# Introduction to Machine Learning

# Module 5:
# Decision trees

# Module Checklist:

- ❏   Decision trees
  - ❏   Intuition
  - ❏   The "best" split
  - ❏   Model performance
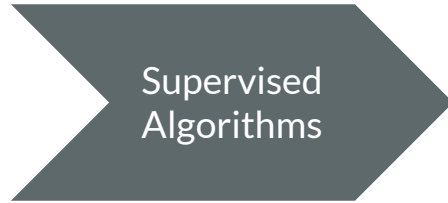  - ❏   Model optimization (pruning)

# Where are we?

We've laid much of the groundwork for machine learning, and introduced our very first algorithm of linear regression.

*Now, we can focus on expanding our algorithm toolkit.* **Decision trees** are another type of algorithm that can accomplish the same objective of prediction that linear regression can. We will also go deeper into the pros and cons of decision trees in this module.

# Where are we?

Supervised Algorithms

Linear regression

Decision tree

Ensemble Algorithms

Let's do a quick intro to decision trees and ensembles!

**Decision Tree**

Today, we will start by looking at a <u>decision tree</u> algorithm.

A decision tree is a set of rules we can use to classify data into categories **(also can be used for regression tasks)**.

Humans often use a similar approach to arrive at a conclusion. **For example, doctors ask a series of questions to diagnose a disease. A doctor's goal is to ask the minimum number of questions needed to arrive at the correct diagnosis.**

What is wrong with my patient?

Help me put together some questions I can use to diagnose patients correctly.
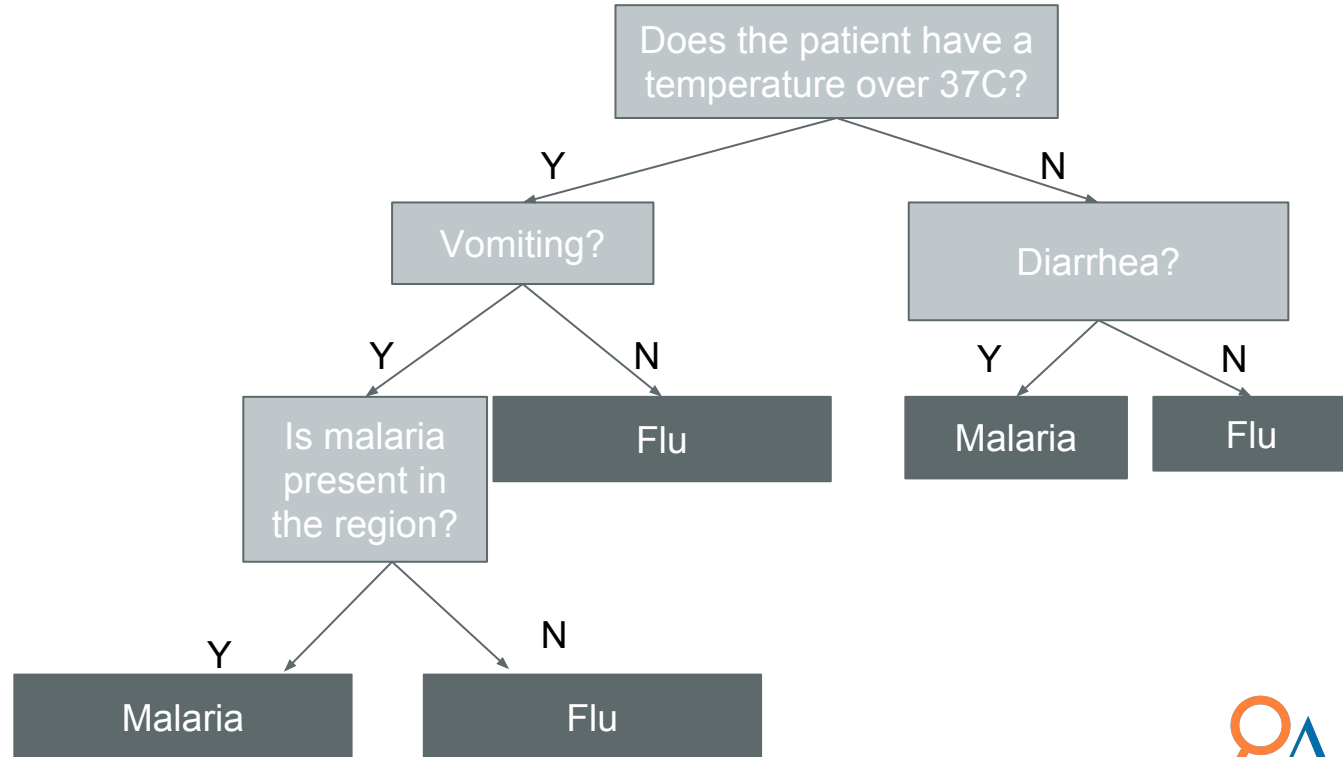
Decision trees are intuitive because they are similar to how we make many decisions.

My mental diagnosis decision tree might look something like this.

How is the way I think about this different from a machine learning algorithm?

Does the patient have a temperature over 37C?

Y — Vomiting?

N — Diarrhea?

Vomiting?
Y — Is malaria present in the region?
N — Flu

Is malaria present in the region?
Y — Malaria
N — Flu

Diarrhea?
Y — Malaria
N — Flu

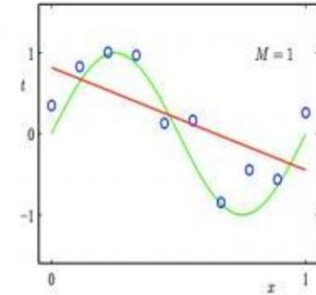Decision trees are intuitive and can handle more complex relationships than linear regression can.

A linear regression is a single global trend line.

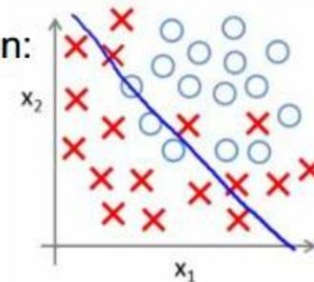This makes it inflexible for more sophisticated relationships.

Sources:
https://www.slideshare.net/ANITALOKITA/winnow-vs-perceptron,
http://www.turingfinance.com/regression-analysis-using-python-statsmodels-and-quandl/
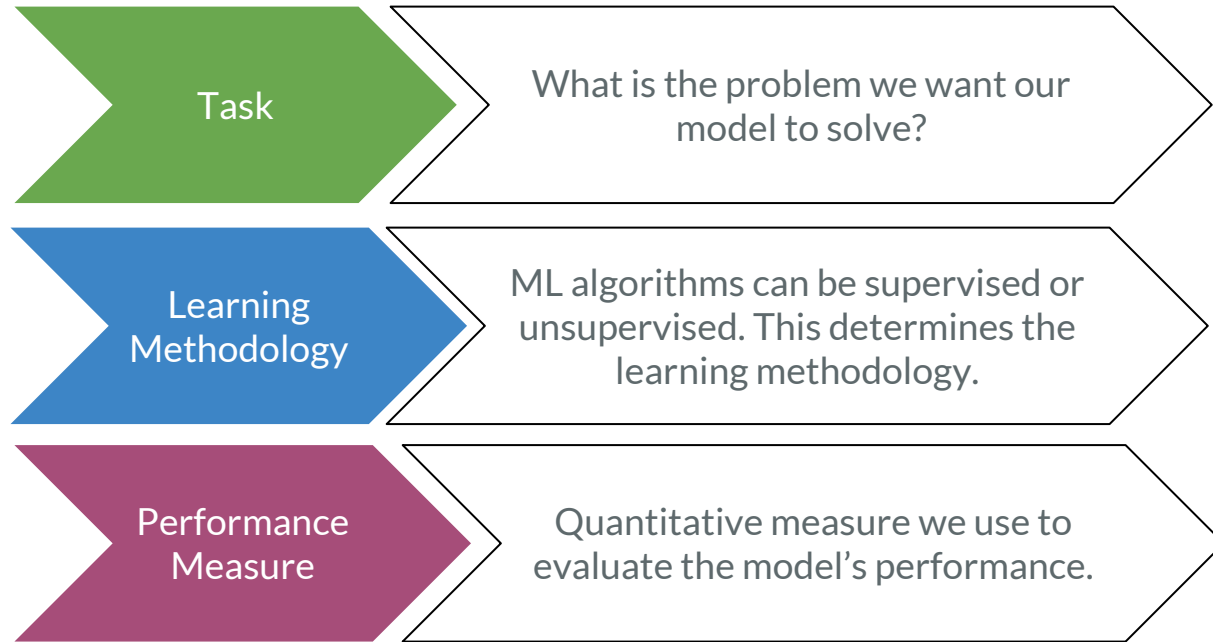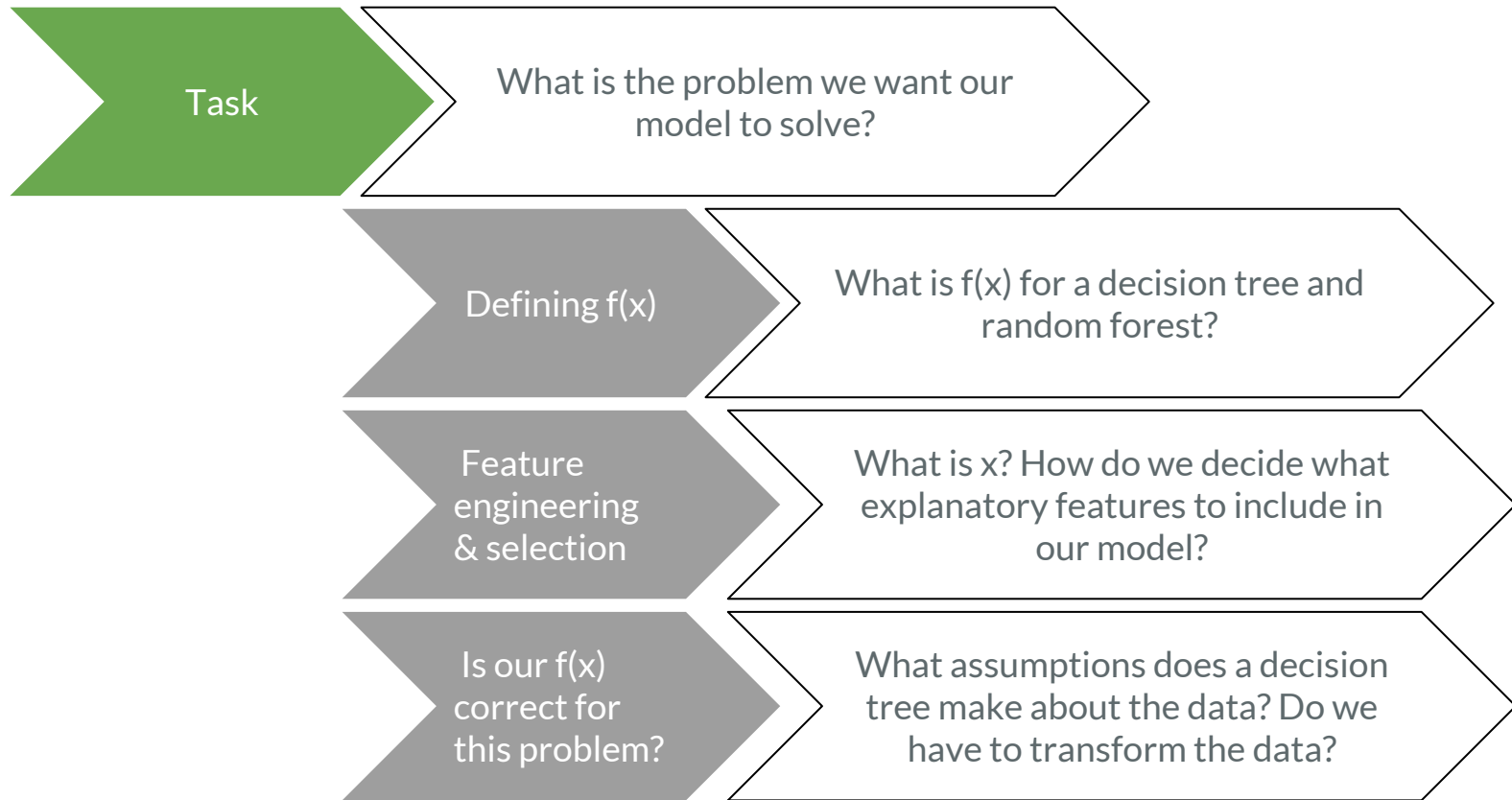
Regression:

predictor too inflexible: cannot capture pattern

Classification:

# We will use our now familiar framework to discuss both decision trees and ensemble algorithms:

| Task | What is the problem we want our model to solve? |
|------|--------------------------------------------------|
| Learning Methodology | ML algorithms can be supervised or unsupervised. This determines the learning methodology. |
| Performance Measure | Quantitative measure we use to evaluate the model's performance. |

| | |
|---|---|
| **Task** | What is the problem we want our model to solve? |
| **Defining f(x)** | What is f(x) for a decision tree and random forest? |
| **Feature engineering & selection** | What is x? How do we decide what explanatory features to include in our model? |
| **Is our f(x) correct for this problem?** | What assumptions does a decision tree make about the data? Do we have to transform the data? |

| Learning Methodology | Decision tree models are supervised. How does that affect the learning processing? |

| How does our ML model learn? | Overview of how the model teaches itself. |

| What is our loss function? | Every supervised model has a loss function it wants to minimize. |

| Optimization process | How does the model minimize the loss function? |

| Performance | Quantitative measure we use to evaluate the model's performance. |
|---|---|
| Measures of performance | Metrics we use to evaluate our model |
| Feature Performance | Determination of feature importance |
| Ability to generalize to unseen data | Overfitting, underfitting, bias, variance |

# Decision Tree

# Decision tree: model cheat sheet

## Pros

- Mimics human intuition closely (we make decisions in the same way!)
- Not *prescriptive* (i.e., decision trees do not assume a normal distribution)
- Can be intuitively understood and interpreted as a flow chart, or a division of feature space.
- Can handle nonlinear data (no need to transform data)

## Cons

- **Susceptible to overfitting** (poor performance on test set)
- High variance between data sets
- Can become unstable: small variations in the training data result in completely different trees being generated

# Task

Task

What are we predicting?

How does a doctor diagnose her patients based on their symptoms?

# A doctor builds a diagnosis from your symptoms.

I start by asking patients a series of questions about their symptoms.

I use that information to build a diagnosis.

| temp X1 | vomiting X2 | Shaking chills X3 | diagnosis Y | predicted diagnosis Y* |
|---------|-------------|-------------------|-------------|------------------------|
| 39.5°C | Yes | Severe | Flu | Flu |
| 37.8°C | No | Severe | Malaria | Malaria |
| 37.2°C | No | Mild | Flu | Malaria |
| 37.2°C | Yes | None | Flu | Flu |

There are some questions the doctor does not need to ask because she learns it from her environment. (For example, a doctor would know if malaria present in this region.) **A machine, however, would have to be explicitly taught this feature.**

A doctor makes a mental diagnosis path to arrive at her conclusion. She is building a decision tree!

How does the way I think about how to decide the value of the split compare to a machine learning algorithm?

Does the patient have a temperature over 37C?

Y

vomiting

N

Severe shaking chills

Y

Is malaria present in the region

N

flu

Y

malaria

N

flu

Y

malaria

N

flu

# Human Intuition

# Decision Tree

Based upon my experience as a doctor, I know there are certain questions whose answers quickly separate flu from malaria.

At each split, we can determine the best question to maximize the number of observations correctly grouped under the right category.

- Both a doctor and decision tree try to arrive at the minimum number of questions (splits) that needs to be asked to correctly diagnose the patient.
- Key difference:  how the order of the questions and the split value are determined. A doctor will do this based upon **experience**, a decision tree will **formalize** this concept using a loss function.

# Machine learning adds the power of data to our doctor's task.

To determine whether a new patient has malaria, a doctor uses her experience and education, and machine learning creates a model using data.

By using a machine learning model, our doctor is able to use both her experience **and** data to make the best decision!

**Task**

Decision trees act in the same manner as human categorization, they (split) the data based upon answers to questions.

Where do I go?

Malaria    No malaria    Malaria

Mr. Model creates a flow chart (the model) by separating our entire dataset into **distinct categories**. Once we have this model, we'll know what category our **new patient** falls into!

To get a sense of how this "splitting" works, let's play a guessing game. **I am thinking of something that is blue.**

You may ask 10 questions to guess what it is.

The first few questions you would probably ask may include:
- Is it alive? (No)
- Is it in nature? (Yes)

Then, as you got closer to the answer, your questions would become more specific.
- Is it solid or liquid? (Liquid)
- Is it the **ocean**? (Yes!)

This process is your mental decision tree. In fact, this strategy captures many of the qualities of our algorithm!

Some winning strategies include:

- Use questions early on that eliminate the most possibilities.
- Questions often start broad and become more granular.

You created your own decision tree **based on your own experience of what you know is blue** to make an educated guess as to what I was thinking of (**the ocean**).

Similarly, a decision tree model will make an educated guess, but instead of using experience, it uses **data**.

Now, let's build a formal vocabulary for discussing decision trees.

Decision Tree · Defining f(x)

Like a linear regression, a decision tree has explanatory features and an outcome. Our f(x) is the decision path from the top of the tree to the final outcome.

outcome

model

feature

model

feature

outcome

| Decision Tree | New vocabulary |

| Split | The "decisions" of the decision tree model. The model decides which features to use to split the data. |
| Root node | Our starting point. The first split occurs here along the feature that the model decides is **most important**. |
| Internal node | Intermediate splits. Corresponds to a subset of the entire dataset. |
| Terminal node | Predicted outcome. Corresponds to a smaller subset of the entire dataset. |

model

Root node

Split

feature

Internal node

outcome

Terminal node

Decision Tree

New vocabulary

A decision tree splits the feature space (the dataset) along features. The order of splits are determined by the algorithm.

Feature 1

Feature 2

Feature 3

Root node

Splits

This visualization of splits in feature space is a different but equally accurate way to conceptualize decision trees as the flow chart in the previous slide.

Splits at the top of the tree have an outsized importance on final outcome.

Root node → Does the patient have a temperature over 37C?

Split →

What is the root node, split, internal node and terminal node in this example?

Y

vomiting

N

Severe shaking chills

Y

Is malaria present in the region

N

flu

Y

malaria

N

flu

Internal node

Y

malaria

N

flu

Terminal node →

Let's look at another question to combine our intuition with formal vocabulary. **What should Sam do this weekend?**

Let's get this weekend started!

$Y$ = Choice of weekend activity

Dancing
Cooking dinner at home
Eating at fancy restaurant
Music concert
Walk in the park

Training data set: You have a dataset of what he has done every weekend for the last two years.

The decision tree f(x) predicts the value of a target variable by learning simple decision rules inferred from the data features.

In this example, we predict Sam's weekend activity using decision rules trained on historical weekend behavior.

Our most important predictive feature is Sam's budget. How do we know this? Because it is the **root node**.

How much $$ do I have?

<$50

>=$50

Raining?

Girlfriend?

Y

N

Y

N

Movie!

Walk in the park!

Concert!

Dancing!

An example of how the algorithm works:

A decision tree splits the dataset ("**feature space**"). Here, the entire dataset = data from 104 weekends. You can see how each split subsets the data into progressively smaller groups.

Now, when we want to predict what will happen **this weekend**, we can!
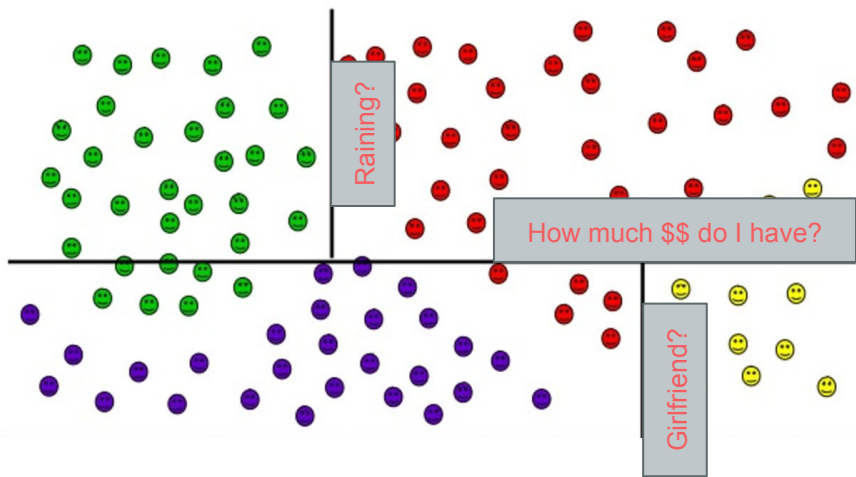


n=104

How much $$ do I have?

< $50          >= $50

n=30    Raining?          n=74    Girlfriend?

Y          N          Y          N

Movie!     Walk in the park!     Concert!     Dancing!

n=17          n=13          n=6          n=68

Decision Tree Task → f(x) as a spatial function

Each of the four weekend options are grouped spatially by our splits.

Visualized in the data:

Raining?

How much $$ do I have?

Girlfriend?

n=104

How much $$ do I have?

< $50          >= $50

n=30    Raining?          Girlfriend?    n=74

Y              N          Y              N

Movie!    Walk in the park!    Concert!    Dancing!

n=17      n=13              n=6         n=68

Now we understand the mechanics of the decision tree task. **But how does a decision tree learn**? How does it know where to split the data, and in what order?

We turn now to examine decision trees' learning methodology.

# Learning Methodology

Remember that the key difference between a doctor and our model is how the <u>order</u> of the questions and the <u>split value</u> are determined.

# Human Intuition

# Decision Tree

Based upon my experience as a doctor, I know there are certain questions whose answers quickly separate flu from malaria.

At each split, we can determine the best question to ask to maximize the number of observations correctly grouped under the right category.

key difference between a doctor and our model is how the <u>order</u> of the questions and the <u>split value</u> are determined.

The two key levers that can be changed to improve accuracy of our medical diagnosis are:
- The order of the questions
- The split value at each node. (For example the temperature boundary)

A doctor will make these decisions based upon experience, a decision tree will set these parameters to minimize our loss function by learning from the data.

# Recap: Loss function

A loss function quantifies how unhappy you would be if you used f(x) to predict Y* when the correct output is y. It is the object we want to minimize.

Linear Regression

Sound familiar? We just went over a similar optimization process for linear regression (our loss function there was MSE).

Remember: All supervised models have a loss function (sometimes also known as the cost function) they must optimize by changing the model parameters to minimize this function.

Parameters of a model are values that Mr. Model controls to minimize our loss function. Mr. Model sets these parameters by learning from the data.

One of the parameters Mr. Model controls are the split values. For example, Mr. Model learns that $50 is the best split to predict whether Sam eats out or stays at home.

How much $$ do I have?

< $50 → Cook at home

>= $50 → Eat at restaurants

**Parameters**

Values that control the behavior of the model. The model learns what parameters are from data.

Learning Methodology → How does our decision tree learn?

We have two decision tree parameters: Split decision rule (value of b) and the ordering of features

How much $$ do I have?

<b → Cook at home

>=b → Girlfriend?

b=0 → Concert!

b=1 → Dancing!

Mr. Model controls B and the ordering of the features.

Mr. Model learns what best fits the data by moving b up or down and by trying many combinations of decision rule ordering.

Central problem: How do we learn the "*best*" split?

Imagine that is a game. We start with a **random** ordering of features and set b to completes random values.

Our **random** initialization of parameters give us a Unsurprisingly high initial error

Mr. Model's job is to change the parameters so that every time he updates the loss goes down.

The game is over when Mr.Model is able to reduce the error to e, the irreducible error.

## Squared Error

# Recap of tasks:

Decision trees also can be used for two types of tasks!

## Regression

Continuous variable

A regression problem is when we are trying to predict a numerical value given some input, such as "dollars" or "weight".

## Classification

Categorical variable

A classification problem is when are trying to predict whether something belongs to a category, such as "red" or "blue" or "disease" and "no disease".

Source: Andrew Ng , Stanford CS229 Machine Learning Course

Our decision tree loss function depends on the type of task.

The most common loss functions for decision trees include:

I can minimize all types of loss functions!

*For classification trees:*
1. Gini Impurity
2. **Entropy**

*For regression trees:*
1. Mean Squared Error

For example, the loss function for a regression decision tree should feel familiar. It is the same loss function we used for linear regression!

For classification trees, we can use information gain!

How would you rank the entropy of these circles?



A

B

C

Source

Information gain attempts to minimize **entropy** in each subnode. Entropy is defined as a degree of disorganization.

**If the sample is completely homogeneous, then the entropy is zero.**

If the sample is an equally divided (50% – 50%), it has entropy of one.
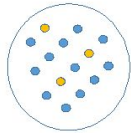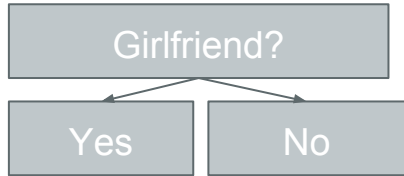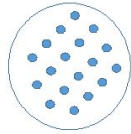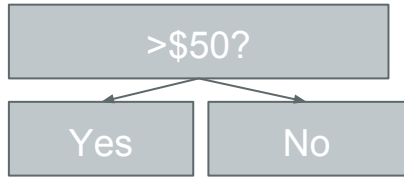
So, how does it work in our decision tree model?

I'm going to randomly generate three splits...

Split

entropy

>$50?

Yes    No

Girlfriend?

Yes    No

Raining?

Yes    No

The **Information Gain** algorithm:

1. Calculates entropy of parent node
2. Calculates entropy of each individual node of split, and calculates weighted average of all sub-nodes available in split.

*The algorithm then chooses the split that has the lowest entropy.*

**Regression trees use Mean Squared Error.**

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \widehat{Y}_i)^2$$

| | |
|---|---|
| $Y - Y_*$ | For every point in our dataset, measure the difference between true Y and predicted Y. |
| $^\wedge 2$ | Square each Y-Y* to get the absolute distance, so positive values don't cancel out negative ones when we sum. |
| Sum | Sum across all observations so we get the total error. |
| mean | Divide the sum by the number of observations we have. |

*This may look familiar - look back at our discussion of linear regression! Decision trees also use MSE.*

The split with lower MSE is selected as the criteria to split the population.

Regression trees use Mean Squared Error.

I'm going to randomly generate three splits…

Split | RMSE

Amount of $
Yes | No
3.83

# of girlfriends
Yes | No
7.28

Likelihood of rain
Yes | No
10.3

The RMSE algorithm proceeds by:
1. Calculating variance for each node.
2. Calculating variance for each split as **weighted average** of each node variance.

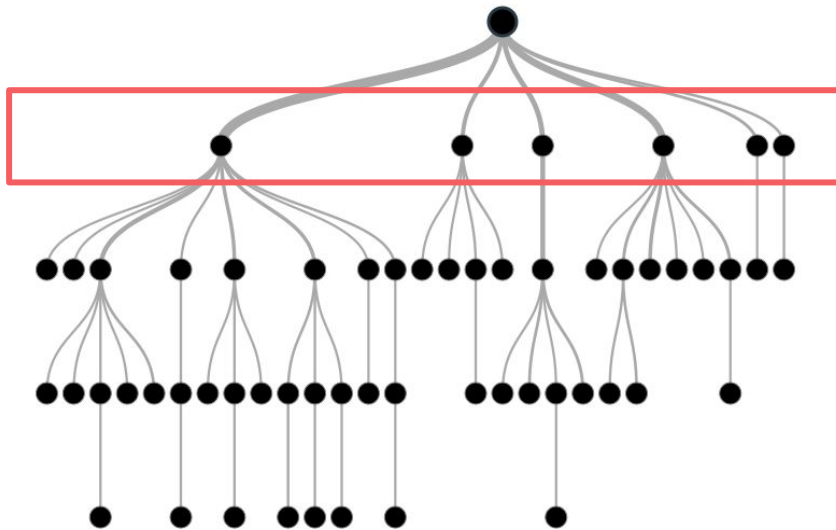*The algorithm then selects the split that has the lowest variance.*

# Model Performance

Decision trees provide us with understanding of what features are most important for predicting Y*

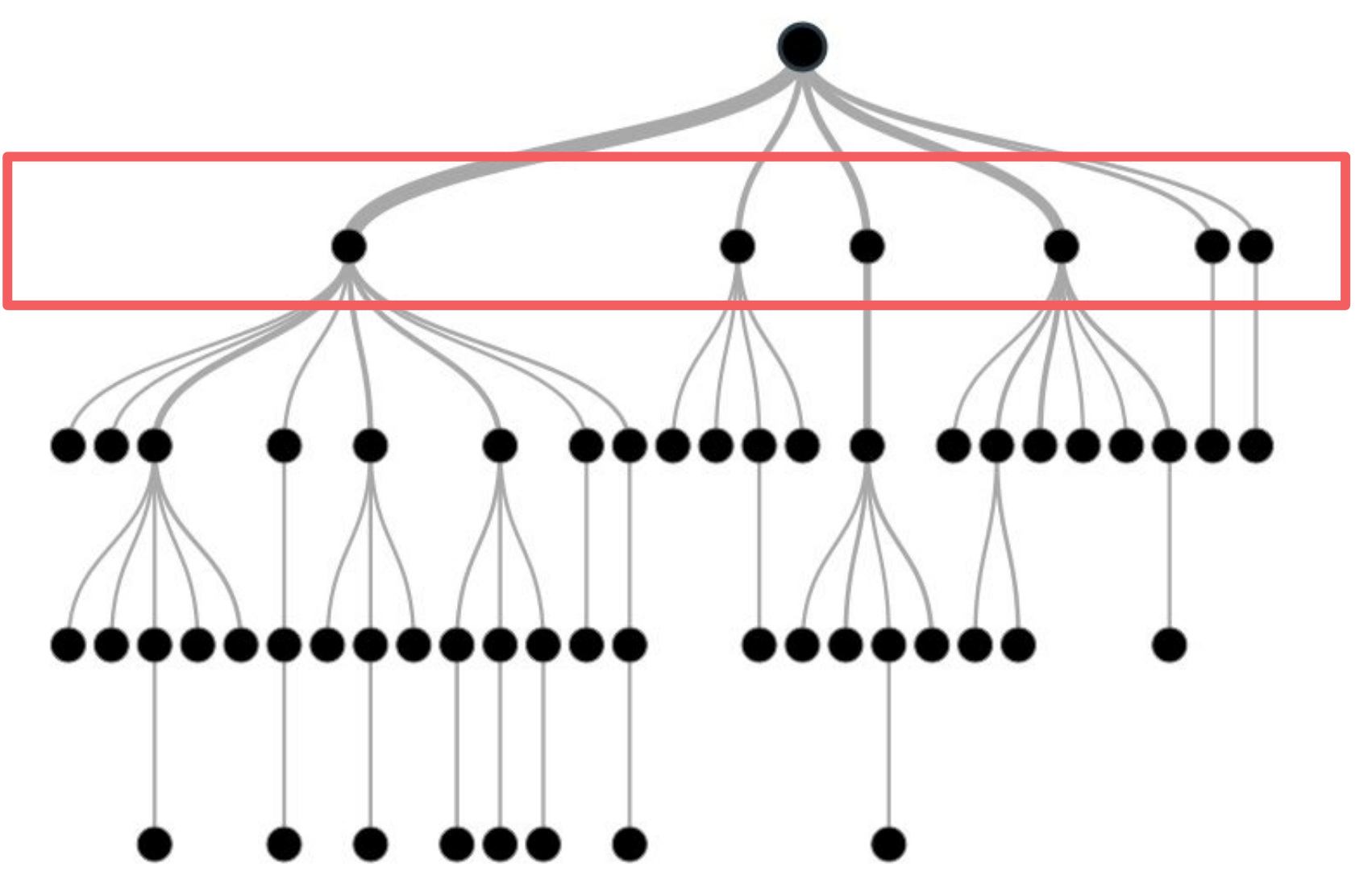

*The intuition is provided by the understanding that the most important splits are at the first nodes.*

Recall our intuition: Important splits happen closer to the root node.

*The decision tree tells us what the important splits are (which features cause splits that reduce cost function the* ***most****).*
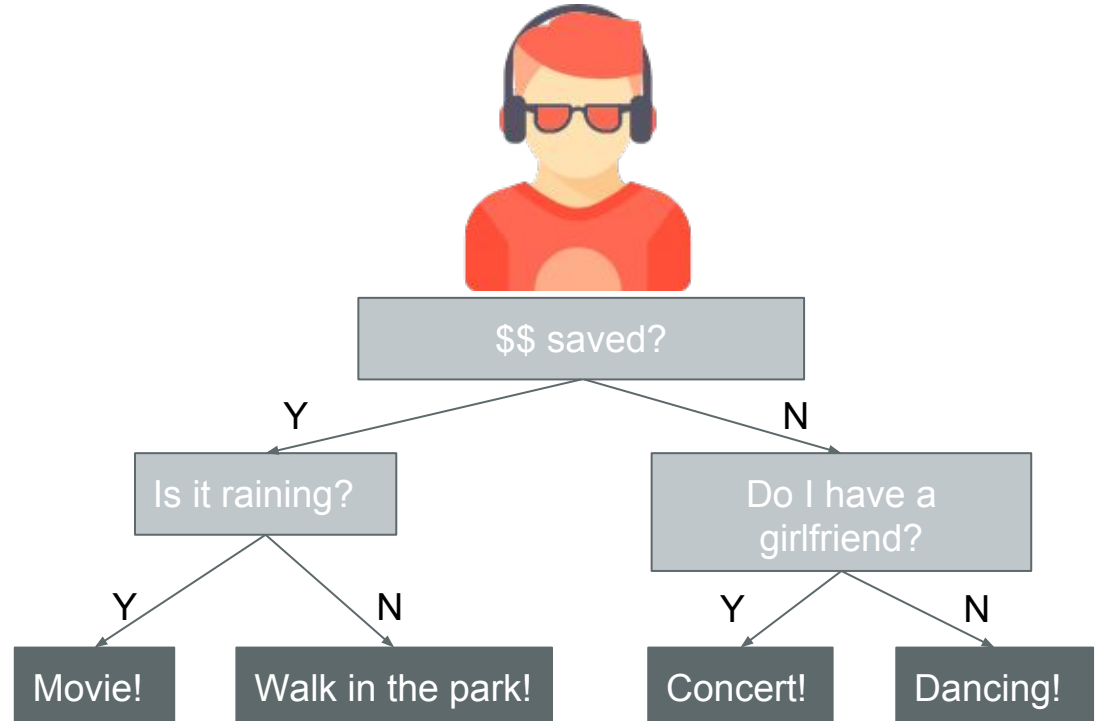
Recall our intuition: Important splits happen earlier

In Sam's case, our algorithm determined that Sam's budget was the most important feature in predicting his weekend plans.
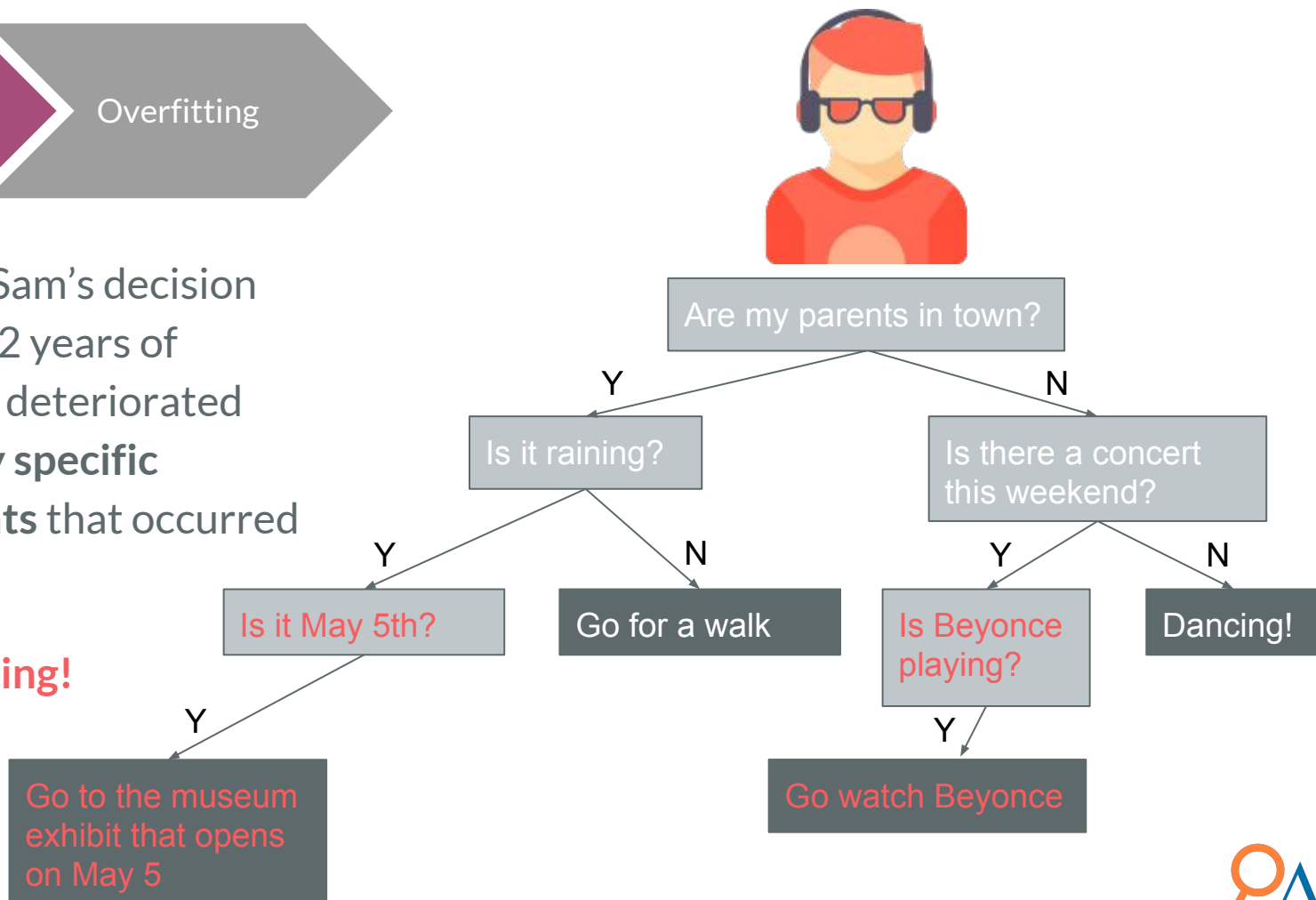
This makes a lot of sense!

$$ saved?

Y — Is it raining?

N — Do I have a girlfriend?

Is it raining?
Y — Movie!
N — Walk in the park!

Do I have a girlfriend?
Y — Concert!
N — Dancing!

Imagine if the Sam's decision tree, based on 2 years of weekend data, deteriorated into **extremely specific weekend events** that occurred in the data.
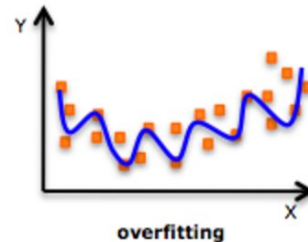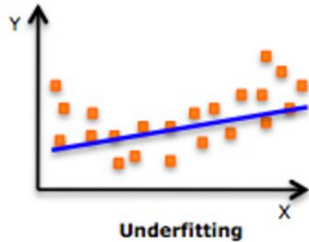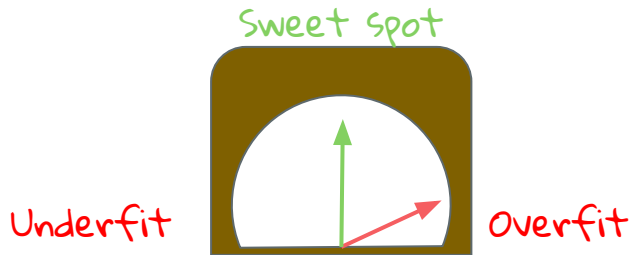
**This is overfitting!**

Are my parents in town?

Y — Is it raining?

N — Is there a concert this weekend?

Is it raining?
Y — Is it May 5th?
N — Go for a walk

Is it May 5th?
Y — Go to the museum exhibit that opens on May 5

Is there a concert this weekend?
Y — Is Beyonce playing?
N — Dancing!

Is Beyonce playing?
Y — Go watch Beyonce

Remember that the most important goal we have is to build a model that will generalize well to unseen data.



Underfitting · Just right! · overfitting

If our train data set overfits, it will not generalize to our test set (unseen data) well.

However, if it underfits, we are not capturing the complexity of the relationship and our loss will be high!
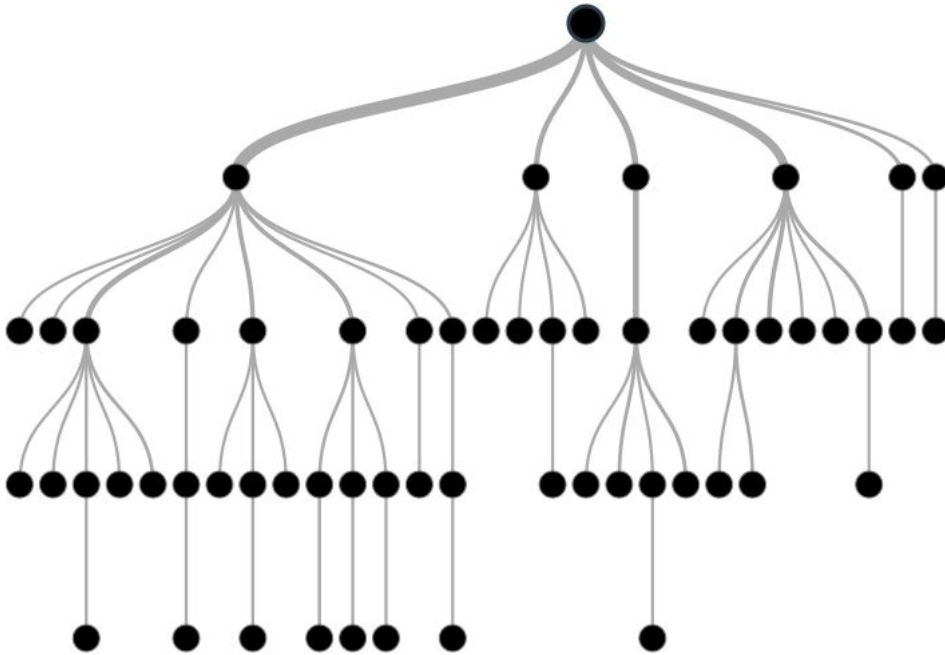


Sweet spot

Underfit · Overfit

Our goal in evaluating performance is to find a sweet spot between overfitting and underfitting.

Performance → Ability to generalize to unseen data

Unfortunately, decision trees are very prone to overfitting. We have to be very careful in evaluating this.

Each time we add a node, we fit additional models to subsets of the data . This means we start to get to know our train data really well but it will impair our ability to generalize to our test data.

Overfitting: If each terminal node is an individual observation, it is overfit.

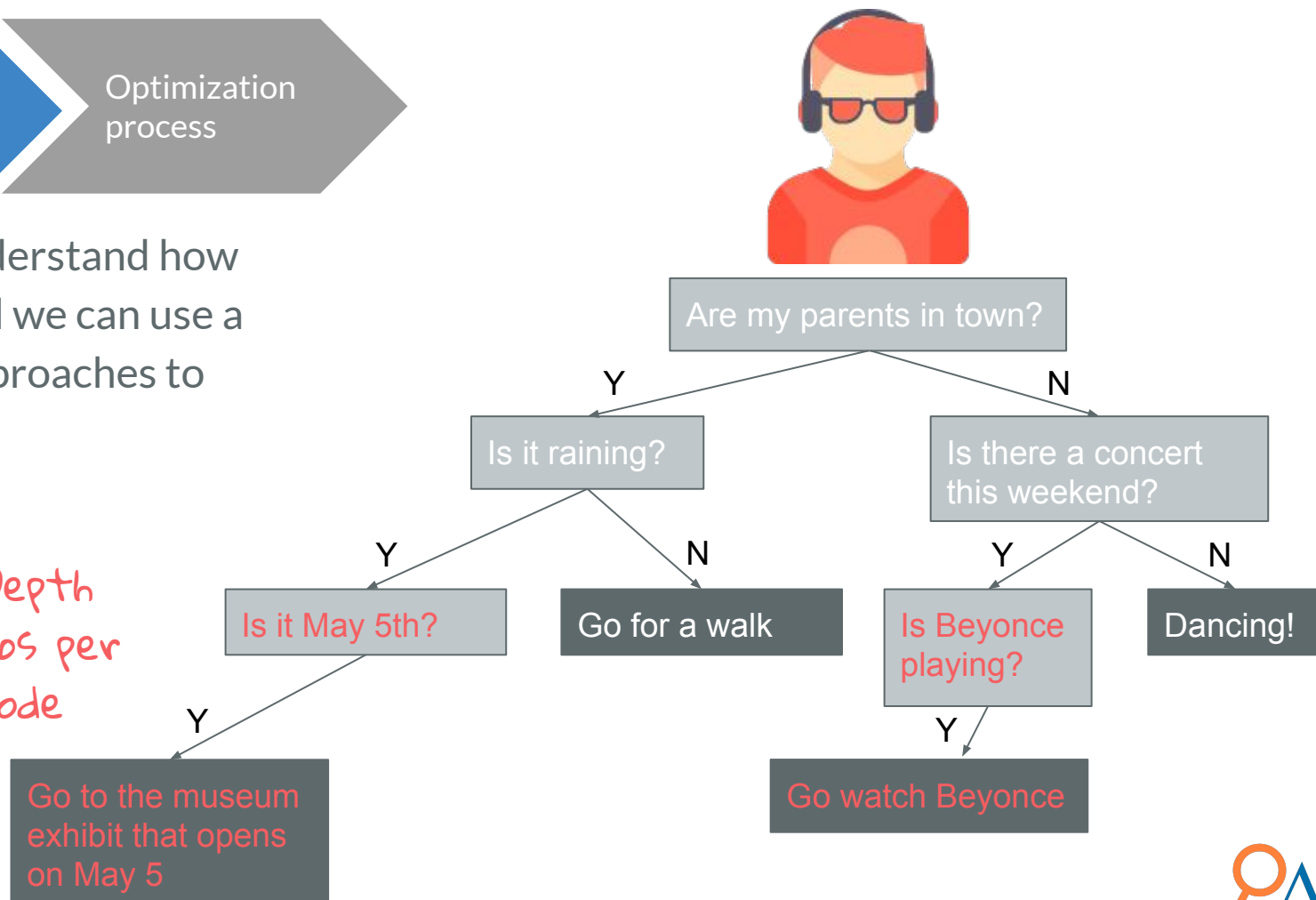Let's take a look at a concrete example of optimizing our model!

# Model Optimization

Now that we understand how overfitting is bad we can use a few different approaches to avoid it:

- Pruning
- Maximum Depth
- Minimum obs per terminal node

Are my parents in town?

Y

N

Is it raining?

Is there a concert this weekend?

Y

N

Y

N

Is it May 5th?

Go for a walk

Is Beyonce playing?

Dancing!

Y

Go to the museum exhibit that opens on May 5

Y

Go watch Beyonce

Pruning, max depth and n_ obs in a terminal node are all examples of hyperparameters.

Hyperparameters are higher level settings of a model that are fixed before training begins.

**Pruning, max depth and n_ obs in a terminal node are all decision tree hyperparameters set before training.**

Their values are **not learned** from the data so Mr. Model cannot say anything about them.

You, not Mr.Model decide what the hyperparameters are!

- In linear regression, the coefficients were **parameters.**

$$Y = a + b_1x_1 + b_2x_2 + \ldots + e$$     The model decides!

  - Parameters are **learned** from the training data using the chosen algorithm.

- Hyperparameters cannot be learned from the training process. They express "higher-level" properties of the model such as its **complexity** or **how fast it should learn**.

We decide!

Source: Quora

Learning Methodology → Decision tree Hyperparameters

Pruning → Top-down

Bottom-up

Maximum depth

Minimum samples per leaf
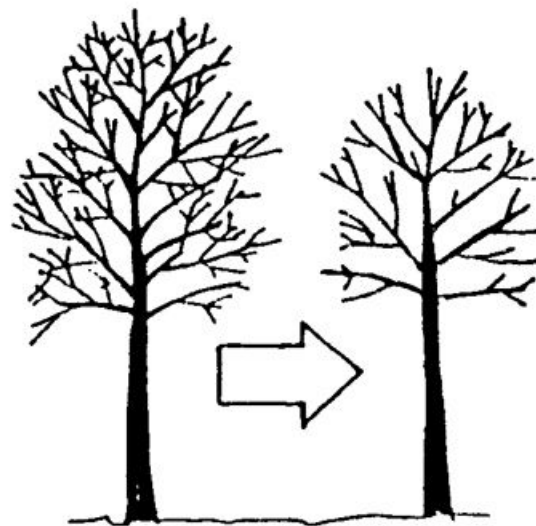
Maximum number of terminal nodes

Pruning reduces the number of splits in the decision tree.

Limiting the number of splits minimizes the problem of **overfitting**.

Two approaches:
- Pre-pruning (top-down)
- Post-pruning (bottom-up)

Read more: Notre Dame's Data Mining class, CSE 40647

Pre-pruning is a top down approach where the tree is limited in depth before it fully grows.

# Pre-pruning slows the algorithm before it becomes a fully grown tree. This prevents irrelevant splits.

- E.g. In the malaria diagnosis example, it probably wouldn't make sense to include questions about a person's favorite color.
    - It might cause a split, but it likely wouldn't be a meaningful split. May lead to overfitting

A useful analogy in nature is a bonsai tree. This is a tree whose growth is slowed starting from when it's a sapling.

Post-pruning is a bottom up approach where the tree is limited in depth after it fully grows.

*Grow the full tree, then prune by merging terminal nodes together.*

1. Split data into training and validation set
2. Using the decision tree yielded from the training set, merge two terminal nodes together
3. Calculate error of tree with merged nodes and tree without merged nodes
4. If tree with merged nodes has lower error, merge leaf nodes and repeat 2-4.

A useful analogy in nature is pruned bushes. Bushes are grown to their full potential, then are cut and shaped.
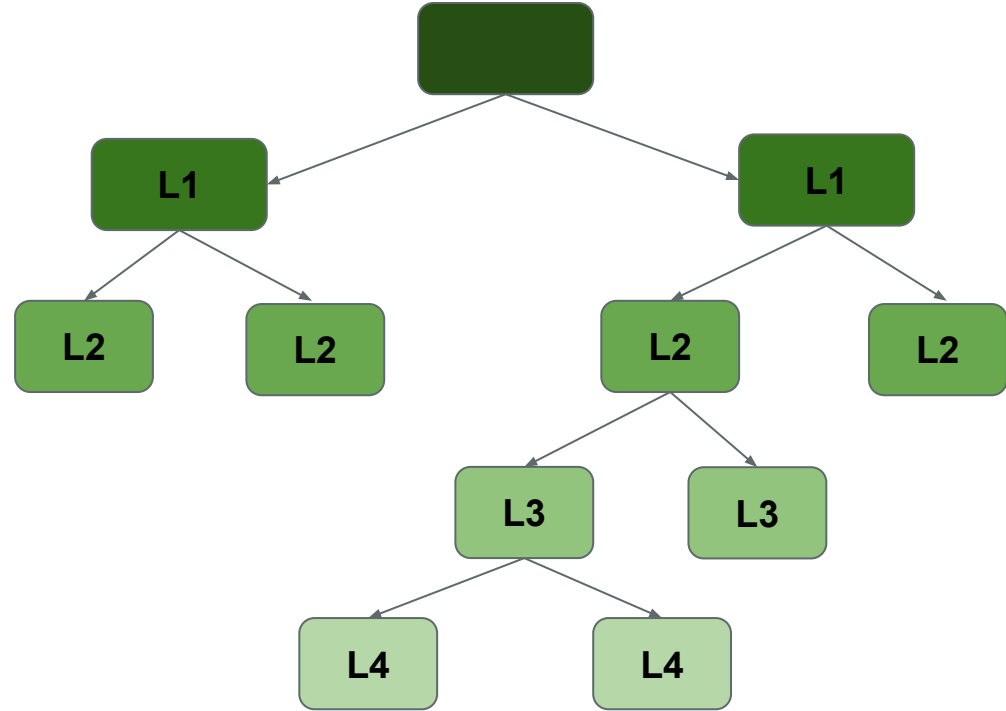
Maximum depth defines the maximum number of layers a tree can have.

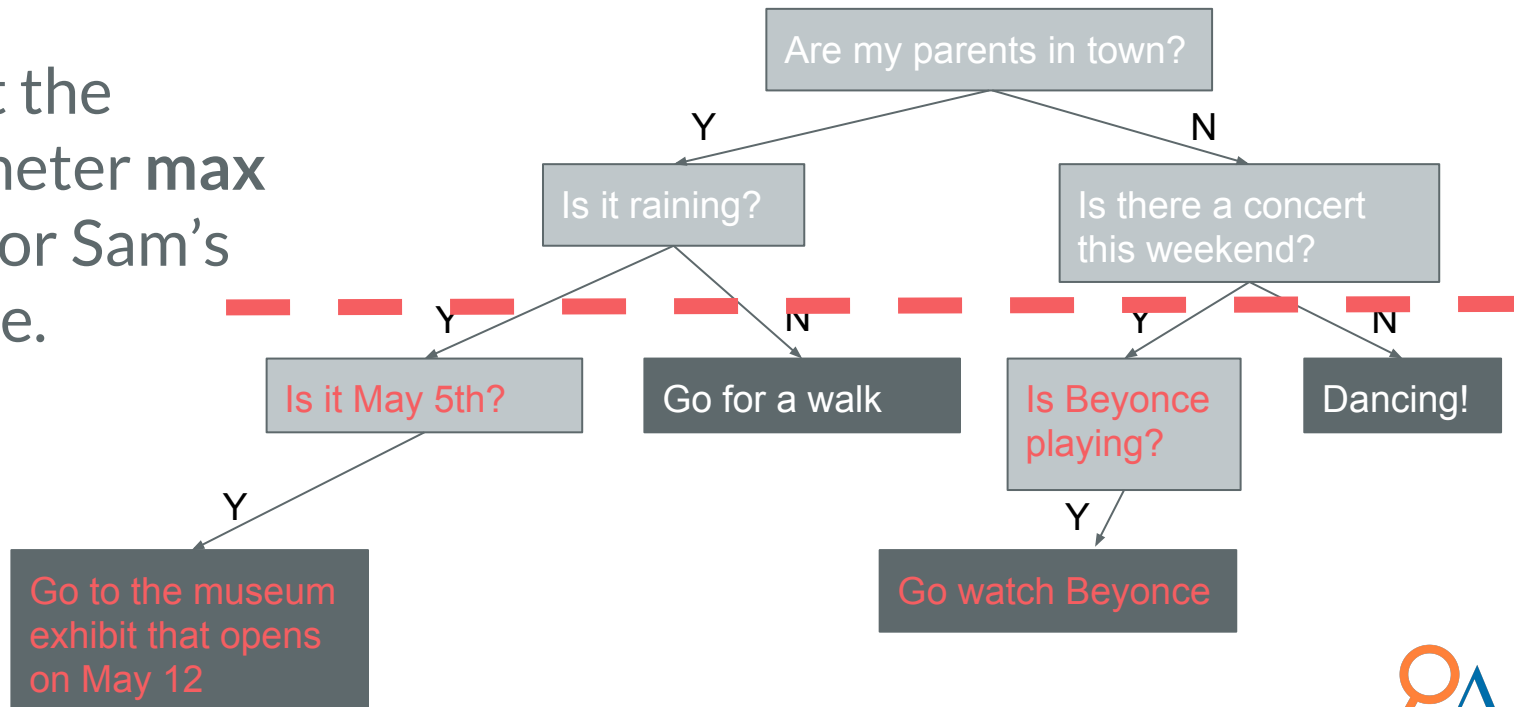For example, a maximum depth of 4 means a tree can be split between 1 to 4 times.



Maximum depth reached, no more splitting

Here we set the hyperparameter **max depth** to 2 for Sam's decision tree.

Are my parents in town?

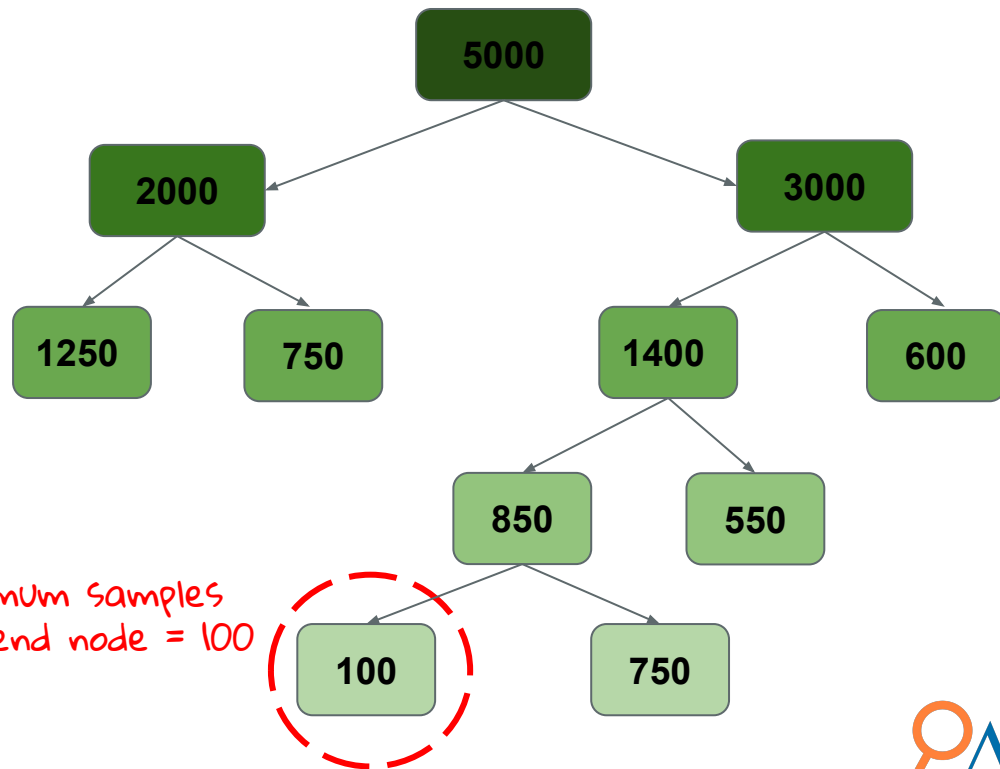Y

N

Is it raining?

Is there a concert this weekend?

Y

N

Y

N

Is it May 5th?

Go for a walk

Is Beyonce playing?

Dancing!

Y

Y

Go to the museum exhibit that opens on May 12

Go watch Beyonce

# Minimum observations per end node

This hyperparameter establishes a minimum number of observations in an end node, which reduces the possibility of modelling data noise.
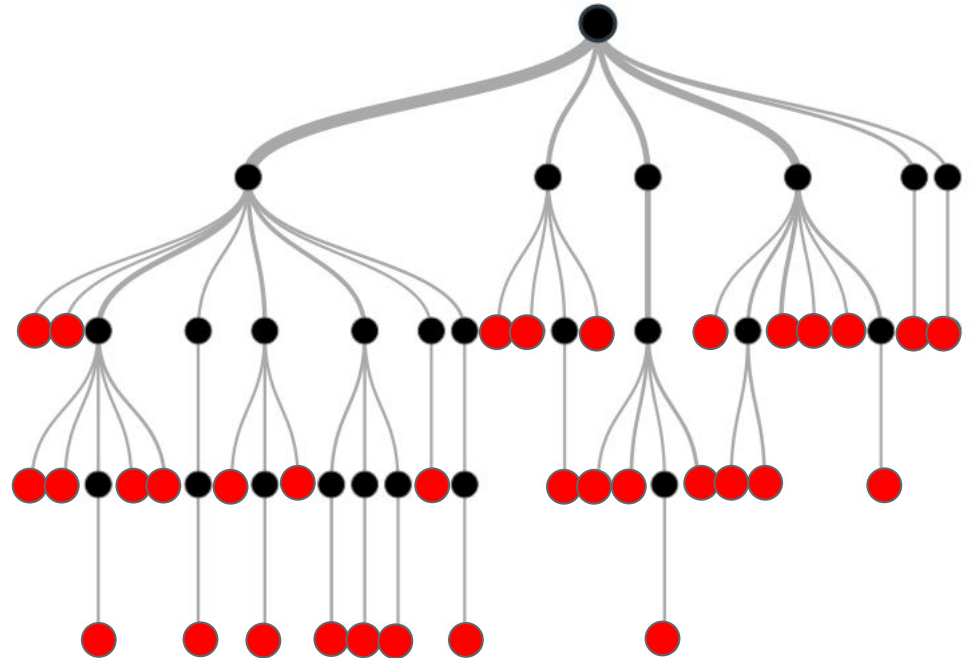


Minimum samples per end node = 100

Source

Maximum number of terminal nodes limits the number of branches in our tree.

The maximum terminal nodes limits the number of branches in our tree. Again, this limits overfitting and reduces the possibility of modelling noise.
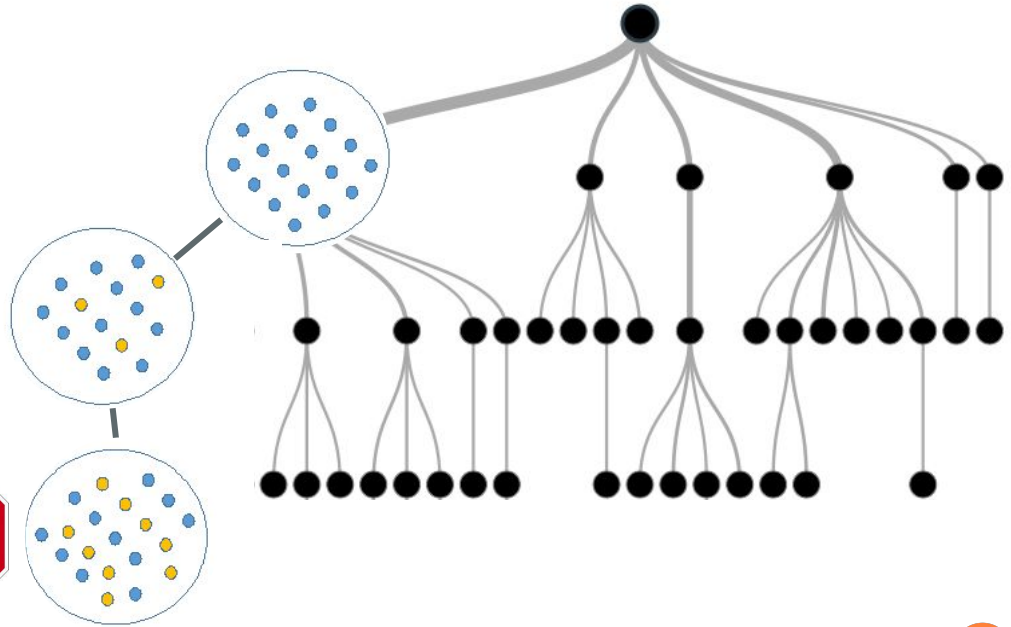
● = terminal nodes

# Minimum impurity

- Recall the Information Gain cost function
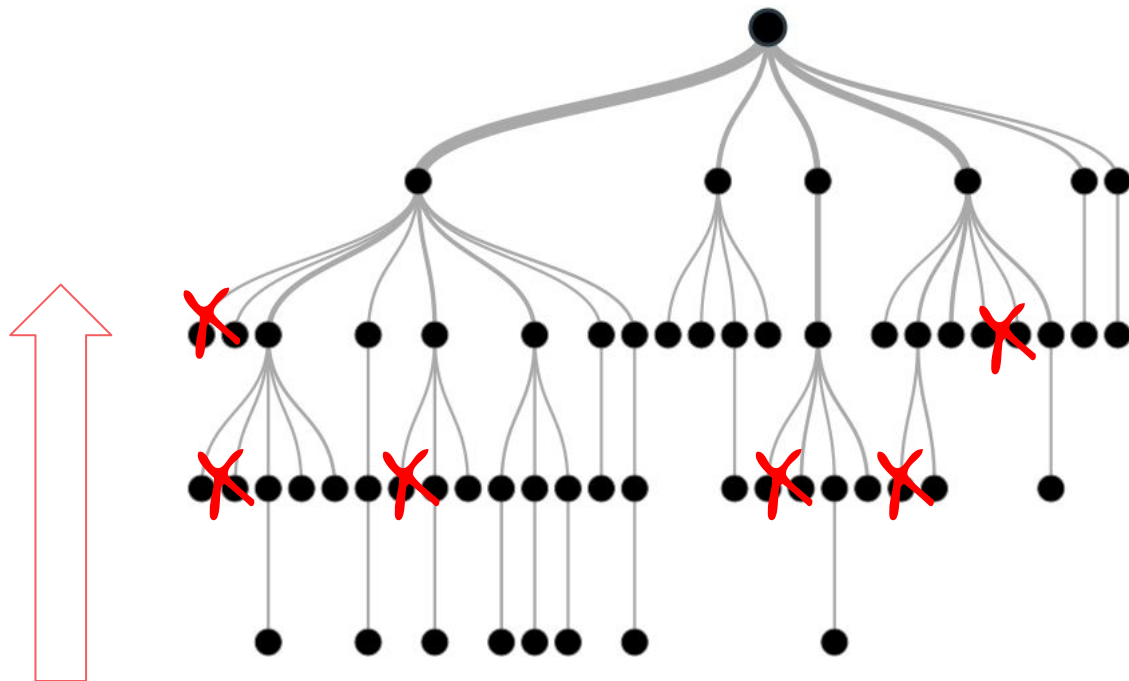- The model iterates until it's reached a certain level of impurity

# Which hyperparameter should I use?

No objectively best hyperparameter for each model. **Use trial and error** - see how each changes your results!

# Bottom-up pruning

Selectively merges some terminal nodes together

# Bottom-up pruning

Let's see how a single merge works:

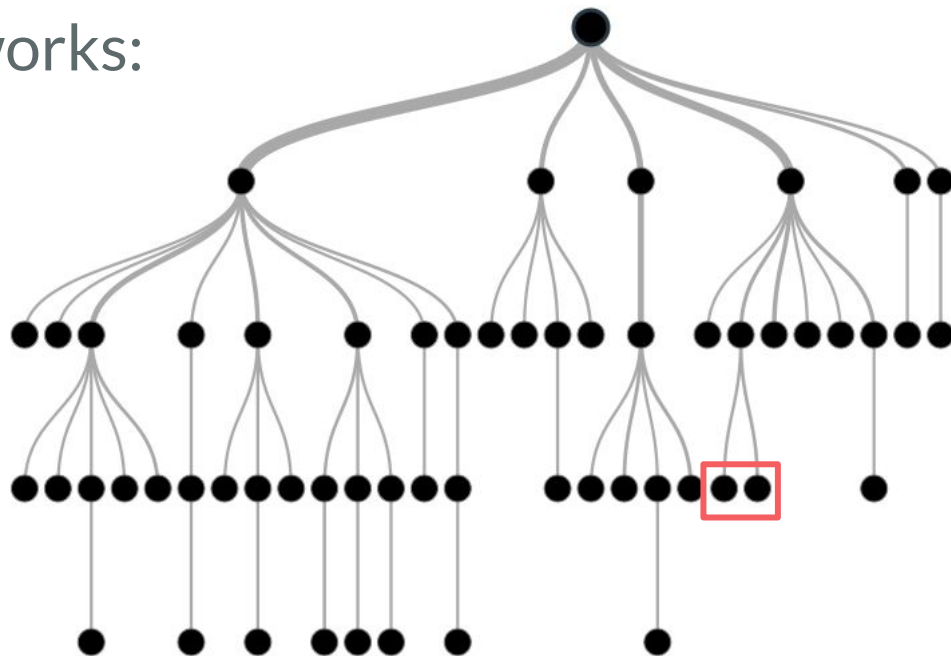The model randomly selects the merge marked in the tree.

The model calculates:
Error of tree without merge = a
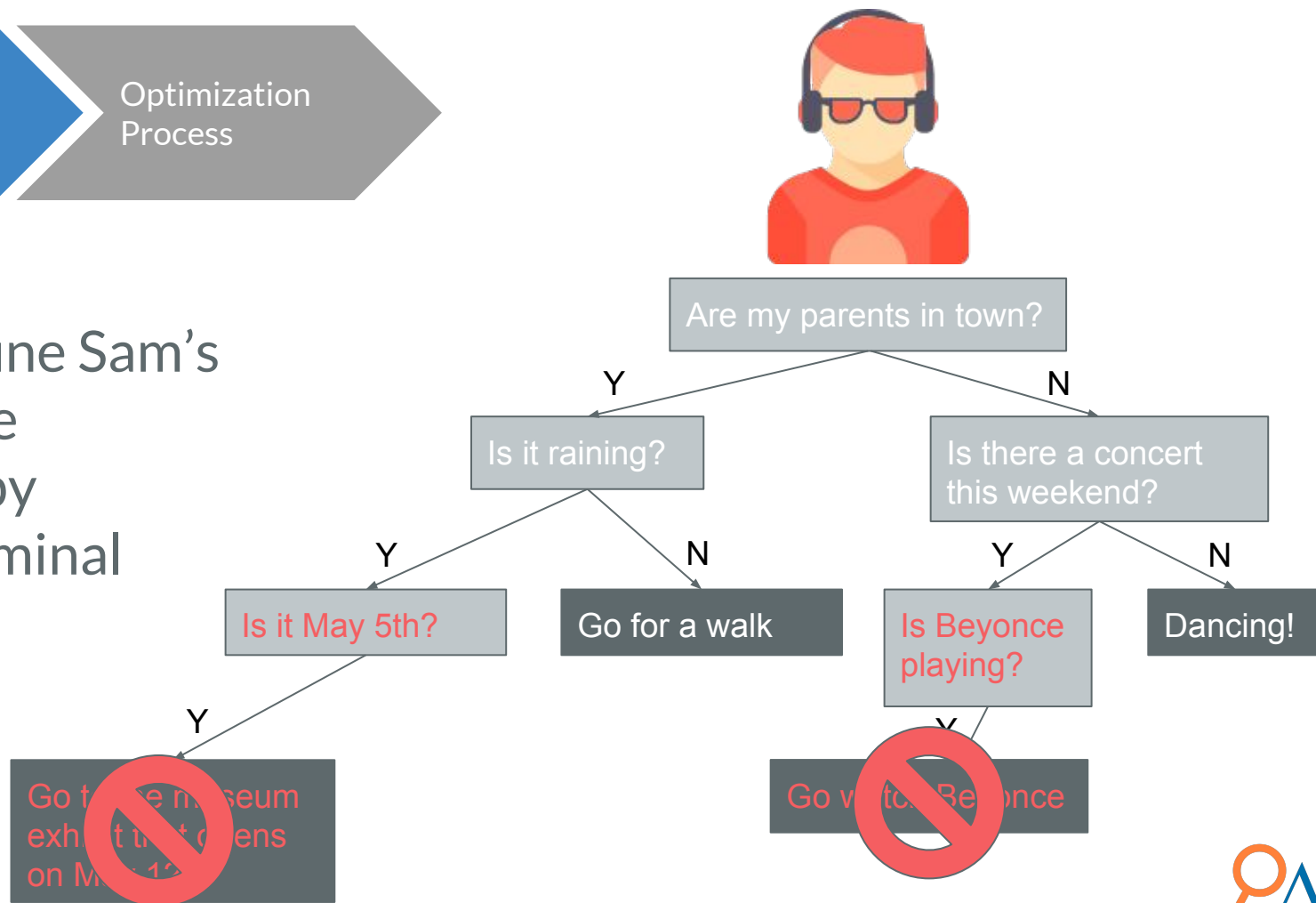Error of tree with merge = b

If b < a, it chooses to merge!

... And repeat

Here we prune Sam's decision tree bottom-up by merging terminal nodes.

Are my parents in town?

Y

N

Is it raining?

Is there a concert this weekend?

Y

N

Y

N

Is it May 5th?

Go for a walk

Is Beyonce playing?

Dancing!

Y

Y

Go to the museum exhibit that opens on May 1?

Go watch Beyonce

# Learning Methodology

Learning Methodology

How does our ML model learn?

What is our loss function?

Optimization process

# Learning Methodology

| Learning Methodology | Supervised Learning |
|---|---|
| How does our ML model learn? | Iteration! Evaluating different splits, deciding which one's best |
| What is our loss function? | Gini Index, Information gain, RMSE |
| Optimization process | Pruning |

# End of theory

# Checklist:

✓    Decision trees
      ✓    Intuition
      ✓    The "best" split
      ✓    Model performance
      ✓    Model optimization (pruning)

# Advanced resources

# Want to take this further? Here are some resources we recommend:

- Textbooks
  - An Introduction to Statistical Learning with Applications in R (James, Witten, Hastie and Tibshirani): Chapter 8
  - The Elements of Statistical Learning: Data Mining, Inference, and Prediction (Hastie, Tibshirani, Friedman): Chapter 9
- Online resources
  - Analytics Vidhya's guide to understanding tree-based methods