

DECISION TREES AND RANDOM FORESTS

Abbas Chokor, Ph.D.

Staff Data Scientist, Seagate Technology

OUR PROGRESS SO FAR

UNIT 1: RESEARCH DESIGN AND EXPLORATORY DATA ANALYSIS

What is Data Science	Lesson 1
Research Design and Pandas	Lesson 2
Statistics Fundamentals I	Lesson 3
Statistics Fundamentals II	Lesson 4
Flexible Class Session	Lesson 5

UNIT 2: FOUNDATIONS OF DATA MODELING

Introduction to Regression	Lesson 6
Evaluating Model Fit	Lesson 7
Introduction to Classification	Lesson 8
Introduction to Logistic Regression	Lesson 9
Communicating Logistic Regression Results	Lesson 10
Flexible Class Session	Lesson 11

UNIT 3: DATA SCIENCE IN THE REAL WORLD

› Decision Trees and Random Forests	Lesson 12
› Natural Language Processing	Lesson 13
› Dimensionality Reduction	Lesson 14
› Time Series Data I	Lesson 15
› Time Series Data II	Lesson 16
› Database Technologies	Lesson 17
› Where to Go Next	Lesson 18
› Flexible Class Session	Lesson 19
› Final Project Presentations	Lesson 20



LAST CLASS

WHAT DID WE LEARN?

- Supervised vs unsupervised algorithms
- Understand and apply k-means clustering
- Density-based clustering: DBSCAN
- Silhouette Metric



You got all objectives?



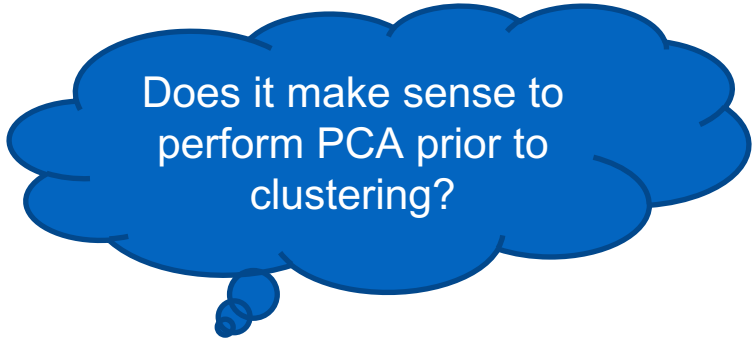
Not all of them...

Let's form groups of 1's and 2's ...

LAST CLASS

ANNOUNCEMENTS

- ❖ You will need to return your parking garage passes.
- ❖ Mid-class survey
- ❖ Final Project 2



Does it make sense to
perform PCA prior to
clustering?

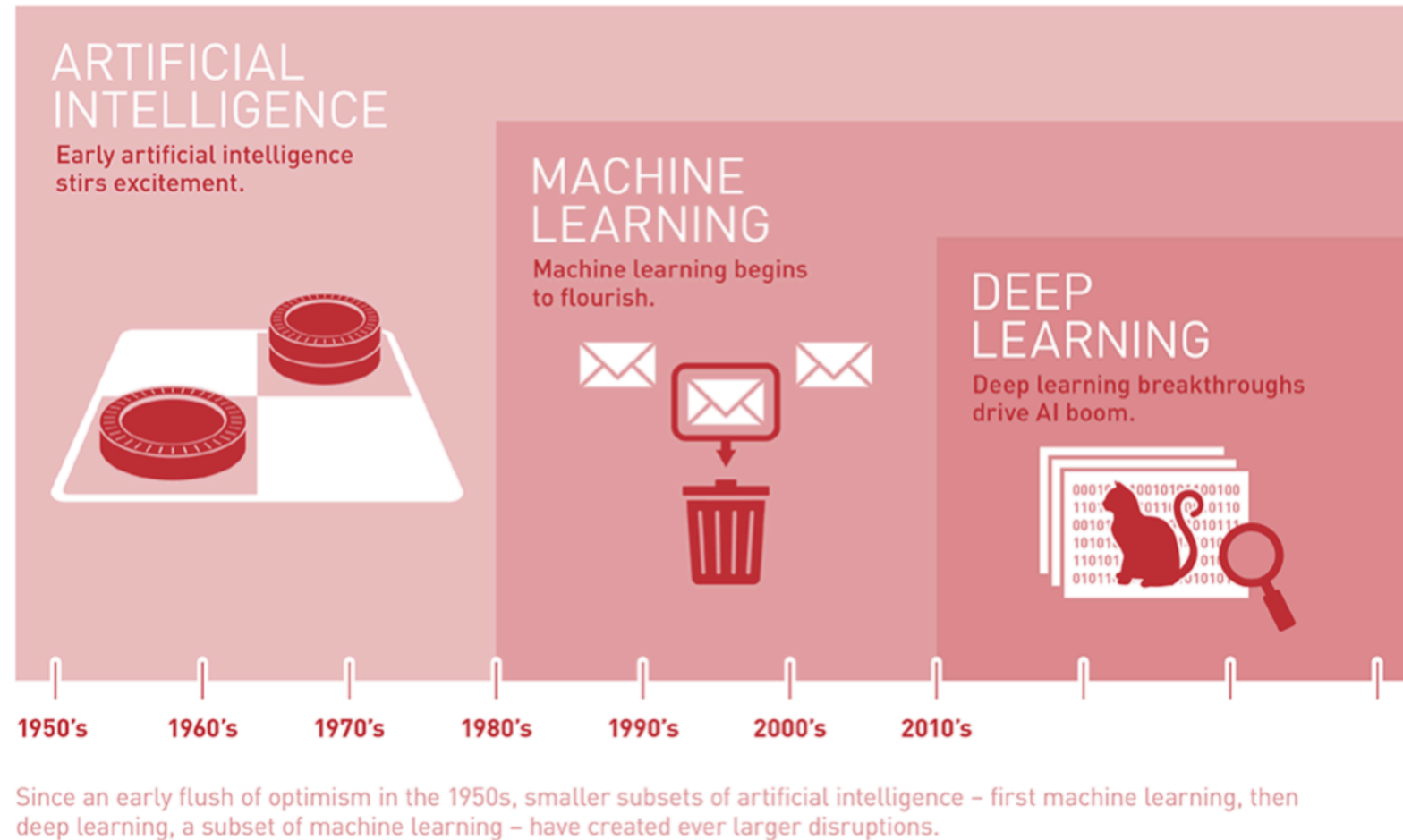


Others?

WHAT IS MACHINE LEARNING

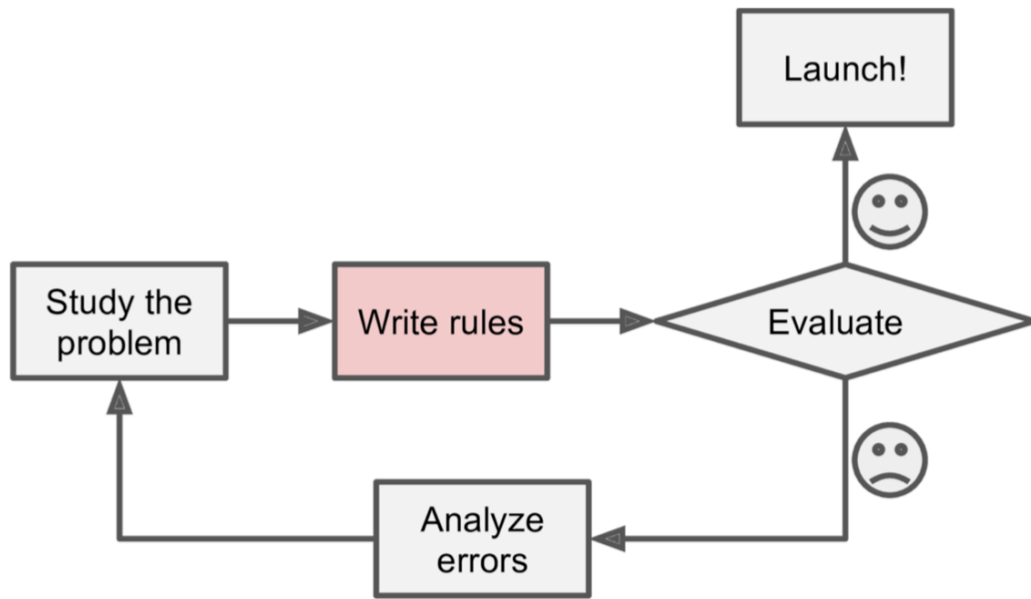
[Machine Learning is the] field of study that gives computers the ability to learn without being explicitly programmed.

Arthur Samuel, 1959



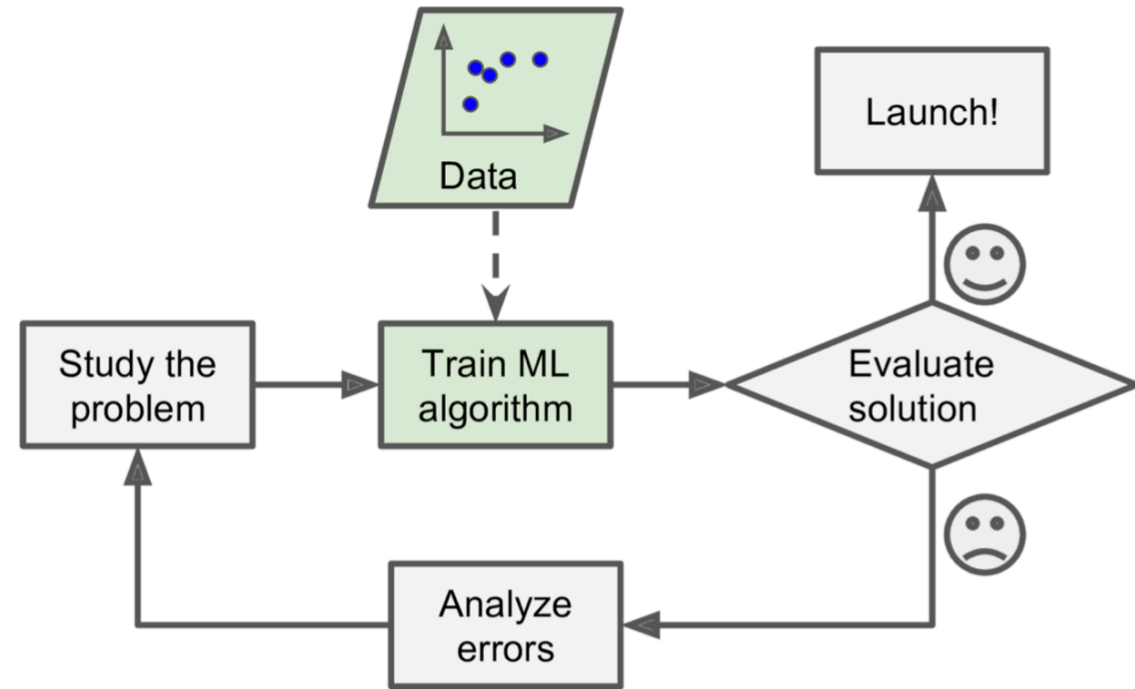
Machine Learning: is the science (and art) of programming computers so they can *learn from data*.

WHY TO USE MACHINE LEARNING?



Traditional approach

Vs.



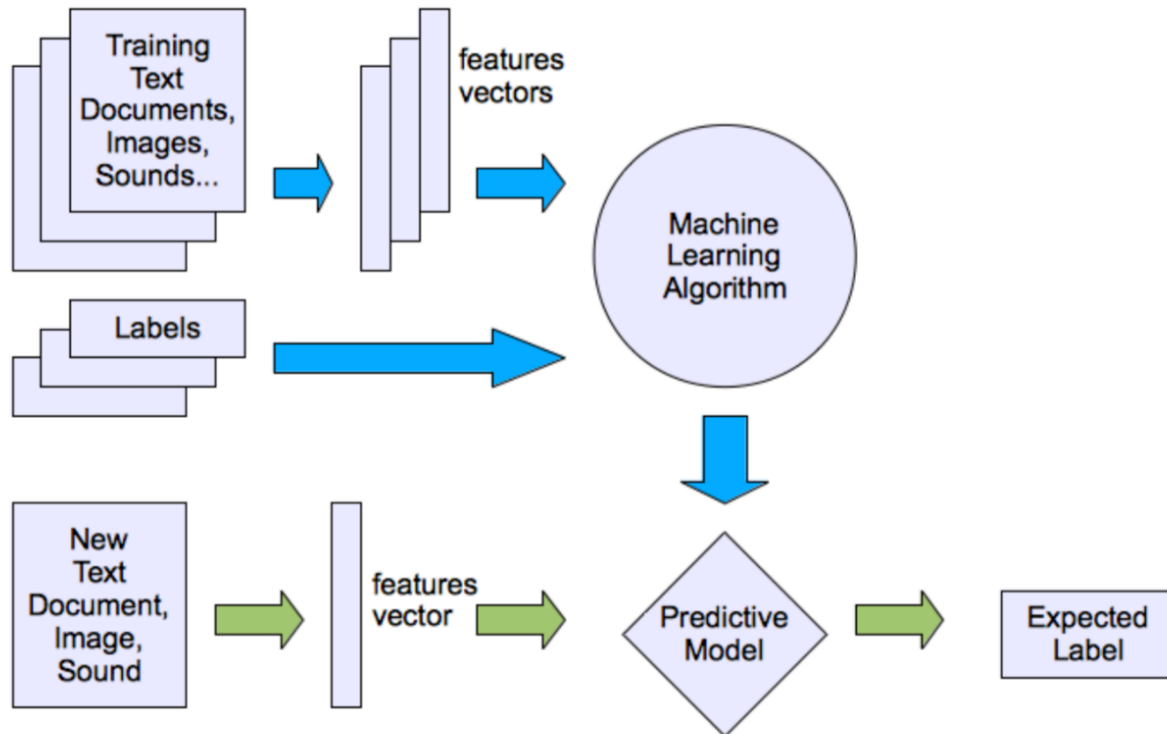
Machine Learning approach

Did you know?

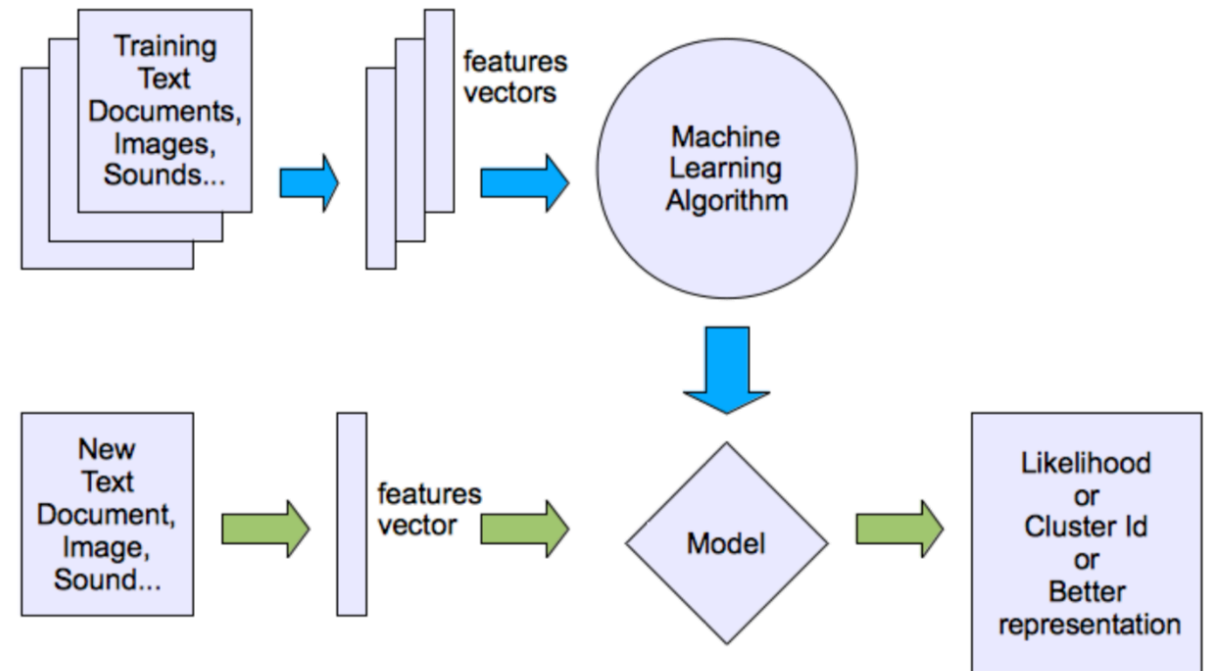
Machine learning algorithms are expected to replace 25% of the jobs across the world, in the next 10 years.

TYPES OF MACHINE LEARNING

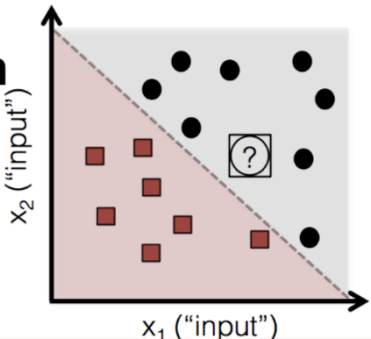
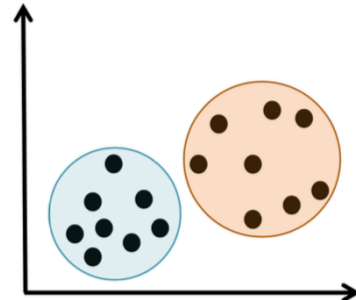
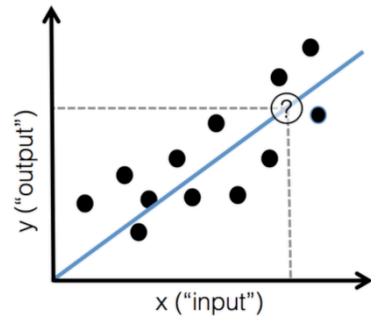
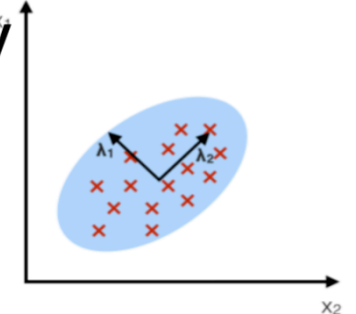
Supervised



Unsupervised



TYPES OF MACHINE LEARNING

	Supervised Working with Labeled Data	Unsupervised Working with Unlabeled Data
Discrete Countable Data	Classification 	Clustering 
Continuous Infinite Data	Regression 	Dimensionality Reduction 

DECISION TREES AND RANDOM FORESTS

LEARNING OBJECTIVES

- Understand and build **decision tree** models for classification and regression
- Understand and build **random forest** models for classification and regression
- Know how to extract the most **important predictors** in a random forest model

GUIDED PRACTICE

EXPLORE THE DATASET

ACTIVITY: EXPLORE THE DATASET



EXERCISE

DIRECTIONS (25 minutes)

We will be using a dataset from StumbleUpon, a service that recommends webpages to users based upon their interests. They like to recommend “evergreen” sites, ones that are always relevant. This usually means websites that avoid topical content and focus on recipes, how-to guides, art projects, etc. We want to determine important characteristics for “evergreen” websites. Follow these prompts to get started:

1. Break into groups.
2. Prior to looking at the data, brainstorm 3-5 characteristics that would be useful for predicting evergreen websites.
3. After looking at the dataset, can you model or quantify any of the characteristics you wanted? See the Notebook for data dictionary and starter code.
4. Does being a news site affect evergreenness? Compute or plot the percent of evergreen news sites.

ACTIVITY: EXPLORE THE DATASET



DIRECTIONS (25 minutes)

5. In general, does category affect evergreenness? Plot the rate of evergreen sites for all Alchemy categories.
6. How many articles are there per category?
7. Create a feature for the title containing “recipe”. Is the percentage of evergreen websites higher or lower on pages that have “recipe” in the title?

Check: Were you able to plot the requested features? Can you explain how you would approach this type of dataset?

DELIVERABLE

Requested features and answers to questions

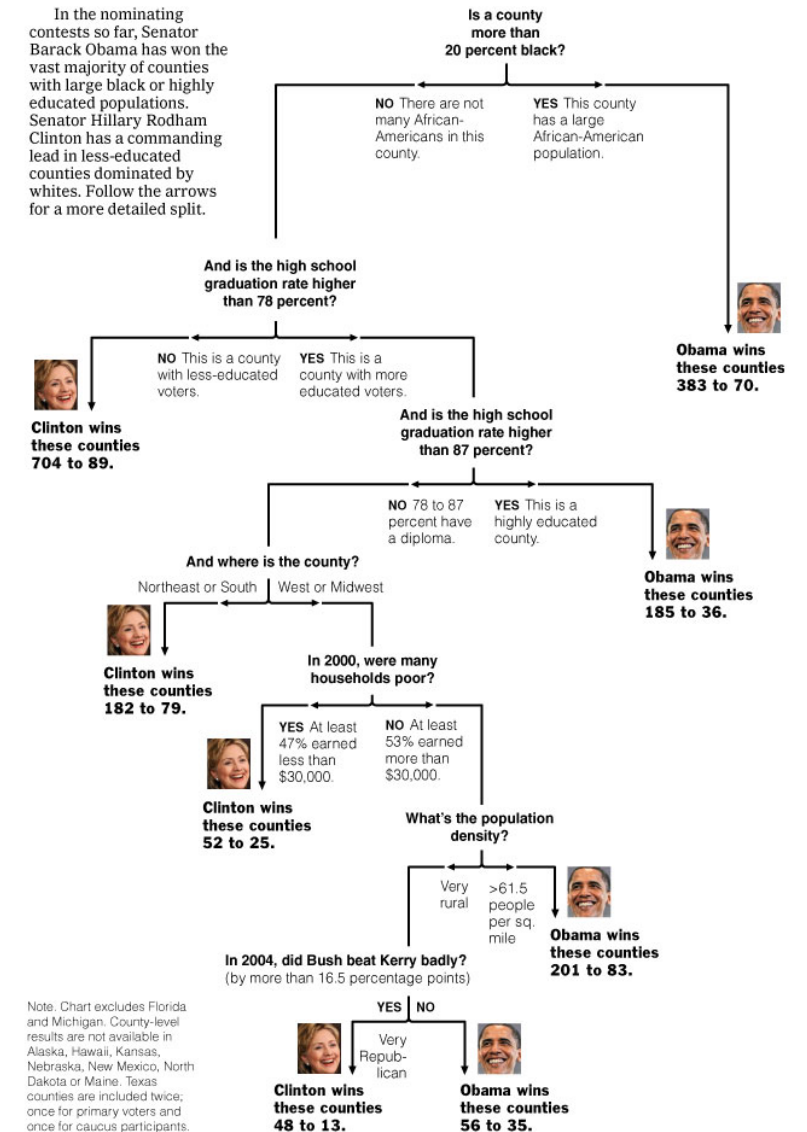
INTRODUCTION

TRAINING DECISION TREES

INTUITION BEHIND DECISION TREES

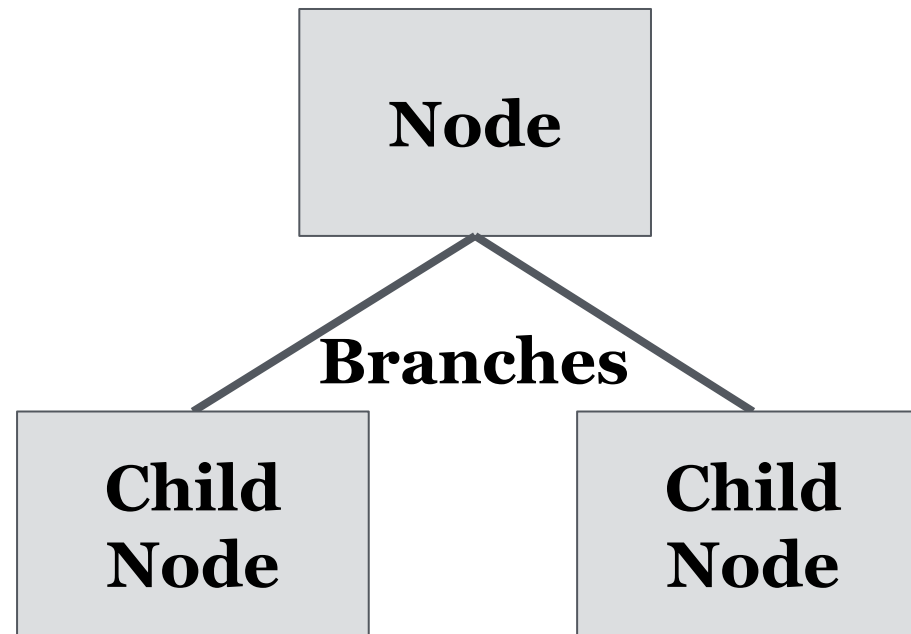
- Decision trees are like the game “20 questions”. They make decision by answering a series of questions, most often binary questions (yes or no).
- We want the smallest set of questions to get to the right answer.
- Each questions should reduce the search space as much as possible.

Decision Tree: The Obama-Clinton Divide



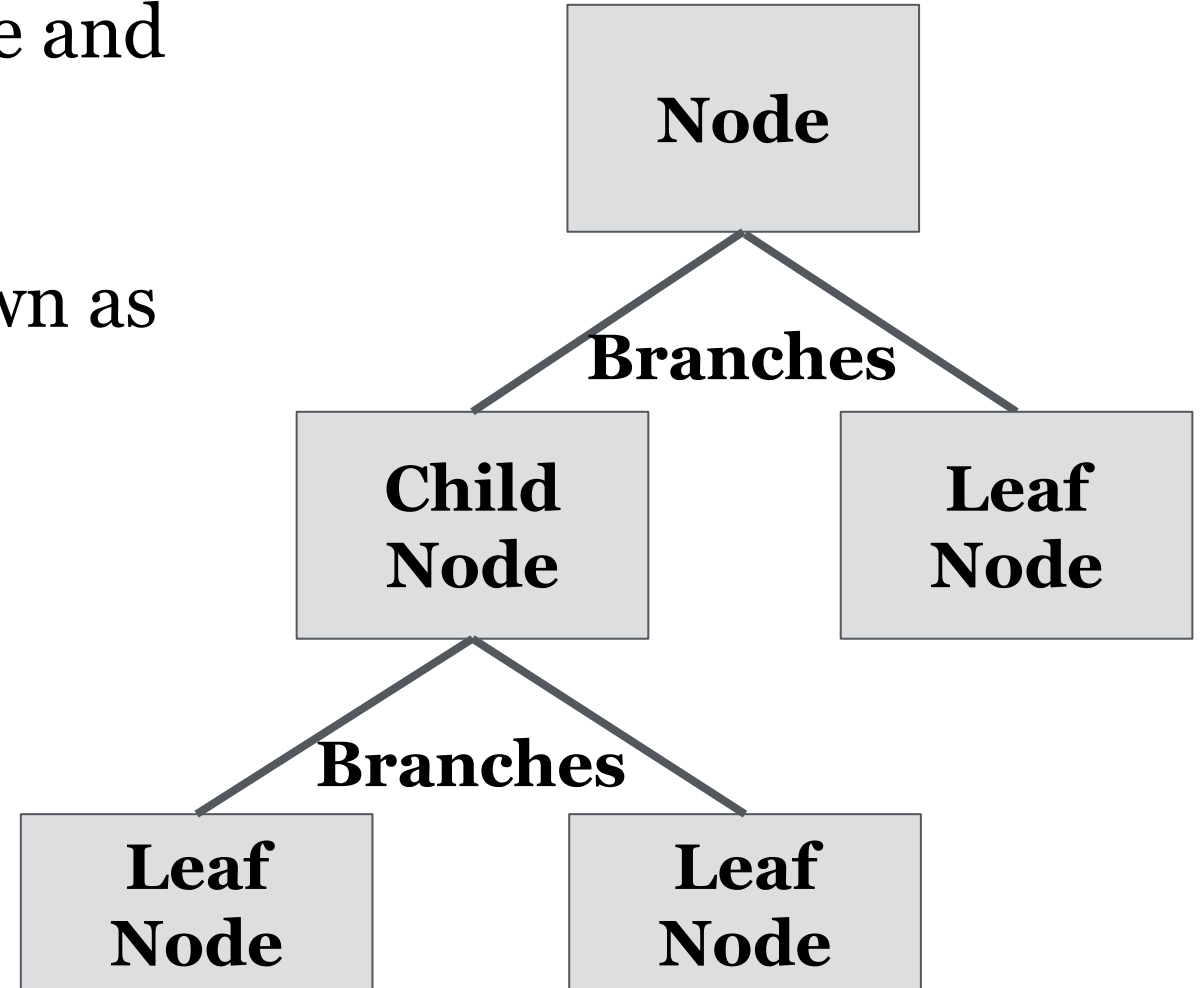
TREES

- Trees are a data structure made up of *nodes* and *branches*.
- Each node typically has two or more branches that connect it to its children.



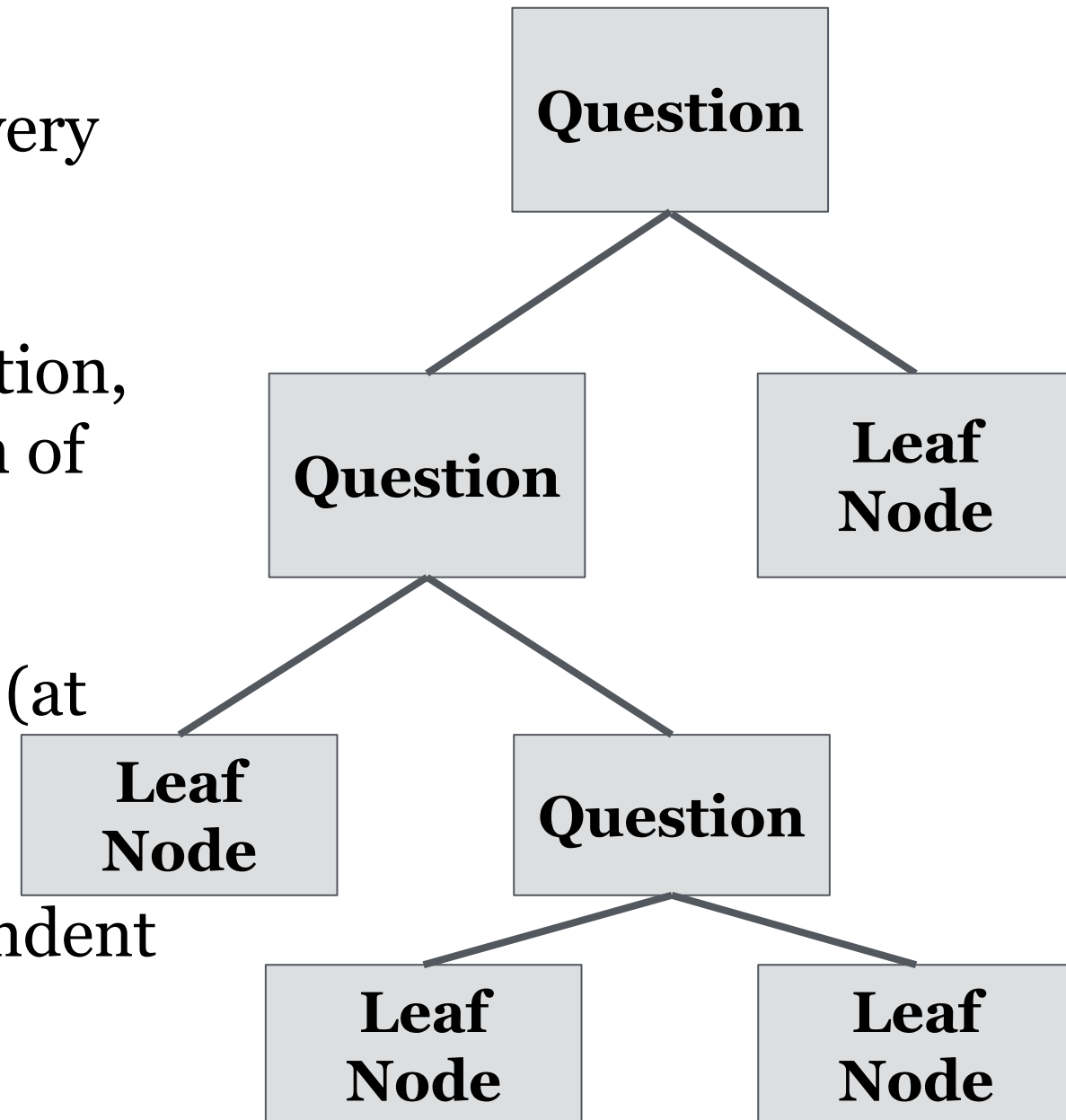
TREES

- Each child is another node in the tree and contains its own *subtree*.
- Nodes without any children are known as *leaf* nodes.



DECISION TREES

- A *decision tree* contains a question at every node.
- Depending upon the answer to the question, we proceed down the left or right branch of the tree and ask another question.
- Once we don't have any more questions (at the *leaf* nodes), we make a prediction.
- Note: The next question is always dependent on the last.



DECISION TREES

- Let's suppose we want to predict if an article is a news article.
- What questions should we ask to make a prediction?
- How many questions should we ask?

DECISION TREES

- We may start by asking: does it mention a President?
- If it does, it must be a news article.
- If not, let's ask another question: does the article contain other political features?
- If not, does the article contain references to political topics?
- We could keep going on in this manner until we were satisfied.

DECISION TREES

Because it isn't computationally feasible to consider every possible partition of the feature space, we take a **top-down, greedy** approach known as recursive binary splitting

- **Top-Down**: The approach begins at the top of the tree and then successively splits the predictor space; each split is indicated via two new branches further down on the tree
- **Greedy**: At each step of the tree-building process, the best split is made at that particular step, rather than looking ahead and picking a split that will lead to a better tree in some future step

ACTIVITY: KNOWLEDGE CHECK

ANSWER THE FOLLOWING QUESTIONS

Let's work as a class to accomplish the following:

1. Using our StumpleUpon dataset, try to predict whether a given article is evergreen.
2. Build a decision tree to determine the above.

DELIVERABLE

Our decision tree



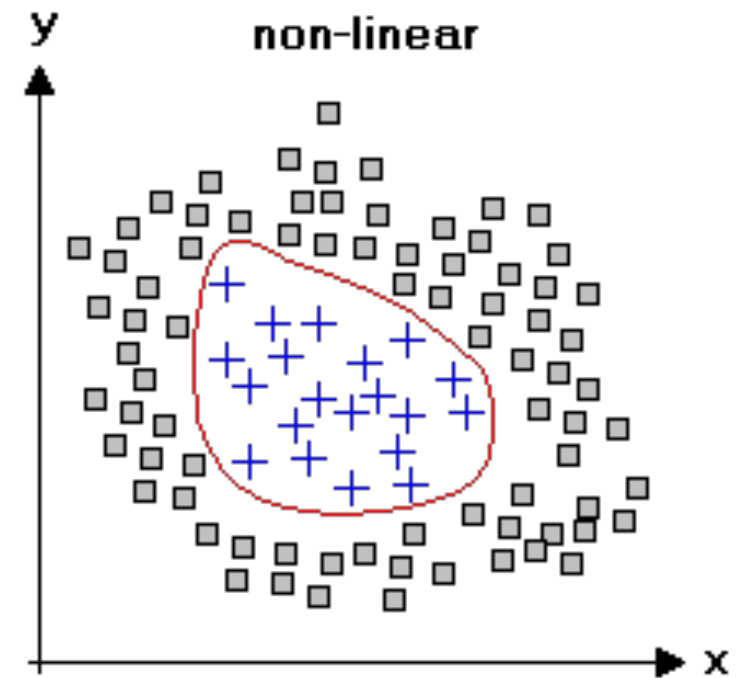
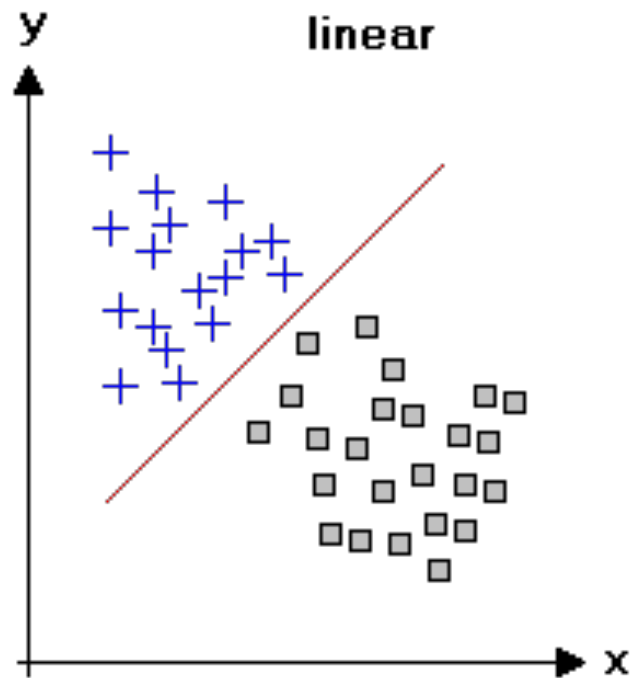
EXERCISE

COMPARISON TO PREVIOUS MODELS

- Decision trees are *non-linear*, an advantage over logistic regression.
- A *linear* model is one in which a change in an input variable has a constant change on the output variable.

COMPARISON TO PREVIOUS MODELS

- Linear vs. non-linear classification models



COMPARISON TO PREVIOUS MODELS

- An example of this difference is the relationship between years of education and salary. In a *linear* model, the increase in salary from 10 to 15 years of education would be the same as the increase in salary from 15 to 20 years of education. In a *non-linear* model, salary can change dramatically for years 0-15 and negligibly from years 15-20.
- Trees automatically contain interaction of features, since each question is dependent on the last.

TRAINING A DECISION TREE MODEL

- Training a decision model is deciding the best set of questions to ask.
- A good question will be one that best segregates the positive group from the negative group and then narrows in on the correct answer.
- For example, in our news article decision tree, the best question is one that creates two groups, one that is mostly news stories and one that is mostly non-news stories.

TRAINING A DECISION TREE MODEL

- We can quantify the *purity* of the separation of groups using Classification Error, Entropy, or Gini Coefficient.
- We want to choose the question that gives us the best *change* in our purity measure. At each step, we can ask, “Given our current set of data points, which question will make the largest change in purity?”
- This is done *recursively* for each new set of two groups until we reach a stopping point.

TRAINING A DECISION TREE MODEL (optional not required)

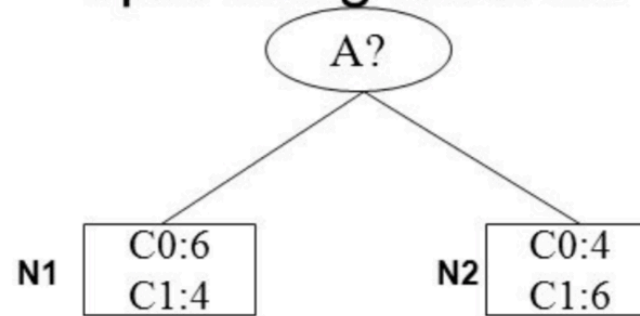
- How does a tree decide where to split?

GINI Index

$$Gini = \sum_{i \neq j} p(i)p(j)$$

i and j are levels of the target variable

• Split using attribute A

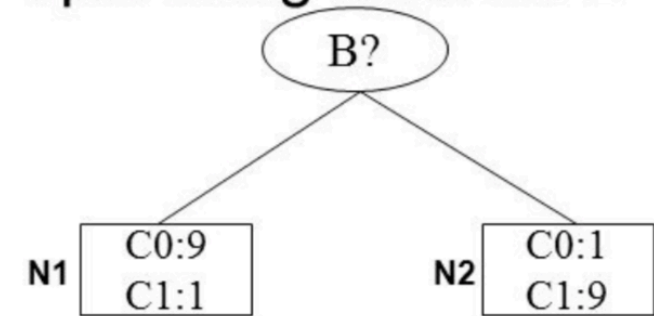


$$Gini(N1) = 1 - 0.6^2 - 0.4^2 = 0.48$$

$$Gini(N2) = 1 - 0.4^2 - 0.6^2 = 0.48$$

$$Gini_{splitA} = 10/20 * 0.48 + 10/20 * 0.48 \\ = 0.48$$

• Split using attribute B



$$Gini(N1) = 1 - 0.9^2 - 0.1^2 = 0.18$$

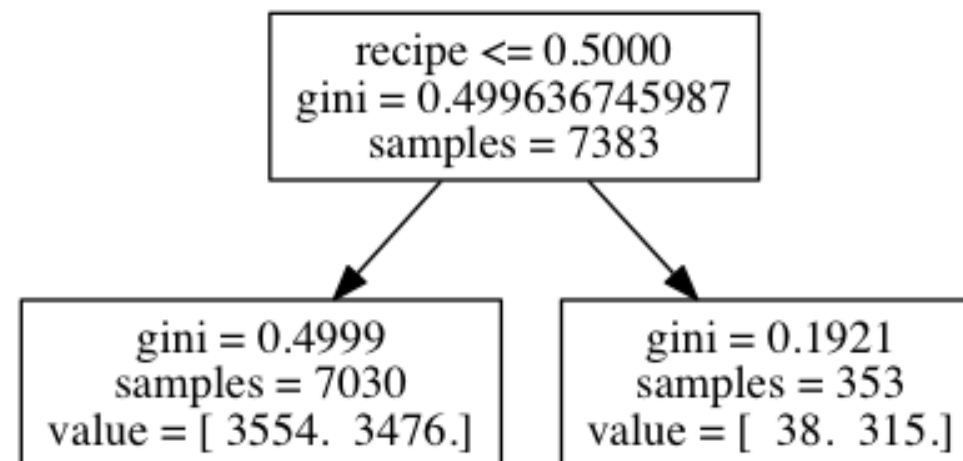
$$Gini(N2) = 1 - 0.1^2 - 0.9^2 = 0.18$$

$$Gini_{splitB} = 10/20 * 0.18 + 10/20 * 0.18 \\ = 0.18$$

$Gini_{splitB} < Gini_{splitA}$ means splitting using attribute B has **lower impurity** in the children nodes. Thus splitting using attribute B is **preferred**.

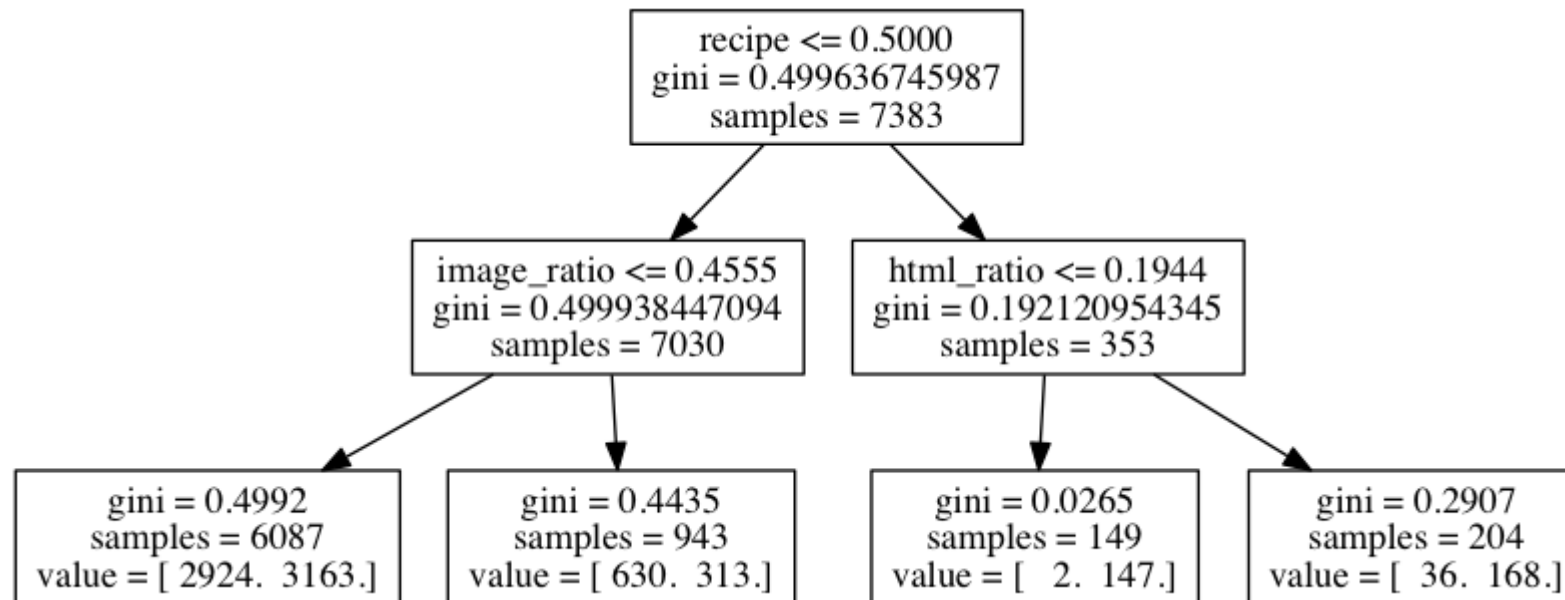
TRAINING A DECISION TREE MODEL

- Let's build a sample tree for our evergreen prediction problem. Assume our features are whether the article contains a recipe, the image ratio, the html ratio.
- First, let's choose the feature that gives us the highest purity, the recipe feature.



TRAINING A DECISION TREE MODEL

- We can take each side of the tree and repeat the process.



- We can continue this process until we have asked as many questions as we want or until our leaf nodes are completely pure.

TRAINING A DECISION TREE MODEL

Pros

- Trees are very easy to explain to people. They are even easier to explain than linear regression
- Decision trees more closely mirror human decision-making than do the regression and classification methods seen so far
- Trees can be displayed graphically and are easily interpreted even by non-experts
- Trees can easily handle qualitative predictors without the need to create binary variables

Cons

- Trees do not generally have the same level of predictive accuracy as some of the other regression and classification methods seen so far. However, by aggregating many decision trees, the predictive performance of trees can be substantially improved

MAKING PREDICTIONS FROM A DECISION TREE

- Predictions are made by answering each of the questions.
- Once we reach a leaf node, our prediction is made by taking the majority label of the training samples that fulfill the questions.
- In our sample tree, if we want to classify a new article, ask:
 - Does the article contain the word recipe?
 - If it doesn't, does the article have a lot of images?
 - If it does, then 630 / 943 article are evergreen.
 - So we can assign a 0.67 probability for evergreen sites.

ACTIVITY: KNOWLEDGE CHECK



ANSWER THE FOLLOWING QUESTIONS

1. How do we classify a new article?
2. How do we make predictions from a decision tree?

DELIVERABLE

Answers to the above questions

GUIDED PRACTICE

DECISION TREES IN SCIKIT-LEARN

ACTIVITY: DECISION TREES IN SCIKIT-LEARN



EXERCISE

DIRECTIONS (15 minutes)

1. In the starter code notebook, work through the exercises titled “Decision Trees in scikit-learn”.
2. In your groups from earlier, work on evaluating the decision tree using cross-validation methods.
3. What metrics would work best? Why?

Check: Are you able to evaluate the decision tree model using cross-validation methods?

DELIVERABLE

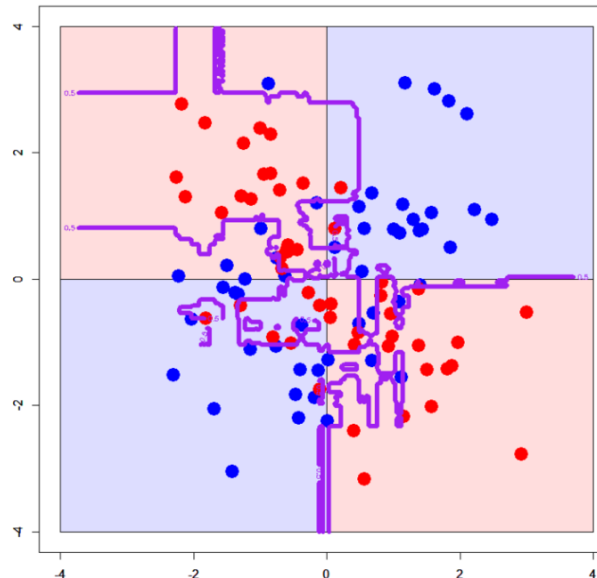
Completed exercises and answer to #3

DEMO

OVERFITTING IN DECISION TREES

OVERFITTING IN DECISION TREES

- Decision trees tend to be weak models because they can easily memorize or overfit to a dataset.
- A model is *overfit* when it memorizes or bends to a few specific data points rather than picking up general trends in the data.



OVERFITTING IN DECISION TREES

- An unconstrained decision tree can learn an extreme tree (e.g. one feature for each word in a news article).
- We can limit our decision trees using a few methods.
 - Limiting the number of questions (nodes) a tree can have.
 - Limiting the number of samples in the leaf nodes.

ACTIVITY: KNOWLEDGE CHECK



ANSWER THE FOLLOWING QUESTIONS

1. Why are decision trees generally thought of as weak models?
2. How can we limit our decision trees?

DELIVERABLE

Answers to the above questions

GUIDED PRACTICE

ADJUSTING DECISION TREES TO AVOID OVERFITTING

ACTIVITY: ADJUSTING DECISION TREES TO AVOID OVERFITTING

DIRECTIONS (15 minutes)

1. You can control for overfitting in decision trees by adjusting one of the following parameters:
 - a. `max_depth`: Control the maximum number of questions.
 - b. `min_samples_in_leaf`: Control the minimum number of records in each node.
2. Test each of these parameters in the starter code notebook.

DELIVERABLE

Code using the above parameters



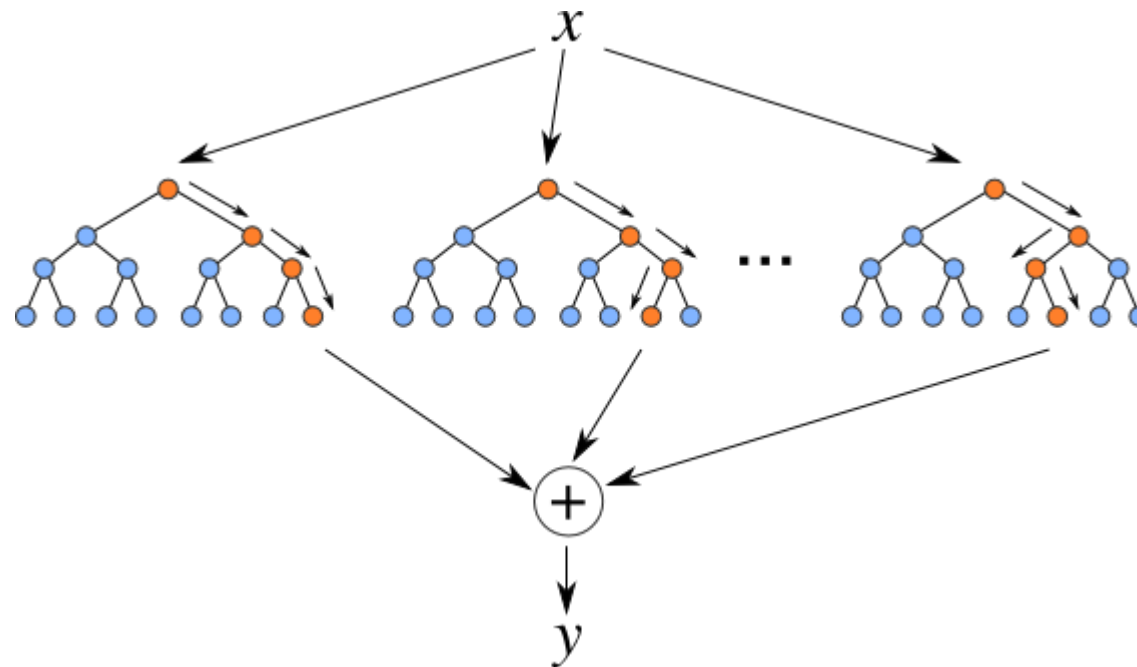
EXERCISE

INTRODUCTION

RUNNING THROUGH THE RANDOM FORESTS

RUNNING THROUGH THE RANDOM FORESTS

- ▶ Random forest models are one of the most widespread classifiers used.
- ▶ They are relatively simple to use and help avoid overfitting.
- ▶ Random Forests are an *ensemble* or collection of individual decision trees.



PROS AND CONS OF RANDOM FORESTS

- Advantages
 - Easy to tune
 - Built-in protection against overfitting
 - Non-linear
 - Built-in interaction effects
- Disadvantages
 - Slow
 - Black-box
 - No “coefficients”
 - Harder to explain

TRAINING A RANDOM FOREST

- Training a random forest model involves training many decision tree models.
- Since decision trees overfit easily, we use many decision trees together and randomize the way they are created.

TRAINING A RANDOM FOREST

- Random Forest Algorithm

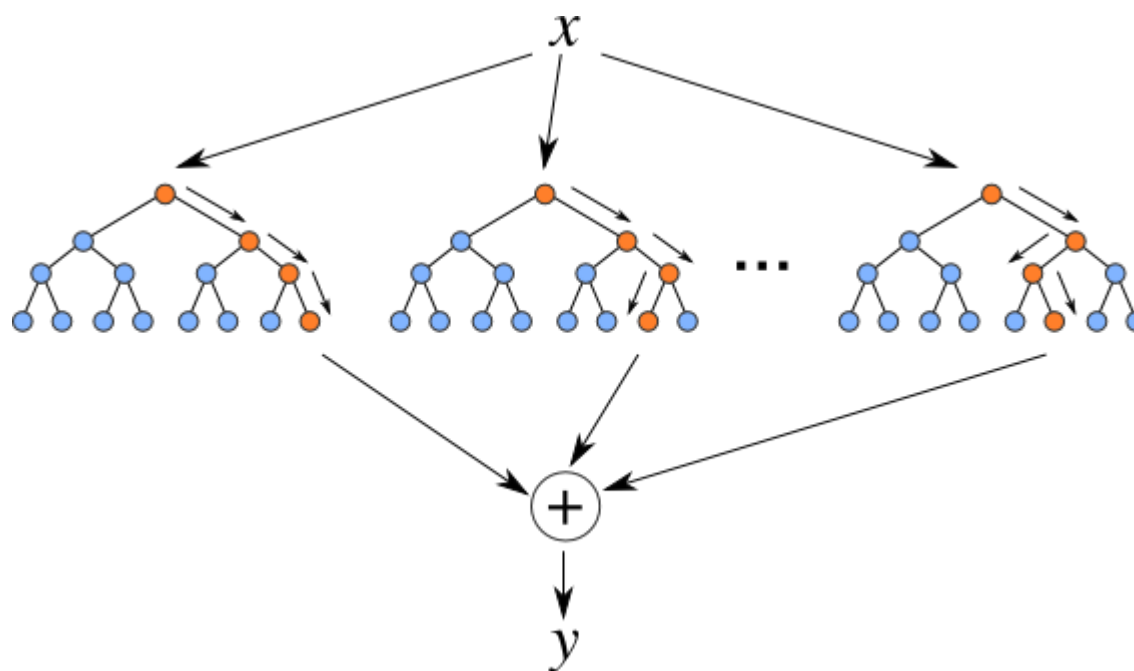
- a. Take a bootstrap sample of the dataset.

- a. Train a decision tree on the bootstrap sample. For each split/feature selection, only evaluate a *limited* number of features to find the best one.

- a. Repeat this for N trees.

PREDICTIONS USING A RANDOM FOREST

- Predictions for a random forest model come from each decision tree.
- Make an individual prediction with each decision tree.
- Combine the individual predictions and take the majority vote.



GUIDED PRACTICE

REGRESSION WITH DECISION TREES AND RANDOM FORESTS

ACTIVITY: REGRESSION WITH DECISION TREES & RANDOM FORESTS



EXERCISE

DIRECTIONS (20 minutes)

1. Build a random forest model to predict the evergreeness of a website. Remember to use the parameter `n_estimators` to control the number of trees used in the model.
2. Take note of the features that give the best splits to determine the most important features.
3. Decision trees and random forests can be used for both classification and regression. In regression, predictions are made by taking the average value of the samples in the leaf node. You can take the average of the individual trees' predictions. Build a regression based random forest model.

DELIVERABLE

The models mentioned above

INDEPENDENT PRACTICE

**EVALUATE
RANDOM FOREST
USING CROSS-
VALIDATION**

ACTIVITY: EVALUATE RANDOM FOREST USING CROSS-VALIDATION



DIRECTIONS (25 minutes)

1. Building upon the previous Guided Practice, add any input variables to the model that you think may be relevant.

1. For each feature:
 - a. Evaluate the model for improved predictive performance using cross-validation.
 - b. Evaluate the importance of the feature.

1. **Bonus:** Just like the 'recipe' feature, add in similar text features and evaluate their performance.

DELIVERABLE

Newly created features and models

CONCLUSION

TOPIC REVIEW

REVIEW Q&A

- What are decision trees?
- What does training involve?
- What are some common problems with decision trees?
- What are random forests?
- What are some common problems with random forests?

OUR PROGRESS SO FAR

UNIT 1: RESEARCH DESIGN AND EXPLORATORY DATA ANALYSIS

What is Data Science	Lesson 1
Research Design and Pandas	Lesson 2
Statistics Fundamentals I	Lesson 3
Statistics Fundamentals II	Lesson 4
Flexible Class Session	Lesson 5

UNIT 2: FOUNDATIONS OF DATA MODELING

Introduction to Regression	Lesson 6
Evaluating Model Fit	Lesson 7
Introduction to Classification	Lesson 8
Introduction to Logistic Regression	Lesson 9
Communicating Logistic Regression Results	Lesson 10
Flexible Class Session	Lesson 11

UNIT 3: DATA SCIENCE IN THE REAL WORLD

Decision Trees and Random Forests	Lesson 12
› Natural Language Processing	Lesson 13
› Dimensionality Reduction	Lesson 14
› Time Series Data I	Lesson 15
› Time Series Data II	Lesson 16
› Database Technologies	Lesson 17
› Where to Go Next	Lesson 18
› Flexible Class Session	Lesson 19
› Final Project Presentations	Lesson 20



LESSON

Q & A

LESSON

EXIT TICKET

DON'T FORGET TO FILL OUT YOUR EXIT TICKET