



# EEC1509 - Machine Learning

## Lesson #10 Random Forest

Ivanovitch Silva  
November, 2018



# Update repository

---

```
git clone https://github.com/ivanovitchm/EEC1509_MachineLearning.git
```

Ou ....

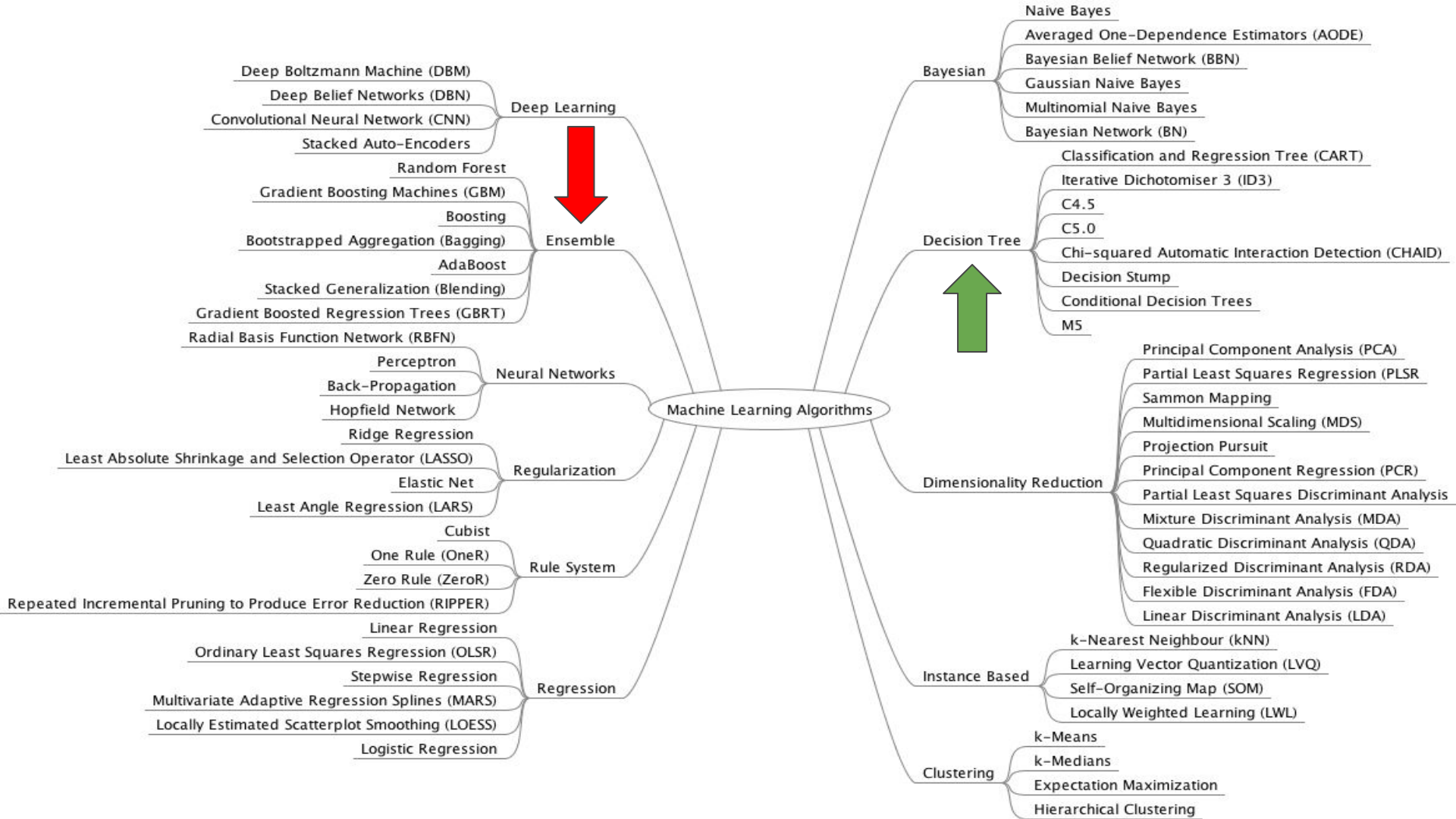
```
git pull
```



# Agenda

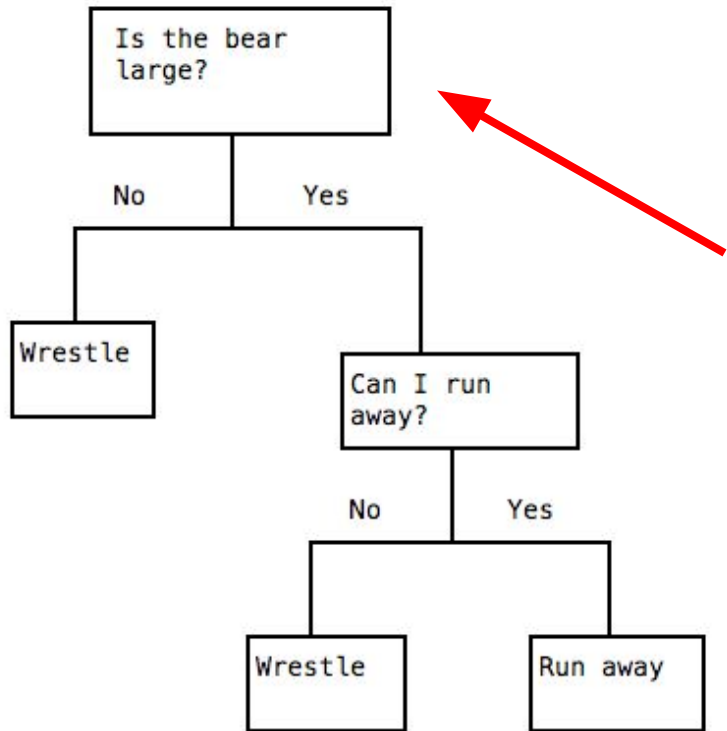
---

- Previously on last class (Decision Trees)
- Ensembles (Random Forest)
- Combining predictions
- Why Ensembling works
- Introduction variation with bagging and random features
- Reducing overfitting using Random Forest



# Decision Tree (classification)

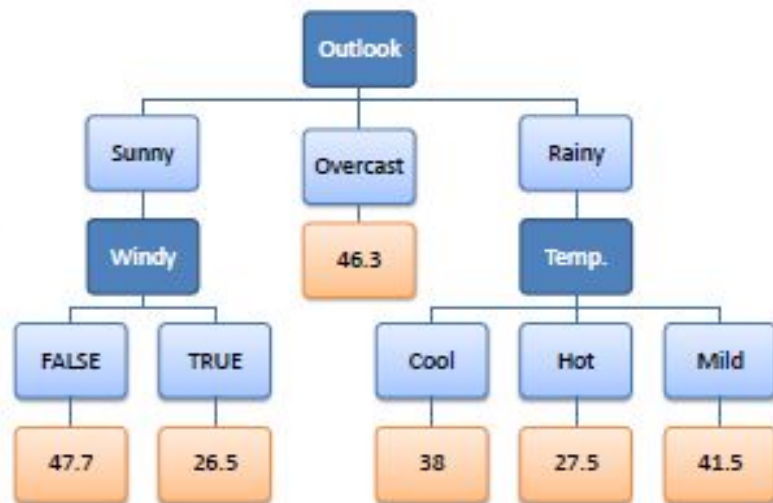
Should I wrestle this bear?



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 (regression)

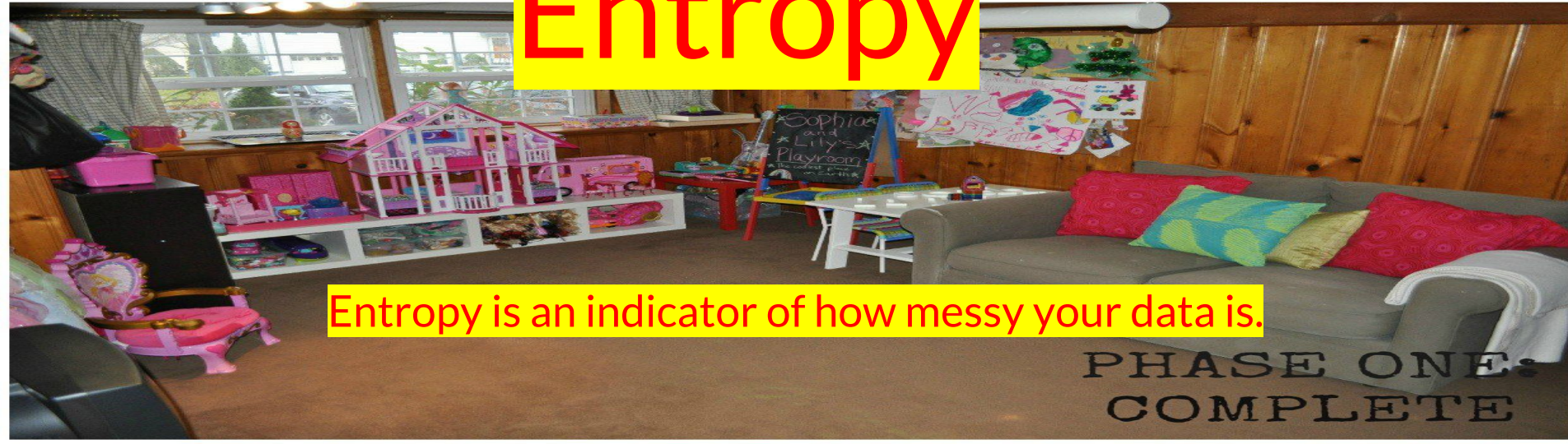
Predictors				Target
Outlook	Temp.	Humidity	Windy	Hours Played
Rainy	Hot	High	False	26
Rainy	Hot	High	True	30
Overcast	Hot	High	False	48
Sunny	Mild	High	False	46
Sunny	Cool	Normal	False	62
Sunny	Cool	Normal	True	23
Overcast	Cool	Normal	True	43
Rainy	Mild	High	False	36
Rainy	Cool	Normal	False	38
Sunny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







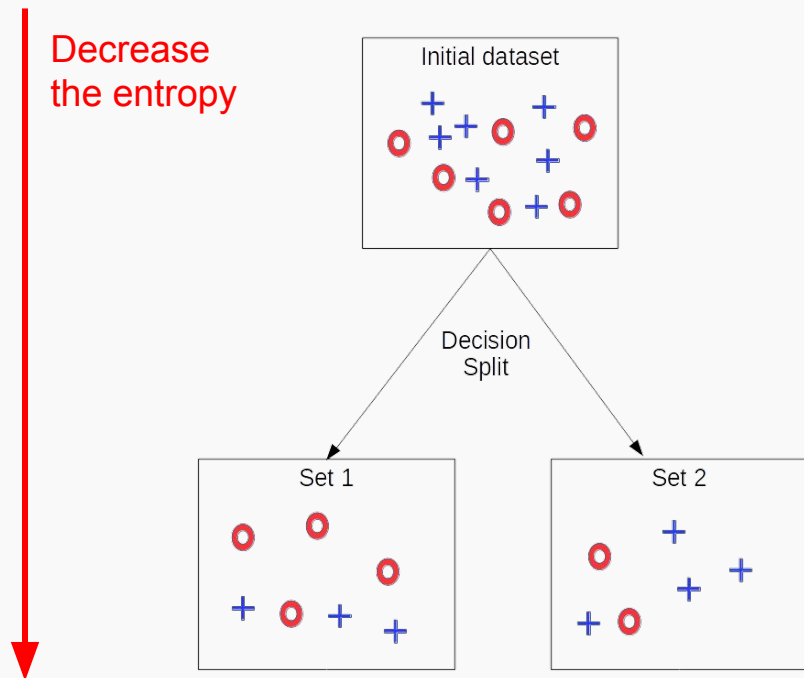
# Entropy



Entropy is an indicator of how messy your data is.

PHASE ONE:  
COMPLETE

# 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.



# 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)
```

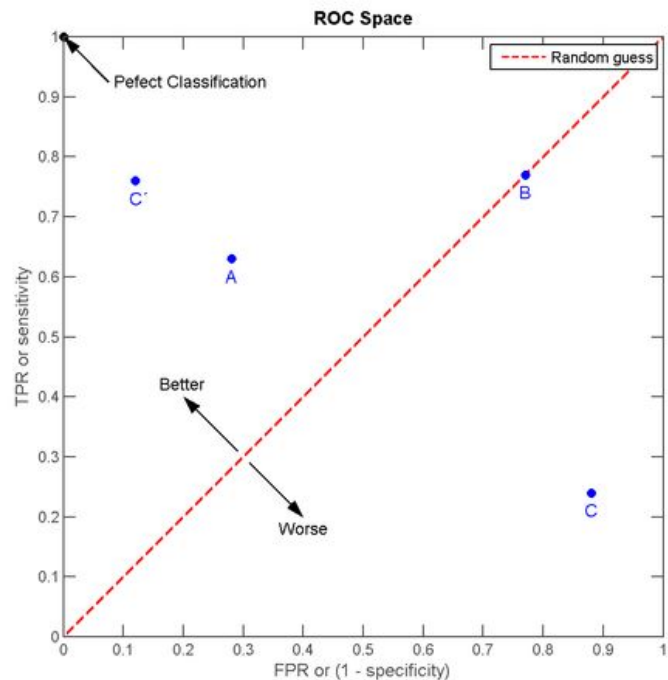
[\[source\]](#)

## `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) ¶
```

[\[source\]](#)

# Receiver Operating Characteristic (ROC)

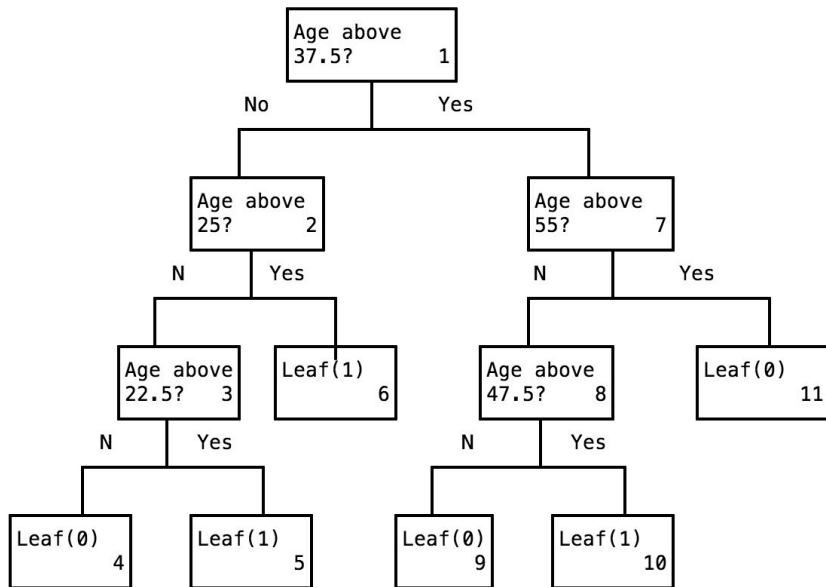


		True condition	
Total population		Condition positive	Condition negative
Predicted condition	Predicted condition positive	<b>True positive, Power</b>	<b>False positive, Type I error</b>
	Predicted condition negative	<b>False negative, Type II error</b>	<b>True negative</b>
		True positive rate (TPR), Recall, Sensitivity, probability of detection $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$

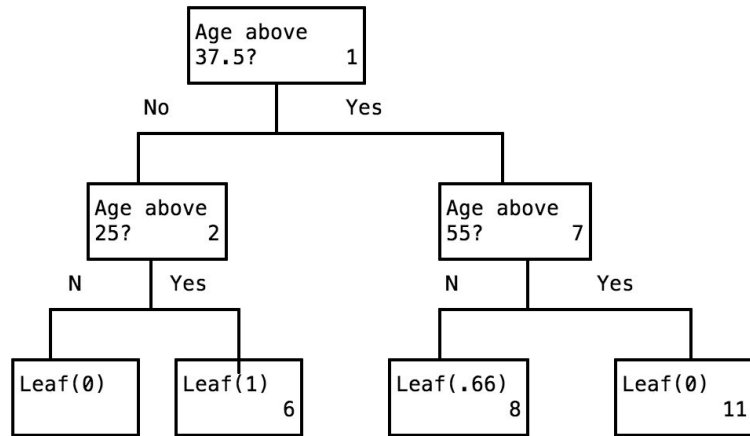
AUC - Area Under Curve

# Decision Tree Overfitting

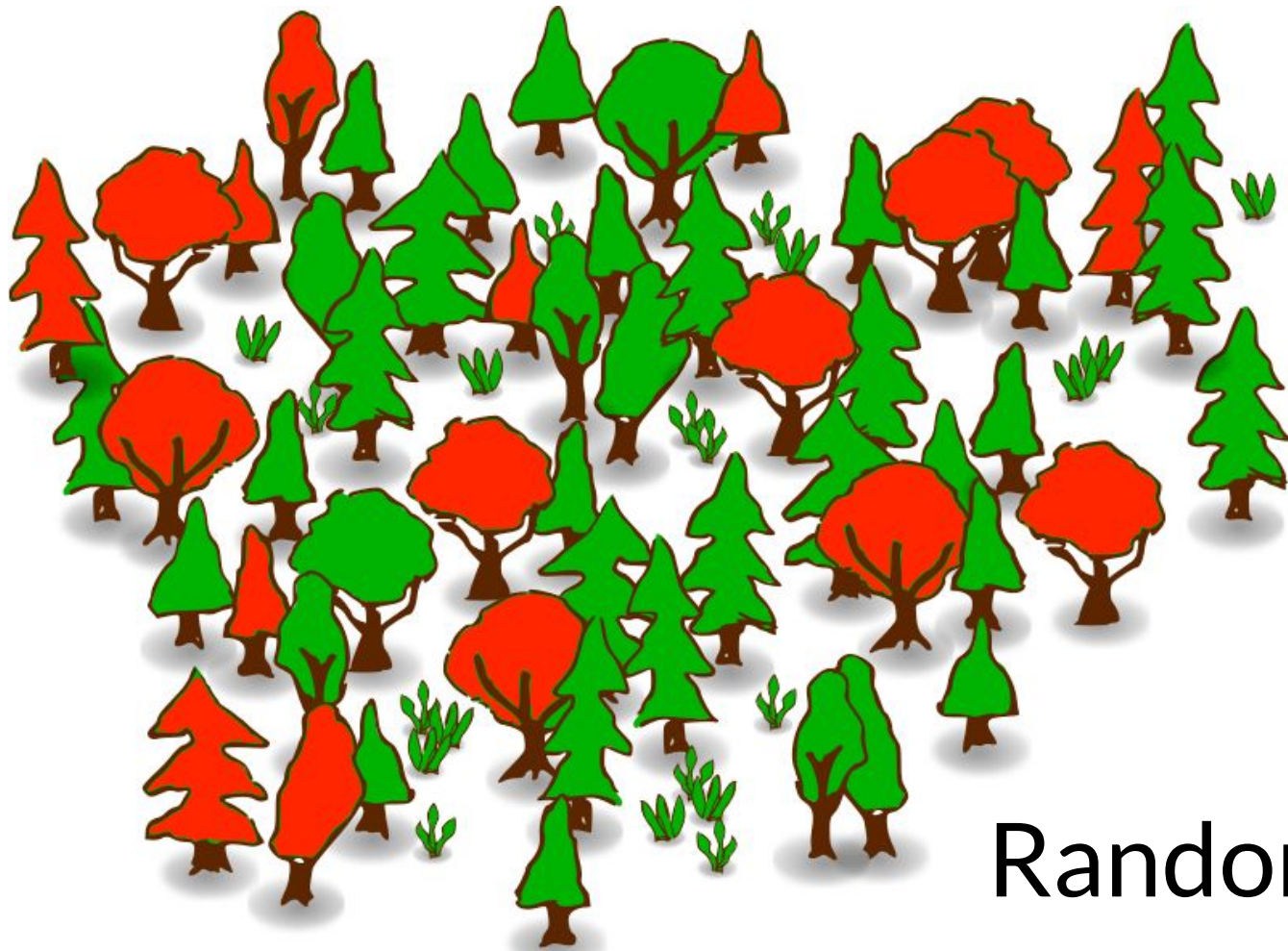
Full tree



Smaller tree



settings	train AUC	test AUC
default (min_samples_split: 2, max_depth: None)	0.947	0.694
min_samples_split: 13	0.842	0.699
min_samples_split: 13, max_depth: 7	0.748	0.743
min_samples_split: 100, max_depth: 2	0.662	0.655



# Random Forest



# The Dataset

---

The target column, or what we want to predict, is whether individuals make less than or equal to 50k a year, or more than 50k a year.

US Census  
1994

	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationship	race	sex	capital_gain
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0

# Combining Model Predictions with Ensembles

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import roc_auc_score
3
4 # features
5 columns = ["age", "workclass", "education_num", "marital_status",
6            "occupation", "relationship", "race", "sex",
7            "hours_per_week", "native_country"]
8
9 # model 1
10 clf = DecisionTreeClassifier(random_state=1, min_samples_leaf=2)
11 clf.fit(train[columns], train["high_income"])
12
13 # model 2
14 clf2 = DecisionTreeClassifier(random_state=1, max_depth=5)
15 clf2.fit(train[columns], train["high_income"])
16
17 # prediction on model 1
18 predictions = clf.predict(test[columns])
19 print(roc_auc_score(test["high_income"], predictions))
20
21 # prediction on model 2
22 predictions = clf2.predict(test[columns])
23 print(roc_auc_score(test["high_income"], predictions))
```

0.6878964226062301

0.6759853906508785

# Combining our predictions

## Majority Voting

DT1	DT2	DT3	Final Prediction
0	1	0	0
1	1	1	1
0	0	1	0
1	0	0	0

settings	test AUC
min_samples_leaf: 2	0.688
max_depth: 2	0.676
combined predictions	0.715

	#1 Model	#2 Model	Mean	Rounded
0	0.166667	0.288507	0.227587	0.0
1	0.000000	0.288507	0.144253	0.0
2	0.000000	0.180918	0.090459	0.0
3	0.000000	0.354167	0.177083	0.0
4	0.000000	0.041009	0.020504	0.0
5	0.000000	0.006875	0.003437	0.0
6	0.000000	0.006875	0.003437	0.0
7	0.333333	0.146179	0.239756	0.0
8	0.000000	0.006875	0.003437	0.0
9	0.666667	0.777431	0.722049	1.0
10	0.000000	0.041009	0.020504	0.0

Probability Voting

# Why Ensembling Works

```
Terminal
  GOAT
  Successfully authenticated!
  ? Enter a name for the repository: example
  ? Optionally enter a description of the repository: Just an example
  ? Public or private: public
  ? Select the files and/or folders you wish to ignore: (Press <space> to select)
  > bower_components
    o index.js
    o lib
    o node_modules
    o package.json
```



# Introduction Variation with Bagging

```

1 # We'll build 10 trees
2 tree_count = 10
3
4 # Each "bag" will have 60% of the number of original rows
5 bag_proportion = .6
6
7 predictions = []
8 for i in range(tree_count):
9     # We select 60% of the rows from train, sampling with replacement
10    # We set a random state to ensure we'll be able to replicate our results
11    # We set it to i instead of a fixed value so we don't
12    # get the same sample in every loop.
13    bag = train.sample(frac=bag_proportion, replace=True, random_state=i)
14
15    # Fit a decision tree model to the "bag"
16    clf = DecisionTreeClassifier(random_state=1, min_samples_leaf=2)
17    clf.fit(bag[colums], bag["high_income"])
18
19    # Using the model, make predictions on the test data
20    predictions.append(clf.predict_proba(test[colums])[:,1])
21 combined = np.sum(predictions, axis=0) / 10
22 rounded = np.round(combined)
23
24 print(roc_auc_score(test["high_income"], rounded))

```

0.7329963297474371

settings	test AUC
min_samples_leaf: 2	0.688
max_depth: 2	0.676
combined predictions	0.715
min_samples_leaf: 2, with bagging	0.732



# Introduction Variation from Random Features

---

```
# Fit a decision tree model to the "bag"
clf = DecisionTreeClassifier(random_state=1, min_samples_leaf=2,
                             splitter="random", max_features="auto")
clf.fit(bag[columns], bag["high_income"])
```

settings	test AUC
min_samples_leaf: 2	0.688
max_depth: 2	0.676
combined predictions	0.715
min_samples_leaf: 2, with bagging	0.732
min_samples_leaf: 2, with bagging and random subsets	0.734

# Put it All Together

---

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 clf = RandomForestClassifier(n_estimators=5, random_state=1,
4                             min_samples_leaf=2)
5 clf.fit(train[columns], train["high_income"])
6
7 predictions = clf.predict(test[columns])
8 print(roc_auc_score(test["high_income"], predictions))
```

0.7347461391939776

# Reducing Overfitting

---

```
1 clf = DecisionTreeClassifier(random_state=1, min_samples_leaf=5)
2
3 clf.fit(train[columns], train["high_income"])
4
5 predictions = clf.predict(train[columns])
6 print(roc_auc_score(train["high_income"], predictions))
7
8 predictions = clf.predict(test[columns])
9 print(roc_auc_score(test["high_income"], predictions)).
```

0.8192570489534683

0.7139325899284541

# Reducing Overfitting

---

```
1 clf = RandomForestClassifier(n_estimators=150, random_state=1,  
2                             min_samples_leaf=5)  
3 clf.fit(train[columns], train["high_income"])  
4  
5 predictions = clf.predict(train[columns])  
6 print(roc_auc_score(train["high_income"], predictions))  
7  
8 predictions = clf.predict(test[columns])  
9 print(roc_auc_score(test["high_income"], predictions))
```

0.7917047295143252

0.7498874343962398

# When to use Random Forest

---

- Strengths of a Random Forest
  - Very accurate predictions
  - Resistance to overfitting
- Weakness
  - They are difficult to interpret
  - They take longer to create



## Lesson #10 - Random Forest.ipynb

