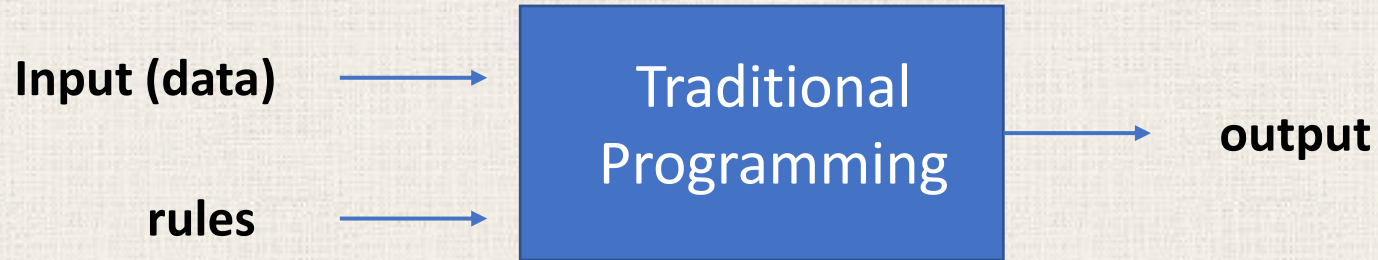
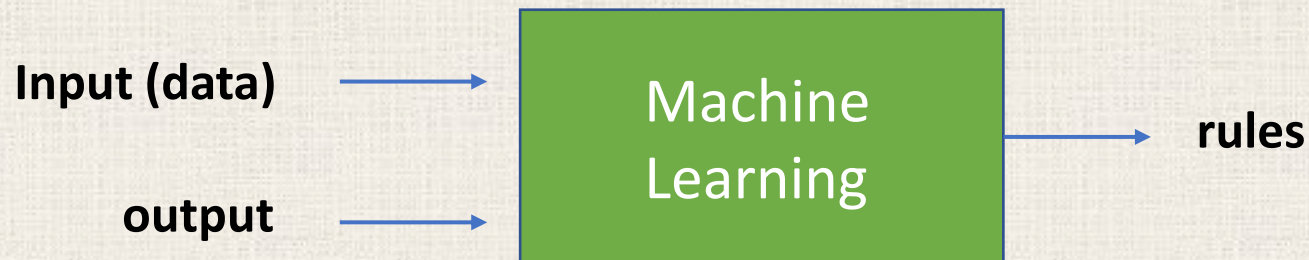


What is Machine Learning?

What is Machine Learning (loan approval example)



- Credit experts identify rules for loan approval
- Build the rules into computer program
- Run the program to approve or deny new applications



- Collect past loan application sample data, approved or denied
- Build a computer program to automatically identify the rules
- Use the rules to process new applications

What is Machine Learning

- Machine learning is *an application of artificial intelligence (AI) that provides systems the ability to automatically **learn and improve from experience** without explicitly programmed*
- Machine learning(ML) is *the scientific study of **algorithms** and **statistical models** that computer systems use to perform a specific task without using explicit instructions, relying on patterns and inference instead.*
- Machine learning algorithms build a **mathematical model** based on **sample data**, known as “training data”, in order to make **predictions** or decisions without being explicitly programmed to perform the task.
- Machine learning is **a tool for turning data into knowledge**

Examples of Machine Learning Applications

- Optical character recognition
- Email filtering
- Face recognition
- Speech recognition
- Product recommendation

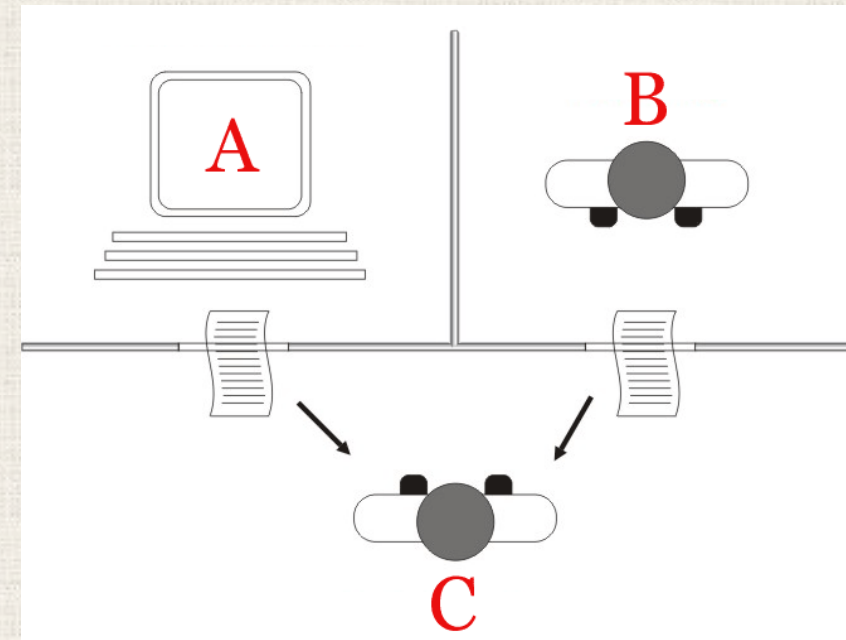
- Loan applications evaluation
- Medical diagnosis
- Outlier (anomaly) detection
- Customer segmentation
- Documents clustering
- Bioinformatics

Machine Learning and Artificial Intelligence (Turing Test)

The Turing Test: proposed by **Alen Turing** (1950) in his paper “*Computer Machinery and Intelligence*”.

Turing is widely considered to be the *father of theoretical computer science and artificial intelligence*.

- A test of a machine’s ability to exhibit intelligent behavior equivalent to a human
- A human evaluator would judge natural language conversation between a human and a machine
- The evaluator would be aware that one of the two partners in conversation is a machine
- All participants would be separated from each other
- The conversation would be limited to a text-only channel
- If the evaluator can not reliably tell the machine from the human, the machine is said to have passed the test.



How about these?

God asked Abraham to sacrifice his son Isaac because he wanted to test his faith. Whose son and whose faith are we talking about?

Would you rather sacrifice one adult to save two children, or two children to save five adults, and why?

The following sentence is true.

The previous sentence is false.

Which of those two sentences is true?

To pass a Turing test, what capabilities do you think a computer program should possess?

To pass the Turing test, what capabilities do you think the computer should possess?

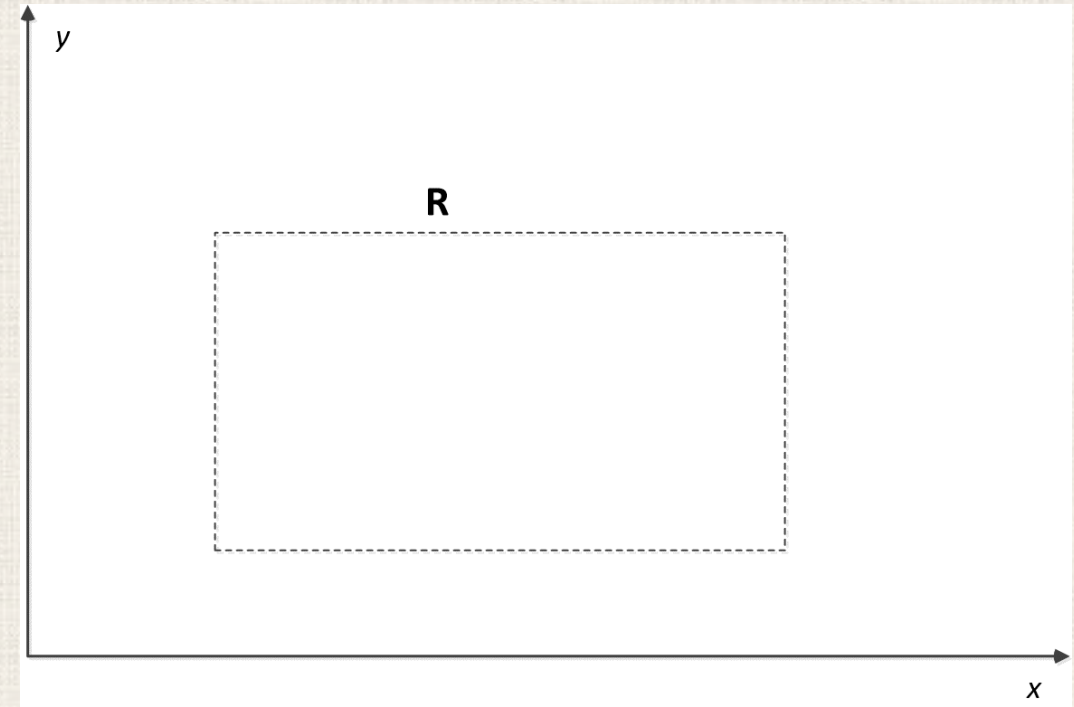
- Natural language processing (to communicate successfully)
- Knowledge representation (to store what is known or learnt)
- Automated reasoning (to use stored knowledge to answer questions and to draw new conclusions)
- Machine learning (to adapt to new situations, to detect new patterns and make prediction)

- **Artificial Intelligence (AI)** traditionally refer to an artificial creation of human-like agents that can learn, reason, plan, perceive or process natural language
- **“general AI”** : hypothetical and not domain specific, but can learn and perform tasks anywhere
- **“narrow AI”** : perform specific tasks within a domain (e.g., language translation, playing chess, prove mathematical theorems, diagnose diseases, auto driving, etc.)
- **AI** includes **acquiring knowledge** and **applying knowledge** to make decision
- **ML** only handles how to *acquire knowledge*: finding hidden patterns behind massive data

Key components of a Machine Learning Problem

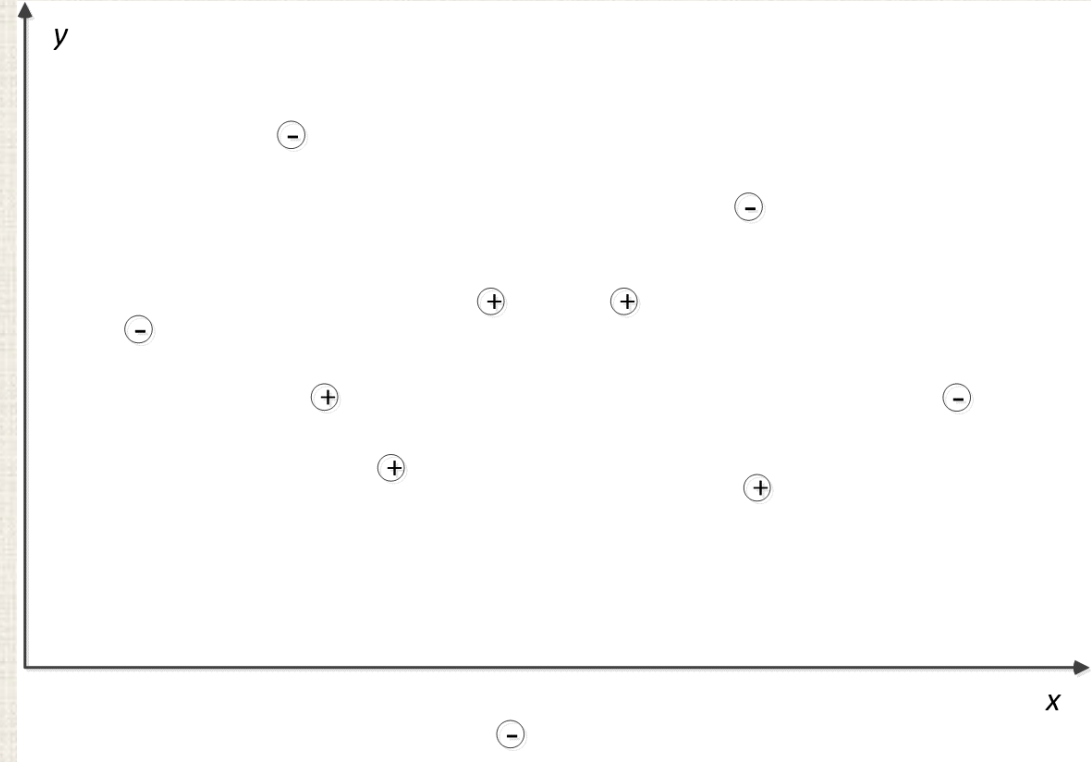
A Rectangle Learning Game

- The objective of the game is to learn an unknown axis-aligned rectangle R – **target** rectangle that separates points on the plane into two classes, **positive** (inside the rectangle) and **negative** (outside of the rectangle).



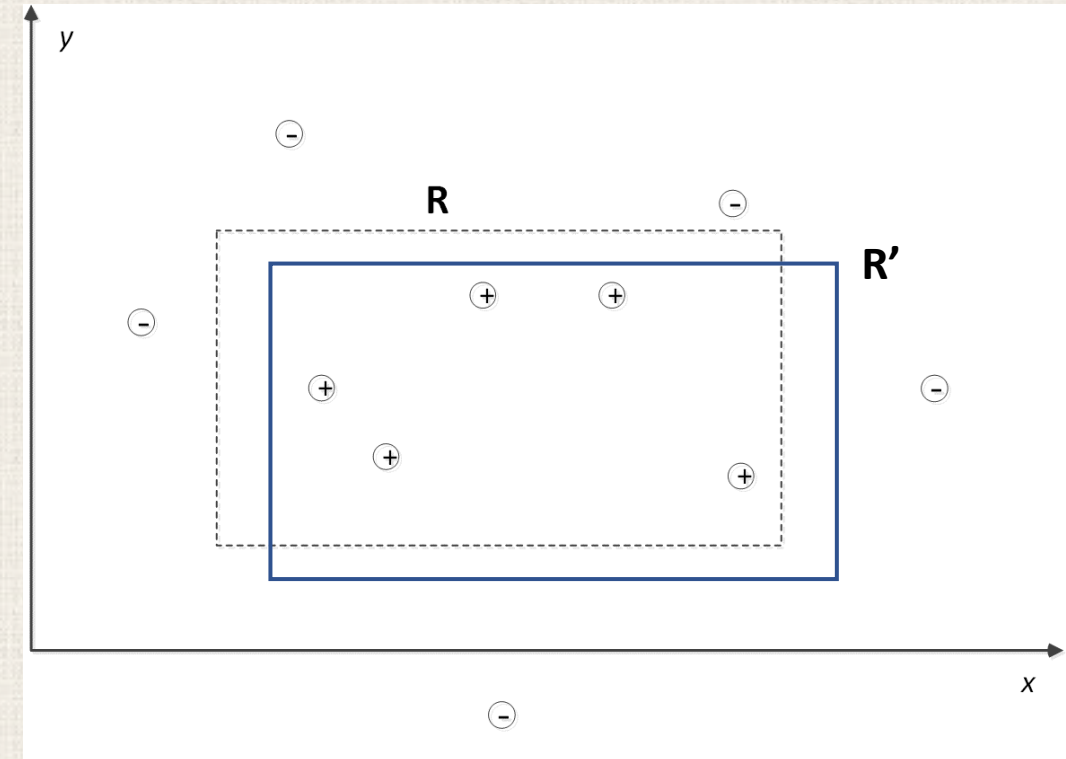
A Rectangle Learning Game

- The learner receives information about R through the following process (repeat n times):
 - a random point p is chosen in the plane according to some fixed probability distribution \mathbb{D} .
 - The learner is given the point p together with a label indicating whether p is contained in R (a positive example) or not contained in R (a negative example).



A Rectangle Learning Game

- The goal of the learner is to use as few examples as possible, and as little computation as possible, to pick a ***hypothesis*** rectangle R' that is a close approximation to R .
- The learner will choose R' from those not only separate given points but will also perform (predict) well on unseen data points (good **generalization**)



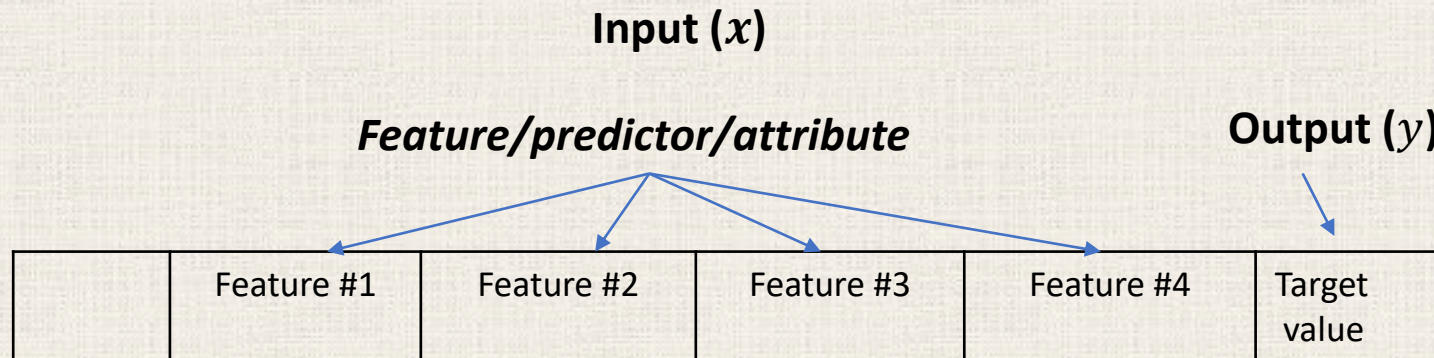
Key components of a Machine Learning Problem

- **A true model (unknown)**
- **A set of data to learn from (a random sample governed by the true model)**
- **A set of hypothesis models**
- **An objective function to measure the learning error**

The iris dataset

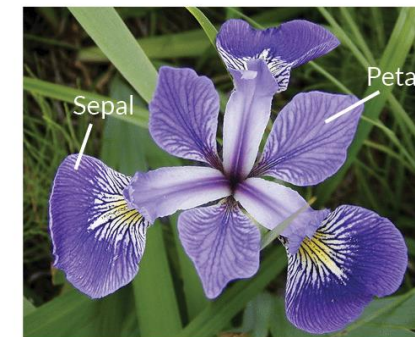
In UCI Machine Learning Repository

<https://archive.ics.uci.edu/ml/index.php>



| | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) | species |
|-----|-------------------|------------------|-------------------|------------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | 0 |
| 1 | 4.9 | 3 | 1.4 | 0.2 | 0 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | 0 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | 0 |
| 4 | 5 | 3.6 | 1.4 | 0.2 | 0 |
| ... | | | | | |
| 50 | 7 | 3.2 | 4.7 | 1.4 | 1 |
| 51 | 6.4 | 3.2 | 4.5 | 1.5 | 1 |
| 52 | 6.9 | 3.1 | 4.9 | 1.5 | 1 |
| 53 | 5.5 | 2.3 | 4 | 1.3 | 1 |
| 54 | 6.5 | 2.8 | 4.6 | 1.5 | 1 |
| ... | | | | | |
| 144 | 6.7 | 3.3 | 5.7 | 2.5 | 2 |
| 145 | 6.7 | 3 | 5.2 | 2.3 | 2 |
| 146 | 6.3 | 2.5 | 5 | 1.9 | 2 |
| 147 | 6.5 | 3 | 5.2 | 2 | 2 |
| 148 | 6.2 | 3.4 | 5.4 | 2.3 | 2 |
| 149 | 5.9 | 3 | 5.1 | 1.8 | 2 |

← Instance/example/sample/observation



Iris Versicolor



Iris Setosa



Iris Virginica

Formalize the Machine Learning Problems

- **The unknown truth:** In machine learning problems, our goal is to extract an **unknown** relationship $y = f(\mathbf{x})$ from data.
 - where the output $y \in \mathcal{R}$ (**target value**) is some quantity that can be predicted from an **input** $\mathbf{x} \in \mathcal{R}^d$,
 - where, $\mathbf{x} = [x_1, x_2, \dots, x_d]^T$ is the **feature vector** and d is the dimension of the feature vector.
 - $x_i, i = 1, 2, \dots, d$ are the **features** (attributes) of the **instances** (examples, objects).
- **The given data set:** We assume that we have access to a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where (\mathbf{x}_i, y_i) is the i^{th} point in the data set and is an example or instance (probably noisy) of the unknown mapping f to be learned.

$$\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}]^T$$

- **The hypothesis model set:** Since learning arbitrary function is intractable, we restrict ourselves to some ***hypothesis class*** H of allowable functions. Specifically, we typically employ a parametric model:

$$h_{\mathbf{w}}(\mathbf{x}) = g(\mathbf{x}, \mathbf{w})$$

Where $\mathbf{w} = [w_1, w_2, \dots, w_d]^T \in \mathcal{R}^d$, whose elements are known as parameters or ***weights***.

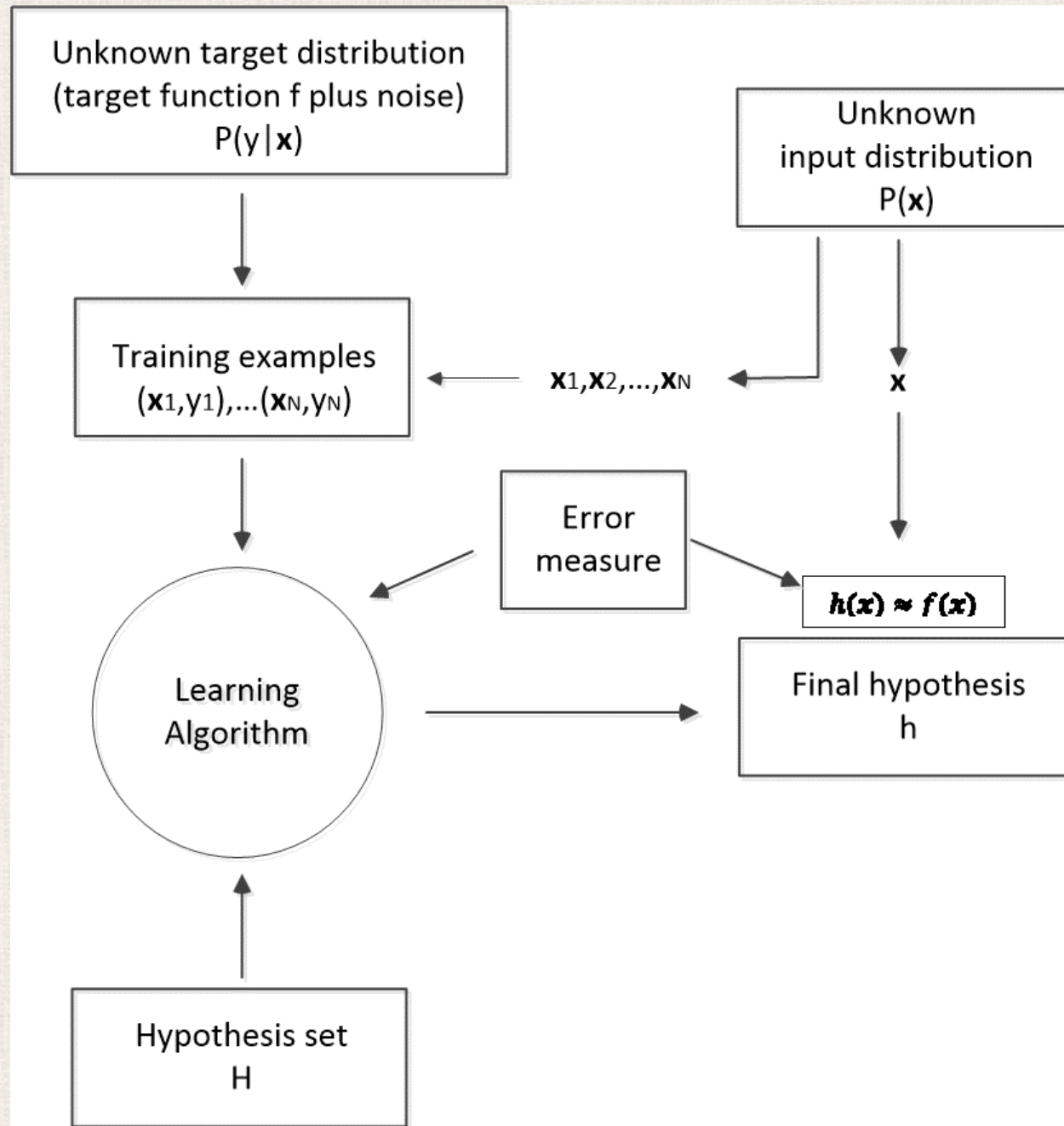
The ***hypothesis class*** is then the set of all functions induced by the possible choice of the parameter \mathbf{w} :

$$H = \{h_{\mathbf{w}} | \mathbf{w} \in \mathcal{R}^d\}$$

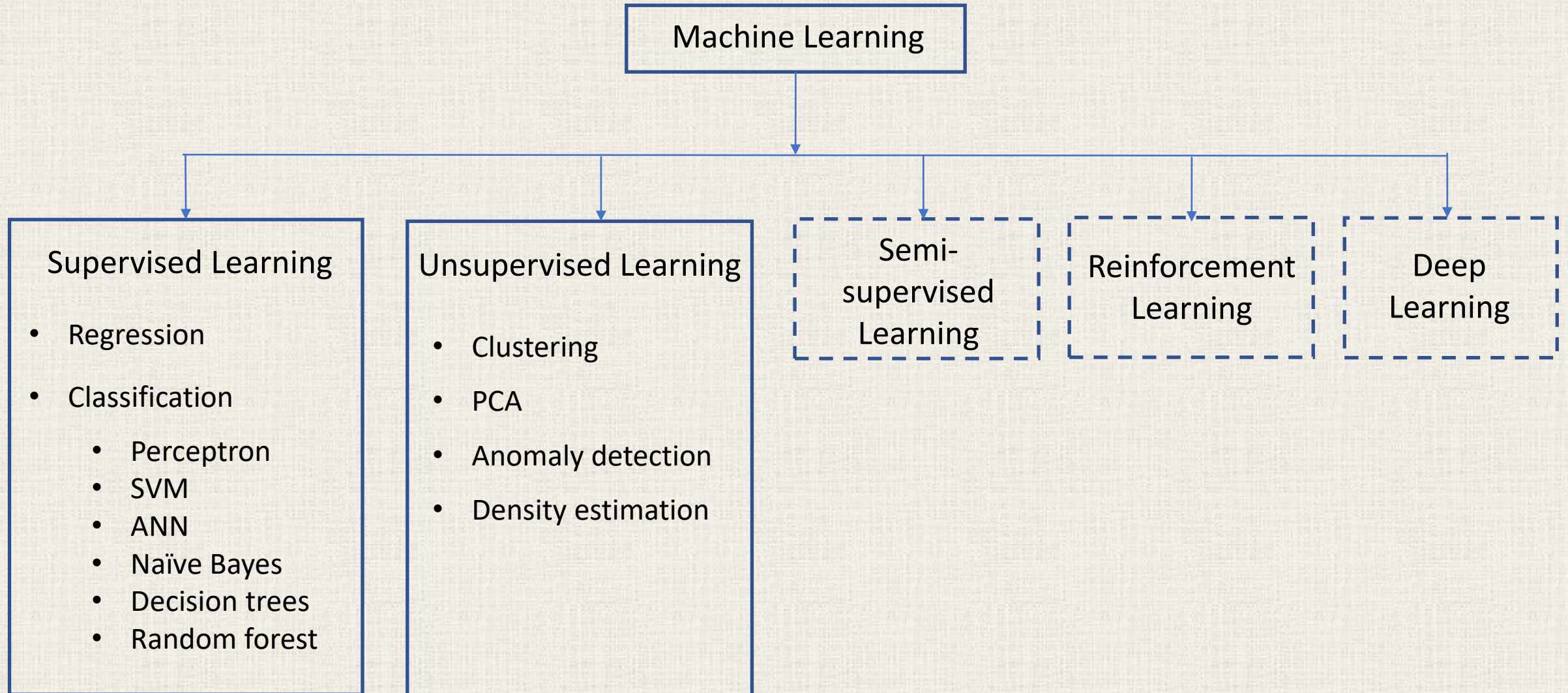
- After designating a ***loss (cost) function*** L , which measures how poorly the prediction \hat{y} of the hypothesis matches the true output y , we can proceed to search for the parameter that best fit the data by minimizing this function:

$$\mathbf{w}^* = \underset{\mathbf{w}}{arg \min} \{L(\mathbf{w})\}$$

Machine Learning Problems



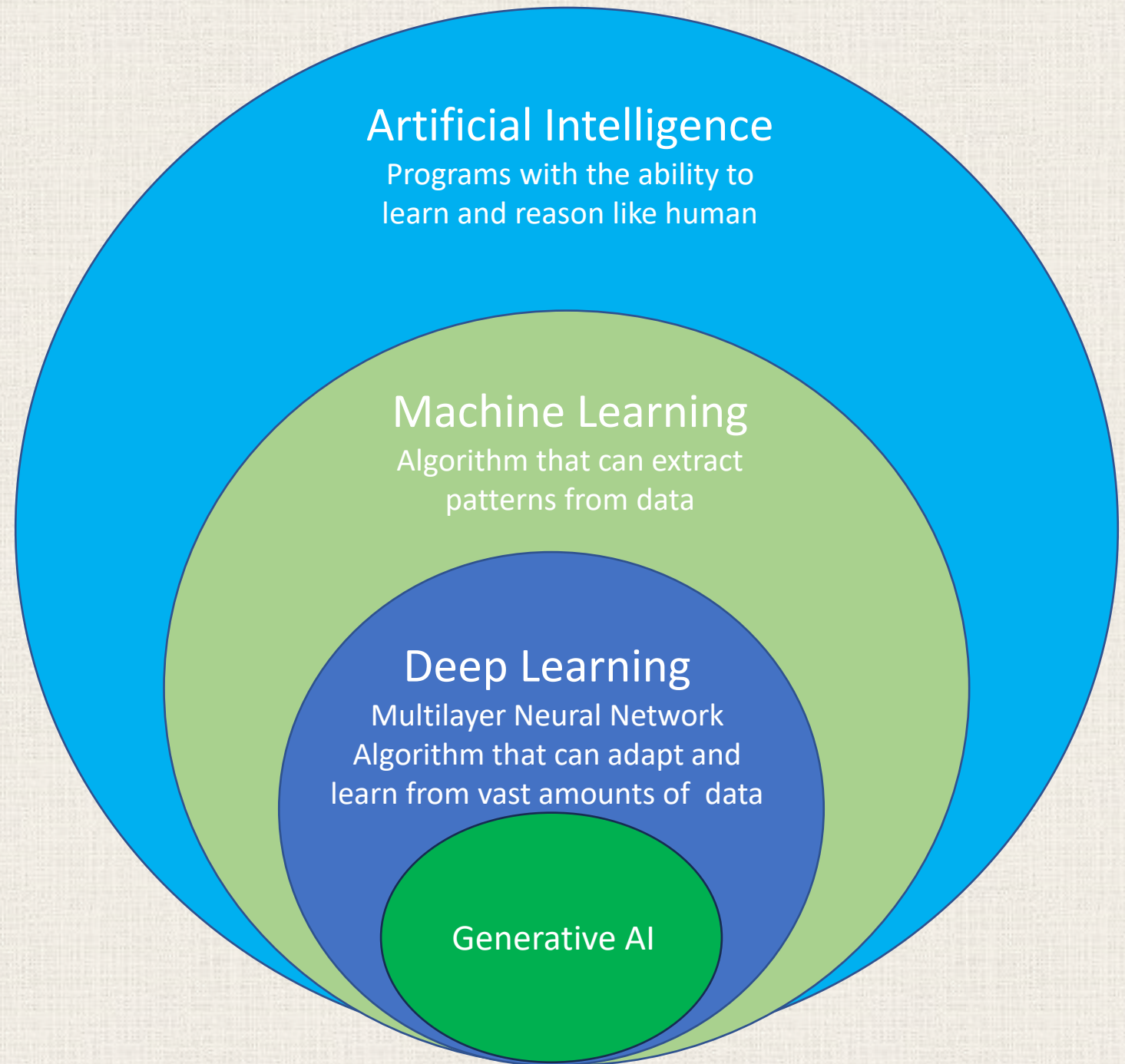
Types of Machine Learning Problems



- ***Supervised learning***: the training data comprises examples of the input feature vectors along with their corresponding target values (labels).
 - When the target values are only from a finite number of discrete numbers, the problem is called ***classification***
 - When the target values are one or more continuous variables, the problem is called **regression**
- ***Unsupervised learning***: the training data consists of a set of input feature vectors without any corresponding target values.
 - ***Clustering***: discover groups of similar examples within the data
 - ***Anomaly detection***
 - ***Density estimate***: to determine the distribution of data within the input space
 - ***Visualization***: project the data from a high dimensional space down to the two or three dimensions (e.g., ***PCA***) for the purpose of visualization

- ***Semi-supervised learning***: only a portion of the training data set are labeled, a large amount of the data are unlabeled. Unlabeled data can be used to improved the learning performance
- ***Reinforcement learning***: there is no desired target values in the training data set. The algorithm interact with a dynamic environment that provides feedback in terms of rewards and punishments. The algorithm discover the optimal output through a process of trial and error.
- ***Deep learning*** is a class of artificial neural networks algorithm that uses multiple layers to progressively extract higher level features from the raw input and eventually make decisions.
 - Typical deep neural network types include *convolutional neural networks* (CNN), *deep belief network* and *recurrent neural networks* (RNN), etc.
 - Deep learning has been used in computer vision, speech recognition, natural language processing, machine translation, drug design, medical image analysis, and board game programs, etc.

Artificial Intelligence, Machine Learning and Deep Learning



Review of Vector Calculus and Linear Algebra

Consider two column vectors $\mathbf{x}, \mathbf{y} \in \mathcal{R}^n$, a matrix $\mathbf{A} \in \mathcal{R}^{m \times n}$, and a real number $\alpha \in \mathcal{R}$:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{A}_{m \times n} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

Then,

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

$$\alpha \mathbf{x} = \begin{bmatrix} \alpha x_1 \\ \alpha x_2 \\ \vdots \\ \alpha x_n \end{bmatrix}$$

Consider two column vectors $\mathbf{x}, \mathbf{y} \in \mathcal{R}^n$, a matrix $\mathbf{A} \in \mathcal{R}^{m \times n}$, and a real number $\alpha \in \mathcal{R}$:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{A}_{m \times n} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

Then,

$$\mathbf{x}^T = [x_1 \quad x_2 \quad \cdots \quad x_n]$$

$$\mathbf{A}_{n \times m}^T = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ a_{12} & a_{22} & \cdots & a_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{mn} \end{bmatrix}$$

Consider two column vectors $\mathbf{x}, \mathbf{y} \in \mathcal{R}^n$, a matrix $\mathbf{A} \in \mathcal{R}^{m \times n}$, and a real number $\alpha \in \mathcal{R}$:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{A}_{m \times n} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

Then,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = [x_1 \quad x_2 \quad \cdots \quad x_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i \quad (\text{inner product})$$

$$\mathbf{z} = \mathbf{A}\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n a_{1i} x_i \\ \sum_{i=1}^n a_{2i} x_i \\ \vdots \\ \sum_{i=1}^n a_{mi} x_i \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix} \quad (\text{linear transformation})$$

■ A **norm** of a real vector is a function $\|\cdot\|: \mathcal{R}^n \rightarrow \mathcal{R}$ that satisfies:

- $\|\mathbf{x}\| \geq 0$, with equality if and only if $\mathbf{x} = \mathbf{0}$
- $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$
- $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (**the triangle inequality**)

■ We will typically be concerned with a few specific norms on \mathcal{R}^n :

- 1-norm: $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ (*Taxicab norm, **L1 norm***)
- 2-norm: $\|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\sum_{i=1}^n x_i^2}$ (*Euclidean norm, **L2 norm***)
- p-norm: $\|\mathbf{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}$ ($p \geq 1$)
- Infinity-norm: $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$

- We will typically be concerned with a few specific norms on \mathcal{R}^n :

- 1-norm: $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ (Taxicab norm, **L1 norm**)

- 2-norm: $\|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\sum_{i=1}^n x_i^2}$ (**Euclidean norm, L2 norm**)

- p-norm: $\|\mathbf{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}$ ($p \geq 1$)

- Infinity-norm: $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$

- Note that 1-norm and 2-norm are special cases of the p-norm, and infinite-norm is the limit of the p-norm as $p \rightarrow \infty$
- A **distance** metric between two vectors \mathbf{x} and \mathbf{y} is then defined as: $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|$

Determinant: (Laplace expansion)

- For a square matrix $\mathbf{A} \in \mathcal{R}^{n \times n}$, the determinant $|\mathbf{A}|$ is a scalar function of \mathbf{A} .

$$|\mathbf{A}| = \sum_{j=1}^n a_{kj} C_{kj}$$

- Where, k is an arbitrary row, and C_{kj} is the ***cofactor*** of the element a_{kj} .

$$C_{kj} = (-1)^{k+j} M_{kj}$$

- Where, M_{kj} is the ***minor*** associated with the element a_{kj} which is the ***determinant*** of the $(n - 1) \times (n - 1)$ matrix formed from \mathbf{A} by crossing out the k th row and the j th column.
- Same method can be applied using a column instead of a row.

Determinant: (Laplace expansion)

- ***Example:*** calculate $|A|$ of $A = \begin{bmatrix} 2 & 4 & 1 \\ 0 & 1 & 0 \\ 2 & 2 & 3 \end{bmatrix}$ using *Laplace Expansion*

$$|A| = 2 \times (-1)^2 \times \begin{vmatrix} 1 & 0 \\ 2 & 3 \end{vmatrix} + 4 \times (-1)^3 \times \begin{vmatrix} 0 & 0 \\ 2 & 3 \end{vmatrix} + 1 \times (-1)^4 \times \begin{vmatrix} 0 & 1 \\ 2 & 2 \end{vmatrix} = 4$$

$$|A| = 0 \times (-1)^3 \times \begin{vmatrix} 4 & 1 \\ 2 & 3 \end{vmatrix} + 1 \times (-1)^4 \times \begin{vmatrix} 2 & 1 \\ 2 & 3 \end{vmatrix} + 0 \times (-1)^3 \times \begin{vmatrix} 2 & 4 \\ 2 & 2 \end{vmatrix} = 4$$

Properties of Determinant: For a square matrix $\mathbf{A} \in \mathcal{R}^{n \times n}$

- $|\mathbf{A}^T| = |\mathbf{A}|$
- $|\mathbf{AB}| = |\mathbf{A}||\mathbf{B}|$
- $|\mathbf{A}^{-1}| = |\mathbf{A}|^{-1}$
- $|\alpha \mathbf{A}| = \alpha^n |\mathbf{A}|$
- $|\mathbf{A}| = \prod_i \lambda_i(\mathbf{A})$, λ_i are the eigenvalues of matrix \mathbf{A}

Eigenvalue and eigenvector:

- For a square matrix $\mathbf{A} \in \mathcal{R}^{n \times n}$, we say that a nonzero vector $\mathbf{x} \in \mathcal{R}^n$ is an eigenvector of \mathbf{A} corresponding to eigenvalue λ if $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$. Then,
 - for any $\gamma \in \mathcal{R}$, \mathbf{x} is an eigenvector of $\mathbf{A} + \gamma\mathbf{I}$ with eigenvalue $\lambda + \gamma$
 - If \mathbf{A} is invertible, then \mathbf{x} is an eigenvector of \mathbf{A}^{-1} with eigenvalue λ^{-1}
 - $\mathbf{A}^k \mathbf{x} = \lambda^k \mathbf{x}$ for any $k \in \mathcal{Z}$ (where $\mathbf{A}^0 = \mathbf{I}$ by definition)

Example:

- Find the eigenvalue of $A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$

$$\begin{aligned} |A - \lambda I| = 0 &\Rightarrow \begin{vmatrix} 2 - \lambda & 1 \\ 1 & 2 - \lambda \end{vmatrix} = 0 \Rightarrow (2 - \lambda)^2 - 1 = 0 \\ &\Rightarrow \lambda^2 - 4\lambda + 3 = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 3 \end{aligned}$$

▪ **Orthogonal matrices:**

- Two vectors \mathbf{x} and \mathbf{y} are said to be **orthogonal** if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$
- A square matrix $\mathbf{Q} \in \mathcal{R}^{n \times n}$ is said to be **orthogonal** if its columns are pairwise orthogonal unit (**orthonormal**) vectors. This definition implies that

$$\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbf{I} \Rightarrow \mathbf{Q}^T = \mathbf{Q}^{-1}$$

- Orthogonal matrices preserve inner products:

$$(\mathbf{Q}\mathbf{x})^T (\mathbf{Q}\mathbf{y}) = \mathbf{x}^T \mathbf{Q}^T \mathbf{Q} \mathbf{y} = \mathbf{x}^T \mathbf{y}$$

- They also preserve 2-norm:

$$\|\mathbf{Q}\mathbf{x}\|_2 = \sqrt{(\mathbf{Q}\mathbf{x})^T (\mathbf{Q}\mathbf{x})} = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2$$

▪ ***Eigendecomposition:***

Given a squared matrix $A \in \mathcal{R}^{n \times n}$ with n linearly independent eigenvectors (diagnosable), then,

$$A = Q\Lambda Q^{-1}$$

Where,

- $Q = [\mathbf{q}_1 \quad \cdots \quad \mathbf{q}_n]$ is a squared matrix with $\mathbf{q}_1, \dots, \mathbf{q}_n$, as its columns
- $\mathbf{q}_1, \dots, \mathbf{q}_n$ are the eigenvectors, $\lambda_1, \dots, \lambda_n$ are the eigenvalues.
- $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$

▪ ***Symmetric matrices:***

- A squared matrix A is said to be ***symmetric*** if it is equal to its own transpose ($A = A^T$)
- **Eigendecomposition:** Given a symmetric matrix $A \in \mathcal{R}^{n \times n}$, then,

$$A = Q\Lambda Q^T$$

Where,

- $Q = [\mathbf{q}_1 \quad \cdots \quad \mathbf{q}_n]$ is an orthogonal matrix with $\mathbf{q}_1, \dots, \mathbf{q}_n$, as its columns
- $\mathbf{q}_1, \dots, \mathbf{q}_n$ are the orthonormal basis of eigenvectors, $\lambda_1, \dots, \lambda_n$ are the eigenvalues.
- $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$

- ***Eigendecomposition: Example***

Let's find the eigendecomposition of the squared matrix:

$$A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

We already find the eigenvalues of this matrix as $\lambda_1 = 1, \lambda_2 = 3$.

To find the eigenvector corresponding to eigenvalue $\lambda_1 = 1$, we substitute $\lambda_1 = 1$ to $A\mathbf{x} = \lambda\mathbf{x}$,

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Which is equivalence to:

$$2x_1 + x_2 = x_1 \Rightarrow x_1 = -x_2$$

$$x_1 + 2x_2 = x_2 \Rightarrow x_1 = -x_2$$

Hence, a possible eigenvector corresponding to the eigenvalue $\lambda_1 = 1$ is:

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

The normalized eigenvector is:

$$\mathbf{x}_1 = \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix}$$

To find the eigenvector corresponding to eigenvalue $\lambda_2 = 3$, we substitute $\lambda_2 = 3$ to $A\mathbf{x} = \lambda\mathbf{x}$,

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 3 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Which is equivalence to:

$$2x_1 + x_2 = 3x_1 \Rightarrow x_1 = x_2$$

$$x_1 + 2x_2 = 3x_2 \Rightarrow x_1 = x_2$$

Hence, a possible eigenvector corresponding to the eigenvalue $\lambda_2 = 3$ is:

$$\mathbf{x}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

The normalized eigenvector is:

$$\mathbf{x}_2 = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

Since matrix A is symmetric, according to the eigendecomposition result, we have,

$$A = Q\Lambda Q^T$$

i.e.,

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

▪ ***Singular value decomposition:***

- The ***singular value decomposition*** (SVD) is a factorization of a real or complex matrix. It is the generalization of the eigendecomposition of a squared matrix to any $m \times n$ matrix.
- Specifically, every matrix $A \in \mathcal{R}^{m \times n}$ can be decomposed as:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

Where,

- $\mathbf{U} \in \mathcal{R}^{m \times m}$ and $\mathbf{V} \in \mathcal{R}^{n \times n}$ are orthogonal matrices and $\mathbf{\Sigma} \in \mathcal{R}^{m \times n}$ is a diagonal matrix with non-negative real numbers on the diagonal;
- The diagonal entries σ_i of $\mathbf{\Sigma}$ are the singular values of \mathbf{A} ;
- The columns of \mathbf{U} are called the left-singular vectors of \mathbf{A} ;
- The columns of \mathbf{V} are called the right-singular vectors of \mathbf{A} ;

- Observe that the SVD factors provide eigendecompositions for $A^T A$ and AA^T :

$$A^T A = (U \Sigma V^T)^T U \Sigma V^T = V \Sigma^T U^T U \Sigma V^T = V \Sigma^T \Sigma V^T$$

$$AA^T = U \Sigma V^T (U \Sigma V^T)^T = U \Sigma V^T V \Sigma^T U^T = U \Sigma^T \Sigma U^T$$

Thus,

- The columns of V (*the right-singular vectors* of A) are eigenvectors of $A^T A$
- The columns of U (*the left-singular vectors* of A) are eigenvectors of AA^T
- The singular values σ_i of A are the square roots of the eigenvalues of $A^T A$ (or equivalently, of AA^T)

- **Example: Eigendecomposition of the following matrix (linear algebra package in Numpy)**

$$A = \begin{bmatrix} 1 & 0 & 1 \\ -1 & -2 & 0 \\ 0 & 1 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.9258 & 0.1892 & -0.5066 \\ -0.3304 & 0.7660 & 0.3506 \\ -0.1834 & -0.6134 & 0.7877 \end{bmatrix} \begin{bmatrix} 0.8019 & 0 & 0 \\ 0 & -2.2470 & 0 \\ 0 & 0 & -0.5550 \end{bmatrix} \begin{bmatrix} 1.3182 & 0.2611 & 0.7315 \\ 0.3155 & 1.0246 & -0.2530 \\ 0.5529 & 0.8598 & 1.2425 \end{bmatrix}$$

- This matrix is not symmetric. The positive definiteness is not defined. The eigendecomposition is in the form:

$$A = Q\Lambda Q^{-1}$$

■ **Quadratic form:**

- Let $A \in \mathcal{R}^{n \times n}$ be a ***symmetric*** matrix. The expression $\mathbf{x}^T A \mathbf{x}$ is called a *quadratic form*
- ***For example,***

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, A = \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix}$$

$$\mathbf{x}^T A \mathbf{x} = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = a_{11}x_1^2 + 2a_{12}x_1x_2 + a_{22}x_2^2$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix}$$

$$\mathbf{x}^T A \mathbf{x} = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$= a_{11}x_1^2 + a_{22}x_2^2 + a_{33}x_3^2 + 2a_{12}x_1x_2 + 2a_{23}x_2x_3 + 2a_{13}x_1x_3$$

- When $A = I$, these two quadratic forms become $\mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2$

■ ***Positive (semi-)definite matrices:***

- A symmetric matrix A is ***positive semi-definite*** if for all $x \in \mathcal{R}^n$, $x^T A x \geq 0$. denoted as $A \succcurlyeq 0$.
- A symmetric matrix A is ***positive definite*** if for all nonzero $x \in \mathcal{R}^n$, $x^T A x > 0$. denoted as $A \succ 0$.
- ***Proposition:*** A symmetric matrix is ***positive semi-definite*** if and only if all of its eigenvalues are nonnegative, and ***positive definite*** if and only if all of its eigenvalues are positive.

Optimization

$$\min_{\mathbf{x}} f(\mathbf{x})$$

Subject to: $g_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m$

$$h_j(\mathbf{x}) = 0, j = 1, 2, \dots, p$$

- Suppose $f(\mathbf{x}): \mathcal{R}^n \rightarrow \mathcal{R}$. **Optimization** is about finding **extrema** (minima or maxima) (\mathbf{x}^*). It is necessary to consider the set of inputs over which we are optimizing. This set $\mathcal{X} \subseteq \mathcal{R}^n$ is called the **feasible set**.
- If \mathcal{X} is the entire domain of the function being optimized, the problem is **unconstrained** optimization. Otherwise, the problem is **constrained** and maybe much harder to solve, depending on the nature of the feasible set.

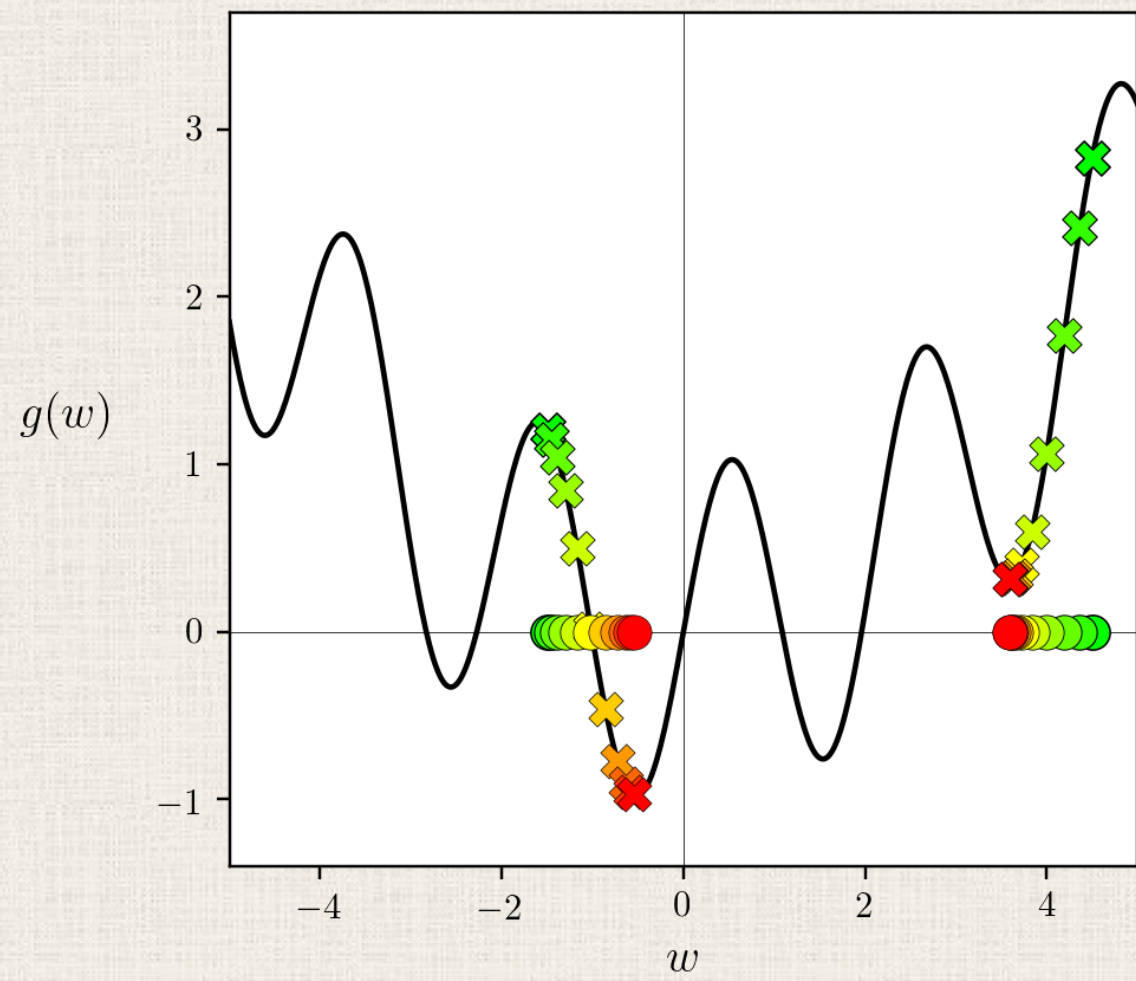
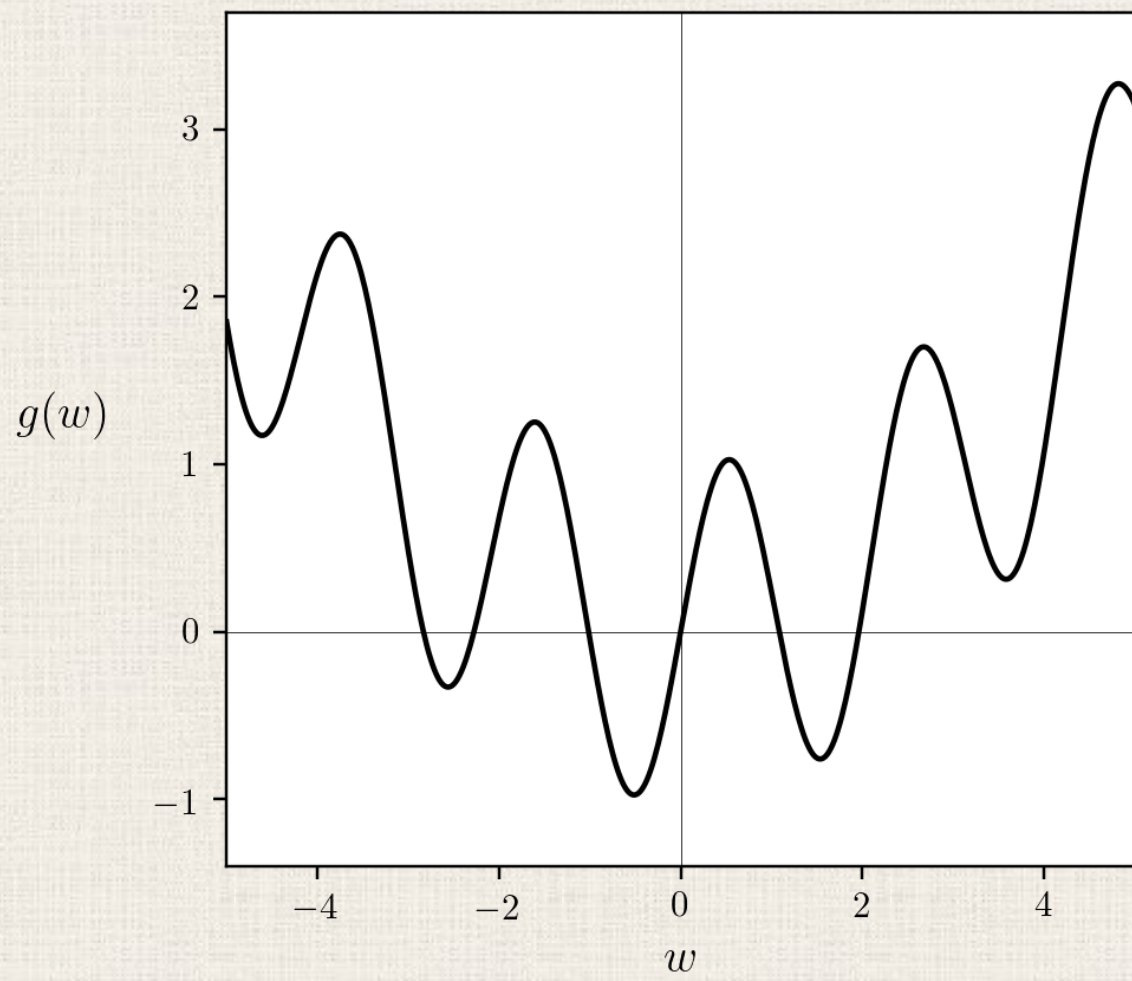
$$\min_{\mathbf{x}} f(\mathbf{x})$$

Subject to: $g_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m$

$$h_j(\mathbf{x}) = 0, j = 1, 2, \dots, p$$

- A point \mathbf{x} is said to be a **local minimum** (resp. local maximum) of f in \mathcal{X} if $f(\mathbf{x}) \leq f(\mathbf{y})$ (resp. $f(\mathbf{x}) \geq f(\mathbf{y})$) for all \mathbf{y} in some neighborhood $\mathcal{N} \subseteq \mathcal{X}$ about \mathbf{x} . Furthermore, if $f(\mathbf{x}) \leq f(\mathbf{y})$ for all $\mathbf{y} \in \mathcal{X}$, then \mathbf{x} is a **global minimum** (similarly for global maximum).
- Observe that maximizing a function f is equivalent to minimizing $-f$, so optimization problems are typically phrased in terms of **minimizing** without loss of generality.

$$g(w) = \sin(3w) + 0.1w^2$$



- **Gradients:** generalize derivatives to scalar functions of several variables.

- The **gradient** of $f: \mathcal{R}^n \rightarrow \mathcal{R}$ denoted ∇f , is given by

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

∇ is called “Del” or “Nabla”

- The gradient has the following very important property: $\nabla f(\mathbf{x})$ points in the direction of steepest **ascent** of f from \mathbf{x} .

- **Gradients:** generalize derivatives to scalar functions of several variables.
 - The **gradient** of $f: \mathcal{R}^n \rightarrow \mathcal{R}$ denoted ∇f , is given by

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

- **Example:** Let $f(x_1, x_2) = 5e^{2x_1} + 10x_2^2$. Find $\nabla f(0,1)$.

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 10e^{2x_1} \\ 20x_2 \end{bmatrix} \Rightarrow \nabla f(0,1) = \begin{bmatrix} 10 \\ 20 \end{bmatrix}$$

■ **The Hessian:**

- The **Hessian** of $f: \mathcal{R}^n \rightarrow \mathcal{R}$ is a matrix of second order partial derivatives:

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}, \quad [\nabla^2 f]_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

- When $n = 2, 3$, we have:

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}, \quad \nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_1 \partial x_3} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \frac{\partial^2 f}{\partial x_2 \partial x_3} \\ \frac{\partial^2 f}{\partial x_3 \partial x_1} & \frac{\partial^2 f}{\partial x_3 \partial x_2} & \frac{\partial^2 f}{\partial x_3^2} \end{bmatrix}$$

■ **The Hessian:**

- The **Hessian** of $f: \mathcal{R}^n \rightarrow \mathcal{R}$ is a matrix of second order partial derivatives:

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}, \quad [\nabla^2 f]_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

- **Example:** Let $f(x_1, x_2) = 5x_2 e^{2x_1} + 10x_2^2$. Find $\nabla^2 f$

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 20x_2 e^{2x_1} & 10e^{2x_1} \\ 10e^{2x_1} & 20 \end{bmatrix}$$

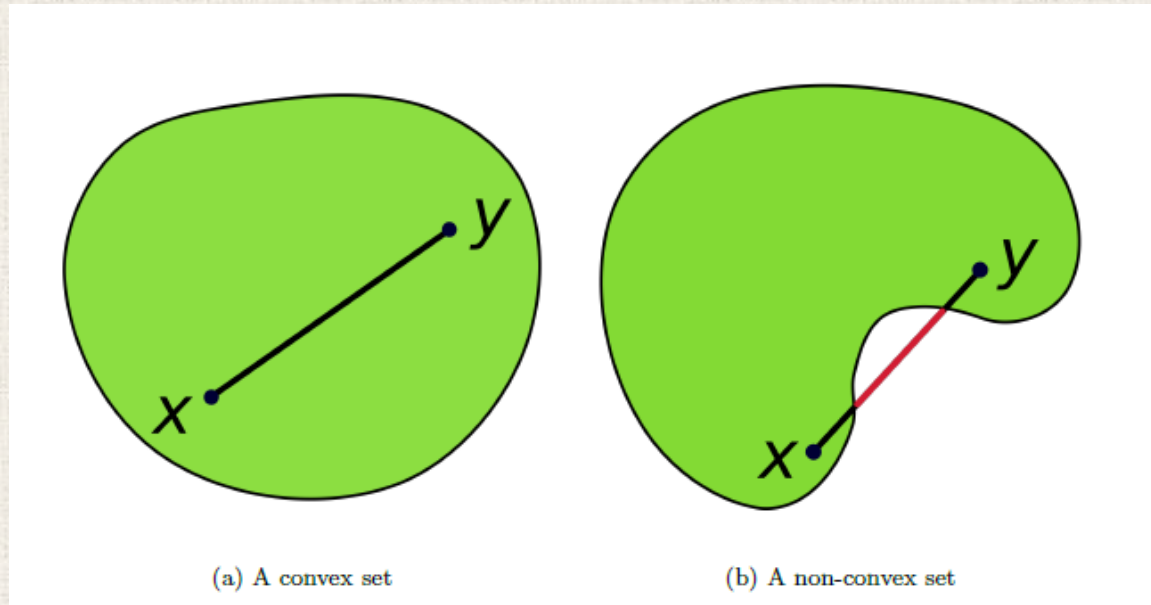
- **Proposition:** if \mathbf{x}^* is a local *extrema* (minimum or maximum) of f and f is continuously differentiable in a neighborhood of \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = 0$
- **Proposition:** if \mathbf{x}^* is a local *minimum* of f and f is twice continuously differentiable in a neighborhood of \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is **positive semi-definite**.
- **Proposition:** Suppose f is twice continuously differentiable with $\nabla^2 f$ positive semi-definite in a neighborhood of \mathbf{x}^* , and that $\nabla f(\mathbf{x}^*) = 0$. Then \mathbf{x}^* is a local minimum of f . Furthermore, if $\nabla^2 f(\mathbf{x}^*)$ is **positive definite**, then \mathbf{x}^* is a **strict** local minimum.

Convexity: Convexity of a function related closely to its minima: do they exist, are they unique?

Convex Sets:

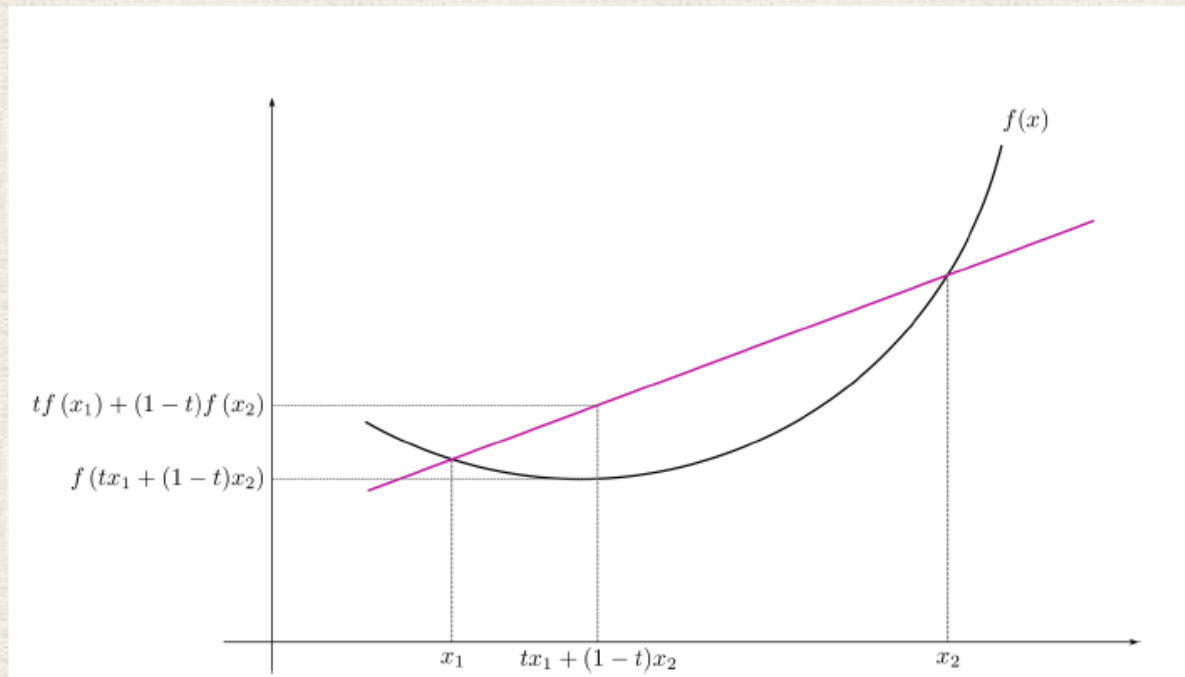
A set $\mathcal{Z} \subseteq \mathcal{R}^d$ is convex if $[tx + (1 - t)y] \in \mathcal{Z}$ for all $x, y \in \mathcal{Z}$ and all $t \in [0,1]$.

Geometrically, this means that all the points on the line segment between any two points in \mathcal{X} are also in \mathcal{X} .



Convex functions:

- A function $f: \mathcal{R}^d \rightarrow \mathcal{R}$ is **convex** if $f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$ for all $x_1, x_2 \in \text{dom}(f)$ and all $t \in [0,1]$.
- We say that f is **strictly convex** if the inequality holds strictly for all $t \in [0,1]$ and $x_1 \neq x_2$.
- Geometrically, **convexity means that the line segment between two points on the graph of f lies on or above the graph itself.**



- Examples of ***convex functions***:

$$f(x) = ax + b$$

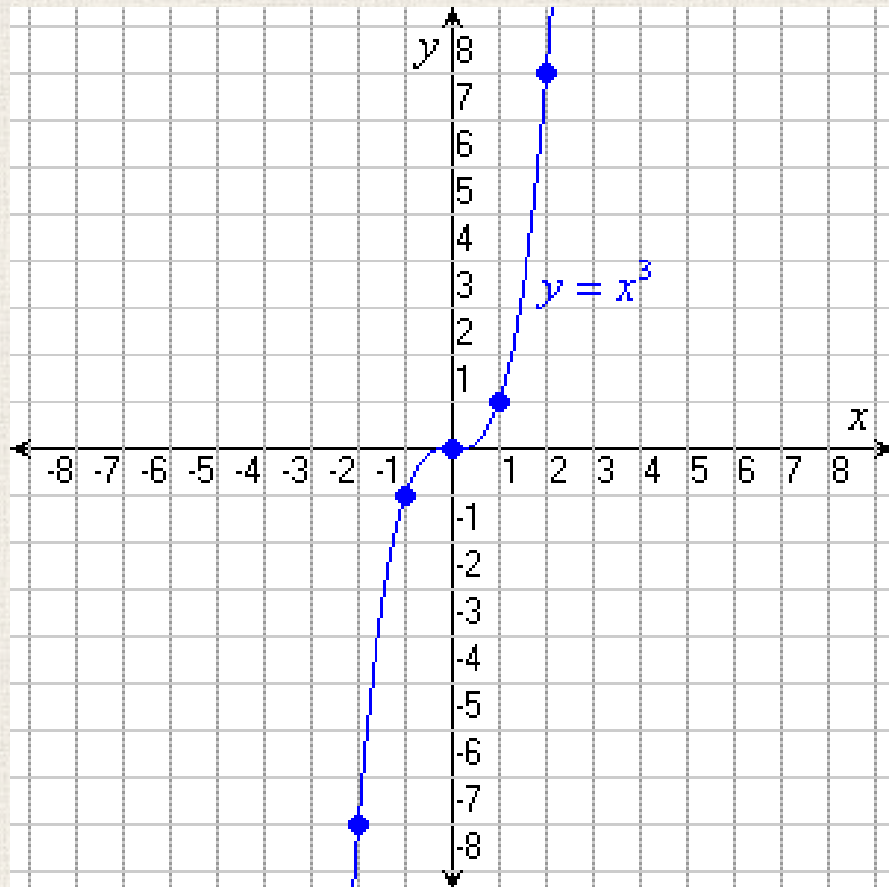
$$f(x) = x^2$$

$$f(x) = |x|$$

$$f(x) = e^x$$

- Every norm is a convex function

- How about x^3 ?



x^3 is neither convex nor concave on entire \mathcal{R} .

However, it is convex when $x \geq 0$

Consequences of convexity:

- **Proposition:** Let \mathcal{X} be a convex set. If f is convex, then any *local minimum* of f in \mathcal{X} is also a *global minimum*.
- **Proposition:** Let \mathcal{X} be a convex set. If f is strictly convex, then there exists at most one local minimum of f in \mathcal{X} . Consequently, if it exists, it is the *unique global minimum* of f in \mathcal{X} .

- **Convex functions:**

- **Proposition:** Suppose f is twice differentiable. Then,
 - f is convex if and only if $\nabla^2 f(\mathbf{x}) \succcurlyeq 0$ for all $\mathbf{x} \in \text{dom}(f)$.
 - If $\nabla^2 f(\mathbf{x}) \succ 0$ for all $\mathbf{x} \in \text{dom}(f)$, then f is strictly convex.
- **Proposition:** norms are convex

- **Proposition:** if f is convex and $\alpha \geq 0$, then αf is convex.
- **Proposition:** if f and g are convex, then $f + g$ is convex.
- **Proposition:** if f_1, \dots, f_n are convex and $\alpha_1, \dots, \alpha_n \geq 0$, then, $\sum_{i=1}^n \alpha_i f_i$ is convex.
- **Proposition:** if f is convex, then $g(\mathbf{x}) \equiv f(\mathbf{A}\mathbf{x} + \mathbf{b})$ is convex for any appropriately-sized \mathbf{A} and \mathbf{b}

$$\|\mathbf{x}\|_2^2? \quad \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2?$$

- **Proposition:** if f and g are convex, then $h(\mathbf{x}) \equiv \max\{f(\mathbf{x}), g(\mathbf{x})\}$ is convex

Gradient Descent Search Method to Find Local Minima

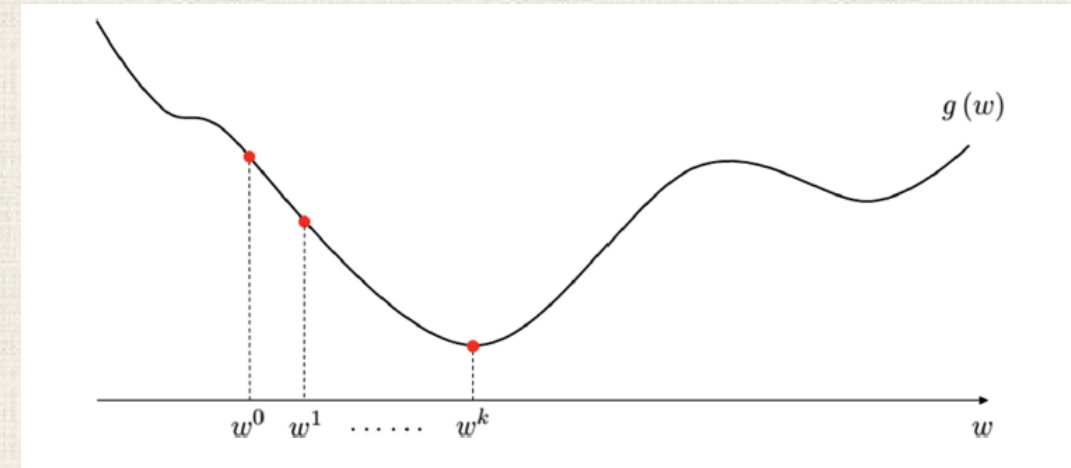
- Given an objective function $g(\mathbf{w})$, to find the value of the parameter \mathbf{w} that minimizes the objective function, we start from an initial point \mathbf{w}^0 , the local optimization method refines this initial point sequentially, pulling it downhill towards points that are lower and lower on the objective function, eventually reaching a minimum.
- This process yields a sequence of $K + 1$ points (**minimizing sequence**):

$$\mathbf{w}^0, \mathbf{w}^1, \dots, \mathbf{w}^K$$

So that,

$$g(\mathbf{w}^0) > g(\mathbf{w}^1) > \dots > g(\mathbf{w}^K)$$

$$\underset{\mathbf{w}}{\operatorname{argmin}} g(\mathbf{w})$$



Gradient descent search in one dimension

- According to **Taylor series expansion** of $g(w)$, we have,

$$g(w + \Delta w) \approx g(w) + g'(w)\Delta w$$

- Substitute Δw with $-\eta g'(w)$, where $\eta > 0$ is a constant (called **learning rate**), we have,

$$g(w - \eta g'(w)) \approx g(w) - \eta (g'(w))^2$$

- Hence, we have, when $g'(w) \neq 0$,

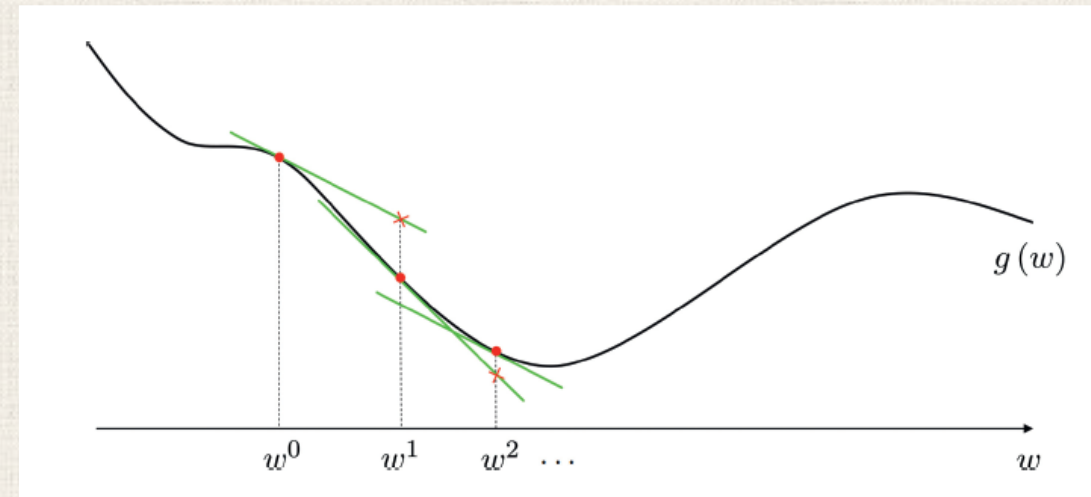
$$g(w - \eta g'(w)) < g(w)$$

- i.e., If w is updated as:

$$w^k = w^{k-1} - \eta g'(w^{k-1})$$

the minimizing sequence can be generated.

$$g(w + \Delta w) = g(w) + g'(w)\Delta w + \frac{g''(w)}{2!}(\Delta w)^2 + \frac{g'''(w)}{3!}(\Delta w)^3 + \dots$$



Gradient descent search in multi-dimension case:

- In higher dimensional case, the Taylor series of the objective function becomes:

$$g(\mathbf{w} + \Delta\mathbf{w}) \approx g(\mathbf{w}) + \nabla g(\mathbf{w})\Delta\mathbf{w}$$

Where,

$$\nabla g(\mathbf{w}) = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

is the gradient of the objective function.

- The gradient descent update of the variable becomes:

$$\mathbf{w}^k = \mathbf{w}^{k-1} - \eta \nabla g(\mathbf{w}^{k-1})$$

Example of Gradient descent Search in one-dimension:

- Consider the following objective function:

$$g(w) = (w)^2 + 2w + 2$$

Find the value of w that minimizes $g(w)$.

- Consider the following objective function:

$$g(w) = (w)^2 + 2w + 2$$

To find the value of w that minimizes $g(w)$, we take the derivative of $g(w)$ and let it equal to zero:

$$g'(w) = 2w + 2 = 0 \Rightarrow w = -1$$

$$g''(w) = 2 > 0$$

Hence, we have, $g(w)$ is minimized at $w^* = -1$ and $g(w^*) = 1$

- Consider the following objective function:

$$g(w) = (w)^2 + 2w + 2$$

Let's solve this problem using the gradient descent search method with:

$$w^0 = 0, \eta = 0.1$$

Where, $w^0 = 0$ is the value of w at step 0 ($k = 0$), i.e., the initial value of w .

In 1D case, the gradient becomes the derivative of the objective function:

$$g'(w) = 2w + 2$$

- ($k = 1$) At $w^0 = 0$, we have,

$$g'(w^0) = 2w^0 + 2 = 2; \quad g(w^0) = (w^0 + 1)^2 + 1 = 2$$

Update w :

$$w^1 = w^0 - \eta g'(w^0) = 0 - 0.1 \times 2 = -0.2$$

- ($k = 2$) At $w^1 = -0.2$, we have,

$$g'(w^1) = 2w^1 + 2 = 1.6$$

$$g(w^1) = (-0.2 + 1)^2 + 1 = 1.64$$

Update w :

$$w^2 = w^1 - \eta g'(w^1) = -0.2 - 0.1 \times 1.6 = -0.36$$

- ($k = 3$) At $w^2 = -0.36$, we have,

$$g'(w^2) = 2w^2 + 2 = 1.28$$

$$g(w^2) = (w^2 + 1)^2 + 1 = 1.41$$

Update w :

$$w^3 = w^2 - \eta g'(w^2) = -0.36 - 0.1 \times 1.28 = -0.49$$

- ($k = 4$) At $w^3 = -0.49$, we have,

$$g'(w^3) = 2w^3 + 2 = 1.02$$

$$g(w^3) = (w^3 + 1)^2 + 1 = 1.26$$

Update w :

$$w^4 = w^3 - \eta g'(w^3) = -0.49 - 0.1 \times 1.02 = -0.59$$

- ($k = 5$) At $w^4 = -0.59$, we have,

$$g'(w^4) = 2w^4 + 2 = 0.82$$

$$g(w^4) = (w^4 + 1)^2 + 1 = 1.17$$

Update w :

$$w^5 = w^4 - \eta g'(w^4) = -0.59 - 0.1 \times 0.82 = -0.67$$

- ($k = 6$) At $w^5 = -0.67$, we have,

$$g'(w^5) = 2w^5 + 2 = 0.66$$

$$g(w^5) = (w^5 + 1)^2 + 1 = 1.11$$

Update w :

$$w^6 = w^5 - \eta g'(w^5) = -0.67 - 0.1 \times 0.66 = -0.74$$

- ($k = 7$) At $w^6 = -0.74$, we have,

$$g'(w^6) = 2w^6 + 2 = 0.52$$

$$g(w^6) = (w^6 + 1)^2 + 1 = 1.07$$

Update w :

$$w^7 = w^6 - \eta g'(w^6) = -0.74 - 0.1 \times 0.52 = -0.79$$

- ($k = 8$) At $w^7 = -0.79$, we have,

$$g'(w^7) = 2w^7 + 2 = 0.42$$

$$g(w^7) = (w^7 + 1)^2 + 1 = 1.04$$

Update w :

$$w^8 = w^7 - \eta g'(w^7) = -0.79 - 0.1 \times 0.42 = -0.83$$

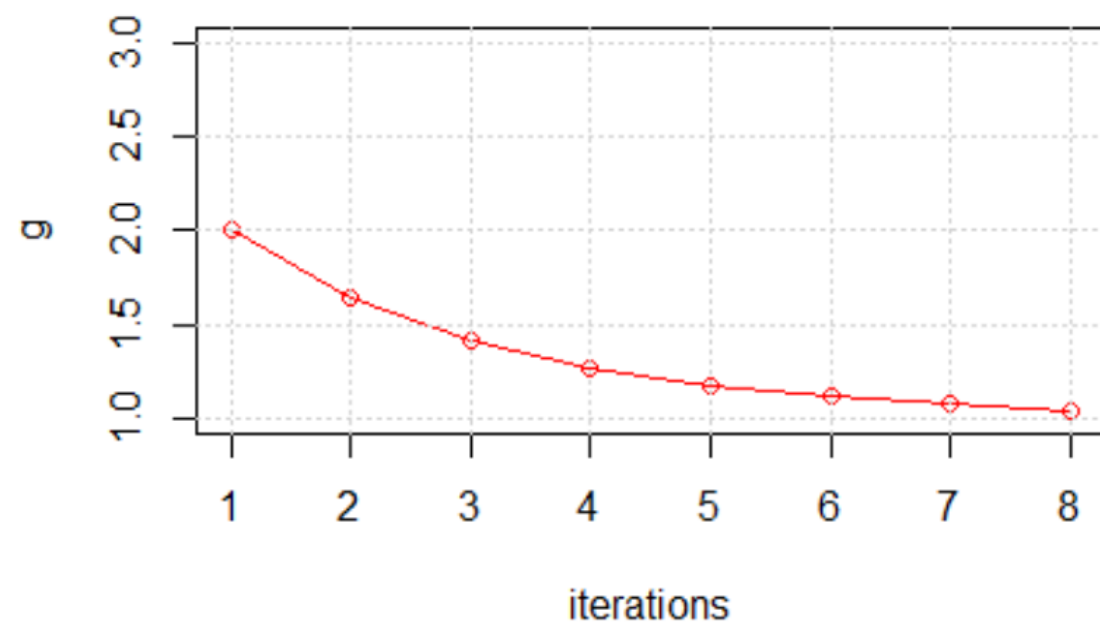
- ($k = 8$) At $w^8 = -0.83$, we have,

$$g'(w^8) = 2w^8 + 2 = 0.34$$

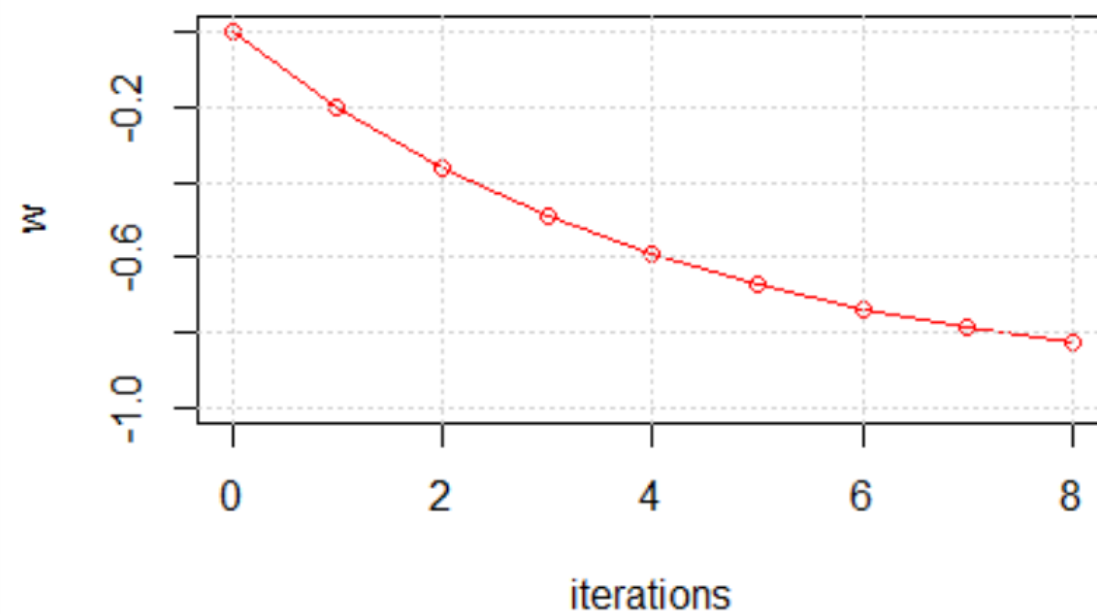
$$g(w^8) = (w^8 + 1)^2 + 1 = 1.03$$

\vdots

convergence curve



convergence curve



Example of Gradient descent Search in 2D:

- Consider the following objective function:

$$g(\mathbf{w}) = 0.5w_1^2 + w_2^2, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

Find the value of \mathbf{w} that minimizes $g(\mathbf{w})$.

- It is easy to find out that the objective function is minimized at:

$$\mathbf{w}^* = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \text{ and } g(\mathbf{w}^*) = 0$$

- The gradient of the objective function is:

$$\nabla g(\mathbf{w}) = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \end{bmatrix} = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix}$$

- The gradient descent update of the variable becomes:

$$\mathbf{w}^k = \mathbf{w}^{k-1} - \eta \nabla g(\mathbf{w}^{k-1})$$

Let's carry out the gradient descent search with:

$$\mathbf{w}^0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \text{ and } \eta = 0.2$$

- ($k = 1$), we have,

$$\nabla g(\mathbf{w}^0) = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad g(\mathbf{w}^0) = 1.5$$

$$\mathbf{w}^1 = \mathbf{w}^0 - \eta \nabla g(\mathbf{w}^0) = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 0.2 \times \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix}$$

- ($k = 2$), at $\mathbf{w}^1 = \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix}$ we have,

$$\nabla g(\mathbf{w}^1) = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 1.2 \end{bmatrix}, \quad g(\mathbf{w}^1) = 0.68$$

$$\mathbf{w}^2 = \mathbf{w}^1 - \eta \nabla g(\mathbf{w}^1) = \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix} - 0.2 \times \begin{bmatrix} 0.8 \\ 1.2 \end{bmatrix} = \begin{bmatrix} 0.64 \\ 0.36 \end{bmatrix}$$

- ($k = 3$), at $\mathbf{w}^2 = \begin{bmatrix} 0.64 \\ 0.36 \end{bmatrix}$ we have,

$$\nabla g(\mathbf{w}^2) = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 0.64 \\ 0.72 \end{bmatrix}, \quad g(\mathbf{w}^2) = 0.3344$$

$$\mathbf{w}^3 = \mathbf{w}^2 - \eta \nabla g(\mathbf{w}^2) = \begin{bmatrix} 0.64 \\ 0.36 \end{bmatrix} - 0.2 \times \begin{bmatrix} 0.64 \\ 0.72 \end{bmatrix} = \begin{bmatrix} 0.512 \\ 0.216 \end{bmatrix}$$

- ($k = 4$), at $\mathbf{w}^3 = \begin{bmatrix} 0.512 \\ 0.216 \end{bmatrix}$, we have,

$$\nabla g(\mathbf{w}^3) = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 0.512 \\ 0.432 \end{bmatrix}, \quad g(\mathbf{w}^3) = 0.178$$

$$\mathbf{w}^4 = \mathbf{w}^3 - \eta \nabla g(\mathbf{w}^3) = \begin{bmatrix} 0.512 \\ 0.216 \end{bmatrix} - 0.2 \times \begin{bmatrix} 0.512 \\ 0.432 \end{bmatrix} = \begin{bmatrix} 0.41 \\ 0.13 \end{bmatrix}$$

- ($k = 5$), at $\mathbf{w}^4 = \begin{bmatrix} 0.41 \\ 0.13 \end{bmatrix}$, we have,

$$\nabla g(\mathbf{w}^4) = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 0.41 \\ 0.26 \end{bmatrix}, \quad g(\mathbf{w}^4) = 0.101$$

$$\mathbf{w}^5 = \mathbf{w}^4 - \eta \nabla g(\mathbf{w}^4) = \begin{bmatrix} 0.41 \\ 0.13 \end{bmatrix} - 0.2 \times \begin{bmatrix} 0.41 \\ 0.26 \end{bmatrix} = \begin{bmatrix} 0.328 \\ 0.078 \end{bmatrix}$$

- ($k = 6$), at $\mathbf{w}^5 = \begin{bmatrix} 0.328 \\ 0.078 \end{bmatrix}$, we have,

$$\nabla g(\mathbf{w}^5) = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 0.328 \\ 0.156 \end{bmatrix}, \quad g(\mathbf{w}^5) = 0.06$$

$$\mathbf{w}^6 = \mathbf{w}^5 - \eta \nabla g(\mathbf{w}^5) = \begin{bmatrix} 0.328 \\ 0.078 \end{bmatrix} - 0.2 \times \begin{bmatrix} 0.328 \\ 0.156 \end{bmatrix} = \begin{bmatrix} 0.262 \\ 0.047 \end{bmatrix}$$

- ($k = 7$), at $\mathbf{w}^6 = \begin{bmatrix} 0.262 \\ 0.047 \end{bmatrix}$, we have,

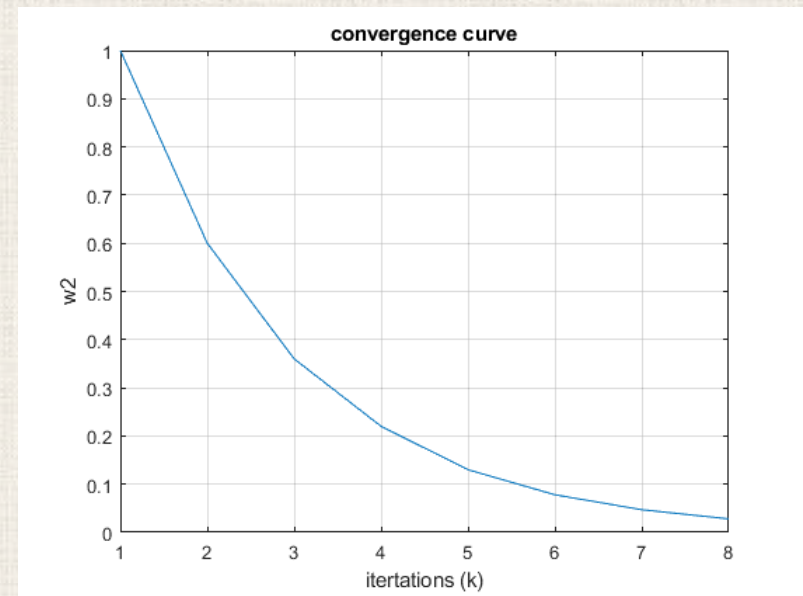
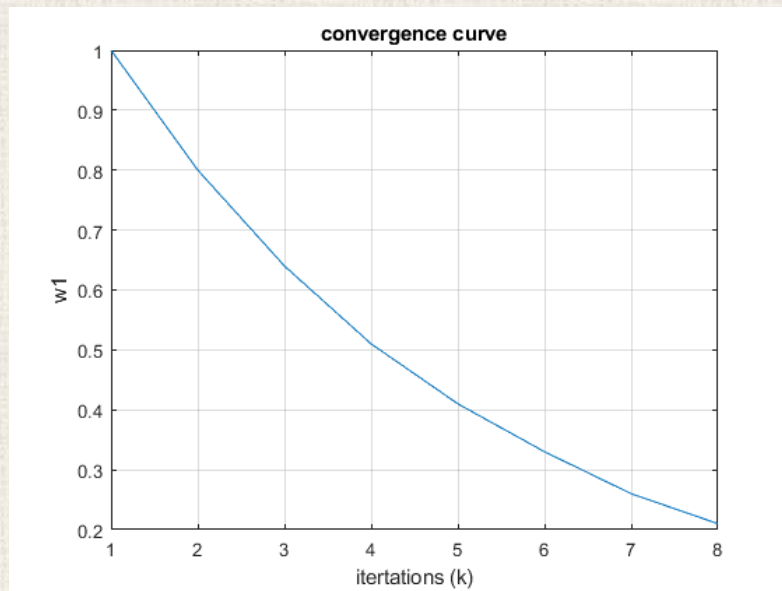
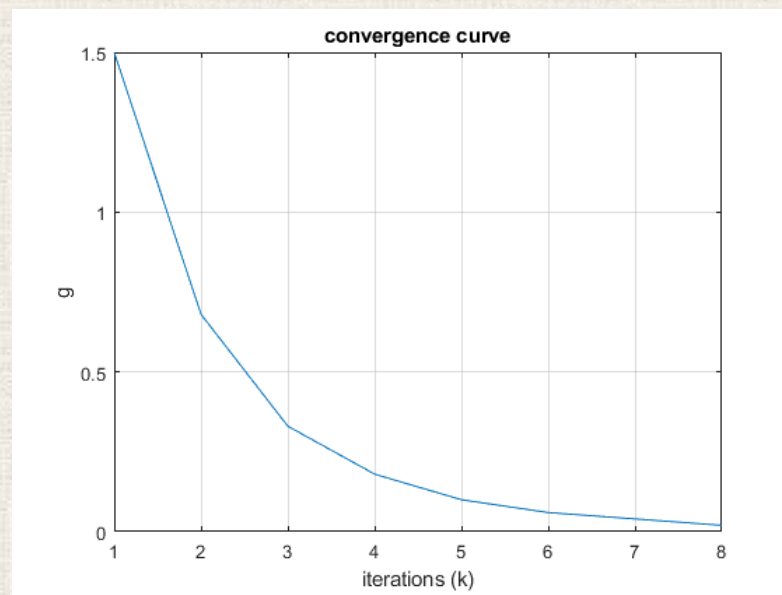
$$\nabla g(\mathbf{w}^6) = \begin{bmatrix} w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 0.262 \\ 0.094 \end{bmatrix}, \quad g(\mathbf{w}^6) = 0.037$$

$$\mathbf{w}^7 = \mathbf{w}^6 - \eta \nabla g(\mathbf{w}^6) = \begin{bmatrix} 0.262 \\ 0.047 \end{bmatrix} - 0.2 \times \begin{bmatrix} 0.262 \\ 0.094 \end{bmatrix} = \begin{bmatrix} 0.21 \\ 0.028 \end{bmatrix}$$

at $\mathbf{w}^7 = \begin{bmatrix} 0.21 \\ 0.028 \end{bmatrix}$, we have,

$$g(\mathbf{w}^7) = 0.023$$

\vdots



Convergence Criterion of Gradient Descent Search Method

Technically, if the learning rate is chosen wisely, the algorithm will halt near stationary points of the objective function. Hence, we can come up with the following criteria to end the algorithm:

- When the magnitude of the gradient is sufficiently small, i.e., for any small number $\epsilon > 0$

$$\|\nabla g(\mathbf{w}^{k-1})\|_2 \leq \epsilon$$

- When steps no longer make significant progress,

$$\|\mathbf{w}^k - \mathbf{w}^{k-1}\|_2 \leq \epsilon$$

- When corresponding evaluations no longer differ substantially,

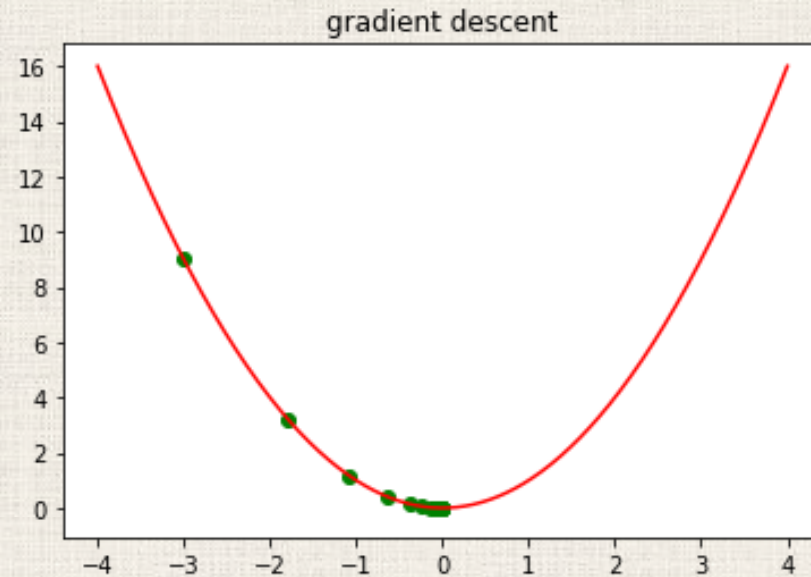
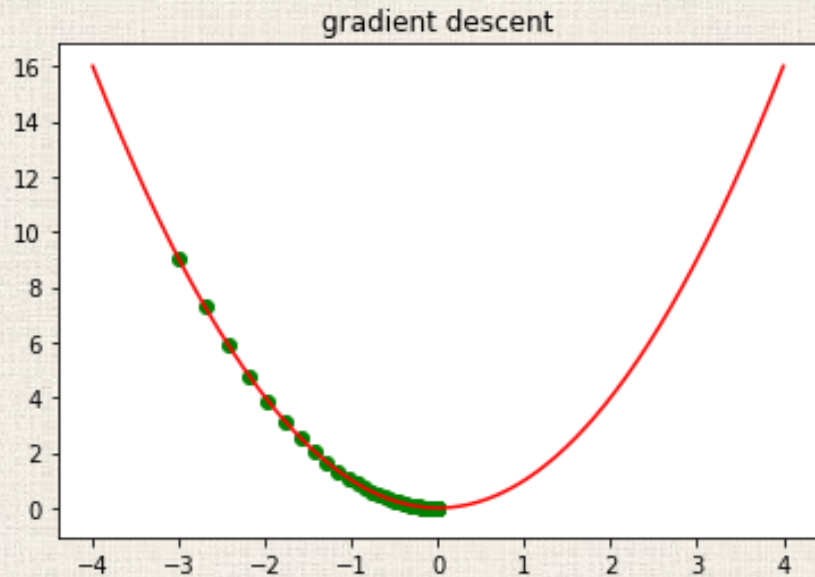
$$|g(\mathbf{w}^k) - g(\mathbf{w}^{k-1})| \leq \epsilon$$

- Run the algorithm for a fixed number of maximum iterations.

Learning rate for the Gradient Descent Search Method

➤ Fixed value learning rate

$$\eta = 10^{-n}, n = 1, 2, 3$$



➤ Diminishing learning rate

$$\eta = \frac{1}{k}, k \text{ is the number of iterations}$$

Learning rate for the Gradient Descent Search Method

- What is the consequence that the learning rate is too small?
- What is the consequence that the learning rate is too large?

Practice Example of Gradient descent Search in two-dimension:

- Consider the following objective function:

$$g(\mathbf{w}) = (w_1)^2 + 2w_1 + 2(w_2)^2 + 2 = (w_1 + 1)^2 + 2(w_2)^2 + 1$$

Use gradient descent search method to find the values of $\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ that minimizes $g(\mathbf{w})$.
You may assume different initial values of \mathbf{w} and different learning rates.