

# Ultimate Cheatsheets for ML/DL

Mattia Bortolaso

October 2025

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Artificial Intelligence (AI)	3
1.1.1	AI Introduction	3
1.1.2	Types of AI	4
1.2	Machine Learning (ML):	5
1.2.1	ML Introduction:	5
1.2.2	Types of Machine Learning:	6
1.2.3	Supervised Learning	7
1.2.4	Unsupervised Learning	7
1.2.5	Reinforcement Learning	8
<b>2</b>	<b>Algorithms</b>	<b>10</b>
2.1	Regression Methods	10
2.2	Classification Methods	11
2.3	Clustering	12
2.4	Dimensionality Reduction	13
<b>3</b>	<b>Models</b>	<b>16</b>
3.1	Neural Network Architectures	16
3.1.1	Feedforward Networks	16
3.1.2	Convolutional Neural Networks (CNN)	16
3.1.3		16
3.2	Techniques	16
3.2.1	Activation Functions (ReLU, GeLU, etc.)	16
3.2.2	Normalization Layers (BatchNorm, LayerNorm)	16
3.2.3	Dropout & Regularization	16
3.3	Training and Optimization	16
3.3.1		16
3.3.2		16
<b>4</b>	<b>Metrics &amp; how to read it</b>	<b>17</b>
4.1	Accuracy & Loss	17
4.1.1	Accuracy	17
4.1.2	Loss	18
4.2	Precision & Recall	18
4.2.1	Precision	19
4.2.2	Recall	19
4.2.3	Confusion Matrix	20
4.2.4	Precision–Recall Trade-off and Curve Interpretation	21

4.3	F1-Score . . . . .	21
<b>5</b>	<b>Implementation &amp; Code Patterns</b>	<b>23</b>
<b>6</b>	<b>Advanced Topics</b>	<b>24</b>

# Chapter 1

## Introduction

This document serves as a comprehensive cheat sheet for studying Machine Learning and Deep Learning. It provides a structured summary of key algorithms, models, technologies, and optimization techniques, covering both foundational and state-of-the-art approaches.

The goal is to offer a concise yet rigorous reference for students, researchers, and practitioners seeking to review or deepen their understanding of modern artificial intelligence methodologies, including theoretical principles, implementation strategies, and code-level insights.

### 1.1 Artificial Intelligence (AI)

In this section we talk about Artificial Intelligence, in this chapter we talk in general about AI as we dive more deep later.

#### 1.1.1 AI Introduction

**Artificial Intelligence (AI)** refers to the ability of machines or computers to mimic the way humans think and make decisions. AI enables machines to think like humans or to replicate certain functions of the human brain. In simple words, **AI is when we enable computers to think.**

Artificial Intelligence allows computers to understand, analyze data, and make decisions without human help or interaction. We use AI systems every day. A few common examples are *Siri*, *Alexa*, and *ChatGPT*.

Other categories where AI is used include:

- Virtual Assistants
- Social Media Algorithms
- Online Shopping Recommendations
- Predictive Text and Autocorrect
- Healthcare Diagnostics
- Language Translation Services
- Fraud Detection in Banking

AI enables computers to become intelligent with numbers and rules, performing calculations at incredible speed and with perfect accuracy. On the other hand, humans possess not only reasoning but also emotions, creativity, and the ability to adapt to complex and unpredictable situations.

### 1.1.2 Types of AI

Artificial Intelligence is divided based on two main categorization — based on capabilities and based on functionality of AI.

The following image illustrates these types of AI:

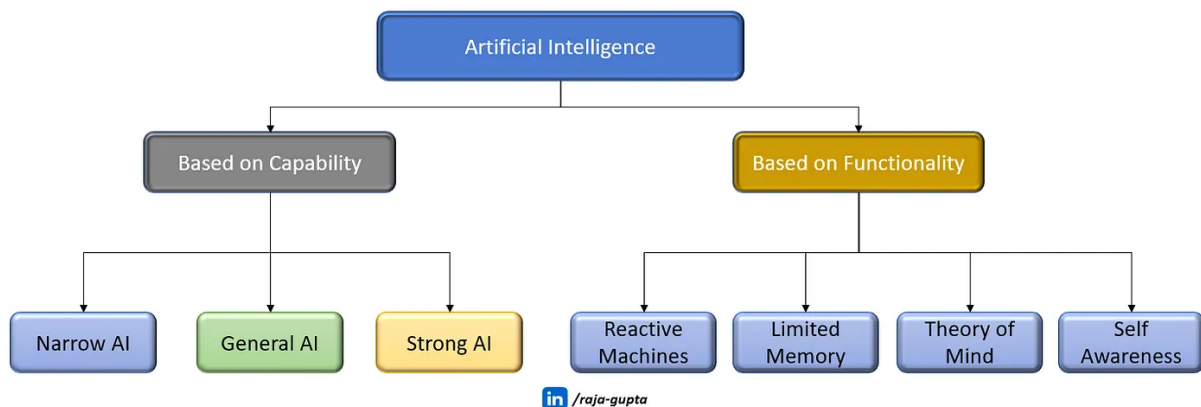


Figure 1.1: Types of AI

Now let's dive in the details of Narrow AI, General AI and Super AI.

- **Narrow AI:** Narrow AI, also known as Weak AI, refers to **artificial intelligence systems that are designed and trained for a specific task** or narrow set of tasks. Narrow AI is focused on performing a single task extremely well, but it cannot perform beyond its field or limitations.
- **General AI:** General AI, also known as Strong AI or artificial general intelligence (AGI), **can understand and learn any intellectual task that a human being can.**
- **Super AI:** Super AI represents a degree of intelligence in systems where **machines have the potential to exceed human intelligence**, outperforming humans in tasks and exhibiting cognitive abilities.

We can also divide types of AI based on functionality:

1. **Reactive Machines:** Reactive machines are AI **systems that have no memory**. These systems operate solely based on the present data, taking into account only the current situation. They can perform a narrowed range of pre-defined tasks.
2. **Limited Memory:** As the name indicates, Limited Memory AI can take informed and improved decisions by looking at its past experiences stored in a temporary memory.

This AI doesn't remember everything forever, but it uses its short-term memory to learn from the past and make better decisions for the future.

**A good example of Limited Memory AI is Self-driving cars.** The AI system in self-driving car utilizes recent past data to make real-time decisions. For instance, they employ sensors to recognize pedestrians, steep roads, traffic signals, and more, enhancing their ability to make safer driving choices. This proactive approach contributes to preventing potential accidents.

3. **Theory of Mind:** This is similar to Super AI — We should pray that we don't reach the state of AI, where **machines have their own consciousness and become self-aware.**

Self-aware AI systems will be super intelligent, and will have their own consciousness, sentiments, and self-awareness. They will be smarter than human mind.

4. **Self-Aware AI:** This is similar to Super AI — We should pray that we don't reach the state of AI, where machines have their own consciousness and become self-aware.

Self-aware AI systems will be super intelligent, and will have their own consciousness, sentiments, and self-awareness. They will be smarter than human mind.

As shown in movie "I, Robot," an AI system named VIKI becomes self-aware and starts making decisions to "protect" humanity in a controversial way.

Similar to Theory of Mind, Self-aware AI also does not exist in reality. Many experts, for example Elon Musk and Stephen Hawking have consistently warned us about the evolution of AI.

## 1.2 Machine Learning (ML):

In this chapter we are talking about Machine Learning, we first start explaining machine learning from a kids prospective and later we add more details.

### 1.2.1 ML Introduction:

Imagine we want to **enable the robot to identify several animals.**

To do so, we will show him pictures of various dogs, cats, bunnies and other animals and **label each picture with the name of the animal.** We train the robot to identify animals based on size, color, body shape, sound etc.

Once the training is completed, the robot will be able to identify these animals we trained him for.

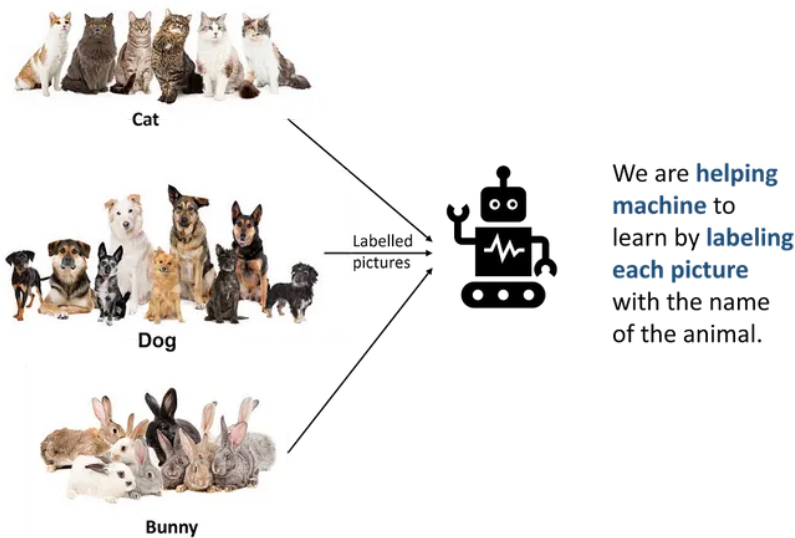
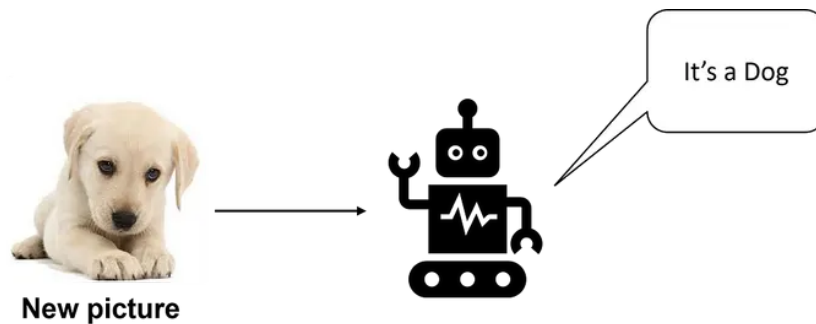


Figure 1.2: Machine Learning scheme



The machine identified the dog even though this picture is not exactly same as any of the dog's picture shown earlier.

**The machine learned to identify a dog.**

Figure 1.3: After training scheme

All dogs do not look alike. However, once robot has seen many pictures of dogs, **it can identify any dog even if it does not exactly look like a specific picture**. We need to show lots of pictures of dog to the robot. More pictures it sees, more efficient it will be.

In simple words Machine Learning is teaching a robot or a machine by giving a lot of examples to work with and learn.

### 1.2.2 Types of Machine Learning:

Machine Learning can be categorized into three main types:

1. Supervised Learning.
2. Unsupervised Learning.

### 3. Reinforcement Learning.

Each type serves different purposes and involves different approaches to learning from data. Let's have a close look into all these types.

#### 1.2.3 Supervised Learning

When we trained our robot by showing pictures of animals, we labelled each picture with the name of the animal. So, **we acted as a teacher to him**. We first told him how does a dog or a cat look like and then only he was able to identify them.

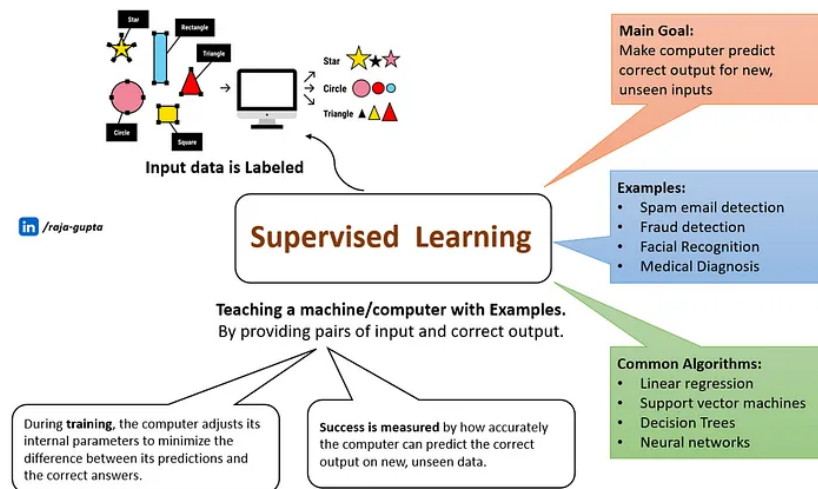


Figure 1.4: Supervised Learning summarized

Supervised Learning is widely used in this field: *Email Spamming Filtering, Image Classification, Facial Recognition, Financial Fraud Detection, Speech Recognition*.

#### 1.2.4 Unsupervised Learning

Let's understand this from a kid's school example. When kids go to their class first day, they meet lots of classmates. At first all classmates are same to them. But with time, they themselves categorized them in different groups:

- They find some classmates very good and want to be friend with them.
- They find some rude or irritating and want to avoid them.
- They find some very good in sports and want to be in the same team as they are.
- and so on ...

**When kids categorized their classmates, nobody told them how to do that.** They did that without anyone's help. — This is how unsupervised learning works.

Let's take a proper machine learning example. Imagine we showed lots of pictures of dogs, cats, bunnies etc. **without any label** to our robot and told him — “*I'm not going to tell you which one is which. Go explore and figure it out*”.

The robot starts to look at these animals, noticing things such as their fur, size, and how they move. It doesn't know their names yet, but it's trying to find patterns and differences on its own.



After exploring, the robot might notice that:

- Some animals have long ears (bunnies)
- Some animals have soft fur and tail (cats)
- Some animals have wagging tails (dogs)

In the end, the robot **might not know the names of the animals**, but it can say that **“These animals are similar in some ways, and those are different in other ways.”** — This is Unsupervised Learning.

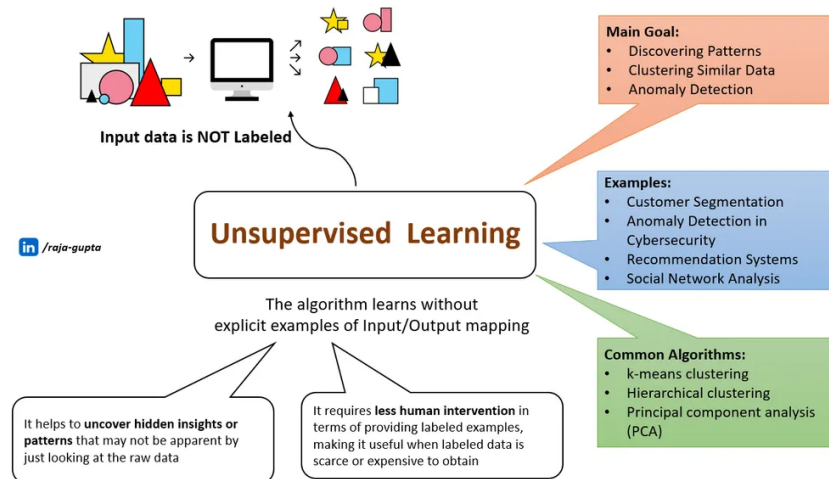


Figure 1.5: Unsupervised Learning summarized

Unsupervised Learning is widely used in this field: *Clustering Customer Segmentation, Anomaly Detection in Cybersecurity, Recommendation System.*

### 1.2.5 Reinforcement Learning

Imagine teaching a dog a new trick — you **reward it with a treat when it does the trick correctly** and give no treat when it doesn't. Over time, the dog **learns to perform the trick to get more treats**.

Similarly, Reinforcement Learning is:

- Training a computer to make a decision.
- By rewarding good choices and punishing bad ones.
- Just as you might train a dog with treats for learning tricks.

In reinforcement learning, there's an agent (for example a robot or computer program) that interacts with an environment. Let's take an example of teaching a computer program to play a game, for example chess.

- In this case, computer program is agent and chess game is the environment.
- The computer program can make different moves in the game, such as moving a chess piece.

- After each move, it receives feedback (reward or penalty) based on the outcome of the game.
- If it loses the game, it receives a negative reward, or a 'penalty'.
- Through trial and error, the program learns which moves lead to the best rewards, helping it figure out the best sequence of moves that leads to winning the game.

Reinforcement learning is powerful because it **allows machines to learn from their experiences** and make decisions in complex, uncertain environments — similar to how we learn from trial and error in the real world.

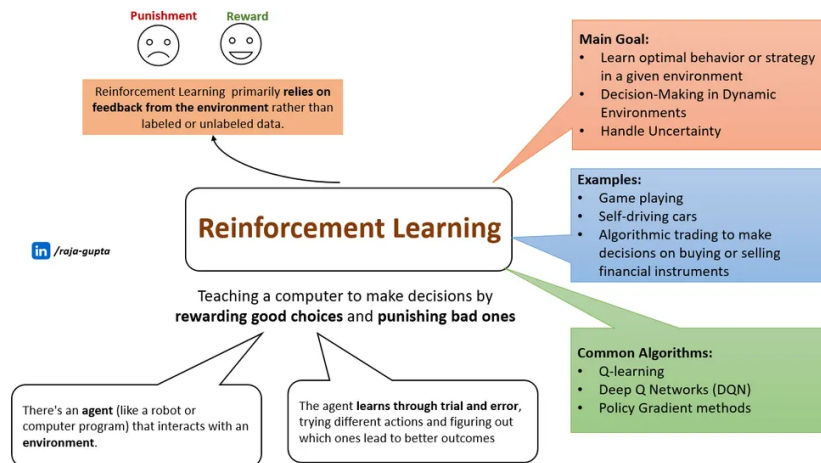


Figure 1.6: Reinforcement Learning summarized

# Chapter 2

## Algorithms

### 2.1 Regression Methods

Regression is one of the fundamental tasks in supervised learning. Its goal is to predict a continuous-valued output based on one or more input variables. Formally, given a dataset of input–output pairs  $\{(x_i, y_i)\}_{i=1}^N$  where  $x_i \in \mathbb{R}^d$  represents the input features and  $y_i \in \mathbb{R}$  represents a continuous target variable, a regression model seeks to learn a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  that approximates the underlying relationship between inputs and outputs.

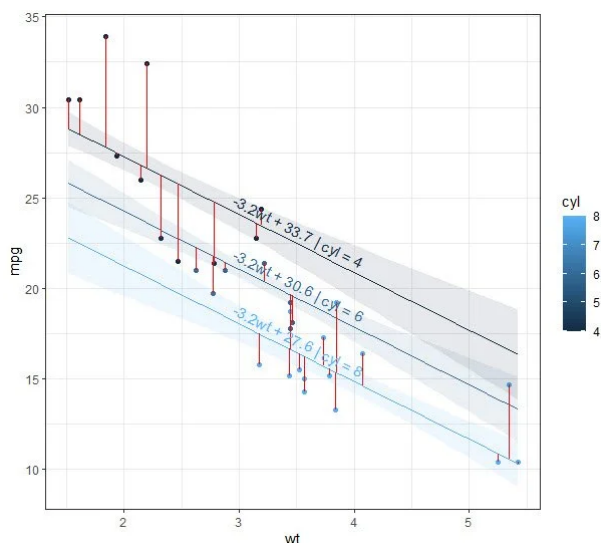


Figure 2.1: Regression examples

In a typical regression problem, the input can represent any structured or unstructured data — such as physical measurements, numerical attributes, or encoded representations of images, audio signals, or text. The model processes these inputs and outputs a single real number or, in the case of multivariate regression, a vector of continuous values.

The simplest and most widely known regression technique is *linear regression*, where the model assumes a linear relationship between the input features and the target:

$$\hat{y} = w^\top x + b,$$

where  $w$  is a weight vector and  $b$  is a bias term. More complex regression models, including polynomial regression, decision trees, and neural networks, can capture nonlinear dependencies between inputs and outputs.

The performance of a regression model is commonly evaluated using metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), or Mean Absolute Error (MAE), which quantify how close the predicted values  $\hat{y}_i$  are to the true targets  $y_i$ .

In summary, regression models map numerical representations of real-world phenomena to continuous outcomes, forming the basis for many predictive modeling tasks, from estimating housing prices and forecasting stock trends to predicting physical properties in scientific domains.

## 2.2 Classification Methods

Classification is another central task in supervised learning, where the objective is to assign each input to one of several discrete categories. Formally, given a dataset of labeled examples  $\{(x_i, y_i)\}_{i=1}^N$  where  $x_i \in \mathbb{R}^d$  represents the feature vector and  $y_i \in \{1, \dots, K\}$  denotes the class label, a classification model learns a function  $f : \mathbb{R}^d \rightarrow \{1, \dots, K\}$  or, equivalently, a probability distribution over classes  $P(y|x)$ .

The model ingests an input vector—which may represent text, an image, an audio signal, or any numerical encoding—and outputs either a class label (hard classification) or a vector of probabilities indicating the model’s confidence for each class (soft classification). The most common approach for probabilistic outputs is to apply a softmax transformation to the model’s final layer, ensuring that all predicted probabilities sum to one.

Several families of models can perform classification:

- **Linear Models:** Logistic Regression and Linear Discriminant Analysis (LDA) are among the simplest and most interpretable classifiers. They model class boundaries using linear decision surfaces and can provide calibrated probability estimates.
- **Support Vector Machines (SVM):** These models seek an optimal separating hyperplane that maximizes the margin between classes. Kernelized SVMs can capture nonlinear decision boundaries.
- **Tree-Based Models:** Decision Trees, Random Forests, and Gradient Boosting Machines (such as XGBoost or LightGBM) partition the feature space hierarchically. They often yield strong predictive performance and can handle heterogeneous data types.
- **Probabilistic Models:** Naïve Bayes and Gaussian Mixture Models (GMMs) represent classes using explicit probability distributions, allowing inference via Bayes’ rule.
- **Neural Network Classifiers:** Deep architectures such as Multilayer Perceptrons (MLPs), Convolutional Neural Networks (CNNs), and Transformers are capable of learning highly nonlinear decision functions directly from raw data.

The quality of a classification model is typically assessed through metrics such as accuracy, precision, recall, F1-score, and the area under the Receiver Operating Characteristic curve (ROC–AUC). These measures quantify how effectively the model distinguishes among different classes and are essential for evaluating performance, particularly in the presence of class imbalance.

In summary, classification methods aim to map structured or unstructured input data to discrete output categories. From simple linear models to large-scale deep networks, the central challenge remains the same: to generalize from observed examples to unseen data, capturing the underlying structure that defines each class.

## 2.3 Clustering

Clustering is a fundamental technique in *unsupervised learning*. Its purpose is to discover inherent groupings within a dataset by identifying patterns or similarities among unlabeled examples. Unlike classification, where each data point is associated with a predefined label, clustering attempts to infer structure directly from the data itself. In this sense, it can be viewed as a way to explore the natural organization of data without prior supervision.

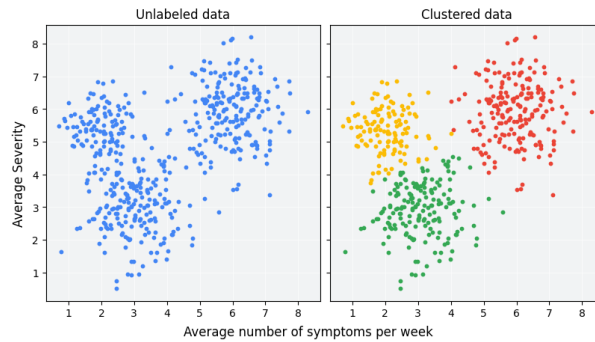


Figure 2.2: Unlabeled data vs Clustered data

Formally, given a set of examples  $\{x_i\}_{i=1}^N$  with  $x_i \in \mathbb{R}^d$ , a clustering algorithm aims to partition the dataset into  $K$  groups, or *clusters*, such that data points within the same cluster are more similar to each other than to those in different clusters. Similarity is typically measured through distance metrics such as Euclidean, Manhattan, or cosine distance, depending on the data representation and problem context.

To illustrate the concept, consider a clinical study designed to evaluate a new treatment protocol. During the experiment, patients report both the frequency and severity of their symptoms per week. Since these observations are continuous and not labeled, researchers can apply clustering analysis to group patients with similar responses to the treatment. For example, one cluster may contain patients who show rapid improvement, another may represent those with moderate progress, and a third may include patients with minimal or no response.

Figure ?? shows a simulated example in which patient data are grouped into three distinct clusters based on their symptom profiles. Such analyses allow researchers to uncover hidden patterns in clinical outcomes and to tailor future interventions according to the characteristics of each patient subgroup.

## Clustering Methods

A variety of clustering algorithms have been developed, each suited to different data structures and assumptions:

- **K-Means Clustering:** A centroid-based algorithm that partitions data into  $K$  clusters by minimizing the within-cluster sum of squared distances. It is simple and efficient, but assumes spherical cluster shapes and requires the number of clusters to be specified in advance.
- **Hierarchical Clustering:** Builds a nested hierarchy of clusters either by iteratively merging smaller clusters (agglomerative) or by splitting larger ones (divisive). The results are often visualized through a dendrogram, which provides insights into the data's structure at multiple levels of granularity.
- **Density-Based Clustering (DBSCAN, HDBSCAN):** Groups data points based on regions of high density, allowing the discovery of arbitrarily shaped clusters and the identification of noise or outliers. It does not require specifying the number of clusters a priori.
- **Gaussian Mixture Models (GMM):** A probabilistic approach in which each cluster is represented as a Gaussian distribution. Unlike K-Means, GMM provides *soft assignments*, meaning that each data point can belong to multiple clusters with different probabilities.
- **Spectral Clustering:** Utilizes the eigenvalues of a similarity matrix to perform clustering in a lower-dimensional space, making it suitable for non-convex or complex data manifolds.

The choice of clustering method depends on the data distribution, dimensionality, and the desired interpretability of results. In practice, clustering is often used as an exploratory tool for data analysis, feature engineering, or as a preprocessing step in larger machine learning pipelines.

In summary, clustering serves as a powerful approach to reveal latent structures within unlabeled datasets. It is widely applied in fields such as bioinformatics, market segmentation, image analysis, and recommender systems, where uncovering hidden patterns can lead to more meaningful understanding and informed decision-making.

## 2.4 Dimensionality Reduction

Dimensionality reduction is a key technique in machine learning and data analysis that seeks to represent high-dimensional data in a lower-dimensional space while preserving as much relevant information as possible. It serves multiple purposes: improving computational efficiency, mitigating the curse of dimensionality, reducing noise, and enhancing interpretability or visualization of complex datasets.

Formally, given a dataset  $\{x_i\}_{i=1}^N$  where each example  $x_i \in \mathbb{R}^D$  is described by  $D$  features, the goal is to find a transformation function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^d$  with  $d \ll D$  such that the lower-dimensional representation  $z_i = f(x_i)$  retains the most informative characteristics

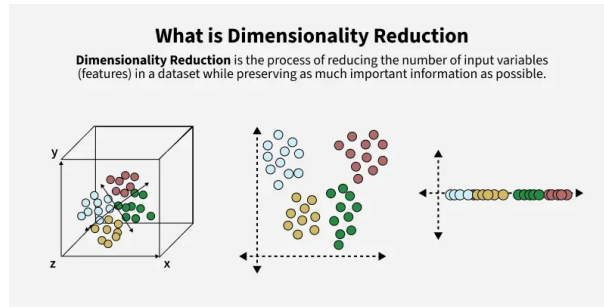


Figure 2.3: Dimensionality Reduction

of the original data. The challenge lies in compressing the representation without losing the essential structure or relationships between data points.

Dimensionality reduction techniques can be broadly divided into two categories: *feature selection* and *feature extraction*. Feature selection involves choosing a subset of the most relevant original variables, while feature extraction constructs new features as combinations or projections of the existing ones.

Several methods have been developed to perform dimensionality reduction, each based on different assumptions and optimization objectives:

- **Principal Component Analysis (PCA):** A linear technique that projects data onto the directions (principal components) of maximum variance. PCA identifies orthogonal axes that best explain the variability in the dataset and is often used for visualization and preprocessing.
- **Linear Discriminant Analysis (LDA):** Although originally designed for classification, LDA can also be used for dimensionality reduction by projecting data onto directions that maximize class separability.
- **Independent Component Analysis (ICA):** Aims to decompose multivariate data into statistically independent components, often used in signal processing and blind source separation tasks.
- **t-Distributed Stochastic Neighbor Embedding (t-SNE):** A nonlinear technique that preserves local similarities among data points, making it particularly useful for visualizing high-dimensional data in two or three dimensions.
- **Uniform Manifold Approximation and Projection (UMAP):** A more recent nonlinear method that approximates the topological structure of data manifolds, often providing better scalability and global structure preservation than t-SNE.
- **Autoencoders:** Neural network-based models that learn a compressed latent representation of data by training to reconstruct the input from a lower-dimensional code. They can capture complex nonlinear relationships and are widely used in deep learning pipelines.

To illustrate, consider a dataset describing patients with hundreds of medical features, such as laboratory values and genetic markers. Applying PCA or UMAP can reveal low-dimensional embeddings that group patients with similar clinical profiles, enabling researchers to visualize disease subtypes or predict treatment outcomes more effectively.

In summary, dimensionality reduction transforms complex, high-dimensional datasets into compact and informative representations. It facilitates efficient learning, improved generalization, and deeper insight into the intrinsic structure of data—serving as both a practical preprocessing step and a powerful tool for exploratory data analysis.



# Chapter 3

## Models

### 3.1 Neural Network Architectures

#### 3.1.1 Feedforward Networks

#### 3.1.2 Convolutional Neural Networks (CNN)

#### 3.1.3

### 3.2 Techniques

#### 3.2.1 Activation Functions (ReLU, GeLU, etc.)

#### 3.2.2 Normalization Layers (BatchNorm, LayerNorm)

#### 3.2.3 Dropout & Regularization

### 3.3 Training and Optimization

#### 3.3.1

#### 3.3.2

# Chapter 4

## Metrics & how to read it

### 4.1 Accuracy & Loss

During the training of a machine learning model, *accuracy* and *loss* are two key metrics used to assess how well the model is learning and generalizing. Together, they provide complementary insights: accuracy measures the proportion of correct predictions, while loss quantifies the numerical difference between the predicted and true outputs. Monitoring these metrics over time allows practitioners to diagnose issues such as underfitting, overfitting, or poor optimization behavior.

Accuracy and loss are typically visualized as curves across training epochs. The shape and trends of these curves reveal how the model's performance evolves during training and validation, helping to identify whether the model is improving, stagnating, or degrading.

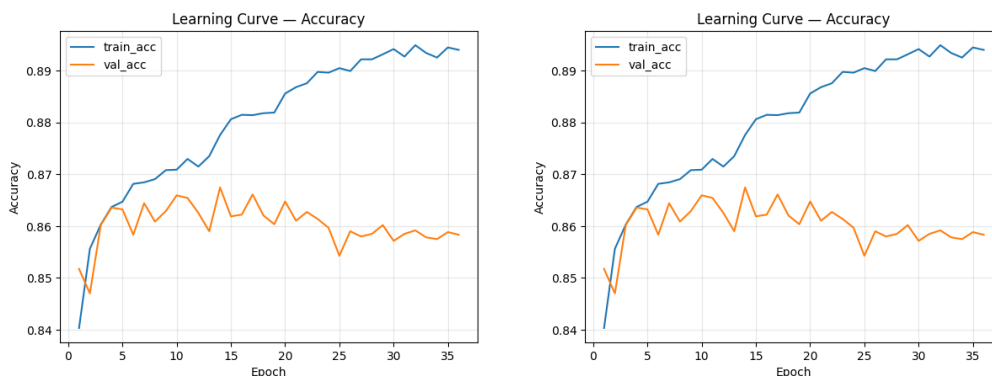


Figure 4.1: Learning Curve for Loss and Accuracy

#### 4.1.1 Accuracy

Accuracy measures the fraction of correctly predicted samples relative to the total number of samples. For a classification task, it is formally defined as:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}.$$

A high accuracy indicates that the model correctly classifies most examples in the dataset. However, accuracy alone can be misleading—especially in imbalanced datasets where some

classes are more frequent than others. In such cases, complementary metrics such as precision, recall, and F1-score provide a more balanced evaluation.

The *accuracy curve* plots model accuracy over successive training epochs for both training and validation sets. A steadily increasing training accuracy accompanied by stable or improving validation accuracy indicates effective learning. Conversely, if validation accuracy stagnates or declines while training accuracy continues to rise, the model may be overfitting.

### 4.1.2 Loss

Loss quantifies how far the model's predictions are from the true target values. It is computed using a *loss function* (or *cost function*) that assigns a penalty to incorrect or inaccurate predictions. The goal of training is to minimize this loss function through iterative optimization, typically via gradient descent or one of its variants.

For example, in regression tasks, the most common loss function is the *Mean Squared Error (MSE)*, defined as:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2,$$

where  $y_i$  represents the true value and  $\hat{y}_i$  the predicted output. In classification tasks, loss functions such as *Cross-Entropy Loss* or *Negative Log-Likelihood* are widely used to measure the divergence between the predicted probability distribution and the true labels.

The *loss curve* typically decreases as training progresses, indicating that the model is minimizing its prediction errors. However, a widening gap between training and validation loss often signals overfitting, whereas consistently high loss on both sets may indicate underfitting or inappropriate model complexity.

In summary, accuracy and loss are fundamental diagnostic tools in machine learning. Their joint analysis provides valuable feedback on model convergence, learning efficiency, and generalization performance, guiding practitioners in fine-tuning architectures and optimization strategies.

## 4.2 Precision & Recall

While accuracy provides a general indication of model performance, it can be misleading in situations where data are imbalanced—when certain classes occur much more frequently than others. In such cases, metrics like *precision* and *recall* offer a more informative and nuanced evaluation of classification models, especially in binary or multi-class settings where the cost of false positives and false negatives differs.

Precision and recall are often analyzed together to assess how well a model identifies positive instances without introducing excessive misclassifications. Their trade-off is frequently visualized through a *precision-recall curve*, which illustrates the relationship between these two quantities as the decision threshold varies.

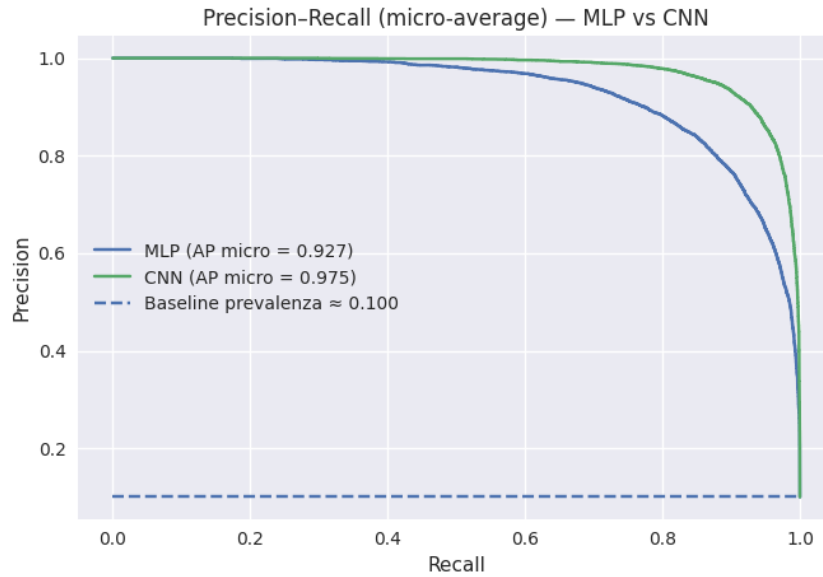


Figure 4.2: Precision Recall Curve

### 4.2.1 Precision

Precision measures the proportion of positive predictions that are actually correct. It answers the question: *Of all instances predicted as positive, how many are truly positive?* Mathematically, precision is defined as:

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}.$$

High precision indicates that the model makes few false positive errors—it is “cautious” in labeling samples as positive. However, a model with very high precision may miss many actual positives if it becomes too conservative, leading to low recall.

In the precision–recall curve, precision typically decreases as recall increases. This occurs because, as the model becomes more inclusive (labeling more instances as positive), it also tends to introduce more false positives.

### 4.2.2 Recall

Recall, also known as *sensitivity* or *true positive rate*, measures the proportion of actual positives that the model successfully identifies. It answers the question: *Of all true positive instances, how many did the model correctly detect?* Formally, recall is defined as:

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}}.$$

A high recall value indicates that the model captures most of the true positives but may include more false positives. In many real-world applications—such as medical diagnosis or fraud detection—high recall is critical, as missing a positive case can be costly or dangerous.

### 4.2.3 Confusion Matrix

The *confusion matrix* is a fundamental tool for evaluating the performance of classification models. It provides a detailed breakdown of the model's predictions compared to the actual ground truth labels, allowing practitioners to understand not only how many predictions were correct, but also the types of errors the model made.

In a binary classification setting, the confusion matrix is a  $2 \times 2$  table that summarizes the counts of true and false predictions for both classes:

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP)	False Negative (FN)
Actual Negative	False Positive (FP)	True Negative (TN)

Here:

- **True Positives (TP):** Instances correctly classified as belonging to the positive class.
- **True Negatives (TN):** Instances correctly classified as belonging to the negative class.
- **False Positives (FP):** Instances incorrectly classified as positive (Type I error).
- **False Negatives (FN):** Instances incorrectly classified as negative (Type II error).

From the confusion matrix, a variety of performance metrics can be derived, including:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}, \quad \text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}.$$

These relationships make the confusion matrix an essential foundation for most evaluation metrics in classification problems.

For *multi-class classification*, the confusion matrix extends to an  $N \times N$  table, where  $N$  is the number of classes. Each row represents the actual class, and each column represents the predicted class. Diagonal entries correspond to correct predictions, while off-diagonal entries indicate misclassifications between specific classes. Visualizing the confusion matrix as a heatmap can help reveal systematic biases—for instance, when the model consistently confuses certain categories due to similarity in features or insufficient training examples.

Confusion Matrix (counts)

T-shirt/top	862	1	14	35	2	1	79	0	6	0
Trouser	1	973	4	16	4	0	2	0	0	0
Pullover	20	0	865	10	65	0	40	0	0	0
Dress	12	6	9	928	21	0	24	0	0	0
Coat	0	0	26	24	914	0	36	0	0	0
Sandal	0	0	0	0	0	978	0	18	0	4
Shirt	114	3	46	33	87	0	711	0	6	0
Sneaker	0	0	0	0	0	9	0	972	1	18
Bag	0	1	2	3	1	2	3	0	988	0
Ankle boot	0	0	0	1	0	5	1	32	1	960

Predicted

Figure 4.3: Confusion Matrix for fashion MNIST dataset

In practice, confusion matrices are often used alongside precision–recall and ROC curves to provide a comprehensive understanding of a model’s behavior. They allow practitioners to identify which types of errors are most common, assess the balance of performance across classes, and guide future steps in data collection, model design, or hyperparameter tuning.

In summary, the confusion matrix is not merely a diagnostic tool but a window into the decision-making process of a classifier. It enables a granular analysis of model performance, transforming raw prediction counts into actionable insights for model refinement and validation.

#### 4.2.4 Precision–Recall Trade-off and Curve Interpretation

Precision and recall are inherently linked: improving one often comes at the expense of the other. By adjusting the model’s decision threshold, one can balance the two metrics depending on the specific requirements of the application. The *precision–recall curve* plots recall on the x-axis and precision on the y-axis, summarizing performance across all possible thresholds.

A model that perfectly separates classes would achieve a precision of 1.0 across all recall levels, resulting in a flat line at the top of the plot. In contrast, a random classifier produces a curve close to the baseline, indicating that precision falls rapidly as recall increases. The area under the precision–recall curve (AUC–PR) provides a scalar summary of overall performance, especially useful when comparing different models on imbalanced datasets.

In summary, precision and recall offer complementary perspectives on model performance. Precision emphasizes correctness among positive predictions, while recall emphasizes completeness in capturing all positive instances. Together, they form the foundation for the *F1-score*, a harmonic mean that provides a single, balanced measure of both precision and recall—allowing for a more holistic evaluation of classification quality.

### 4.3 F1-Score

The *F1-score* is a widely used metric that combines *precision* and *recall* into a single measure of a model’s effectiveness. It is particularly valuable when dealing with imbalanced datasets, where accuracy alone can be misleading, and when both false positives and false negatives carry significant consequences.

Precision measures how many of the predicted positive instances are truly positive, while recall measures how many of the actual positive instances were correctly identified. Since these two metrics often trade off against each other, the F1-score provides a balanced summary by computing their harmonic mean:

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}.$$

This formulation ensures that a high F1-score can only be achieved when both precision and recall are reasonably high. If either precision or recall is low, the F1-score will also decrease sharply, making it more sensitive to performance imbalances than a simple arithmetic mean.

The F1-score ranges between 0 and 1, where 1 indicates perfect precision and recall, and 0 indicates the worst possible performance. It is especially useful when the cost of false positives and false negatives is roughly equal, or when the dataset exhibits skewed class distributions.

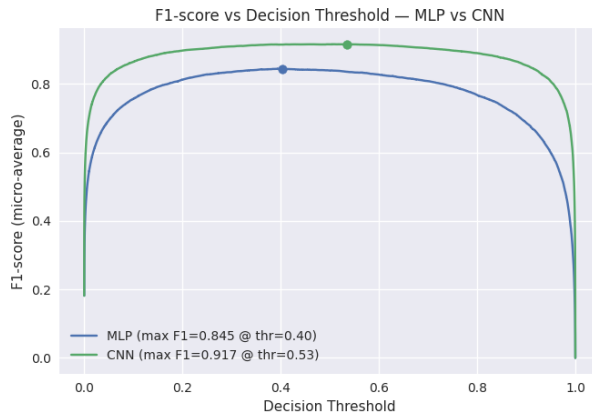


Figure 4.4: F1 score graph with threshold

In a binary classification context, a single F1-score value summarizes model performance. However, for multi-class problems, it can be generalized using different averaging strategies:

- **Macro F1-score:** Computes the F1-score independently for each class and then takes the unweighted mean. It treats all classes equally, regardless of their frequency.
- **Micro F1-score:** Aggregates all true positives, false positives, and false negatives across classes before computing the F1-score. It gives more weight to frequent classes.
- **Weighted F1-score:** Computes the class-specific F1-scores and averages them proportionally to the number of true instances in each class. This approach balances performance evaluation according to class importance.

The *F1-score curve* can also be plotted against varying classification thresholds to visualize the trade-off between precision and recall. By adjusting the decision threshold, practitioners can tune the balance between sensitivity (recall) and specificity (precision) to best fit the requirements of a specific application—such as maximizing recall in medical diagnosis or maximizing precision in spam detection.

In summary, the F1-score provides a robust and interpretable single-number metric for evaluating classification performance, especially in the presence of imbalanced data. It captures the delicate equilibrium between precision and recall, offering a more truthful representation of model effectiveness than accuracy alone. For this reason, it is a standard metric in research benchmarks and real-world machine learning evaluations alike.

## Chapter 5

# Implementation & Code Patterns



# Chapter 6

## Advanced Topics