

Lecture 1: Introduction to Pattern Recognition

Course: 2110573 Pattern Recognition

2026-01-13
version 1.0

1 What is Pattern Recognition?

Pattern recognition is a fundamental branch of machine learning. A common definition is:

“Pattern recognition is a branch of machine learning that focuses on the recognition of patterns and regularities in data, although it is in some cases considered to be nearly synonymous with machine learning.” — *Wikipedia*

This field is closely related to, and often overlaps with, other domains such as:

- Artificial Intelligence (AI)
- Data Mining (DM)
- Knowledge Discovery in Databases (KDD)
- Statistics
- Data Science

At their core, all these fields are concerned with the same fundamental question: **how do we learn from data?**

1.1 A Note on AI

- **Classical Definition:** A system that appears intelligent.
- **Modern Context:** The field is currently focused on “Narrow AI,” which involves creating specialized systems that are very good at a single task (e.g., image recognition, language translation). This is largely synonymous with Machine Learning (ML).
- **Artificial General Intelligence (AGI):** This is the hypothetical ability of an agent to understand or learn any intellectual task a human can. It remains a topic of philosophical debate and long-term research.

2 Course Philosophy

The main goal of this course is to go beyond treating models as “black boxes.” In this course, you will:

- Understand models on a deeper level.
- Implement algorithms from scratch to solidify your understanding.

As François Chollet states, “if you understand something clearly, you should be able to describe it in precise algorithmic terms to a computer.”

3 Types of Machine Learning

There are three main categories of machine learning:

1. **Supervised Learning:** Learn a model from labeled data, i.e., pairs of (input, output).
2. **Unsupervised Learning:** Discover hidden structures in unlabeled data (input only, no output).
3. **Reinforcement Learning:** Train an agent to take actions in an environment to maximize a cumulative reward.

There is also a precursor to machine learning known as **rule-based systems**. An example is a 7-segment display decoder, where explicit ‘IF-THEN’ rules map the state of the 7 segments to a specific digit.

4 The Machine Learning Workflow

The goal of machine learning is to find the best function, $F(x)$, that maps an input x to an output y automatically from data. The typical workflow involves several key stages.

4.1 Feature Extraction

This is the process of transforming raw data into a numerical **feature vector** (x) that is informative and digestible for a model. This is one of the most critical steps, as the quality of the features directly impacts the model’s performance. This adheres to the principle of “Garbage in, Garbage out.”

4.2 Modeling (Training)

During the training phase, a learning algorithm uses a **training set** of feature vectors and their corresponding desired outputs (y) to learn a model (h).

$$\text{Training Data } (x, y) \rightarrow \text{Learning Algorithm} \rightarrow \text{Model } (h)$$

4.3 Evaluation (Testing)

In the testing phase, the trained model (h) is used to make predictions on new, unseen input data.¹

$$\text{New Input } (x) \rightarrow \text{Model } (h) \rightarrow \text{Predicted Output } (\hat{y})$$

The model’s performance is then evaluated by comparing its predictions (\hat{y}) to the actual correct answers, or **ground truth** (y).

5 Model Evaluation

To compare different models, we need metrics to quantify their performance.

5.1 Ground Truth

We evaluate a model by comparing its output to the “correct answer” or ground truth. However, establishing ground truth can be difficult, especially for subjective tasks (e.g., rating a machine translation) or in cases with ambiguous data (e.g., labeling complex road signs or lane lines).

5.2 Metrics for Classification

For a binary classification or detection problem, there are four possible outcomes, often summarized in a **confusion matrix**:

		Detector Prediction	
Actual	Positive	Negative	
Positive	True Positive (TP)	False Negative (FN)	
Negative	False Positive (FP)	True Negative (TN)	

¹We call the act of using a model on some input data x as “inference.”

From this, we derive several common metrics:

- **Accuracy:** $\frac{TP+TN}{TP+TN+FP+FN}$. Note: Accuracy can be misleading on biased datasets.
- **Precision:** $\frac{TP}{TP+FP}$. Of all the positive predictions, how many were actually correct? (High precision is critical for tasks like spam detection).
- **False Alarm Rate:** $\frac{FP}{FP+TN}$ measures how often the model incorrectly predicts a negative case as positive
- **Specificity (True Negative Rate):** $\frac{TN}{TN+FP}$ measures how well the model correctly identifies negative cases.
- **Recall (Sensitivity, True Positive Rate):** $\frac{TP}{TP+FN}$. Of all the actual positive cases, how many did the model correctly identify? (High recall is critical for tasks like cancer screening).
- **F1-Score:** $\frac{2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}}{\frac{1}{Precision} + \frac{1}{Recall}}$. The harmonic mean of precision and recall, providing a single score that balances both.²

5.3 High Precision or High Recall

Choosing an appropriate metric is important because different types of model errors lead to different levels of impact in real-world applications. This also affects how we define the positive class, which is usually the class we care about the most. In general the positive class would be rarer than the negative class. **Recall** focuses on finding as many positive samples as possible, while **Precision** focuses on making sure that the predicted positive samples are truly positive. For example,

- cancer detection (**Recall**) : the goal is to identify as many patients as possible; missing a patient is more harmful than a false alarm, and the positive class is cancer patients. the positive class is cancer patients
- face recognition for bank identity verification (**Precision**) : misidentifying one person as another causes serious damage, while asking a user to retry is acceptable; the positive class is the verified user
- In COVID screening ATK vs PCR
 - ATK (**Recall**) : screening and finding as many infected people as possible in the first stage
 - PCR (**Precision**) : confirming true COVID cases in the second stage

the positive class is COVID patients.

5.4 Example Calculations

Suppose we have a model that predicts whether it will rain or not, with the results shown in the table. In this case, the positive class is rain, because we are more interested in weather conditions when it rains than when it does not.

		Detector Prediction	
		Rain	not Rain
Actual	Rain	15	3
	not Rain	6	39

- The positive class : Rain
- Precision = $\frac{15}{15+6} = \frac{5}{7}$
- Recall = $\frac{15}{15+3} = \frac{5}{6}$
- FAR = $\frac{6}{6+39} = \frac{2}{15}$
- TNR = $\frac{39}{6+39} = \frac{13}{15}$
- F1-score = $2 \cdot \frac{\frac{5}{7} \cdot \frac{5}{6}}{\frac{5}{7} + \frac{5}{6}} = \frac{10}{13}$

²Although F-1 Score captures both precision and recall, it is not recommended to use it as a substitute for both. Typical use cases will focus on either precision or recall, so it is better to report both rather than the average.

5.5 Other Considerations

Besides predictive accuracy, other factors are important when evaluating models:

- Training and testing time
- Memory requirements
- Latency (time to make a single prediction)
- Parallelizability

6 Course Walkthrough

The course is structured into two main parts:

- **Traditional Machine Learning:** Covers foundational topics like K-Means, Regression, Naive Bayes, GMM, SVM, and basic Neural Networks (NN).
- **Deep Learning:** Moves into advanced topics including CNNs, Transformers, Generative Models (GANs, VAE, Diffusion), Self-supervised learning, and Reinforcement Learning.

The schedule includes regular homework assignments, a midterm, a final project, and a final exam.

7 Introduction to Unsupervised Learning

Unsupervised learning is a type of machine learning where the goal is to discover hidden patterns and structures in unlabeled data (i.e., data without a predefined target variable ' y ').

Key applications include:

- Customer/product segmentation
- General data analysis and exploration
- Identifying the number of speakers in a meeting recording
- Assisting supervised learning tasks

Clustering is a fundamental task in unsupervised learning that tries to automatically discover natural groupings within the data.

8 Nearest Neighbour Classification

Before diving into K-means, it's helpful to understand classification using nearest neighbors.

8.1 K-Nearest Neighbour (KNN) Classification

The core idea is to classify a new data point based on the labels of its neighbors in the training data.

1. **Given a query data point:**
2. For every point in the training data, compute the distance to the query point.
3. Identify the 'K' nearest data points (neighbors).
4. Assign a label to the query point by taking a majority vote among its K neighbors.

A simple version (where $K=1$) is called Nearest Neighbour classification, but it is highly susceptible to noise. Using a larger K and a voting scheme makes the method more robust. For weighted k-NN, the votes of closer neighbors can be given more importance (e.g., weighted by the inverse of their distance).

8.2 Distance Metrics

To find the “closest” neighbors, we need a distance or similarity measure. Common measures include:

- **Euclidean Distance:** The straight-line distance between two points.

$$d(X_1, X_2) = \sqrt{\sum_{i=1}^n (x_{1,i} - x_{2,i})^2}$$

- **Cosine Similarity:** Measures the cosine of the angle between two vectors, indicating their orientation similarity.

$$\text{sim}(X_1, X_2) = \frac{X_1 \cdot X_2}{\|X_1\| \|X_2\|} = \frac{\sum_{i=1}^n x_{1,i} x_{2,i}}{\sqrt{\sum_{i=1}^n x_{1,i}^2} \sqrt{\sum_{i=1}^n x_{2,i}^2}}$$

Other measures like Jaccard distance and Earth Mover’s distance also exist.

8.3 KNN Runtime

The runtime complexity for classifying J query points against a training set of N points is **O(JN)**. This can be computationally expensive for large datasets. Techniques like using centroids can speed this up.

9 K-means Clustering

K-means is an iterative algorithm that partitions a dataset into K distinct, non-overlapping clusters.

9.1 The Algorithm

The algorithm aims to minimize the distance between data points and their assigned cluster’s centroid.

1. **Initialization:** Randomly select K data points from the dataset to serve as the initial centroids.
2. **Assignment Step:** Assign each data point in the dataset to the nearest centroid. The “nearest” is determined using a distance metric, typically Euclidean distance.
3. **Update Step:** Recalculate the centroid of each cluster by taking the mean of all data points assigned to it.
4. **Repeat:** Repeat the Assignment and Update steps until the centroids no longer change significantly, meaning the cluster assignments have stabilized.

9.2 Characteristics and Assumptions

- The number of clusters, k , must be specified in advance.
- The algorithm is guaranteed to converge, but it may converge to a **local minimum**.
- **Bad initializations** can lead to poor clustering results. A common solution is to run the algorithm multiple times with different random initializations and select the best outcome (e.g., the one with the lowest overall variance).
- The model implicitly assumes that clusters are spherical, evenly sized, and have similar density.

10 Selecting the Number of Clusters (K)

Choosing the right value for K is a critical step in K-means clustering.

10.1 The Elbow Method

One of the most common methods is the Elbow Method.

- Run the K-means algorithm for a range of K values (e.g., K=1 to 10).

- For each K, calculate the **fraction of explained variance**.

- Fraction of explained variance = $\frac{\text{Between-cluster variance}}{\text{All-data variance}}$
- Between-cluster variance = $\sum_{i=1}^K n_i \frac{\|M_i - M\|^2}{N-1}$ ($\|\cdot\|$ denotes the norm of the vector.)
- All-data variance = $\sum_{j=1}^N \frac{\|x_j - M\|^2}{N-1}$

where M_i is the mean of cluster i , M is the overall mean, n_i is the number of data points in cluster i , x_j is a data point, and N is the total number of data points.

- Plot the fraction of explained variance as a function of K.

- The “elbow” of the curve—the point where the rate of increase sharply declines—is considered the optimal K. This represents the point of diminishing returns, where adding more clusters doesn’t significantly explain more variance.

10.2 Example application of K-means: Real Estate Analysis

This section presents an example of using K-means to cluster real estate data in Thailand, as shown in Figure 1³. We use users’ website viewing behavior (cookies) of each real estate project as features and calculate the distance between users based on the projects they view. The real estate we analyze are shown in Figure 2. The idea is that if users tend to view similar projects, those projects are likely to be similar. The results in Figure 3 reveal interesting patterns, such as clusters forming along BTS (Skytrain) lines (the blue cluster) and clusters based on residential zones. This shows that K-means can help us discover meaningful data groups; however, it does not explain why the zones are separated in this way or what other factors influence the clustering. Therefore, this analysis mainly provides initial insights, which help guide further investigation into what additional data should be collected and how to collect it. The real estate website might use this information to plan advertisements or property development. In general, **clustering, as an unsupervised algorithm, is not an end goal by itself, but a starting point for further investigations or actions by the user.**

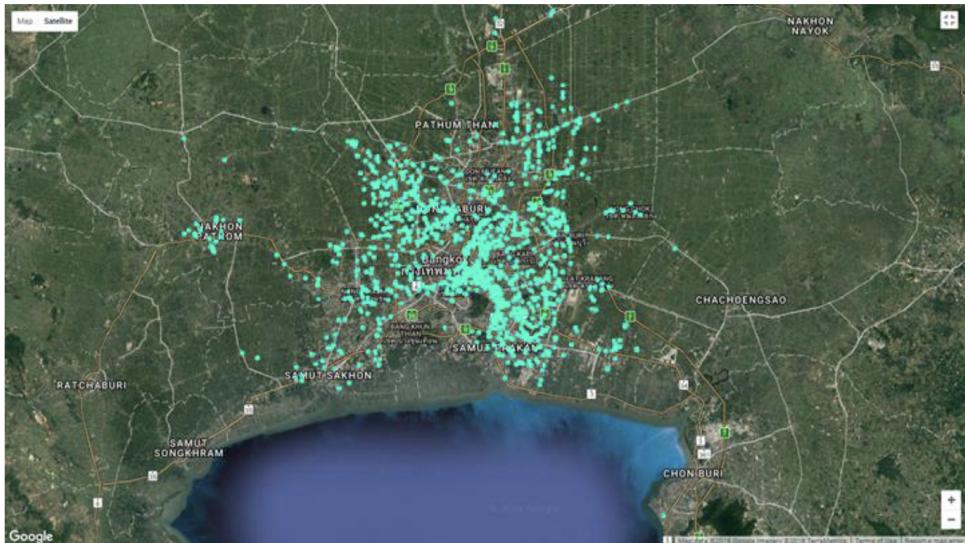


Figure 1: the location data of real estate properties in Thailand

³The data was provided by Home.co.th, a popular real estate website in Thailand.

Real estate project and view activities by each account

Project/Visitor	Visitor 1	Visitor 2	Visitor 3	Visitor 4	Visitor 5	...	Visitor 1,000,000
Project 1	y	n	y	n	n	...	n
Project 2	n	n	n	y	y	...	n
Project 3	n	n	n	n	n	...	n
Project 4	y	n	y	y	n	...	n
Project 5	y	n	y	n	n	...	n
...
Project 4000	n	y	n	n	n	...	y

(a) User-project interactions, feature creation, project similarity analysis, and K-means clustering in a real estate dataset.

features							
Project/Visitor	Visitor 1	Visitor 2	Visitor 3	Visitor 4	Visitor 5	...	Visitor 1,000,000
Project 1	y	n	y	n	n	...	n
Project 2	n	n	n	y	y	...	n
Project 3	n	n	n	n	n	...	n
Project 4	y	n	y	y	n	...	n
Project 5	y	n	y	n	n	...	n
...
Project 4000	n	y	n	n	n	...	y

(b) We use the views as the feature vector for each real estate project.

These two projects garner similar interests							
Project/Visitor	Visitor 1	Visitor 2	Visitor 3	Visitor 4	Visitor 5	...	Visitor 1,000,000
Project 1	y	n	y	n	n	...	n
Project 2	n	n	n	y	y	...	n
Project 3	n	n	n	n	n	...	n
Project 4	y	n	y	y	n	...	n
Project 5	y	n	y	n	n	...	n
...
Project 4000	n	y	n	n	n	...	y

(c) The goal of clustering is to help identify items that has similar looking features.

Figure 2: The process of analyzing a real estate dataset, including user interaction analysis, similarity detection, and K-means clustering.

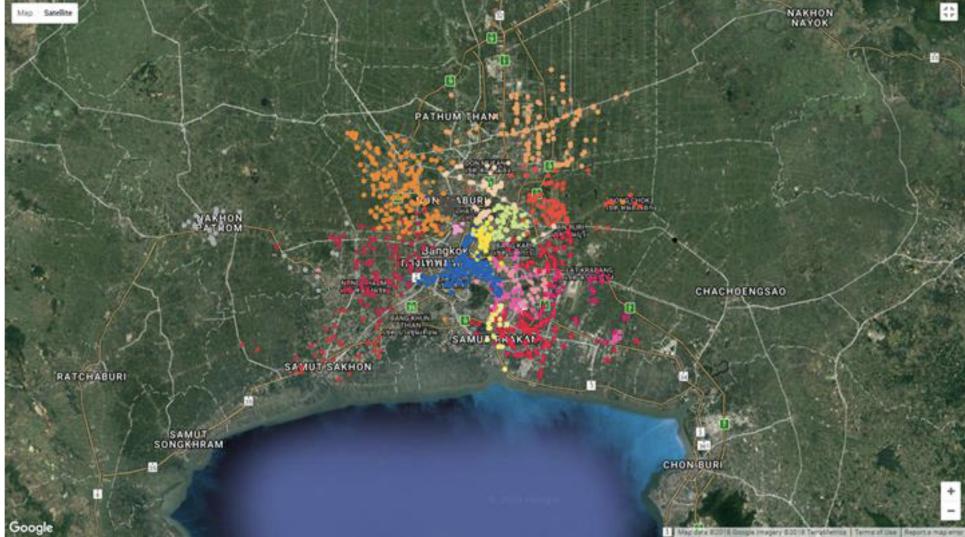


Figure 3: The grouping of real estate properties in Thailand after applying K-means clustering

10.3 Metrics

Clustering is mainly used to gain insight from data so we can use those patterns in later analysis, but we can use metrics to help evaluate how well a model separates the data into meaningful groups, which usually compare how compact each cluster is versus how far apart the clusters are overall. Popular metrics include:

- **Silhouette Score:** Measures how well data points are separated between clusters, with values in the range $[-1, 1]$. A value close to 1 indicates well-separated clusters, 0 indicates overlapping clusters, and negative values suggest incorrect cluster assignment. It is easy to understand but may perform poorly on very dense or highly detailed data. It is based on the average distance of a data point to others in the same cluster, and the minimum average distance to the nearest different cluster.
- **Davies–Bouldin Index (DBI):** Evaluates clustering quality by comparing pairs of clusters using the ratio between within-cluster distance to the centroid and the distance between centroids of different clusters. Lower DBI values indicate better clustering, meaning clusters are compact and well separated.
- **Calinski–Harabasz Index (CHI):** Measures clustering performance by comparing the variance between cluster centroids with the variance of data points within each cluster. A higher CHI value indicates better-defined and more distinct clusters.

There are variants of clustering metrics that require labels. These can also be used as guiding metrics. Consult the scikit-learn website for some examples.⁴

However, as I have mentioned earlier, clustering is just an intermediate step in an overall analysis. The best ways to evaluate them is on a downstream task. For example, if the clustering is used to group customers for creating different product categories. A good metric might be how well a classifier classifies the customers into the clustered groups.

10.4 Other Methods

- **Variance Threshold:** Choose the minimum K that explains a certain percentage (e.g., 95%) of the total variance.
- **Performance Maximization:** If the clustering is part of a larger supervised learning task, choose the K that maximizes the performance (e.g., accuracy) on a validation or test set.

⁴<https://scikit-learn.org/stable/modules/clustering.html#clustering-evaluation>

11 Notes of K-means and Other Clustering Methods

While K-means is powerful and widely used, other clustering methods exist that may be more suitable for different types of data.

- **K-mode and K-median:** Variations of K-means used for categorical data (K-mode) or to reduce sensitivity to outliers (K-median).
- **Spectral Clustering:** A technique that clusters data using the eigenvalues of a similarity matrix, effective for complex and non-convex cluster shapes.
- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** A density-based method that is very robust, can find arbitrarily shaped clusters, and does not require the number of clusters to be specified beforehand.

12 Introduction to Regression

Regression analysis is a set of statistical processes for estimating the relationships between a dependent variable (often called the ‘outcome’ or ‘target’) and one or more independent variables (often called ‘predictors’, ‘covariates’, or ‘features’). The goal is to predict a continuous output value.

For example, we might want to predict the amount of rainfall (our target, y) based on several features (x), such as the type of crops planted, temperature, etc. We are given a training dataset of m examples, like so:

Feature 1	Feature 2	...	Feature n	Target (y)
$x_{1,1}$	$x_{1,2}$...	$x_{1,n}$	y_1
$x_{2,1}$	$x_{2,2}$...	$x_{2,n}$	y_2
\vdots	\vdots	\ddots	\vdots	\vdots
$x_{m,1}$	$x_{m,2}$...	$x_{m,n}$	y_m

13 The Linear Regression Model

In linear regression, we assume a linear relationship between the input features and the output target. Our model, or *hypothesis*, is represented by the function $h_\theta(x)$:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n \quad (1)$$

Here, the θ_i values are the *parameters* (or *weights*) of the model. To simplify this notation, we can introduce $x_0 = 1$. This allows us to write the hypothesis in a more compact vector form:

$$h_\theta(\mathbf{x}) = \sum_{j=0}^n \theta_j x_j = \boldsymbol{\theta}^T \mathbf{x} \quad (2)$$

where $\boldsymbol{\theta}$ and \mathbf{x} are column vectors representing the parameters and the feature values, respectively. Our goal is to find the optimal values for the parameters $\boldsymbol{\theta}$ that make our predictions $h_\theta(\mathbf{x})$ as close to the actual values y as possible.

14 The Cost Function (Mean Squared Error)

To find the best parameters $\boldsymbol{\theta}$, we need a way to measure how well the model is performing. We do this using a *cost function*, also known as a *loss function*. A common choice for regression problems is the **Mean Squared Error (MSE)**.

The MSE measures the average of the squares of the errors, i.e., the average squared difference between the estimated values and the actual value. For easier computation of the gradient, we often use half of the sum of squared errors:

$$J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(\mathbf{x}_i) - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}_i - y_i)^2 \quad (3)$$

where m is the number of training examples. We want to find the $\boldsymbol{\theta}$ that minimizes $J(\boldsymbol{\theta})$.

15 Minimizing the Cost Function

15.1 Gradient Descent

Gradient Descent is an iterative optimization algorithm used to find the minimum of a function. The main idea is to take repeated steps in the opposite direction of the gradient (or derivative) of the function at the current point, because this is the direction of steepest descent.

The update rule for a single parameter θ_j is:

$$\theta_j := \theta_j - \alpha \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} \quad (4)$$

where α is the *learning rate*, a small positive number that controls the step size. We update the parameter θ in Equation 4. In this case, we use the Mean Squared Error (MSE) loss defined in Equation 3. Based on this loss function, we can compute the gradient with respect to θ as follows.

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} = \frac{1}{2m} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} (\theta^T \mathbf{x}_i - y_i)^2 \quad (5)$$

$$= \frac{1}{2m} \sum_{i=1}^m 2(\theta^T \mathbf{x}_i - y_i) \frac{\partial}{\partial \theta_j} (\theta^T \mathbf{x}_i - y_i) \quad (6)$$

$$= \frac{1}{m} \sum_{i=1}^m (\theta^T \mathbf{x}_i - y_i) x_{i,j} \quad (7)$$

$$= \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}_i) - y_i) x_{i,j} \quad (8)$$

So, the complete update rule for each parameter θ_j in the model is:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}_i) - y_i) x_{i,j} \quad (9)$$

This update is performed simultaneously for $\forall j = 0, \dots, n$. This method is known as **Batch Gradient Descent** because it uses all m training examples in each step. Variations include **Stochastic Gradient Descent (SGD)** which uses a small subset of examples (mini-batch) at a time.

15.2 The Normal Equation: An Analytical Solution

Instead of an iterative algorithm, it's also possible to solve for the optimal $\boldsymbol{\theta}$ analytically. This method is called the Normal Equation. We can express the cost function in matrix form. Let \mathbf{X} be the design matrix (with each row being a training example \mathbf{x}_i^T) ($\mathbf{X} = [\mathbf{x}_1^T; \mathbf{x}_2^T; \dots; \mathbf{x}_n^T]$) and \mathbf{y} ($\mathbf{y} = [y_1; y_2; \dots; y_n]$) be the vector of target values.

$$J(\boldsymbol{\theta}) = \frac{1}{2m} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \frac{1}{2m} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \quad (10)$$

To minimize this, we take the gradient with respect to $\boldsymbol{\theta}$ and set it to zero.

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \frac{1}{2m} (\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \\ &= \frac{1}{2m} (2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\mathbf{X}^T \mathbf{y}) \\ &= \frac{1}{m} (\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{y}) \end{aligned}$$

Setting the gradient to zero:

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{y} &= 0 \\ \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} &= \mathbf{X}^T \mathbf{y} \end{aligned}$$

This gives us the final equation for $\boldsymbol{\theta}$:

$$\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (11)$$

Pros and Cons

- **Gradient Descent:** Needs a learning rate α . Works well with a very large number of features. Is iterative.
- **Normal Equation:** No need for a learning rate. It's a one-step calculation. However, calculating the inverse of $\mathbf{X}^T \mathbf{X}$ is computationally expensive, with a complexity of roughly $O(n^3)$ where n is the number of features. It can be slow if n is very large. Also, $\mathbf{X}^T \mathbf{X}$ may not be invertible if features are redundant or if there are more features than training examples.

16 Summary

- Basic concept of learning from data
- We have learned two types of learning methods so far
 1. **Supervised learning – Linear Regression**
 - Learn a model F from pairs of (x, y)
 - Goal – find a model, use the model in applications
 2. **Unsupervised learning – K-mean clustering**
 - Discover the hidden structure in unlabeled data x (no y)
 - Goal - find insights, turn it into actions
- The concept of setting the loss function and optimize for it is a key concept in machine learning that will keep coming up in this course.