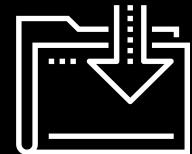


# Introduction to Machine Learning

Fintech  
Lesson 10.1



# Class Objectives

---

By the end of this lesson, you will be able to:



Recognize the differences between supervised and unsupervised machine learning (ML).



Define clustering and how it is used in finance.



Apply the K-means algorithm to identify clusters in a given dataset.



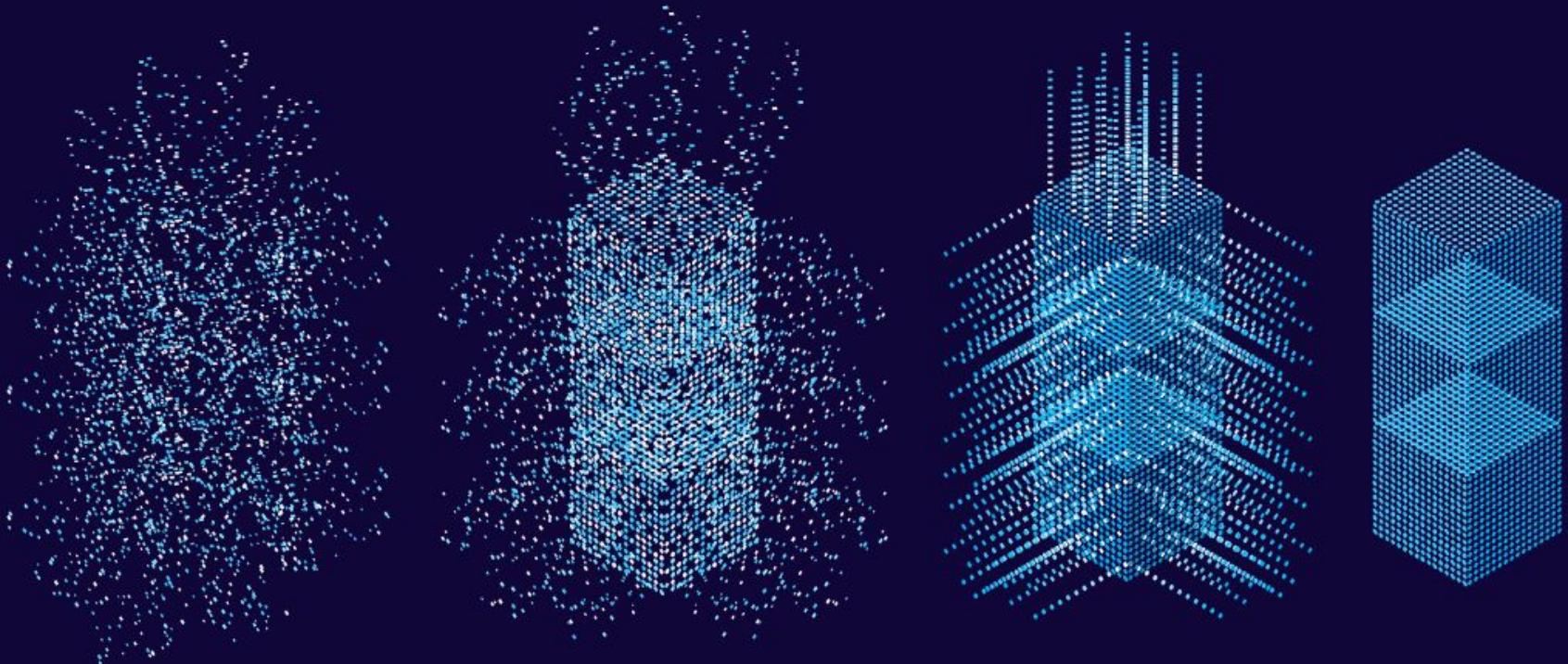
Use the elbow method to determine the optimal number of clusters for a dataset.



**WELCOME**

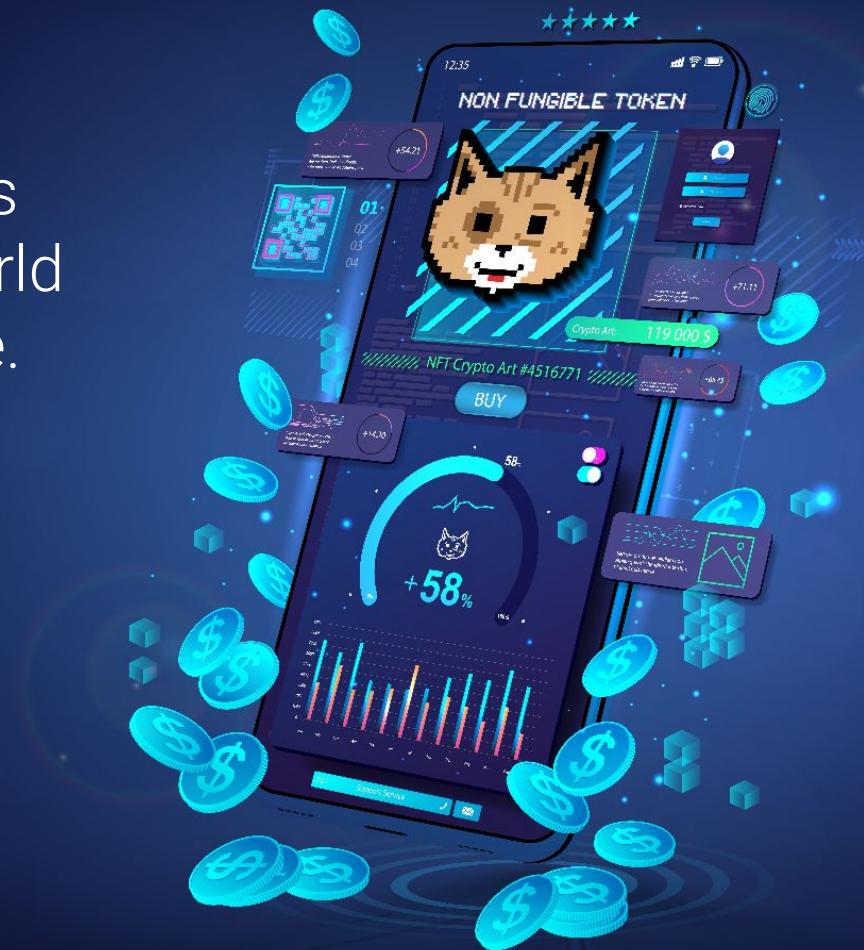
# Introduction to Machine Learning in Finance

Machine learning is a type of artificial intelligence that uses algorithms and statistical models to make decisions or predictions about data.



Today, machine learning is changing the financial world at an unprecedented pace.

Changes are being driven by automation power, which enables decisions to be made more quickly and efficiently than ever before.





## This Week's Challenge Assignment

---

Create a machine learning model that groups cryptocurrencies to assemble investment portfolios that are based on the profitability of those cryptocurrencies.

# The Mysticism of Machine Learning

---

Despite the mainstream use of the term “machine learning,” most people still don’t know what machine learning *really* is.



**Machine learning** is the practice of applying computer algorithms and statistics to create models that can learn from data and then make decisions or predictions about future data.

# Machine Learning

---

Algorithms learn how to make decisions without needing anyone to program all that logic.

They learn the patterns, behavior, and logic on their own directly from the data, and then they use that knowledge to make decisions and predictions.





Here's an example of how machine learning can be useful...

# Machine Learning

---

Imagine that you work as a fraud analyst in a bank, and you want to identify fraudulent transactions.

## Option 1

Create a 5,000-line `if-else` decision structure that evaluates every price range and product category to determine if a transaction counts as fraudulent.

## Option 2

Use machine learning algorithms to review all of the transactions that an account owner has ever made.

Then, you can group the transactions and predict whether the most recent transaction counts as fraudulent.



This is the kind of machine learning solution that you'll learn to build!



# Why is machine learning essential for fintech?

# Machine Learning in Fintech

---



Applications for machine learning vary widely—from the textual analysis of fraudulent contracts to the determination of interest rates.



Machine learning applications have streamlined many operational processes that are associated with finance.



Incorporating machine learning has helped the finance industry increase its responsiveness to customer demands.

# Machine Learning in Fintech

---

As AI gets smarter and machine learning models become more entrenched in the finance industry, understanding how they work will become a necessary skill for all fintech professionals.

**...and you'll have a head start!**



# ML Application: Financial Advising

---

One intriguing application of machine learning: replacing the traditional financial advisor with artificial intelligence (AI).  
For example, the following fintech companies automate retirement saving:



# ML Application: Financial Advising

---

Using machine learning to recommend tailored portfolios and manage the financial advising process means that customers no longer need to pay high annual fees for traditional advisors.

Cost-reducing automation also extends beyond products oriented toward end users (like individuals saving for retirement)...

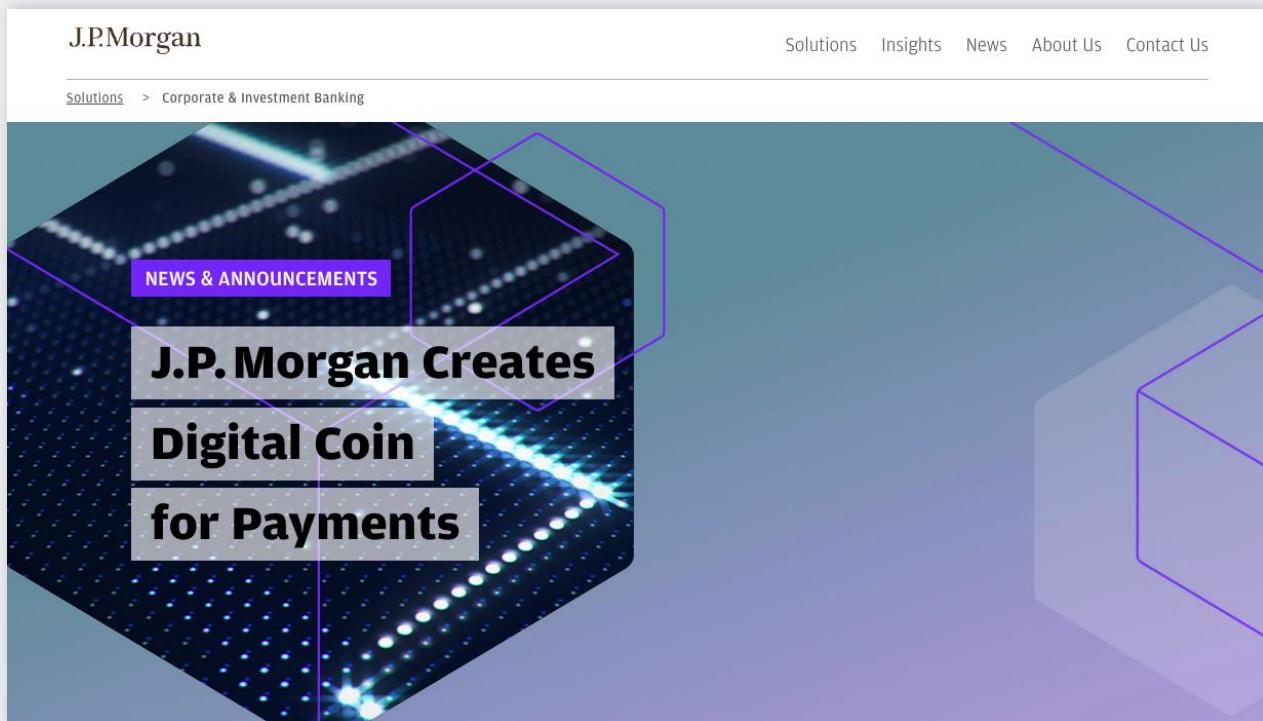


# ML Application: Financial Advising

J.P. Morgan, for example, has developed a program named **Onyx**.

Onyx uses natural language processing (NLP) to automatically run due diligence on its commercial-credit contract agreements.

These lines of code compress the 360,000 hours of work that humans previously did—into mere seconds.



# ML Application: Forecasting Market Results

Machine learning models are used to forecast financial-market results.

These forecasts range from loan evaluation (for example, predicting the default rates on bonds or consumer loans) to high-frequency algorithmic trading in the stock market.



# ML Application: Forecasting Market Results

Algorithmic trading driven by statistical models has been around since the late 1990s.

- One example is the quant hedge fund Two Sigma Securities, which automatically trades more than 300 million shares per day (*about one share per day for every person in the United States*).
- The fund has been in business since the early 2000s.
- Within the last five years, algorithmic trading has grown to account for 75% to 80% of all equity trades in the United States.



The screenshot shows a magazine-style article from **Institutional Investor**. The header includes navigation links: Portfolio, Corner Office, Culture, Premium, Research, Video, and Innovation. A search bar and user options (Search, Subscribe, Sign In, Register) are also present. The main title is **Inside the Geeky, Quirky, and Wildly Successful World of Quant Shop Two Sigma**, written by Stephen Taub on June 28, 2019. Below the title is a caricature illustration of two men, David Siegel and John Overdeck, smiling. The text below the illustration discusses their achievement of the Lifetime Achievement Award and mentions Paul Tudor Jones and Tom Hill.

**CORNERS OFFICE**

## Inside the Geeky, Quirky, and Wildly Successful World of Quant Shop Two Sigma

Co-founders David Siegel and John Overdeck have earned Institutional Investor's Lifetime Achievement Award for hedge fund management this year. Paul Tudor Jones and Tom Hill dish on the \$60 billion quant powerhouse's early days.

By Stephen Taub June 28, 2019

Illustration by II

# ML Application: Forecasting Market Results

Algorithmic trading requires many kinds of decisions, all of which are made using machine learning models.

More cryptocurrency exchanges and trading firms are seeking individuals who can write machine learning code that trades profitably.



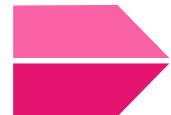
# Additional ML Applications in Fintech

---

There are many applications of machine learning in fintech, including:



Identifying money laundering and legal violations



Recognizing satellites for the real-time awareness of trading opportunities in commodities



Predicting customer churn in financial products



Even the venture capital and private equity industries have begun to experiment with machine learning models to predict the likelihood of start-up success.

An aspect of machine learning called **natural language processing (NLP)** is helping to redefine the investment process.

NLP models are used to pore over social media, news feeds, and annual report documents in search of specific words or phrases.

This can help predict the future direction of a company's stock price, enabling a trading position ahead of the market's move.





What are some examples of  
machine learning models that  
you've heard of?

# Types of ML

---

Examples include:



Regression



Clustering



Neural networks



Deep learning

# Types of ML

---

We can group all of these models into two main buckets:

01

## Supervised learning

The algorithm learns on a **labeled dataset**, where each example in the dataset is tagged with the answer.

This provides an answer key that can be used to evaluate the accuracy of the training data.

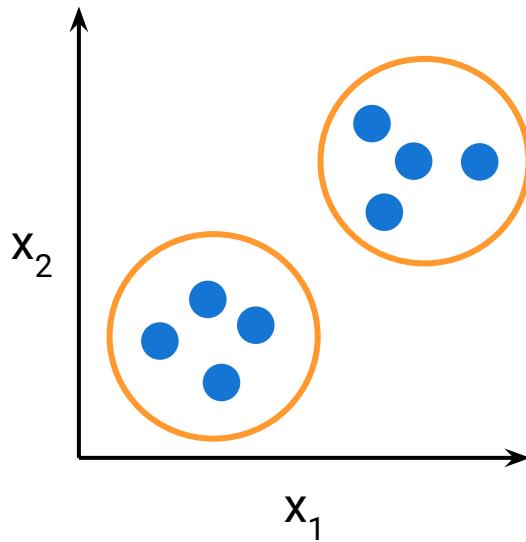
02

## Unsupervised learning

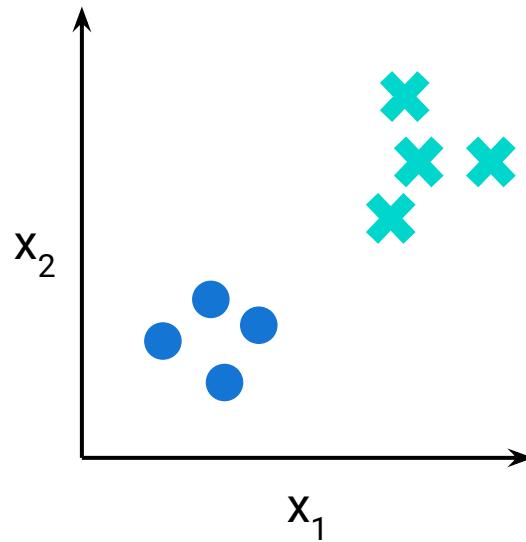
The algorithm tries to make sense of an **unlabeled dataset** by extracting features and patterns on its own.

# Supervised Learning vs. Unsupervised Learning

Supervised Learning



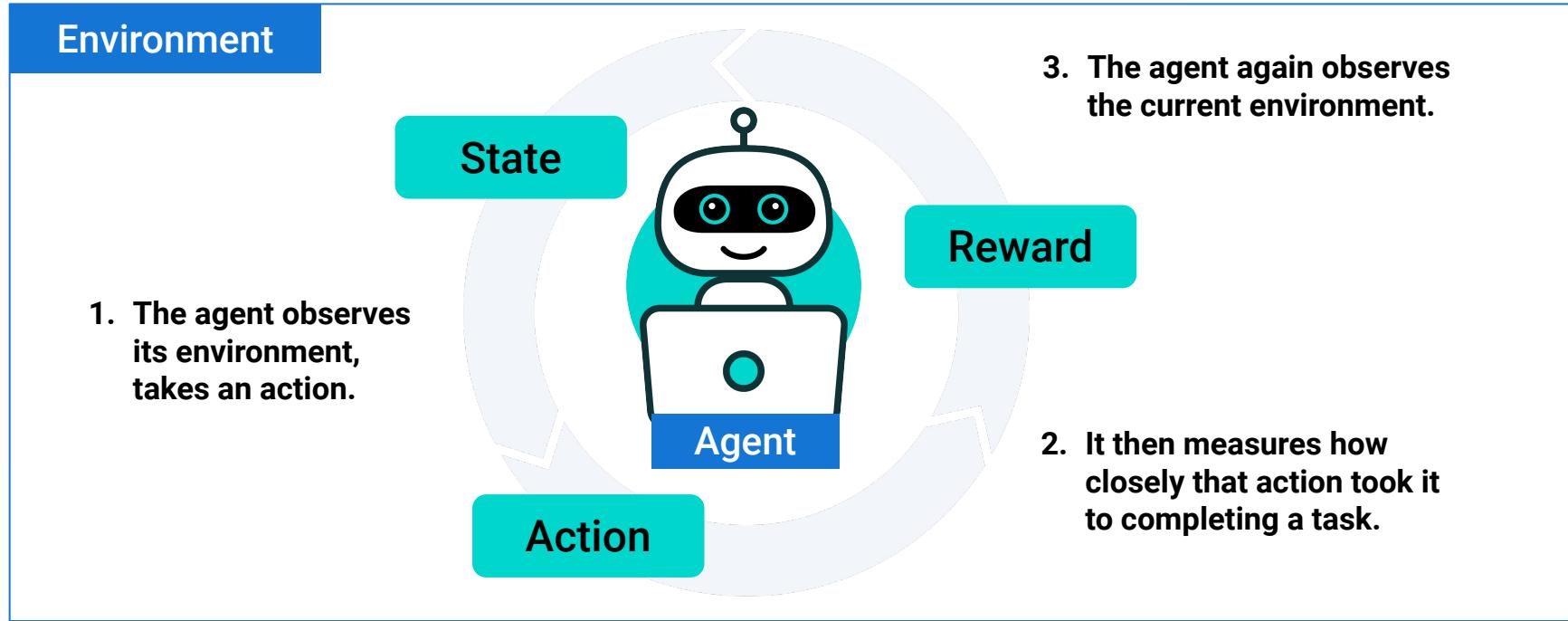
Unsupervised Learning



vs.

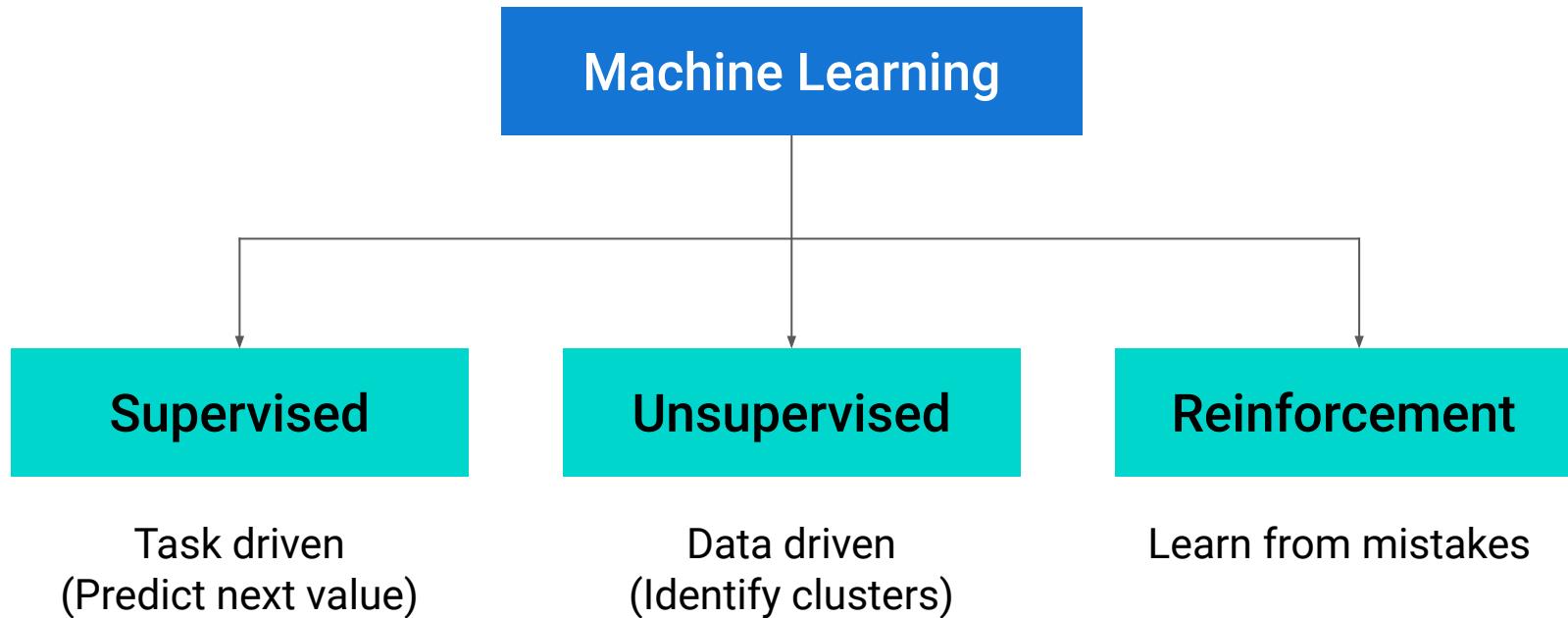
# Reinforcement Learning

This third type of machine learning algorithm is used less frequently but still has important applications in finance.

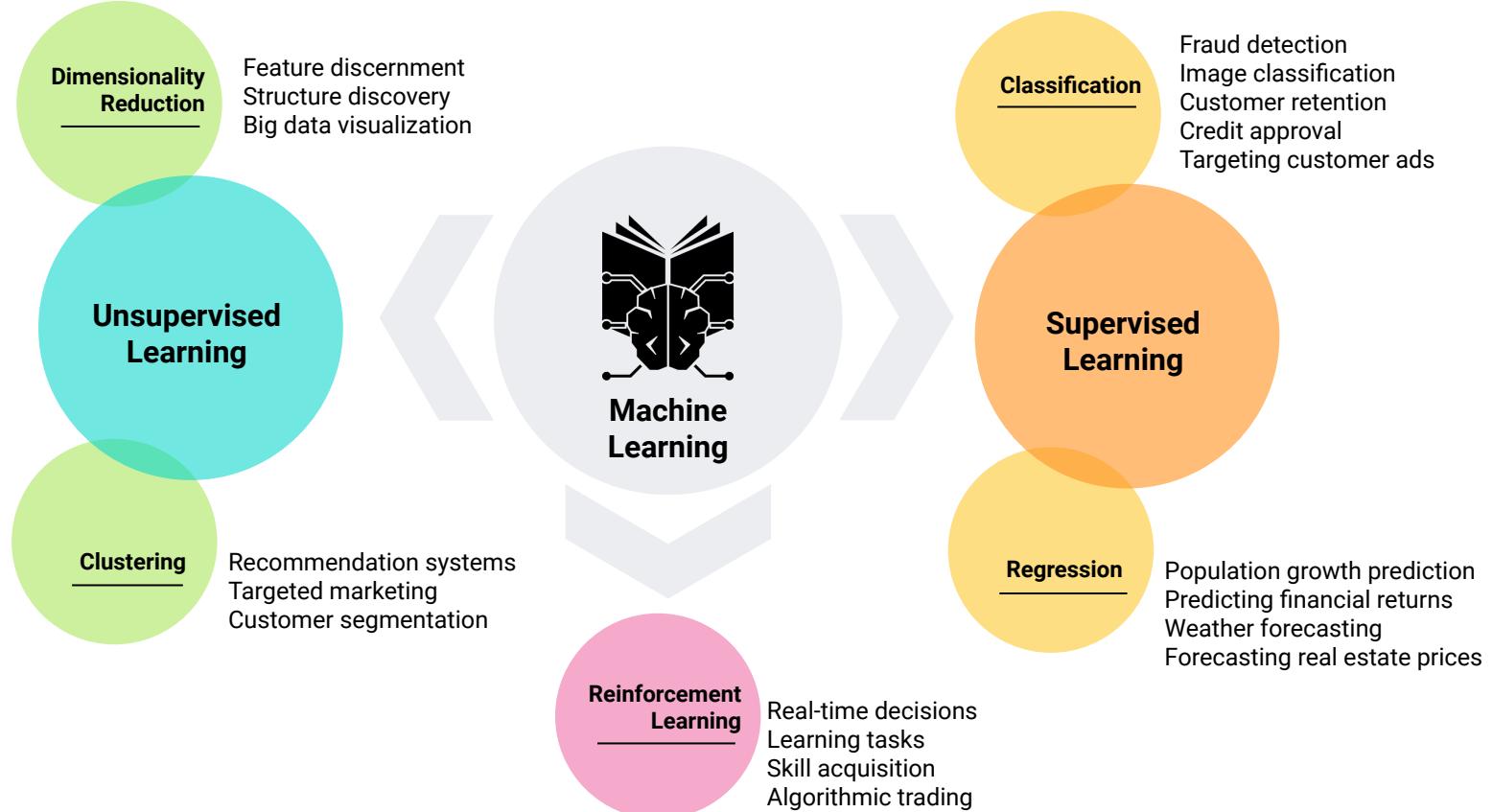


# Three Types of ML

---



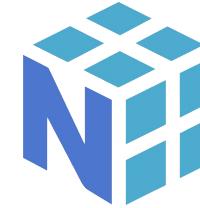
# Types of ML



# Types of ML

---

Most Python libraries for machine learning use a common interface to build and use machine learning models.



# Questions?





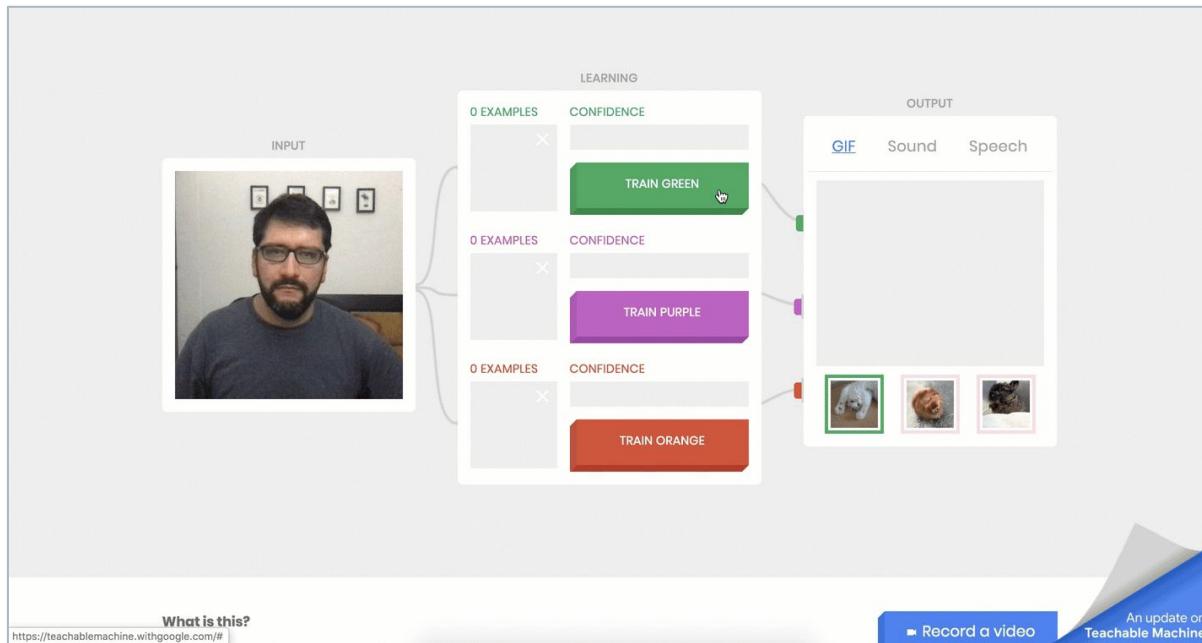
## Instructor Demonstration

---

Machine Learning Is Awesome

# Teachable Machine in Action

The [Teachable Machine project from Google](#) shows the fundamental mechanism of a neural network by training a model that recognizes gestures from your webcam to predict one of three classes.



# Introduction to Unsupervised Learning

**Unsupervised learning algorithms** use test data to construct models that categorize relationships among data points.

# Introduction to Unsupervised Learning

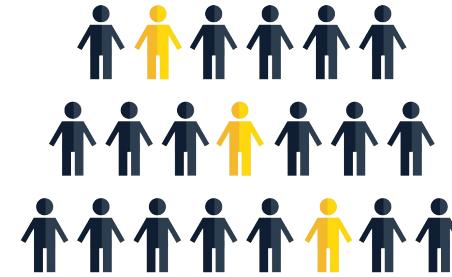
For example, when you're reviewing a particular item for purchase on a website, unsupervised learning algorithms might be used to identify related items that are frequently bought together.

The screenshot shows a product page for the Fenix PD35 TAC LED Flashlight. The main image is a black tactical flashlight. Below it, the product name "PD35 TAC" and "1000 LUMENS" is displayed. To the left is a "Frequently Bought Together" section showing four items: a Fenix ARB-L18-3500 battery, a Fenix ARE-X1 Charging Kit, and a Fenix AER-02 Remote Pressure Switch. An "ADD ALL TO CART" button is visible at the bottom of this section. The page also includes a "Total Price: \$131.80" and a "Save" button.



This power to recognize data patterns has broad applications in finance.

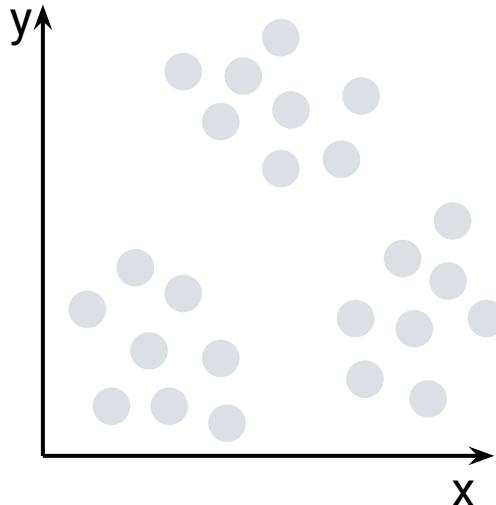
Unsupervised learning can be used to **identify clusters**, or related groups, of clients to target with product offerings or marketing campaigns.



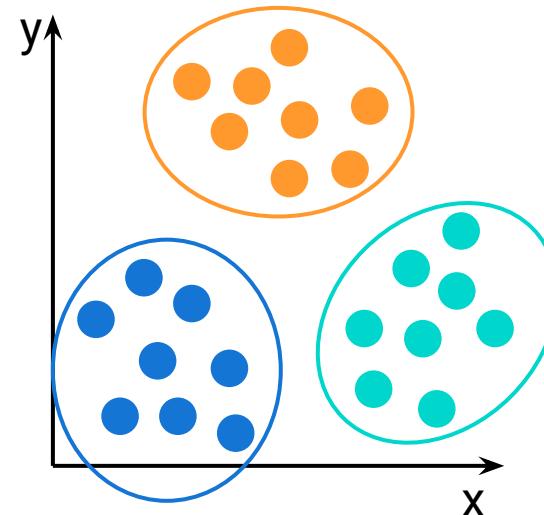
# Introduction to Unsupervised Learning

The **K-means algorithm** is used for marketing use cases because of its ability to segment customers for financial benefits.

**Before K-means**



**After K-means**



# Introduction to Unsupervised Learning

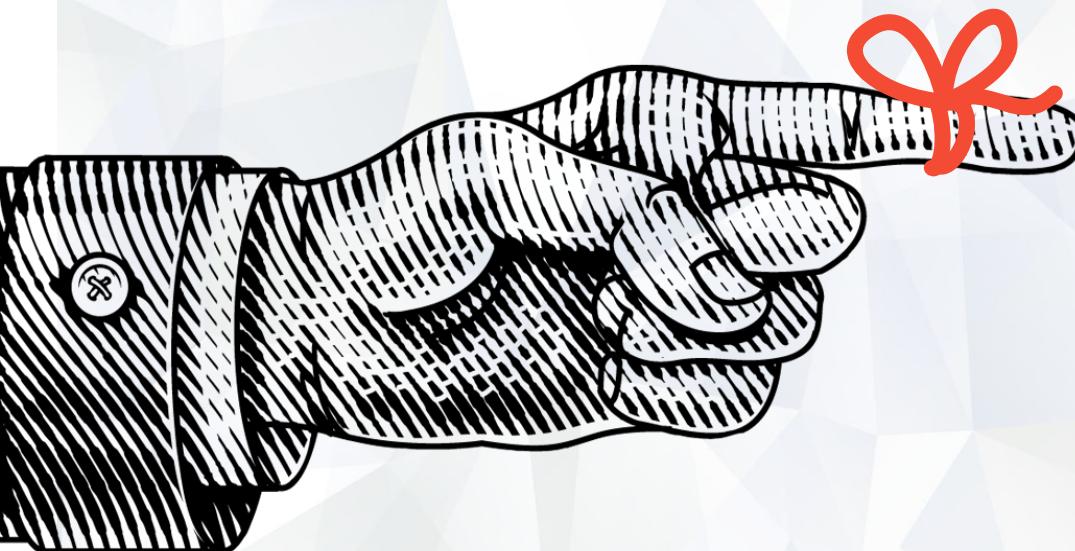
---

## In the lesson:

You'll apply unsupervised learning by using the K-means algorithm to define customer segments, or clusters.

## In the Challenge assignment:

You'll apply the K-means algorithm to cluster data and prepare a cryptocurrency portfolio proposal to the company's board of directors.



*Remember,*

the two most frequently used  
methods of machine learning  
are **supervised learning** and  
**unsupervised learning**.

# Supervised vs. Unsupervised Learning

---

Supervised Learning	Unsupervised Learning
Input data is labeled.	Input data is unlabeled.
Uses training datasets.	Uses input datasets.
<b>Goal:</b> Predict a class or value.	<b>Goal:</b> Determine patterns or group data, called data clusters.

# Challenges of Unsupervised Learning

---

Unsupervised learning comes with challenges:



Because the data isn't labeled, we don't know if the output is correct.



The algorithm is creating its own categories for the data, so an expert is needed to determine if these categories are meaningful.



Even with challenges, unsupervised learning can be useful for a variety of fintech applications, including the following customer segmentation tasks:

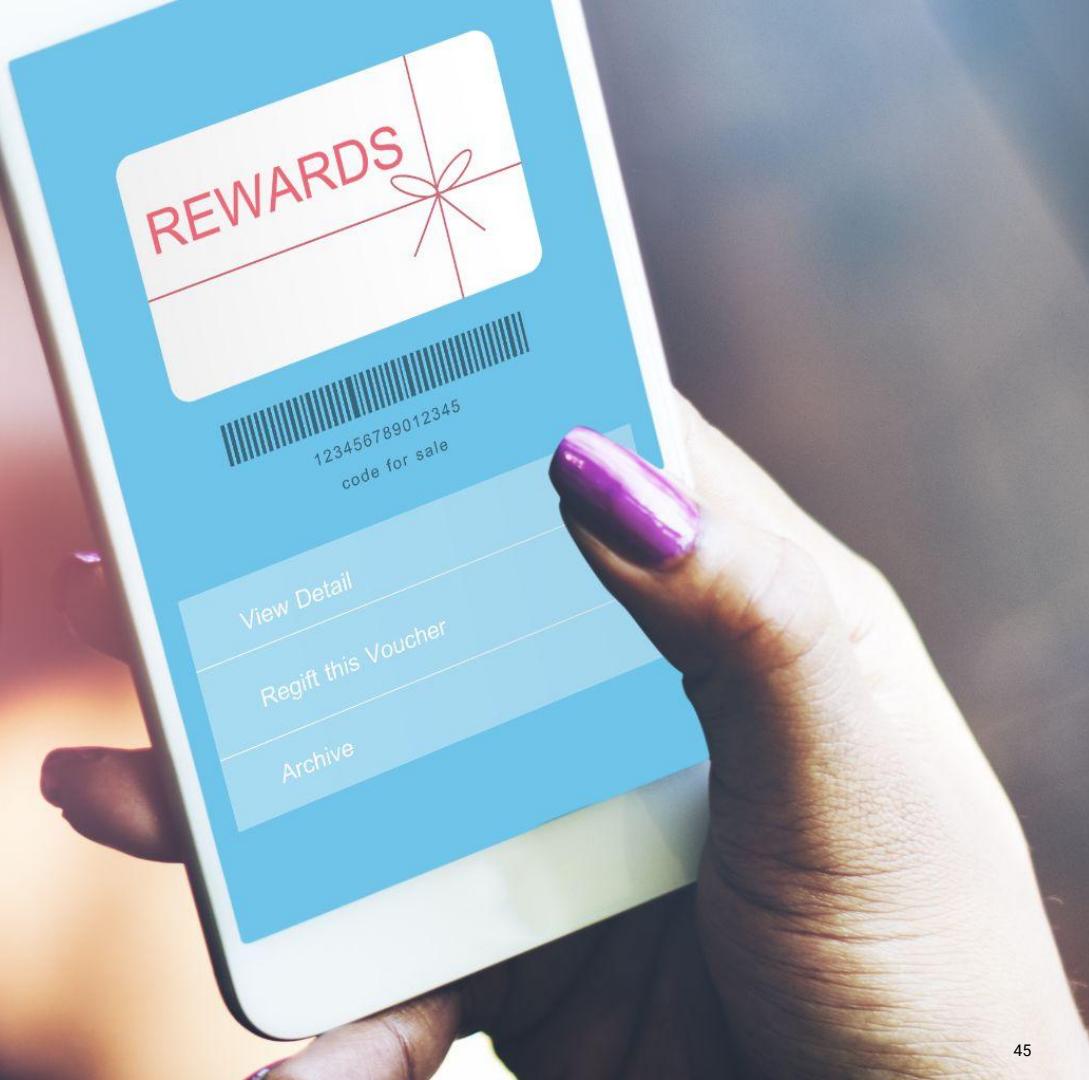
- Grouping customers by spending habits
- Finding fraudulent credit card charges
- Identifying unusual data points (outliers) within the dataset



How might clustering be used  
by fintech businesses?

## One possible answer:

Clustering can be used to group customers by spending habits and create customized offers via email or mobile apps.





# How might anomaly detection be used by credit card companies?



## One possible answer:

Anomaly detection can be used to detect potentially fraudulent credit card transactions by grouping transactions into “normal” or “abnormal.”

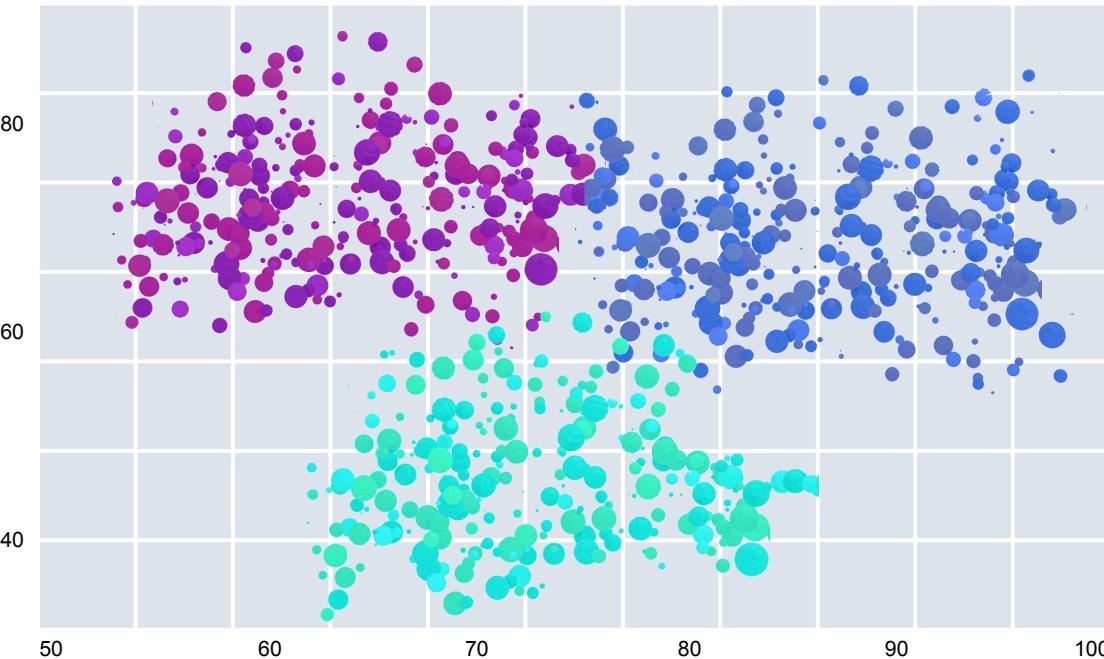
# Clustering Explained

**Clustering** is grouping data together so that every member of that group is similar in some way.

# Clustering Explained

---

Unsupervised learning models are often created using a clustering algorithm.





# Instructor Demonstration

---

## Clustering Explained

# Clustering Explained

---

The process of clustering data points into groups is called **centering**.



In advanced analytics, centering helps to determine the number of classes or groups to create.



Centering improves the performance of logistic regression models by ensuring that all data points share the same starting mean value.



Data points with the same starting mean value are clustered together.

# Questions?





# The K-means Algorithm

Suggested Time:

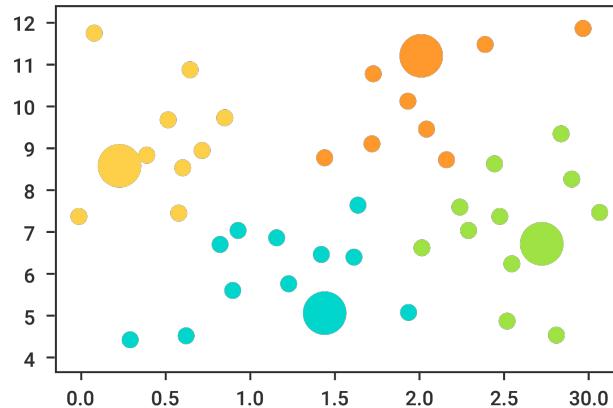
---

15 Minutes

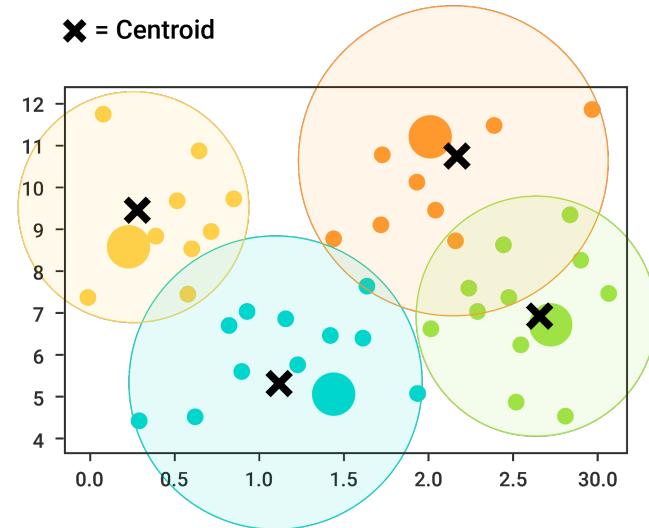
The **K-means algorithm** is the simplest and most common algorithm used to group data points into clusters.

# The K-means Algorithm

K-means takes a predetermined amount of clusters and then assigns each data point to one of those clusters.



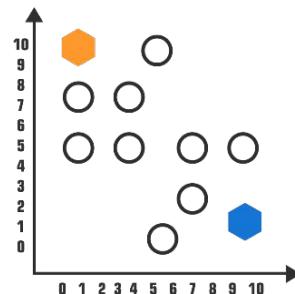
The algorithm assigns points to the closest cluster center.



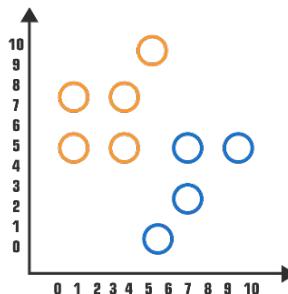
The algorithm readjusts the cluster's center by setting each center as the mean of all the data points contained within that cluster.

# The K-means Algorithm

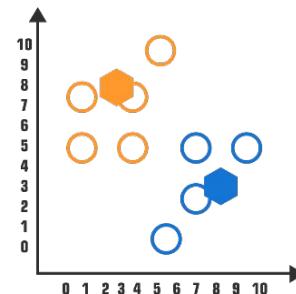
The K-means algorithm then repeats this process, again and again, each time getting a little bit better at separating the data points into distinct groups.



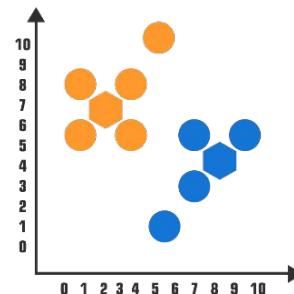
Randomly select K-clusters



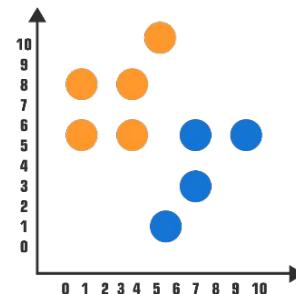
Each object assigned to similar centroid randomly



Cluster centers updating depending on renewed cluster mean



Reassign data points; update cluster centers



Reassign data points

# Questions?





# Activity: Segmenting Customers

In this activity, you will use the K-means algorithm to segment customer data for mobile versus in-person banking service ratings.

Suggested Time:

---

20 Minutes



Time's Up! Let's Review.

# Questions?





## Instructor Demonstration

---

### Review Segmenting Customers

# Questions?





Countdown timer

15:00

(with alarm)

Break



# Introduction to Clustering Optimization

# Introduction to Clustering Optimization

---



The appropriate clustering algorithm and parameter settings depend on the individual dataset and intended use of the results.



Cluster analysis is not an automatic task.



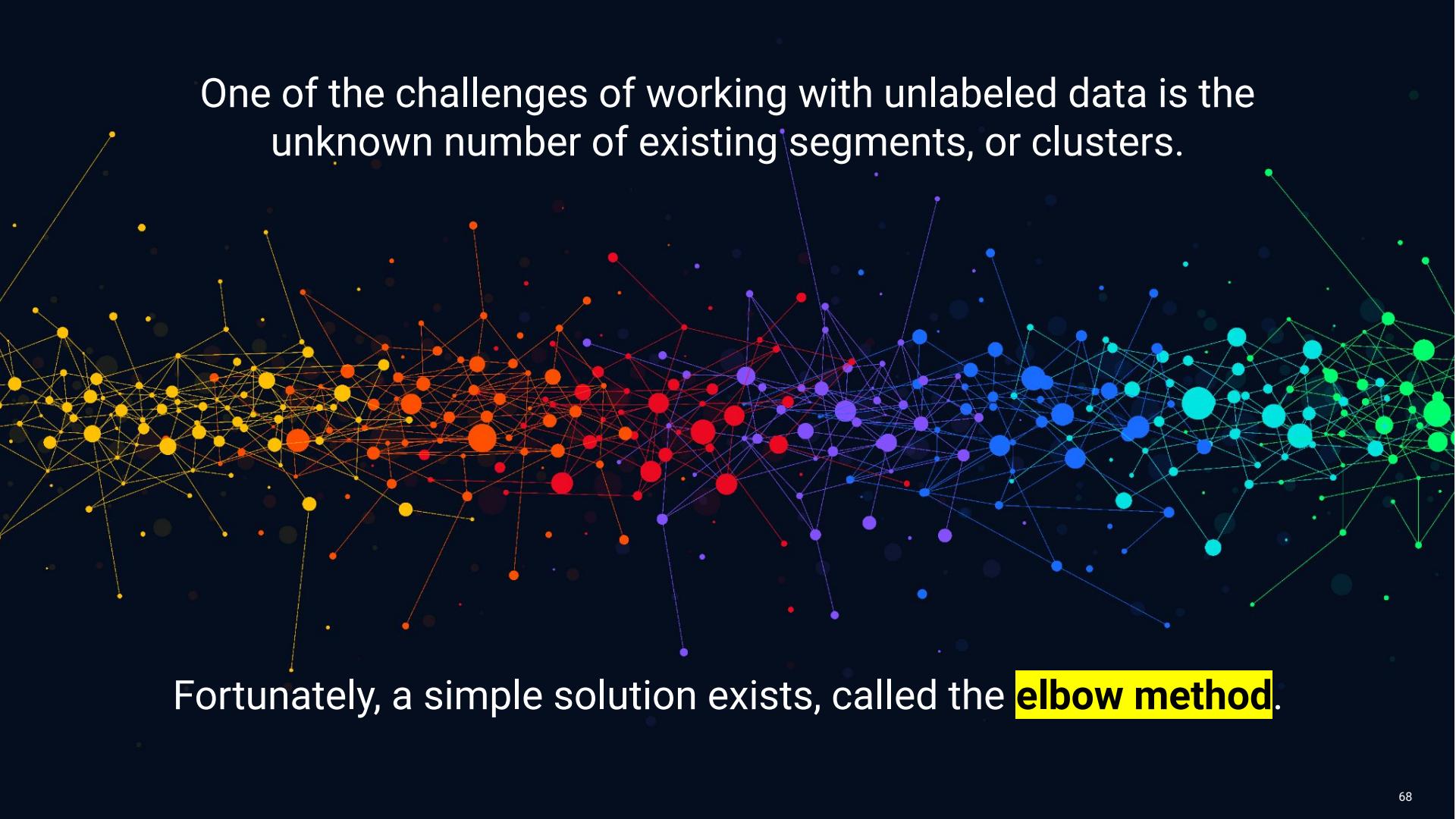
As a fintech professional, you will need to do some trial and error to find the optimal clusters.



This process includes modifying the data preprocessing and model parameters until the result achieves the desired properties.



How do you know the optimal  
number of clusters, or value of k,  
and how do you find it?



One of the challenges of working with unlabeled data is the unknown number of existing segments, or clusters.

Fortunately, a simple solution exists, called the **elbow method**.

# Questions?



# The Elbow Method

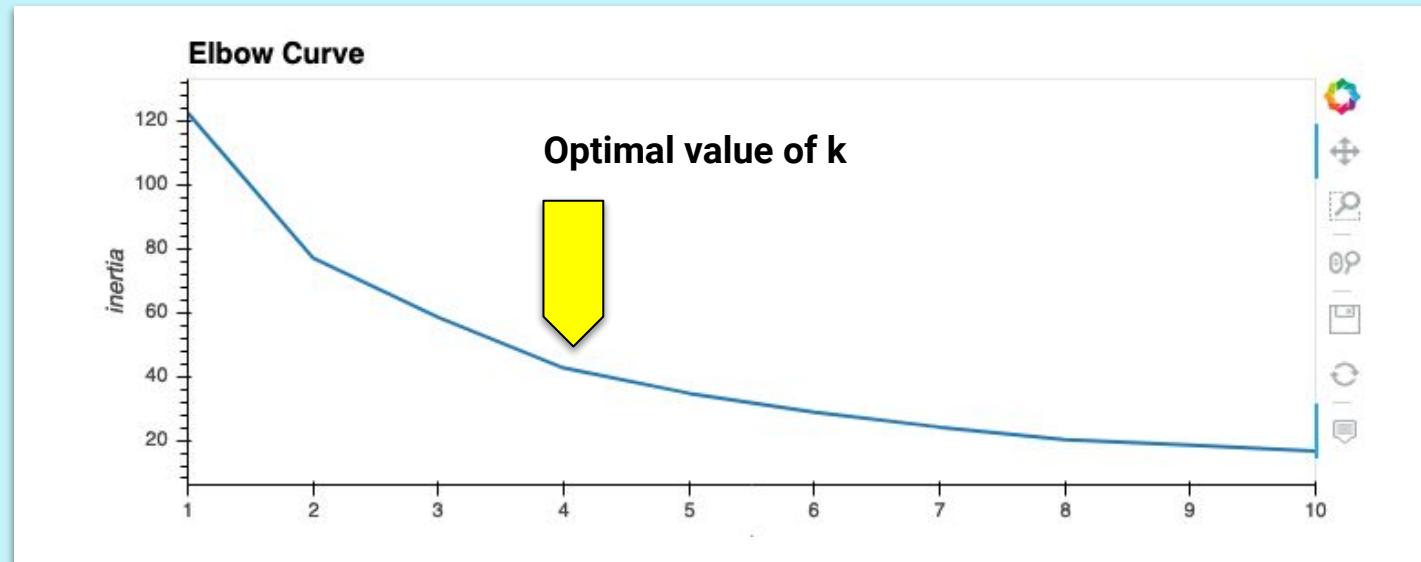


Since the K-means algorithm needs to have the amount of clusters provided ahead of time, how can you be sure that the amount of clusters you chose is correct?

# The Elbow Method

One method for determining the optimal value of  $k$ , or the number of clusters in a dataset, is the **elbow method**.

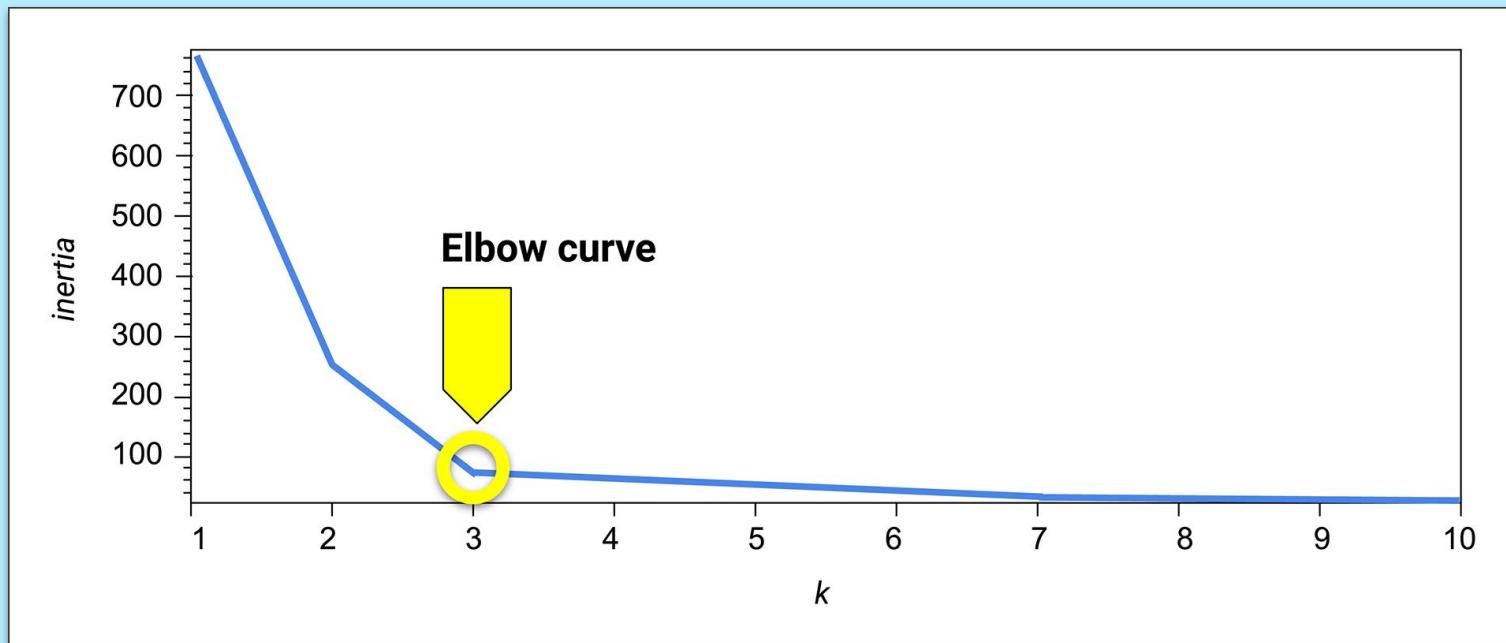
- The elbow method runs the K-means algorithm for a range of possibilities for  $k$ , or the number of clusters.
- The resulting elbow curve plots the number of clusters,  $x$ , versus an objective function called inertia.



# Elbow Curve

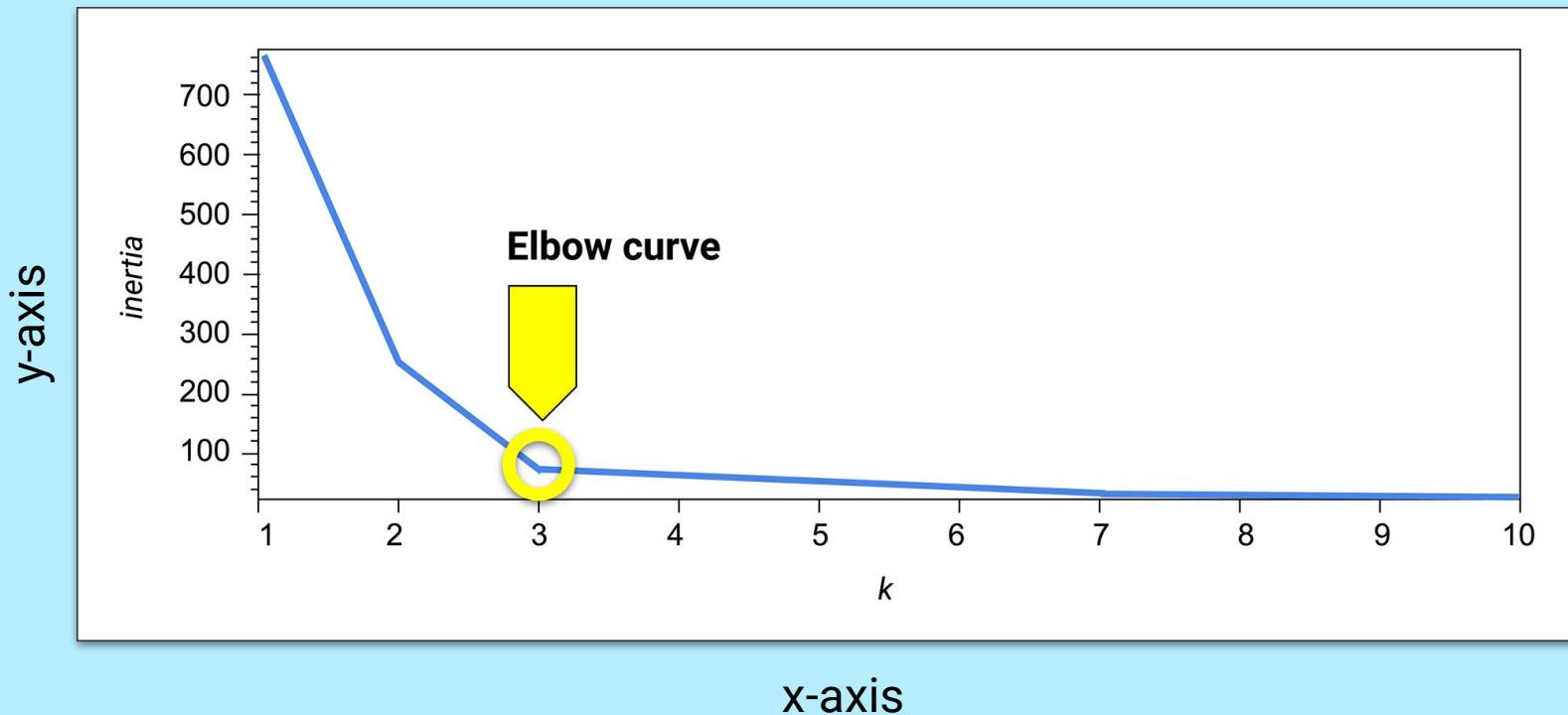
The **elbow curve** is commonly used to figure out the best value of  $k$ .

It is essentially used to determine the number of clusters at which the data points become tightly clustered.



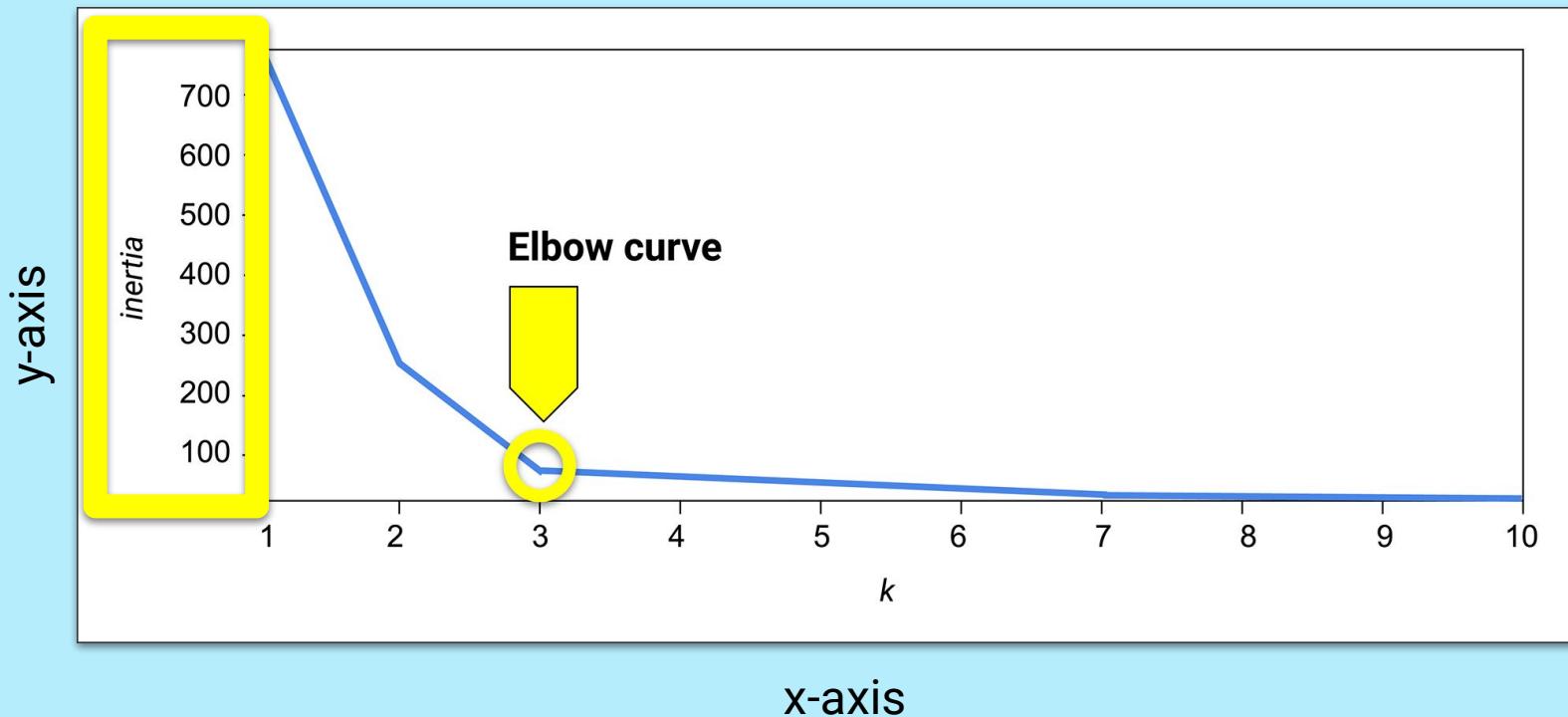
# Elbow Curve

On the elbow curve, the x-axis is the value of clusters, while the y-axis is a metric used to assess the value of  $k$ .



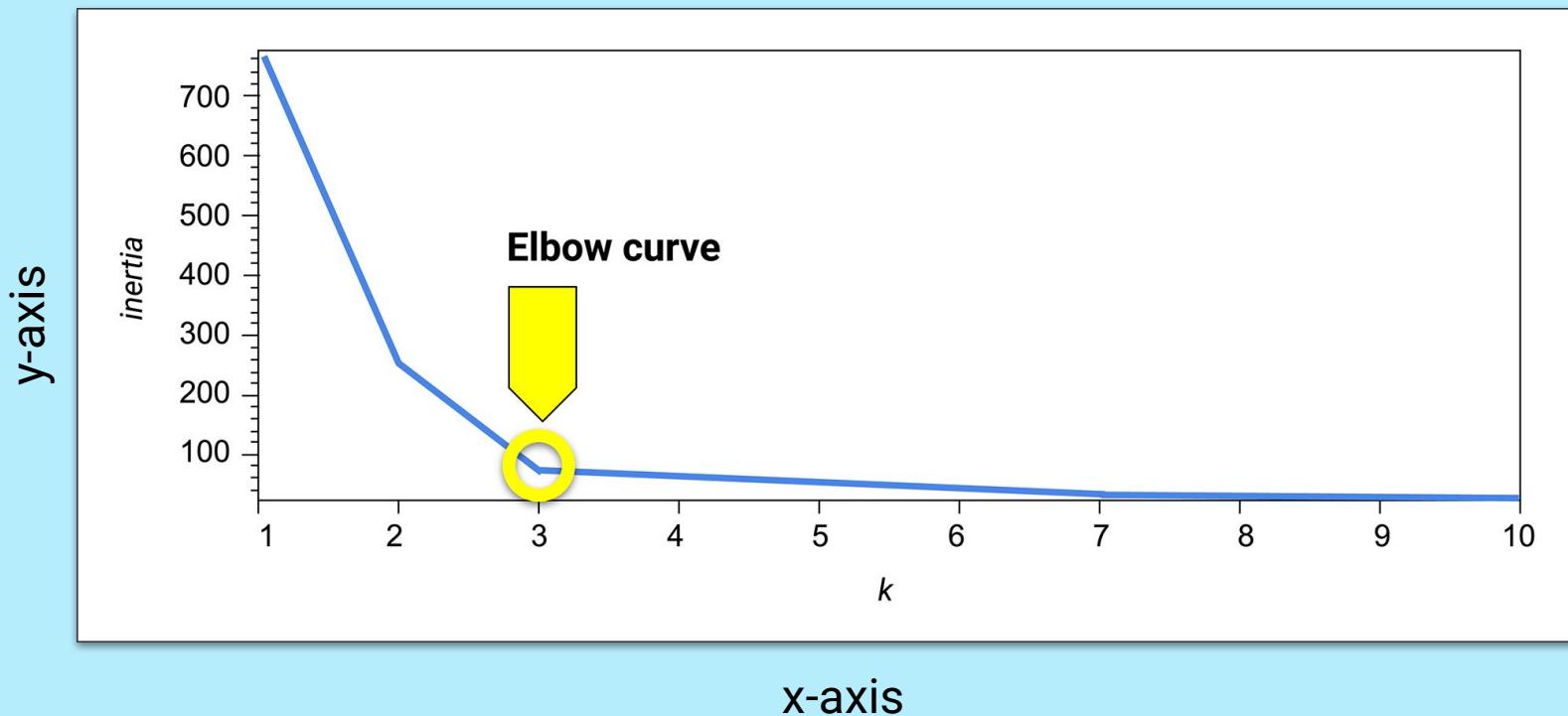
# Elbow Curve

The **inertia** is commonly used as an objective function. It is the sum of the squared distances of samples to their closest cluster center.



# Elbow Curve

A low inertia value means that the data points are tightly clustered around the cluster center.



# Inertia

---

Inertia involves complicated math, but it is basically a measure of how concentrated the elements are in a dataset.

High concentration

Datasets with a high concentration of elements (where elements are tightly grouped together) have a **low** inertia value.

This means that there is a small standard deviation for the elements in the cluster relative to the cluster mean value.

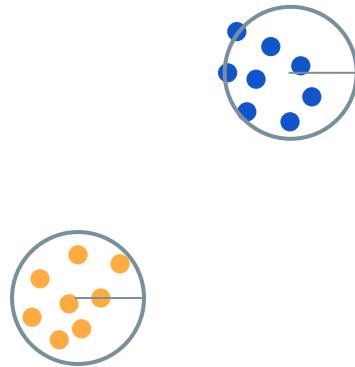
Low concentration

Datasets with a low concentration of elements (where elements are spread out) have a **high** inertia value.

This means that there is a high standard deviation for the elements in the cluster relative to the cluster mean value.

## Low Inertia

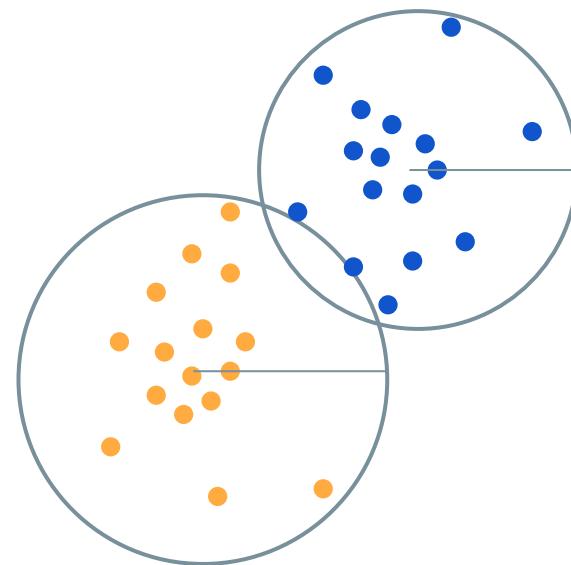
Radius of circle is small =  
small standard deviation from cluster mean



vs.

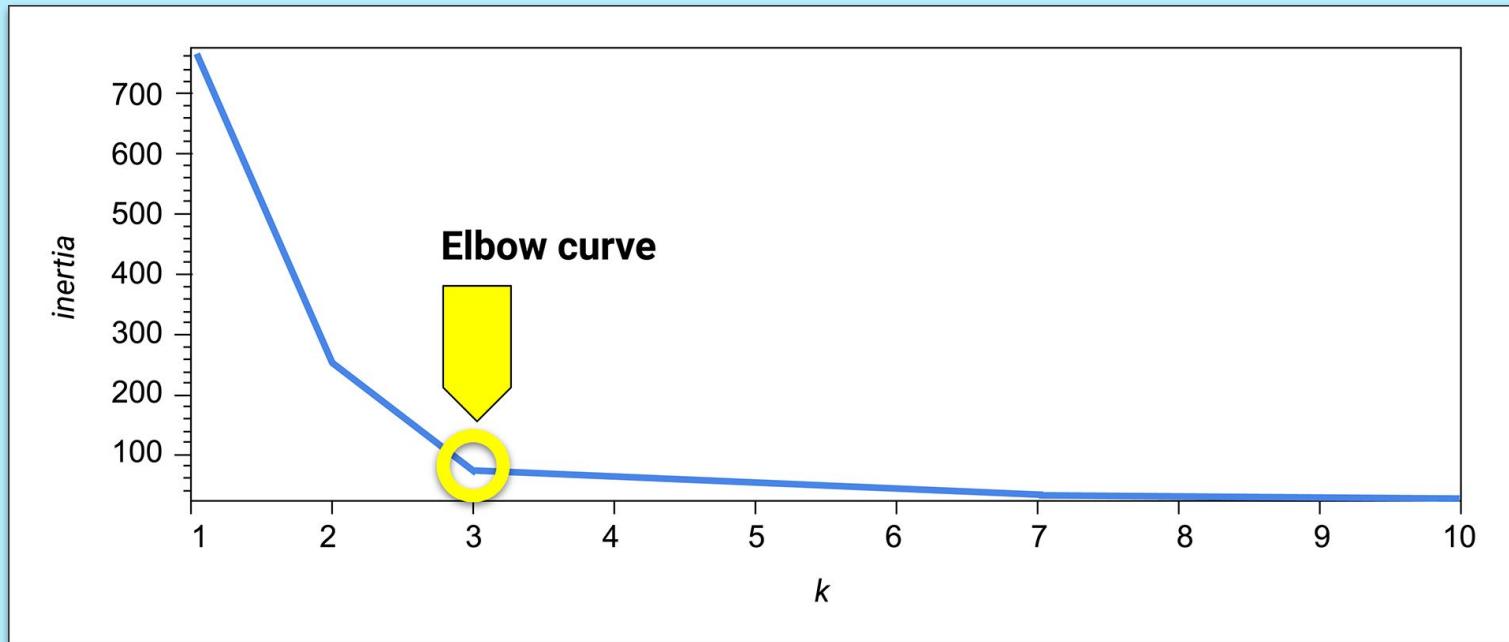
## High Inertia

Radius of circle is large =  
large standard deviation from cluster mean



# The Elbow Method

The goal is to find a value for  $k$  that corresponds to a measure of inertia that shows minimal change for each additional cluster (or value of  $k$ ) that is added to the dataset. **The spot is indicated by the bend in the elbow.**





## The Elbow Method

Suggested Time:

---

20 Minutes

# Questions?





## Activity: Finding k

In this activity, you will use the elbow method to determine the optimal number of clusters that should be used to segment a dataset of stock pricing information.

Suggested Time:

---

20 minutes



Time's Up! Let's Review.

# Questions?





# Recap

---

Congratulations on acquiring your first machine learning skills! You can now:



Explain the differences between supervised and unsupervised machine learning.



Describe the purpose of clustering and how it's used in finance.



Use the K-means algorithm to identify dataset clusters, and how to optimize this algorithm by using the elbow method.



Optimize the K-means algorithm by using the elbow method.



## Next Class

You will learn how to preprocess the data that goes into these types of models, and you'll create models that can adapt and perform better on more complex types of data.

# Questions?



The  
End