

# Classification: Linear separating planes

---

## What is a Perceptron?

A **Perceptron** is one of the simplest types of artificial neural networks and forms the building block for more complex models. It's mainly used for **binary classification problems**, where the goal is to separate data into two classes.

---

### 📌 Key Components of a Perceptron

#### 1. Inputs (features):

- Example:  $x_1, x_2, \dots, x_n$
- These are the features that describe each data point.

#### 2. Weights:

- Each input is multiplied by a corresponding weight:  $w_1, w_2, \dots, w_n$
- Weights represent the importance of each feature.

#### 3. Bias (intercept):

- A constant term  $b$  that allows the decision boundary to be shifted.

#### 4. Activation function:

- The perceptron applies a step function that outputs either 0 or 1 depending on the weighted sum:

$$\text{Output} = \begin{cases} 1 & \text{if } (w_1x_1 + w_2x_2 + \dots + w_nx_n + b) > 0 \\ 0 & \text{otherwise} \end{cases}$$

---

### ✓ How It Works

1. It takes input features.
2. It computes the weighted sum:  $z = w \cdot x + b$
3. It applies the activation function:
  - If  $z$  is greater than 0 → classify as class 1.
  - If  $z$  is less than or equal to 0 → classify as class 0.

The perceptron learns by adjusting the weights and bias during training so that it better classifies the data.

---

## Geometric Interpretation

- The perceptron tries to find a **linear boundary (plane or line)** that separates two classes.
  - It's similar to drawing a line that divides red and blue dots in the graph.
- 

## Perceptron Learning Algorithm

### Step 1 – Generate Data

- You create a dataset  $D$  consisting of points  $x$  in the  $x$ - $y$  plane.
  - Example: Points  $(x, y)$  like  $(1, 2)$ ,  $(3, 4)$ , etc.
  - This dataset will have two classes labeled by  $y = +1$  or  $y = -1$ .
- 

### Step 2 – Pick a Random Separator

- Choose a random line defined by  $ax + by + c = 0$ .
  - This line is used to label points by checking which side they fall on:
    - Points on one side get a label of  $+1$ , others  $-1$ .
  - This line acts as the **ideal separator** for training purposes.
- 

### Step 3 – Initialize Weights

- Start with an initial weight vector  $w_0$ , often  $(0, 0)$ .
- These weights define the decision boundary  $w \cdot x = 0$ .

- The algorithm iteratively adjusts these weights based on classification errors.
- 

#### Step 4 – Identify Misclassified Points

- For each point  $x_n$ , calculate  $w_t^T x_n$ .
  - If  $\text{sign}(w_t^T x_n) \neq y_n$ , then this point is misclassified.
  - Example: The algorithm checks if the current separator correctly classifies all points.
  - If it finds an error, it selects that point to update the weights.
- 

#### Step 5 – Update Weights

- The update rule is:

$$w_{t+1} = w_t + \eta y_n x_n$$

Where:

- $\eta$  is the learning rate.
- $y_n$  is the correct label (+1 or -1).
- $x_n$  is the input vector (coordinates of the point).

#### Intuition:

- If the point is misclassified, the weight vector is nudged in the direction of  $y_n x_n$ .
  - The goal is to rotate the boundary so that it better classifies the point.
- 

#### Step 6 – Repeat Until Convergence

- The algorithm keeps iterating through the dataset:
    1. Find a mistake.
    2. Correct it by updating the weights.
  - It stops when no mistakes are left—meaning all points are classified correctly.
- 

#### Step 7 – Return Final Weights

- The final weight vector  $w_{PLA}$  is the learned separator.
  - It is guaranteed to converge if the data is **linearly separable**, meaning there exists a boundary that separates the two classes without errors.
- 

## ✓ Geometric Understanding

- Each iteration is adjusting the angle between the current weight vector  $w_t$  and the input vector  $x_n$ .
  - The algorithm tries to make  $w_t$  align better with  $x_n$  when it's classified wrongly.
  - Correcting the angle through iterative updates helps to find a boundary that separates the classes.
- 

## ✓ Key Takeaways

1. **Simple yet powerful** – The perceptron learns from mistakes to find a separating boundary.
  2. **Iterative process** – Weights are updated gradually, correcting errors as they occur.
  3. **Converges with linearly separable data** – It only guarantees success when a perfect separating line exists.
  4. **Foundation for other models** – Understanding perceptrons is crucial before diving into more complex models like SVMs or neural networks.
- 

## ✓ Limitations

- Can only solve problems that are **linearly separable**.
  - Cannot solve complex patterns like the XOR problem without multiple layers.
  - Sensitive to learning rate and initial weights.
- 

## ✓ Applications

- Simple classification problems like spam detection, image recognition (basic), and pattern matching.
- 

## ✓ Summary

The perceptron is a foundational algorithm in machine learning:

- ✓ It's a linear classifier
- ✓ It learns by updating weights based on classification errors
- ✓ It works well when the data can be separated by a straight line
- ✓ It introduces core concepts like feature weighting and activation

## Pocket Algorithm – Explanation

The **Pocket Algorithm** is an extension of the Perceptron Algorithm that works even when the data is **not linearly separable**. It was designed to handle cases where no perfect separating line exists, unlike the standard perceptron which assumes perfect separability.

---

### 📌 Why is it needed?

- The regular **Perceptron Algorithm** assumes that a separating hyperplane exists that can classify all points correctly.
  - But in real-world datasets, data may not be perfectly separable — some overlap or noise can prevent a flawless boundary.
  - The **Pocket Algorithm** solves this by keeping track of the best solution seen so far, instead of waiting to find a perfect solution.
- 

### ✓ How it works – Intuition

- It still iterates through the dataset and updates the weights like the perceptron.
- However, whenever an update is made, it checks if this new weight vector results in fewer misclassifications than the best one found so far.
- If it is better, it stores this as the "best solution in the pocket".

- Even if the algorithm continues making updates, it always “remembers” the best separator it has found so far.
- 

## ✓ Steps of the Pocket Algorithm

1. **Initialize weights**  $w$  randomly or as zeros.
  2. Set  $w_{\text{pocket}} = w$ , and calculate the number of misclassified points with this  $w$ .
  3. **Iterate:**
    - Find a misclassified point  $x_n$ .
    - Update  $w$  using the standard perceptron rule:
$$w \leftarrow w + \eta y_n x_n$$
    - Check the new  $w$ :
      - If it has **fewer misclassifications** than  $w_{\text{pocket}}$ , update the pocket:
$$w_{\text{pocket}} \leftarrow w$$
  4. Repeat the process for a fixed number of iterations or until a stopping criterion is met.
  5. **Return**  $w_{\text{pocket}}$  as the final solution.
- 

## ✓ Key Features

- ✓ Works even when the data is **not linearly separable**
  - ✓ Tracks the best solution instead of waiting for perfection
  - ✓ Allows the algorithm to avoid bad solutions caused by noisy data
  - ✓ Makes it practical for real-world scenarios
- 

## ✓ Graphical Intuition

- Imagine a pocket where you keep your best pair of shoes while trying on new ones.
- Even if some new options don't work better, you still have the best pair saved.

- Similarly, the algorithm always stores the best-performing weight vector, protecting against bad iterations.
- 

## Comparison with Perceptron

Feature	Perceptron	Pocket Algorithm
Linearly separable data required	Yes	No
Tracks best solution		
Practical for noisy datasets	No	Yes
Final output	First perfect separator	Best separator found after updates

---

## Use Cases

- When data has noise or overlaps between classes.
  - When it's impossible to find a perfect decision boundary.
  - When you want robustness in classification tasks.
- 

# Using Probability Bounds

---

## Key Elements Explained

### 1. Population (Bin):

- The total collection of items is huge (like all fruits in the basket).
- The probability that a randomly chosen item is an orange is  $\mu$ , and the probability it's an apple is  $1-\mu$ .
- But  $\mu$  is unknown!

### 2. Sampling:

- You randomly pick  $N$  fruits from the basket.
- In the sample, you count how many are oranges and how many are apples.

- From this, you calculate the fraction of oranges  $v$  and apples  $1-v$ .
- Now,  $v$  is known because you observed it in the sample!

### 3. Probability Bounds:

- Even though you only have a sample, you want to understand how close  $v$  is to the true proportion  $\mu$ .
- If the sample size  $N$  is large enough,  $v$  is **probably close to  $\mu$**  — meaning you can trust that your sample reflects the true population.
- The slide introduces the idea of a bound, written as:

$$P(|v - \mu| > \epsilon) \leq \text{bound}$$

- This means that the probability that your sample proportion  $v$  differs from the true proportion  $\mu$  by more than some acceptable error  $\epsilon$  is bounded or limited.

### 4. Inequality Names:

- Chebyshev's inequality, Chernoff bound, and Hoeffding's inequality are mathematical tools that formalize how this bound works.
- They provide a way to say, "Even with uncertainty, here's how far off we might expect to be."

## Intuition

- ✓ You can't know the true population exactly, but by sampling a sufficient number of items, you can estimate it with high confidence.
- ✓ The sample proportion  $v$  acts as your observed estimate of the unknown  $\mu$ .
- ✓ Probability bounds give you mathematical guarantees on how far off your estimate is likely to be.
- ✓ These bounds are essential in machine learning, statistics, and data science to make decisions under uncertainty.

## Example

- Suppose you have a huge basket with oranges and apples.

- You don't know how many are oranges — that's  $\mu$ .
  - You sample  $N = 100$  fruits and find  $v = 0.6$  (60% are oranges).
  - You want to know: "How confident can I be that the real proportion  $\mu$  is close to 0.6?"
  - Using a probability bound, you can say something like: "With 95% confidence, the true proportion is within 0.1 of the sample estimate."
- 

## Generalization Error

- We want to understand how well our model performs on unseen data.
  - Two types of error are used:
    - **In-sample error ( $E_{in}$ )**: How well the model performs on training data.
    - **Out-of-sample error ( $E_{out}$ )**: How well the model performs on unseen data.
  - The goal is to ensure that the difference  $|E_{in} - E_{out}|$  is small.
- 

## Probability Bounds

- For a **single hypothesis** (i.e., a fixed hyperplane in d-dimensional space), the probability that the difference between  $E_{in}$  and  $E_{out}$  is greater than some error  $\epsilon$  is bounded by:

$$P(|E_{in} - E_{out}| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

→ The error decreases exponentially as the number of samples  $N$  increases.

- For **multiple hypotheses** (many possible planes), the bound becomes:

$$P(|E_{in}(g) - E_{out}(g)| > \epsilon) \leq 2M e^{-2\epsilon^2 N}$$

where  $M$  is the number of hypotheses being considered.

---

## Reducing $M$ – The Key Idea

- If you consider every possible separating plane,  $M$  becomes infinitely large → the bound worsens.

- By grouping hypotheses or applying structure, you can reduce M from infinity to something more manageable like  $2^N$ .
  - Further simplification leads to considering  $f(Nd_{vc})$ , where  $d_{vc}$  is the **VC dimension**.
- 

## VC Dimension

- $d_{vc}$  measures the complexity of the model.
  - It's related to the number of features or parameters that the model has.
  - A higher  $d_{vc}$  means a more complex model that can fit more patterns but also risks overfitting.
  - It can be determined through techniques like **cross-validation**.
- 

## How Much Data is Enough?

- The rule of thumb derived from theory is:

$$N \geq 10d_{vc}$$

→ You need at least 10 times more data points than the complexity of the model to ensure low generalization error.

- Example: If  $d_{vc} = 100$ , then  $N \geq 1000N$ .
- 

## Final Equation

The bound on error is expressed as:

$$E_{in}(g) = E_{out}(g) + \sqrt{\frac{8}{N} \log \left( \frac{4f(Nd_{vc})}{\delta} \right)}$$

where:

- N = number of samples
  - $d_{vc}$  = model complexity
  - $\delta$  = probability tolerance (how confident we want to be)
-

## Key Takeaways

- ✓ More complex models (higher  $d_{vc}$ ) require more data to generalize well.
  - ✓ Probability bounds help us understand the relationship between data size and error.
  - ✓ VC dimension is a practical way to measure model complexity.
  - ✓ The rule  $N \geq 10 d_{vc}$  gives a guideline for how much data is needed.
  - ✓ There's a trade-off: increasing complexity reduces training error but risks increasing generalization error.
-