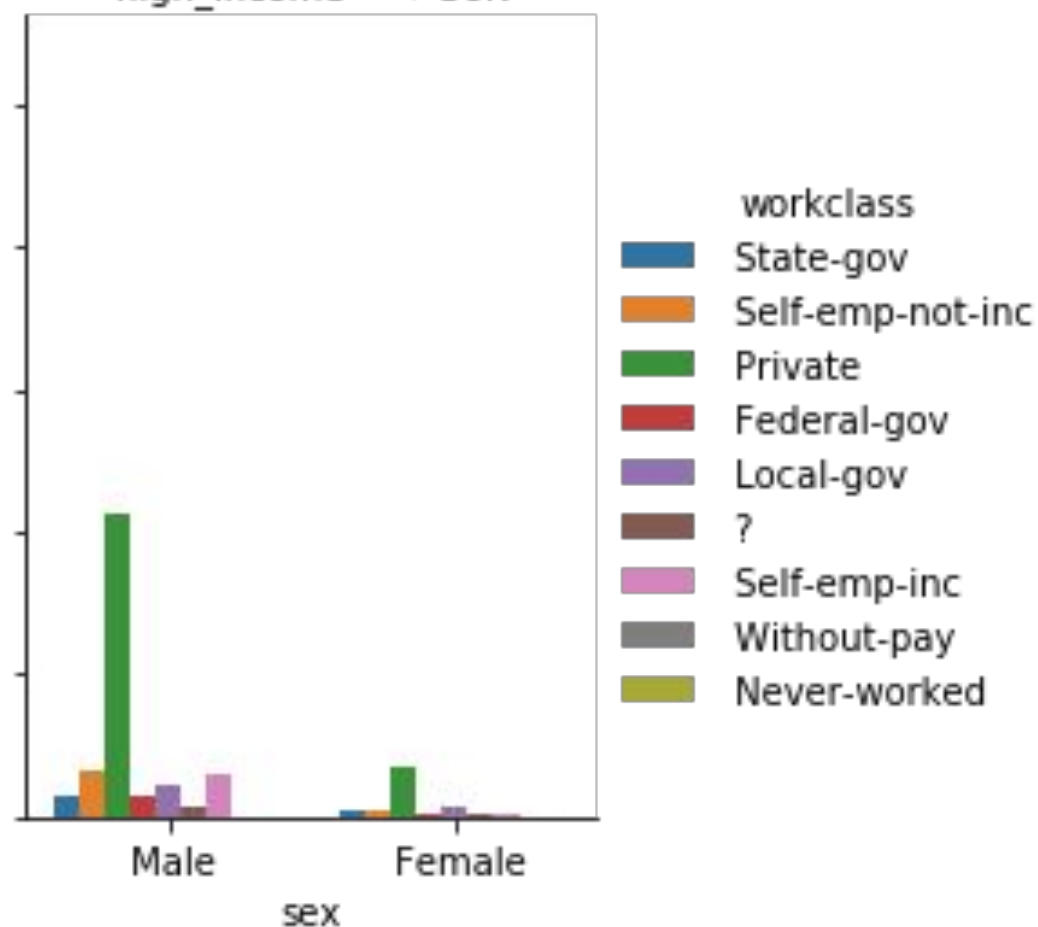
# EXPLORATORY DATA ANALYSIS (EDA)



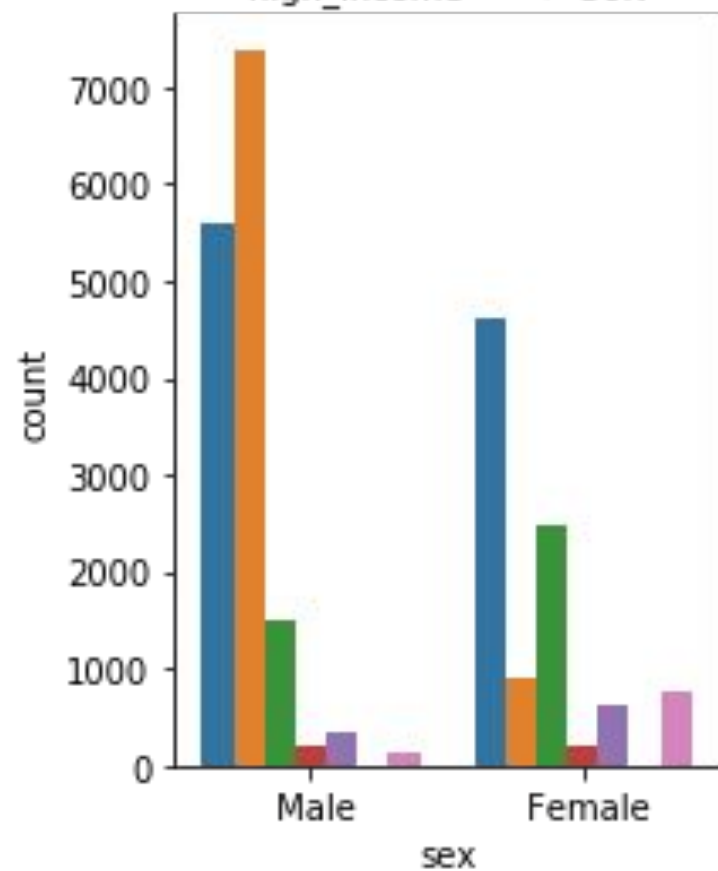
high\_income = &lt;=50K



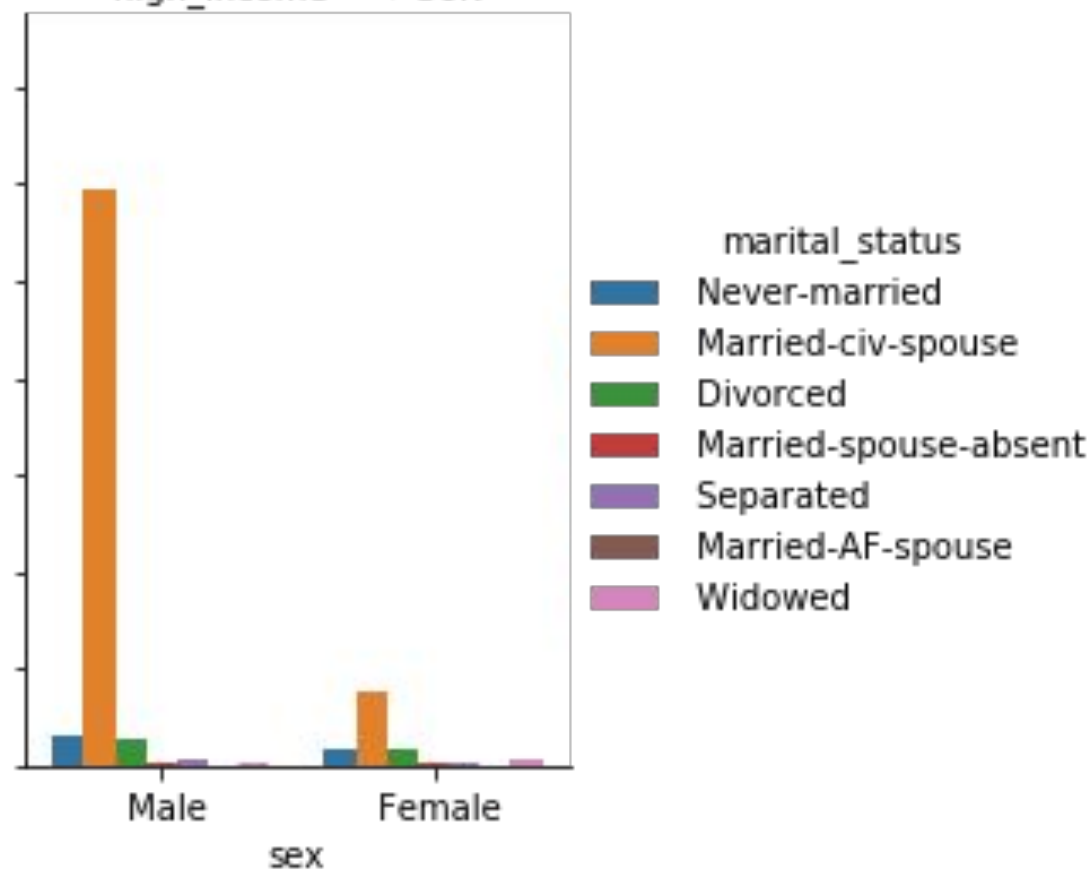
high\_income = &gt;50K



high\_income = &lt;=50K

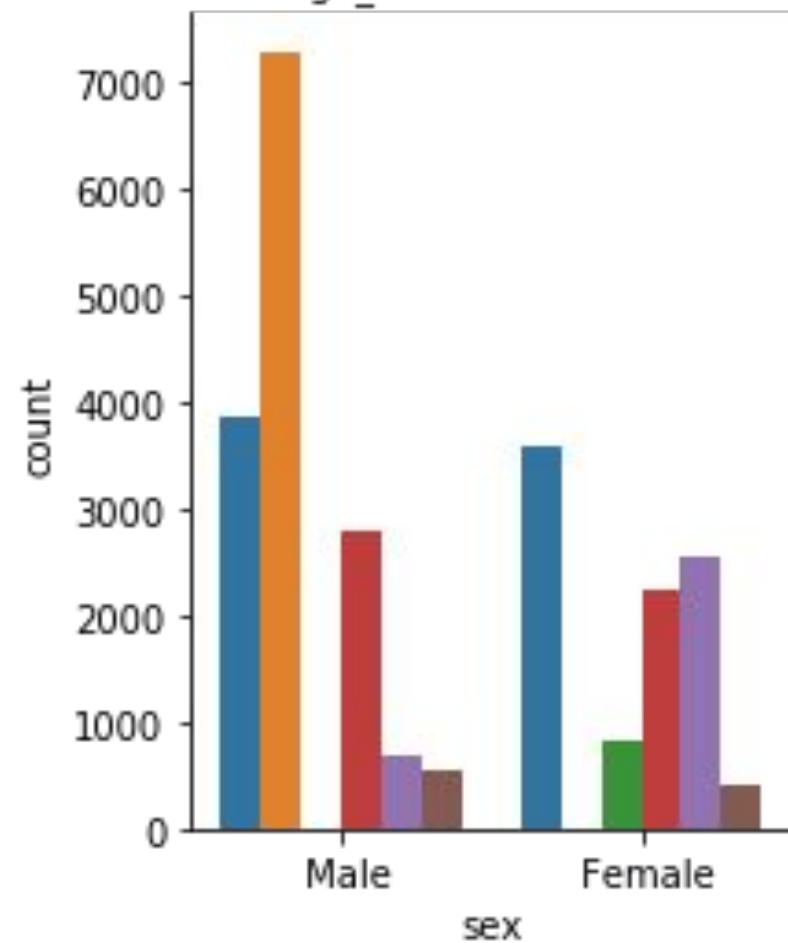


high\_income = &gt;50K

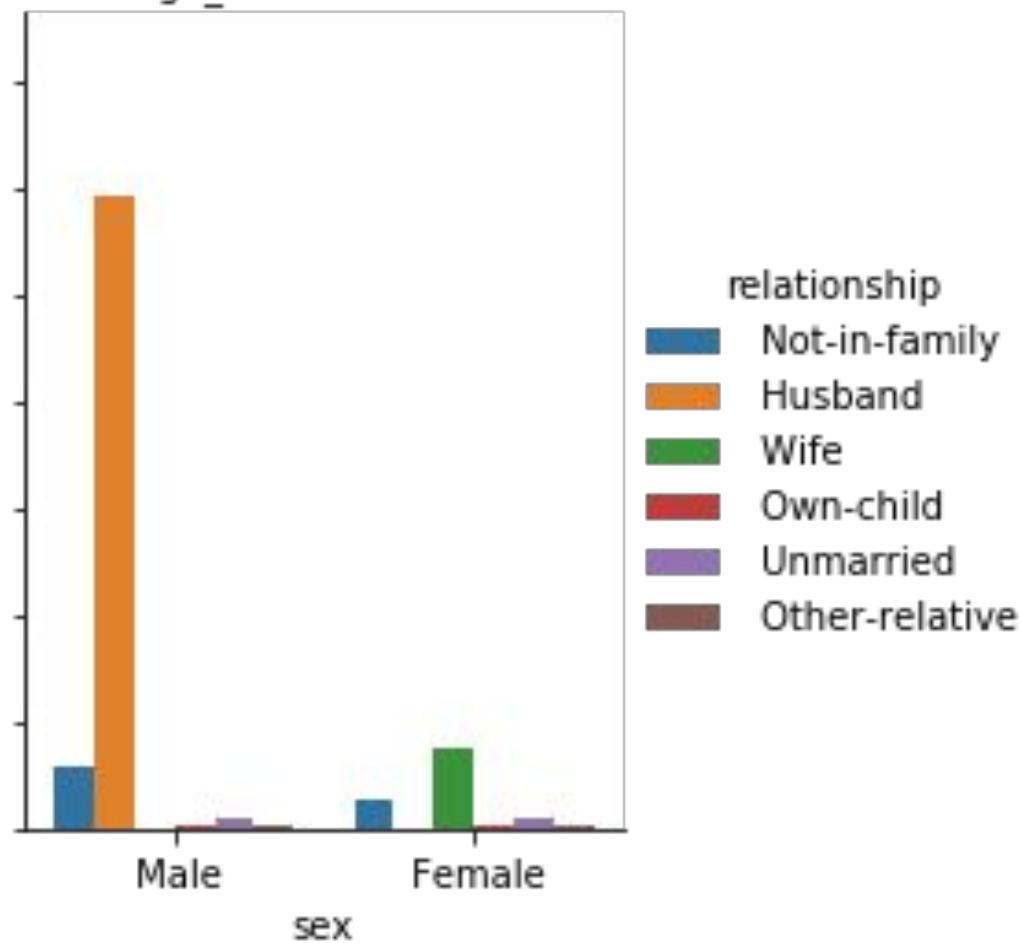


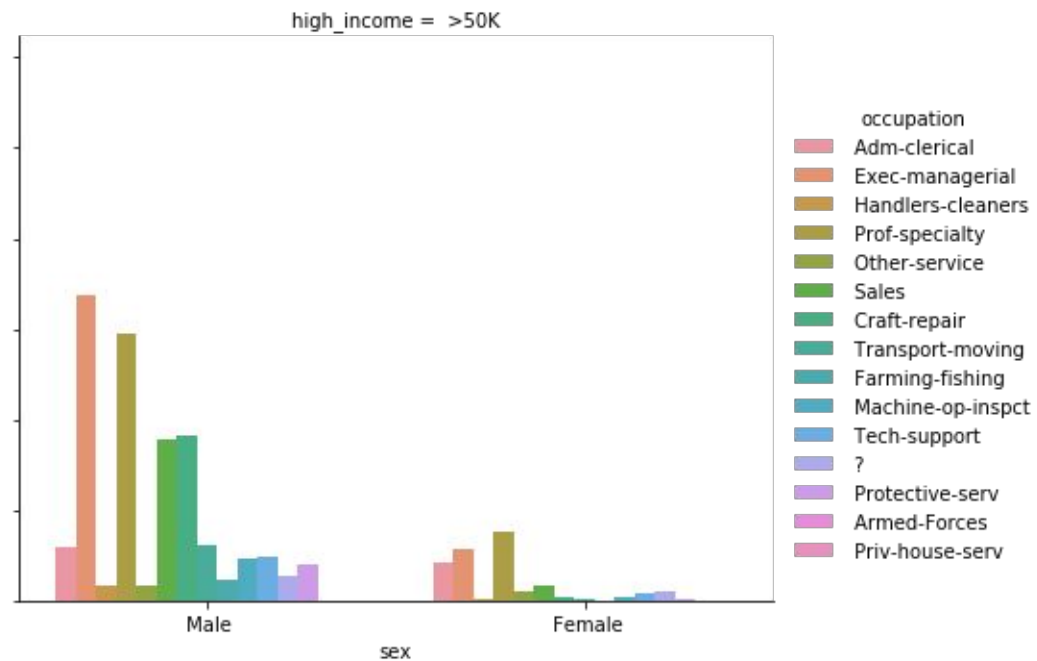
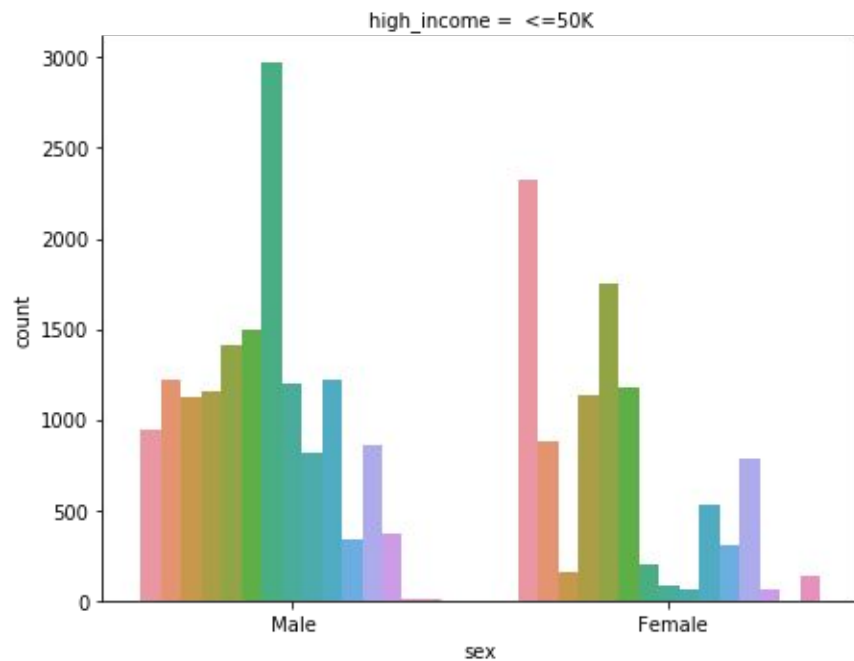


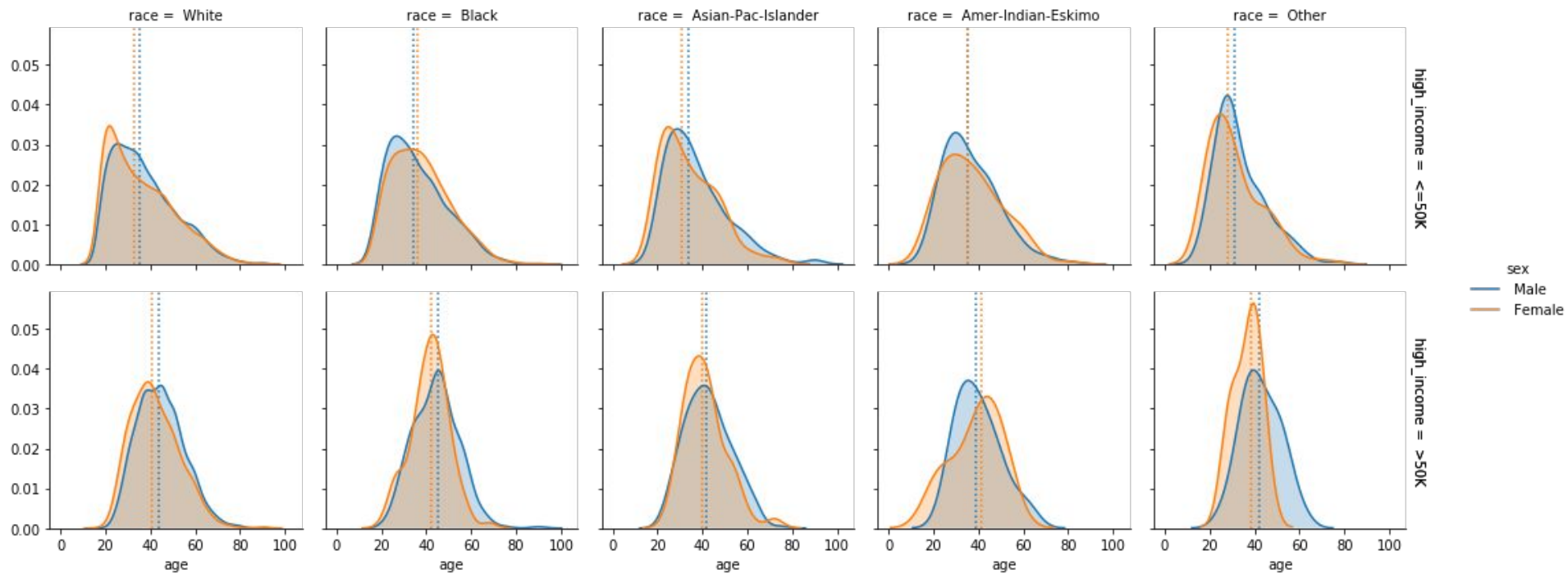
high\_income = &lt;=50K

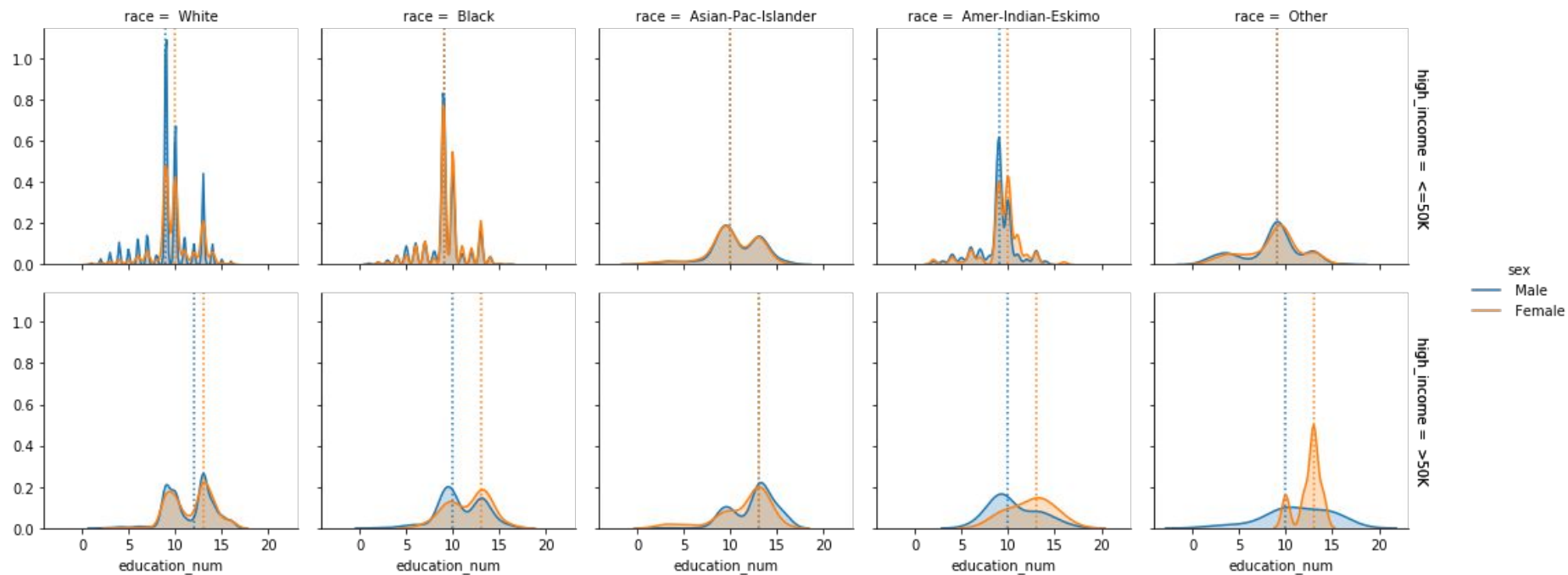


high\_income = &gt;50K











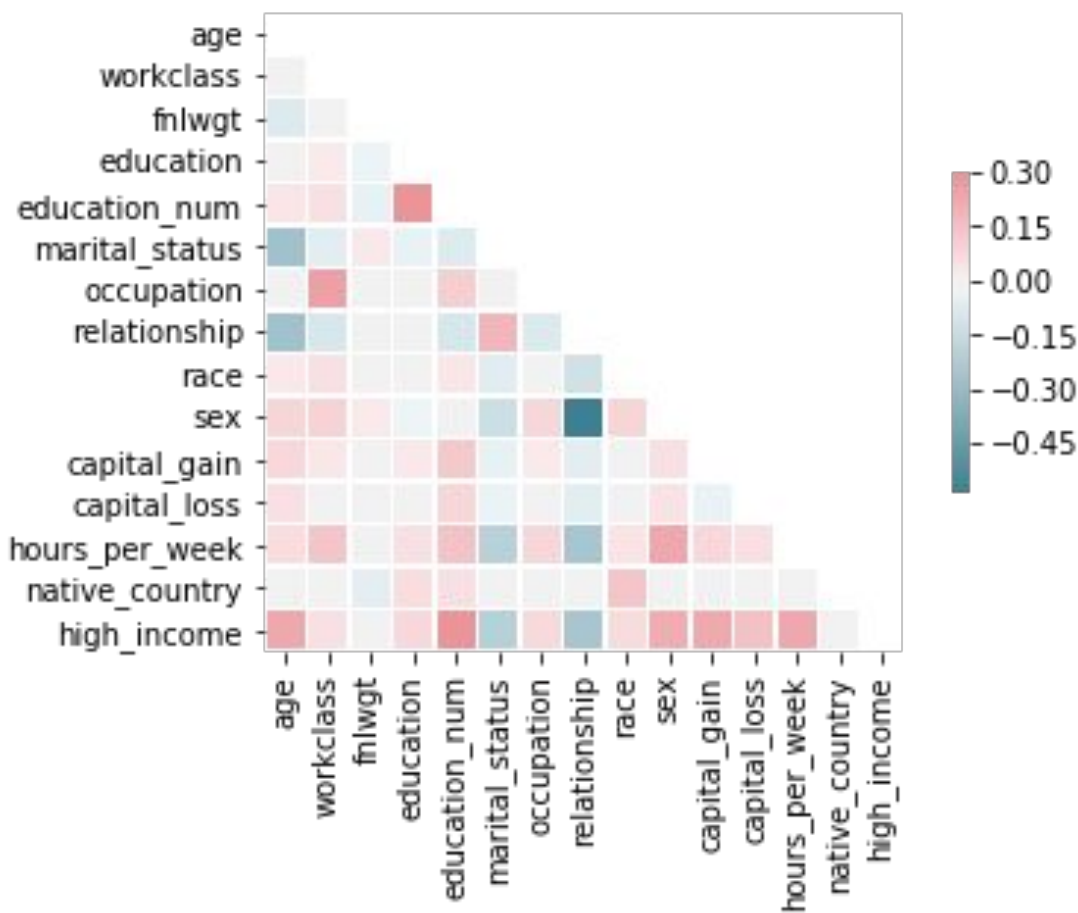
CLEAN, PREPARE AND  
MANIPULATE DATA

	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationship	race	sex	high_income
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	<=50K
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	<=50K
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	<=50K
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	<=50K
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	<=50K



```
for name in income.select_dtypes("object").columns.to_list():
    col = pd.Categorical(income[name])
    income[name] = col.codes
```

	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationship	race	sex	high_income
0	39	7	77516	9	13	4	1	1	4	1	0
1	50	6	83311	9	13	2	4	0	4	1	0
2	38	4	215646	11	9	0	6	1	4	1	0
3	53	4	234721	1	7	2	6	0	2	1	0
4	28	4	338409	9	13	2	10	5	2	0	0



high_income	1.000000
education_num	0.335154
age	0.234037
hours_per_week	0.229689
capital_gain	0.223329
sex	0.215980
capital_loss	0.150526
education	0.079317
occupation	0.075468
race	0.071846
workclass	0.051604
native_country	0.015840
fnlwgt	-0.009463
marital_status	-0.199307
relationship	-0.250918



# Train and Test

[illegible]

# Algorithm Tuning

```
1 # create a pipeline
2 pipe = Pipeline([("classifier",DecisionTreeClassifier())])
3
4 # create a dictionary with the hyperparameters
5 search_space = [{"classifier":[DecisionTreeClassifier()],
6                     "classifier__criterion": ["gini","entropy"]},
7                 {"classifier":[LogisticRegression()],
8                     "classifier__solver": ["liblinear"]},
9                 {"classifier": [KNeighborsClassifier()],
10                    "classifier__n_neighbors": [5,7]}]
11
12 # create grid search
13 kfold = KFold(n_splits=num_folds,random_state=seed)
14 grid = GridSearchCV(estimator=pipe,
15                     param_grid=search_space,
16                     cv=kfold,
17                     scoring=scoring,
18                     n_jobs=-1)
19
20 # fit grid search
21 best_model = grid.fit(X_train,y_train)
```

0.807049 (0.004056) with: {'classifier': DecisionTreeClassifier(class\_weight=None, criterion='entropy', max\_depth=None, max\_features=None, max\_leaf\_nodes=None, min\_impurity\_decrease=0.0, min\_impurity\_split=None, min\_samples\_leaf=1, min\_samples\_split=2, min\_weight\_fraction\_leaf=0.0, presort=False, random\_state=None, splitter='best'), 'classifier\_criterion': 'gini'}

0.814727 (0.008882) with: {'classifier': DecisionTreeClassifier(class\_weight=None, criterion='entropy', max\_depth=None, max\_features=None, max\_leaf\_nodes=None, min\_impurity\_decrease=0.0, min\_impurity\_split=None, min\_samples\_leaf=1, min\_samples\_split=2, min\_weight\_fraction\_leaf=0.0, presort=False, random\_state=None, splitter='best'), 'classifier\_criterion': 'entropy'}

~~0.794725 (0.012523) with: {'classifier': LogisticRegression(C=1.0, class\_weight=None, dual=False, fit\_intercept=True, intercept\_scaling=1, l1\_ratio=None, max\_iter=100, multi\_class='warn', n\_jobs=None, penalty='l2', random\_state=None, solver='warn', tol=0.0001, verbose=0, warm\_start=False), 'classifier\_solver': 'liblinear'}~~

0.774109 (0.008792) with: {'classifier': KNeighborsClassifier(algorithm='auto', leaf\_size=30, metric='minkowski', metric\_params=None, n\_jobs=None, n\_neighbors=5, p=2, weights='uniform'), 'classifier\_n\_neighbors': 5}

0.783745 (0.009897) with: {'classifier': KNeighborsClassifier(algorithm='auto', leaf\_size=30, metric='minkowski', metric\_params=None, n\_jobs=None, n\_neighbors=5, p=2, weights='uniform'), 'classifier\_n\_neighbors': 7}

# Finalize the Model

```
1 # final model
2 predict = best_model.predict(X_test)
3 print(accuracy_score(y_test, predict))
4 print(confusion_matrix(y_test, predict))
5 print(classification_report(y_test, predict))
```

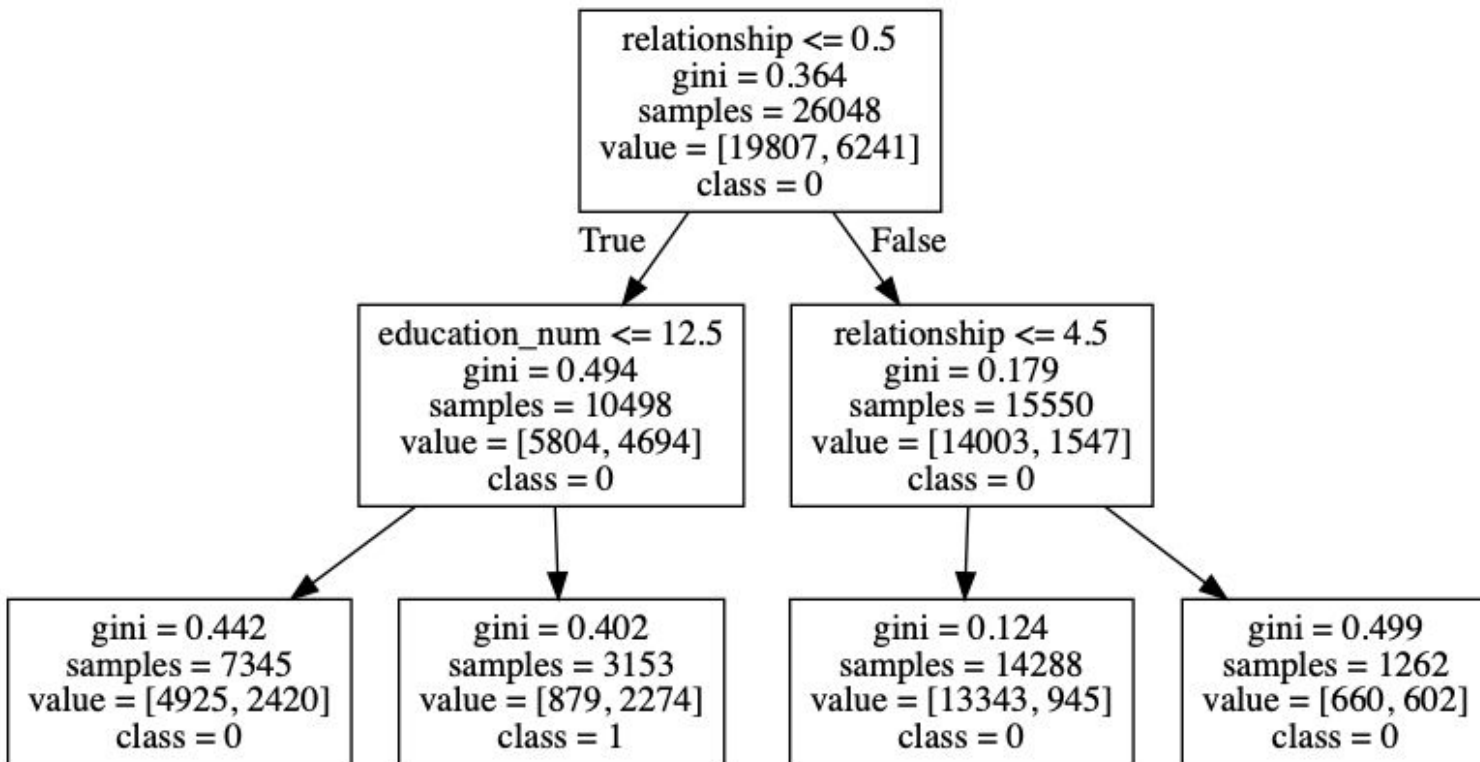
```
0.8191309688315677
```

```
[[4342  603]
```

```
 [ 575  993]]
```

	precision	recall	f1-score	support
0	0.88	0.88	0.88	4945
1	0.62	0.63	0.63	1568
accuracy			0.82	6513
macro avg	0.75	0.76	0.75	6513
weighted avg	0.82	0.82	0.82	6513

# Visualizing a Decision Tree Model



# Knowing when to use decision trees

---

The main advantages of using decision trees is that they're:

- Easy to interpret
- Relatively fast to fit and make predictions
- Able to handle multiple types of data
- Able to pick up nonlinearities in data, and usually fairly accurate

The main disadvantage of using decision trees is their **tendency to overfit**.

# Lesson #06 - Decision Trees.ipynb

