

Introduction to Machine Learning

Praanesh Balakrishnan Nair

January 28, 2025

Contents

1	Introduction	2
2	Well-posed Learning Problem by Tom Mitchell (1998)	2
3	Components of a dataset	2
3.1	Features	2
3.2	Data Points	2
3.3	Feature Vector	2
3.4	Distance and Similarity Matrix	3
3.4.1	Distance Matrix	3
4	Classification	3
4.1	Types of Classification Learning	3
4.1.1	Supervised Learning	3
4.1.2	Unsupervised Learning	3
4.1.3	Semi-supervised Learning	3
4.2	Model Accuracy	3
4.2.1	Error Rate	3
4.2.2	Confusion Matrix	4
4.2.3	Receiver Operating Characteristic Curve (ROC) and Area Under the Curve (AUC)	4
4.3	Model Validation Techniques	4
4.3.1	Internal Validation	4
4.3.2	External Validation	4
5	Regression	5
5.1	Model Accuracy	5
5.2	Types	5
5.2.1	Linear Regression	5
5.2.2	Polynomial Regression	6
5.3	Performance Measures	6
5.3.1	Mean Squared Error:	6
5.3.2	Mean Absolute Error:	6
5.3.3	Root Mean Squared Error:	6
5.3.4	R ² Score	6

1 Introduction

Quote by Herbert Alexander Simon:

Learning is the process by which any system improves its performance from experience

2 Well-posed Learning Problem by Tom Mitchell (1998)

A computer program

- Performs Task T
- Has some Performance P
- Learns from Experience E

Sl. No	Task	Performance	Experience
1.	Classifying Emails as Spam/Not Spam	Number of emails correctly classified	Watching you Label Emails
2.	Playing Chess	Percent of games won	Watch enemy play
3.	Handwriting Recognition	Percent of correct recognitions	Sample images

3 Components of a dataset

3.1 Features

- Individual measurable properties, which are going to be used as input to the machine learning model. Eg. Age of people, Dimensions of a house, etc

3.2 Data Points

- Multiple samples of features.
- Eg:

Sl. No	Age	Height	Weight	BP
1	19	175	68	999
2	25	169	69	0
...

Each of these rows are data points. In each data point, you have different samples of the same features

3.3 Feature Vector

- Features in one data-point is often mathematically represented as a Vector
- Eg:

Sl. No	Age	Height	Weight	BP	Feature Vector
1	19	175	68	999	[19, 175, 68]
2	25	169	69	0	[25, 169, 0]
...

3.4 Distance and Similarity Matrix

3.4.1 Distance Matrix

- Basically Adjacency Matrix
- $d(i, j)$ = distance between i^{th} data point and j^{th} data point.
- It's Symmetric
- Diagonal Elements are 0
- The actual distance between i^{th} data point and j^{th} data point can be measured in many ways:
 1. **Euclidean Distance** = $\sqrt{\sum (x_i - y_i)^2}$
 2. **Manhattan distance** = $\sum |x_i - y_i|$
 3. **Cosine Distance** = $1 - \frac{x \cdot y}{|x||y|}$

4 Classification

Here, you input some data, and the output is a classification of it.

4.1 Types of Classification Learning

4.1.1 Supervised Learning

- Classify features X_i into classes/labels Y_i
- You find the pattern of data that is associated with one label, and use that pattern to classify.

4.1.2 Unsupervised Learning

- You have only features X_i and no labels
- You find patterns, so that similar patterns form one label, and anything different will be given another label.

4.1.3 Semi-supervised Learning

The entire dataset consists of labelled and unlabelled data

1. Perform supervised learning on the labelled data
2. Now you use this to predict the labels of the unlabelled data. The predicted labels are called **psuedo-labels**.
3. Now do supervised learning on the combined data

4.2 Model Accuracy

4.2.1 Error Rate

- In classification, the model accuracy is quantified by the **error rate**.
- **Error Rate** = $\frac{\text{number of misclassifications}}{\text{total number of data points}}$
- **Error Rate** = $\frac{\sum_{i=1}^n I(y_i \neq \hat{y}_i)}{n}$, where $I(y_i \neq \hat{y}_i)$ is 1 if it's a mismatch, and 0 if it's a match

4.2.2 Confusion Matrix

- Matrix where:
 - rows signify Ground truth (Row1: +, Row2: -)
 - columns signify predicted output (Column1: +, Column2: -)
- If row and column have same sign, it means the model has predicted correctly (it's a **true output**).
- If row and column have opposite signs, it means the model has predicted incorrectly (it's a **false output**).

	Predicted +	Predicted -
Actual +	True Positive	False Negative
Actual -	False Positive	True Negative

- Here are things you can derive from the confusion matrix:

	Predicted +	Predicted -	
Actual +	True Positive	False Negative	Sensitivity/Recall
Actual -	False Positive	True Negative	Specificity
	Precision	Negative Predictive Value	

- Sensitivity** = $\frac{\text{Diag. Element of Row 0}}{\text{Row 0}} = \frac{TP}{TP+FN}$
- Specificity** = $\frac{\text{Diag. Element of Row 1}}{\text{Row 1}} = \frac{TN}{FP+TN}$
- Precision** = $\frac{\text{Diag. Element of Column 0}}{\text{Column 0}} = \frac{TP}{TP+FP}$
- Negative Predictive Value** = $\frac{\text{Diag. Element of Column 1}}{\text{Column 1}} = \frac{TN}{FN+TN}$

4.2.3 Receiver Operating Characteristic Curve (ROC) and Area Under the Curve (AUC)

- ROC is the plot between True Positive and False Positive
- AUC is the area under ROC
- $0 \leq AUC \leq 1$
- $AUC = \int_0^1(\text{ROC Curve})$

4.3 Model Validation Techniques

4.3.1 Internal Validation

- Separation between clusters should be high
- Cohesion (distance between points in a cluster) should be low

4.3.2 External Validation

1. Dice Coefficient

- $D(A, B) = \frac{2|A \cap B|}{|A| + |B|}$
- If $D(A, B) = 0$, then there's no overlap. Similarly if $D(A, B) = 1$, they are the same set.
- A could be the data we have and B could be some external data.

2. Jaccard Similarity Index

- $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$

5 Regression

Here, you input some data and you get a quantitative response.

5.1 Model Accuracy

- In regression, it's called the **quality of fit** and it's the quantification of the degree of closeness of predicted response and the true response
- The most commonly used measure for this, is the **Mean Square Error (MSE)**.
- $MSE = \frac{\sum_{i=1}^n (y_i - f(x_i))^2}{n}$, where y_i is the true value, $f(x_i)$ is the predicted value
- Accuracy of the Model $\propto \frac{1}{MSE}$

5.2 Types

5.2.1 Linear Regression

1. What it is

- Given input X_i , predict a quantitative response Y_i .
- You find the line closest to all of the data points.
- The line is given as: $h_{\theta}(x) = \theta_0 + \theta_1(x)$

2. How it works

- You have a cost function given as $J(\theta_0, \theta_1) = \text{Mean Square Error}$
- Simply minimize the cost function i.e. find values for θ_0 and θ_1 such that $J(\theta_0, \theta_1)$ has the smallest value.
- To find those values, you either run a **really large loop** and iterate through all the values of θ_1 and θ_0 possible, or you use something called the **gradient descent**.

3. Gradient Descent

- $\theta_i = \theta_i - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_i}$, where $i = 0$ **or** 1 , and α is the learning rate (user defined)
- In every iteration, you update each parameter by subtracting $\alpha \frac{\partial \theta_i}{\partial J(\theta_0, \theta_1)}$
- $\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_i}$ is the change in $J(\theta_0, \theta_1)$ with respect to θ_i i.e. how much $J(\theta_0, \theta_1)$ changes for a small change in θ_i .
- This change tells you the direction you have to go in the graph, to reduce the cost function. The direction is simply a positive or negative value which should be added to each of the parameters.
- On performing multiple iterations of this method, you finally reach the minimum of this function.
- When you have two or more parameters like this, you shouldn't directly change θ_i , because . Instead, you do:
 - $temp_0 = \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$
 - $temp_1 = \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$
 - $\theta_0 = temp_0$
 - $\theta_1 = temp_1$

4. An extension to this: Multivariate Linear Regression

- So far, we've had one input/variable x , and the line was in a 2D plane.
- Now, we have multiple inputs/variables, and hence the line is in a multidimensional space.
- The line is given as:

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$h_0 = [\theta_0 \quad \theta_1 \quad \theta_2 \dots] \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ \dots \end{bmatrix}$$

$$h_0 = \theta^T X$$

5.2.2 Polynomial Regression

- $h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2$

5.3 Performance Measures

5.3.1 Mean Squared Error:

$$MSE = \frac{\sum_{i=1}^n (y_i - f(x_i))^2}{n}$$

5.3.2 Mean Absolute Error:

$$MAE = \frac{\sum_{i=1}^n |(y_i - f(x_i))|}{n}$$

5.3.3 Root Mean Squared Error:

$$RMSE = \sqrt{MSE}$$

5.3.4 R^2 Score

$$R^2 = 1 - \frac{SS_{Residuals}}{SS_{Total}}$$