

CSE 5819 Assignment #1

by: Aayushi Verma (uef24001)

This assignment is organized as follows:

1. ChatGPT - pg 1–6
2. Problems - pg 7–9

1 ChatGPT

1.1 What are the different learning tasks in machine learning?

In machine learning, there are several different learning tasks, each suited to different types of problems and data. These tasks can be broadly categorized into the following:

1. Supervised Learning

- In supervised learning, the model is trained on labeled data, where each input is paired with a corresponding output. The goal is to learn a mapping from inputs to outputs that can generalize to new, unseen data.
- *Examples:* Classification (e.g., spam detection), Regression (e.g., predicting house prices).

2. Unsupervised Learning

- In unsupervised learning, the model is trained on unlabeled data. The task is to discover patterns, structures, or relationships in the data without any predefined labels.
- *Examples:* Clustering (e.g., customer segmentation), Dimensionality Reduction (e.g., PCA).

3. Semi-Supervised Learning

- Semi-supervised learning uses a combination of a small amount of labeled data and a large amount of unlabeled data to improve learning accuracy. It is especially useful when labeling data is expensive or time-consuming.
- *Examples:* Combining labeled and unlabeled data to improve image classification.

4. Reinforcement Learning

- In reinforcement learning, an agent learns to make decisions by interacting with an environment. The agent receives rewards or penalties based on its actions and learns to maximize cumulative rewards over time.
- *Examples:* Game playing (e.g., AlphaGo), Robotics.

5. Self-Supervised Learning

- Self-supervised learning is a type of unsupervised learning where the model generates labels from the input data itself. The model then learns to predict part of the data from other parts, effectively creating its own supervision.
- *Examples:* Predicting the next word in a sentence, Image inpainting.

6. Transfer Learning

- Transfer learning involves transferring knowledge from one domain (the source domain) to another domain (the target domain). It is especially useful when the target domain has limited data.
- *Examples:* Using a pre-trained image recognition model on a new, but related task.

7. Multi-Task Learning

- Multi-task learning involves training a model to perform multiple related tasks simultaneously. This allows the model to leverage shared information between tasks, leading to better generalization.
- *Examples:* Simultaneous classification of multiple object categories in images.

8. Few-Shot Learning

- Few-shot learning aims to enable a model to learn from a very small number of labeled examples, often by leveraging prior knowledge or learning paradigms.
- *Examples:* Image classification with only a few labeled examples per class.

9. Zero-Shot Learning

- Zero-shot learning aims to classify data into categories that the model has not been explicitly trained on by leveraging relationships between known and unknown classes.

- *Examples:* Classifying images of unseen animal species.

10. Active Learning

- In active learning, the model iteratively selects the most informative data points from which it should learn, reducing the amount of labeled data required while improving learning efficiency.
- *Examples:* Choosing which images to label next in a large dataset.

1.2 What is the definition of supervised learning?

Supervised learning is a type of machine learning where a model is trained on a labeled dataset. In this context, "labeled" means that each training example is paired with an output label or target that the model is expected to predict. The goal of supervised learning is to learn a mapping from inputs (features) to outputs (labels) so that the model can accurately predict the output for new, unseen data.

Key Components:

- **Input Data (Features):** The attributes or variables that describe the data points.
- **Output Data (Labels):** The correct outcomes or categories associated with each input data point.
- **Training Process:** The model learns by finding patterns or correlations between the input features and the corresponding output labels.
- **Prediction:** Once trained, the model can be used to predict the label for new input data.

Examples of Supervised Learning:

- **Classification:** Assigning an email as either "spam" or "not spam."
- **Regression:** Predicting the price of a house based on its features like size, location, and number of bedrooms.

1.3 What are the widely used methods to solve supervised learning problems?

There are several widely used methods to solve supervised learning problems, each suited to different types of data and tasks. Here's a breakdown of some of the most common methods:

1. Linear Models

- **Linear Regression:** Used for regression tasks; it models the relationship between the input features and the output as a linear combination of the features.
- **Logistic Regression:** Used for binary classification tasks; it models the probability that a given input belongs to a particular class.

2. Decision Trees and Ensemble Methods

- **Decision Trees:** A tree-like model of decisions, where each internal node represents a feature, each branch represents a decision rule, and each leaf node represents an outcome (class or value).
- **Random Forest:** An ensemble method that builds multiple decision trees and merges them to get a more accurate and stable prediction.
- **Gradient Boosting Machines (GBM):** An ensemble technique that builds trees sequentially, with each new tree correcting errors made by the previous ones.
 - **XGBoost:** A highly efficient and popular implementation of gradient boosting.
 - **LightGBM:** A faster and more memory-efficient implementation of gradient boosting.
 - **CatBoost:** Another gradient boosting algorithm that handles categorical data efficiently.

3. Support Vector Machines (SVM)

- **Support Vector Machines:** A powerful method for both classification and regression tasks. SVMs find the optimal hyperplane that separates the data into different classes with the maximum margin.

4. k-Nearest Neighbors (k-NN)

- **k-Nearest Neighbors:** A simple, instance-based learning method where the output for a new data point is determined by the majority class (in classification) or average (in regression) of its k-nearest neighbors in the feature space.

5. Neural Networks

- **Artificial Neural Networks (ANNs):** A model inspired by the human brain's structure, capable of capturing complex patterns in data. They consist of layers of interconnected neurons where each layer transforms the input data progressively.
- **Convolutional Neural Networks (CNNs):** Specialized neural networks designed for processing grid-like data such as images. They are particularly effective in image recognition tasks.
- **Recurrent Neural Networks (RNNs):** Neural networks designed for sequential data, such as time series or natural language. Variants like Long Short-Term Memory (LSTM) networks are widely used.

6. Naive Bayes

- **Naive Bayes Classifier:** A probabilistic classifier based on Bayes' theorem, assuming independence between features. Despite its simplicity, it's effective for text classification tasks like spam detection.

7. Bayesian Methods

- **Bayesian Linear Regression:** A version of linear regression where the model parameters are treated as random variables with prior distributions.
- **Bayesian Networks:** Graphical models that represent probabilistic relationships among variables.

8. Discriminant Analysis

- **Linear Discriminant Analysis (LDA):** A method used for classification that finds a linear combination of features that best separates two or more classes.
- **Quadratic Discriminant Analysis (QDA):** Similar to LDA, but allows for quadratic decision boundaries.

9. Gaussian Processes

- **Gaussian Process Regression (GPR):** A non-parametric approach that models the distribution over functions and makes predictions by considering the joint distribution of observed data and unseen data.

10. Ensemble Learning

- **Bagging (Bootstrap Aggregating):** An ensemble method that improves model stability and accuracy by training multiple models on different subsets of the data and averaging their predictions.
- **Stacking:** An ensemble method where multiple models are trained, and their outputs are used as inputs for another model (meta-model), which makes the final prediction.

11. Regularization Techniques

- **Ridge Regression (L2 Regularization):** A linear regression model that penalizes the square of the coefficients, helping to prevent overfitting.
- **Lasso Regression (L1 Regularization):** Similar to ridge regression, but it penalizes the absolute value of the coefficients, which can lead to feature selection by shrinking some coefficients to zero.

12. Deep Learning

- **Deep Neural Networks (DNNs):** Neural networks with many layers, capable of learning high-level abstractions from data. They are widely used in tasks like image recognition, natural language processing, and speech recognition.

1.4 What is the definition of unsupervised learning?

Unsupervised learning is a type of machine learning where the model is trained on a dataset without labeled outputs. In this approach, the model tries to learn the underlying structure, patterns, or distribution in the data without any explicit guidance on what the correct output should be. The goal is to discover hidden patterns, groupings, or features within the data.

Key Characteristics:

- **No Labeled Data:** The data provided to the model does not include labels or target outcomes.
- **Pattern Discovery:** The model identifies patterns, relationships, or structures within the data based on the similarities and differences between data points.
- **Exploratory:** Unsupervised learning is often exploratory in nature, aiming to understand the underlying structure of the data.

Examples of Unsupervised Learning:

- **Clustering:** Grouping similar data points together based on certain criteria.
 - *Example:* Segmenting customers into distinct groups based on purchasing behavior.
- **Dimensionality Reduction:** Reducing the number of features while retaining the essential information.
 - *Example:* Principal Component Analysis (PCA) for reducing the complexity of data.
- **Anomaly Detection:** Identifying outliers or unusual data points that do not fit the general pattern.
 - *Example:* Detecting fraudulent transactions in financial data.
- **Association Rule Learning:** Finding rules that describe large portions of the data.
 - *Example:* Market basket analysis, where the model identifies items that frequently co-occur in transactions.

In summary, unsupervised learning involves analyzing and finding hidden structures in unlabeled data, making it useful for exploratory data analysis, pattern recognition, and feature extraction.

1.5 What are the widely used methods to solve unsupervised learning problems?

Unsupervised learning involves a variety of methods designed to discover patterns, structures, and relationships within unlabeled data. Here are some of the most widely used methods to solve unsupervised learning problems:

1. Clustering

- **k-Means Clustering:** Partitions the data into k distinct clusters based on feature similarity. It assigns each data point to the nearest cluster center and iteratively refines the clusters.
- **Hierarchical Clustering:** Builds a tree-like structure of clusters, where each level represents a different granularity of grouping. It can be agglomerative (bottom-up) or divisive (top-down).
- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** Identifies clusters based on the density of data points, allowing the discovery of clusters of arbitrary shapes and handling noise (outliers) effectively.
- **Gaussian Mixture Models (GMMs):** Assumes that the data is generated from a mixture of several Gaussian distributions, each representing a cluster. It uses the Expectation-Maximization (EM) algorithm to estimate the parameters.

2. Dimensionality Reduction

- **Principal Component Analysis (PCA):** Reduces the dimensionality of the data by transforming it into a set of linearly uncorrelated variables (principal components), capturing the maximum variance.
- **t-Distributed Stochastic Neighbor Embedding (t-SNE):** Reduces dimensions by modeling the probability distribution of pairwise distances between data points, particularly effective for visualizing high-dimensional data in 2D or 3D.

- **Autoencoders:** A type of neural network used to learn a compressed representation of the data. It consists of an encoder that compresses the data and a decoder that reconstructs it.
- **Independent Component Analysis (ICA):** Separates a multivariate signal into additive, independent components, often used in signal processing and source separation.

3. Association Rule Learning

- **Apriori Algorithm:** Identifies frequent itemsets in transactional data and generates association rules, such as "If a customer buys bread, they are likely to buy butter."
- **Eclat Algorithm:** An improvement over Apriori, it focuses on itemset intersections to find frequent itemsets more efficiently.

4. Anomaly Detection

- **Isolation Forest:** Detects anomalies by randomly partitioning data points and isolating those that are fewer in number and different from the rest.
- **One-Class SVM:** A variant of Support Vector Machines that identifies the majority class (normal data) and classifies outliers as anomalies.
- **LOF (Local Outlier Factor):** Measures the local density deviation of a data point with respect to its neighbors, identifying those with significantly lower density as outliers.

5. Latent Variable Models

- **Latent Dirichlet Allocation (LDA):** A generative probabilistic model used for topic modeling, where documents are represented as mixtures of topics, and topics are mixtures of words.
- **Factor Analysis:** Models the data as linear combinations of underlying factors, used for identifying the hidden variables that explain observed data.

6. Self-Organizing Maps (SOMs)

- **Self-Organizing Maps:** A type of artificial neural network trained using unsupervised learning to produce a low-dimensional (typically 2D) representation of the input space, preserving the topological properties.

7. Dimensionality Reduction with Manifold Learning

- **Locally Linear Embedding (LLE):** Captures the local neighborhood structure of the data and projects it into a lower-dimensional space.
- **Multidimensional Scaling (MDS):** Aims to place each data point in a lower-dimensional space while preserving the pairwise distances as much as possible.

8. Clustering-Based Methods for Image and Text Data

- **Bag of Words (BoW):** Represents text data as a bag of words, ignoring grammar and word order, and uses clustering techniques to identify patterns or topics.
- **Word2Vec:** Embeds words into vectors based on their context in the corpus, allowing clustering and pattern recognition in text data.

9. Generative Models

- **Generative Adversarial Networks (GANs):** Consists of a generator and a discriminator; the generator creates synthetic data, and the discriminator evaluates how close the synthetic data is to real data, leading to the generation of realistic data.
- **Variational Autoencoders (VAEs):** A type of autoencoder that learns the distribution of the data in a continuous latent space, allowing for data generation.

10. Matrix Factorization

- **Singular Value Decomposition (SVD):** Factorizes a matrix into three other matrices and is widely used in collaborative filtering for recommendation systems.
- **Non-negative Matrix Factorization (NMF):** Similar to SVD but with the constraint that the resulting matrices have no negative elements, used in topics modeling and image processing.

2 Problems

2.1 Is any of the following problems an example of a supervised learning problem? (You can choose multiple choices). (10 pts)

- ☐ A Given a set of patient data, including various features such as breast images, family history, cell nuclei measurements (for example, radius, texture, perimeter, area, smoothness of the cells), and the biopsy results of benign and malignant, the goal is to classify whether a tumor is benign (non-cancerous) or malignant (cancerous).
- ☐ B Given images of crop (wheat or rice) leaves and biologist's labels in terms of whether the crop is healthy or suffering from a specific disease (for example, powdery mildew, leaf rust, or bacterial blight), the goal is to classify whether any new plant image is healthy.
- ☐ C Given a large dataset of customer behaviors, such as purchase history, browsing patterns, demographics, and engagement metrics, the goal is to group customers into distinct segments based on similarities in their behavior and preferences. There are no predefined labels for the customer segments.
- ☐ D Given historical data of a satellite's position, velocity, and various environmental factors (such as gravitational forces, solar radiation pressure, and atmospheric drag), the goal is to predict the satellite's future position in space at a given time.

Options A, B, and D are supervised learning problems.

2.2 Is any of the following problems an example of an unsupervised learning problem? (You can choose multiple choices). (10 pts)

- ☐ A Given a large dataset of images of faces, each image is represented as a high-dimensional vector (e.g., each pixel in an image is a feature). The goal is to reduce the dimensionality of the data while preserving the most important information, so that the faces can still be recognized or compared with minimal computational cost.
- ☐ B In a manufacturing plant, machines produce data through sensors, such as temperature, vibration, and sound levels. The goal is to detect unusual patterns or behaviors (anomalies) that could indicate machine faults or maintenance needs. The challenge is that there are no pre-observed labels of normal or faulty operation for every possible scenario.
- ☐ C An autonomous vehicle must learn how to drive through complex environments, such as urban roads, highways, and intersections, while interacting with other vehicles, pedestrians, traffic signals, and obstacles. The goal is to enable the vehicle to make safe and efficient driving decisions, such as when to accelerate, brake, change lanes, or stop at traffic lights.
- ☐ D In an online learning platform, some students may frequently revisit course materials and may benefit from more practice problems, while some other students need to be challenged with more advanced content. An algorithm might analyze student data (e.g., study habits, online learning activity, quiz completion times, etc.) and groups students with similar learning patterns.

Options A, B, and D are unsupervised learning problems.

2.3 Is any of the following methods capable of solving a supervised learning problem? (You can choose multiple choices). (10 pts)

- ☐ A Linear regression
- ☐ B Random forests
- ☐ C Support vector machines
- ☐ D K-nearest neighbors

Options A, B, C, and D are all algorithms for solving supervised learning problems.

2.4 For a quadratic form as follows (Eq. 1), please compute the gradient of this function $f(x)$ where x is a vector in R^n , A is a $n \times n$ matrix, B is a vector in R^n , and C is a scalar. (20 pts) (Write your derivation. Do not just copy-paste what GPT provides if you use GPT.)

$$f(x) = x^T A x + B^T x + C \quad (1)$$

To find the gradient of $f(x)$, we have:

$$\begin{aligned} \nabla_x f(x) &= \frac{\partial}{\partial x} (x^T A x + B^T x + C) \\ &= \frac{\partial}{\partial x} x^T A x + \frac{\partial}{\partial x} B^T + \frac{\partial}{\partial x} C \end{aligned} \quad (2)$$

Let's take the derivatives of Eq. 2. First, from the Matrix Cookbook (posted on HuskyCT) Eq. 81, we can use the following identity:

$$\frac{\partial}{\partial x} x^T A x = (A + A^T)x \quad (3)$$

Using normal differentiation rules, we have:

$$\frac{\partial}{\partial x} B^T x = B \quad (4)$$

and:

$$\frac{\partial}{\partial x} C = 0 \quad (5)$$

Therefore, substituting Eq. 3, 4 and 5 back into Eq. 2, we find the gradient of $f(x)$ as:

$$\begin{aligned} \therefore \nabla_x f(x) &= \frac{\partial}{\partial x} x^T A x + \frac{\partial}{\partial x} B^T + \frac{\partial}{\partial x} C \\ &= (A + A^T)x + B + 0 \\ &= (A + A^T)x + B \end{aligned} \quad (6)$$

2.5 The following function (Eq. 7) is a quadratic form where $a_{12} = a_{21}$, can you write out its gradient, and compare with the above gradient formula you derived in Problem 4. (10 pts)

$$f(x) = a_{11}x_1^2 + a_{12}x_1x_2 + a_{21}x_2x_1 + a_{22}x_2^2 + b_1x_1 + b_2x_2 + c \quad (7)$$

Since we are given $a_{12} = a_{21}$, Eq. 7 simplifies to:

$$f(x) = a_{11}x_1^2 + 2a_{12}x_1x_2 + a_{22}x_2^2 + b_1x_1 + b_2x_2 + c \quad (8)$$

We need to find the gradient with respect to x . Since there are multiple x , i.e. x_1 and x_2 , we take the derivative of Eq. 8 with respect to each x_1 and x_2 .

$$\frac{\partial f(x)}{\partial x_1} = 2a_{11}x_1 + 2a_{12}x_2 + b_1 \quad (9)$$

$$\frac{\partial f(x)}{\partial x_2} = 2a_{12}x_1 + 2a_{22}x_2 + b_2 \quad (10)$$

Therefore the gradient of $f(x)$ is:

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2a_{11}x_1 + 2a_{12}x_2 + b_1 \\ 2a_{12}x_1 + 2a_{22}x_2 + b_2 \end{bmatrix} \quad (11)$$

We can re-write Eq. 11 like:

$$\nabla_x f(x) = \begin{bmatrix} 2a_{11}x_1 + 2a_{12}x_2 \\ 2a_{12}x_1 + 2a_{22}x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \quad (12)$$

$$\begin{aligned} &= 2 \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \\ &= 2Ax + B \end{aligned} \quad (13)$$

If we compare Eq. 13 with Eq. 6 derived in Question 4, then given the condition $a_{12} = a_{21}$, both Eq. 13 and Eq. 6 are equivalent to each other, since $a_{12} = a_{21}$ implies a symmetrical matrix, in which case $A = A^T$. Therefore:

$$\begin{aligned}\therefore \nabla_x f(x) &= (A + A^T)x + B \\ &= (2A)x + B \\ &= 2Ax + B\end{aligned}\tag{14}$$

2.6 For a sigmoid function $\sigma(s) = \frac{1}{1+e^{-s}}$, where s is a scalar, prove that its derivative satisfies: $\sigma'(s) = \sigma(s) \cdot (1 - \sigma(s))$. (20pts)

We have:

$$\begin{aligned}\sigma(s) &= \frac{1}{1 + e^{-s}} \\ &= (1 + e^{-s})^{-1}\end{aligned}\tag{15}$$

Using the chain rule:

$$\begin{aligned}\frac{d\sigma(s)}{ds} &= -1 \cdot (1 + e^{-s})^{-2} \cdot \frac{d}{ds}(1 + e^{-s}) \\ &= -(1 + e^{-s})^{-2} \cdot \frac{d}{ds}(1 + e^{-s})\end{aligned}\tag{16}$$

From Eq. 16, we have:

$$\begin{aligned}\frac{d}{ds}(1 + e^{-s}) &= \frac{d}{ds}(1) + \frac{d}{ds}(e^{-s}) \\ &= -e^{-s}\end{aligned}\tag{17}$$

Substituting Eq. 17 back into Eq. 16, we have:

$$\begin{aligned}\frac{d\sigma(s)}{ds} &= -(1 + e^{-s})^{-2} \cdot -e^{-s} \\ &= (1 + e^{-s})^{-2} \cdot e^{-s} \\ &= \frac{e^{-s}}{(1 + e^{-s})^2} \\ &= \frac{e^{-s}}{(1 + e^{-s})(1 + e^{-s})} \\ &= \frac{1}{(1 + e^{-s})} \cdot \frac{e^{-s}}{(1 + e^{-s})} \\ &= \sigma(s) \cdot \frac{e^{-s}}{(1 + e^{-s})}\end{aligned}\tag{18}$$

Note that we can rewrite:

$$\begin{aligned}\frac{e^{-s}}{(1 + e^{-s})} &= 1 - \frac{1}{(1 + e^{-s})} = 1 - \sigma(s) \\ &= \frac{1(1 + e^{-s}) - 1}{(1 + e^{-s})} \\ &= \frac{1 + e^{-s} - 1}{(1 + e^{-s})} \\ &= \frac{e^{-s}}{(1 + e^{-s})}\end{aligned}\tag{19}$$

Therefore substituting Eq. 19 back into Eq. 18, we have:

$$\therefore \frac{d\sigma(s)}{ds} = \sigma(s) \cdot (1 - \sigma(s))\tag{20}$$

Therefore we proved that $\sigma'(s) = \sigma(s) \cdot (1 - \sigma(s))$.