

Chapter - 1

Learning from Data

Aviral Janveja

1 Introduction

Machine Learning is a field concerned with learning patterns from data, without explicitly programming them. For example, we learn to recognize trees not by memorizing their written definition, but by repeatedly looking at them, either in real life or through pictures. Over time, we learn to identify their distinct physical features. Similarly, machine learning systems work by learning patterns directly from example data-points.

2 The Machine Learning Process

Let us look at the main components of the machine learning process, using another simple example from the finance domain : Suppose a customer applies for a loan. The bank must now decide whether granting the loan is a good idea or not, since its goal is to maximize profit while minimizing risk. Naturally, the bank has no magic formula to ascertain whether a customer will make timely repayments or not. However, they may rely on past customer data and how their repayment behavior turned out to be.

The idea here, is to learn the underlying pattern that connects customer information to their repayment behavior. The expectation being that it will help us correctly predict the creditworthiness of fresh applicants with high probability. With this example in mind, let us now outline the main components of the machine learning process.

2.1 Data

As discussed, we are essentially learning patterns from data, therefore without sufficient data-points, the process of machine learning cannot even begin. In line with the above example, the data from N past customers is described as follows :

$$\text{Data} : (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3) \dots (\mathbf{x}_N, y_N)$$

Here \mathbf{x}_n is a column matrix representing customer information such as age, salary, years in residence, current debt and so on. Whereas $y_n \in \{+1, -1\}$ is a binary value representing whether the customer was creditworthy or not in hindsight.

2.2 The Underlying Pattern

In machine learning, our goal is to uncover the underlying pattern from the given set of data-points. In the loan example from above, this underlying pattern is the unknown mechanism that relates customer information like age, salary and so on to their creditworthiness.

We do know, there is some underlying pattern, some relation that connects the input and output in the above data. However, it is not possible to pin it down exactly or write it down as a known formula. The relation is essentially unknown, but we have data on it. That is why we are interested in machine learning or learning from data, in the first place.

But how could this relation be described formally? Could it be a deterministic relation, like a mathematical function (even though unknown) or is it a probabilistic relation? In order to answer that, let us ask ourself the following question, a very practical question based on the example taken above :

Is it not possible that two customers, with identical application information, may end up having very different repayment behavior?

Yes, it is entirely possible and in fact quite common in real world scenarios. Two identical loan applicants may behave differently. It could happen due to unforeseen circumstances or other factors that are not captured in our data. Therefore, we arrive at an important observation. The underlying relation that we are trying to capture is not really deterministic, but rather **probabilistic** in nature.

This probabilistic relation can be described formally via a probability distribution, commonly called the **target distribution** :

$$P(y \mid \mathbf{x})$$

That is, the probability of output y , given input \mathbf{x} . This is actually the underlying pattern that we try to learn.

For example :

When we roll a fair die, we cannot predict the exact outcome, but we can describe the **likelihood of each outcome**, using a probability distribution :

$$P(k) = \frac{1}{6} ; k = 1, 2, 3, 4, 5, 6$$

The key idea to understand is that : a probability distribution tells us how likely each outcome is, not which outcome will occur in a specific trial. Likewise, in the loan example :

$$P(y = +1 \mid \mathbf{x}) = 0.8$$

tells us the likelihood of creditworthiness, given an input \mathbf{x} , not a guaranteed outcome.

2.3 The Learning Model

Error Measures

The learning model defines the **hypothesis set** and the **learning algorithm**. The hypothesis set H is a set of hypothesis functions $\{h_1, h_2, h_3 \dots h_m\}$ out of which the learning algorithm picks the **final hypothesis function** h_f based on the available data-points.

This final hypothesis function is our approximation of the unknown target function. The goal in machine learning is that h_f approximates f well, that is :

$$h_f \approx f$$

An an example, we examine a simple learning model called Perceptron next, to understand how the above process works, continuing with the bank-loan example.

3 Perceptron

Linear classifiers such as the perceptron do not model the entire underlying target distribution. Instead, they optimize a linear decision boundary as per the labels assigned to the past data-points. It can be thought of as, the perceptron approximating the target distribution learning the most probable label, for a given x .

For example, let us say the underlying probability distribution given a point x is as follows :

$$P(y = +1 \mid x) = 0.8$$

and

$$P(y = -1 \mid x) = 0.2$$

Then, final label $y = +1$ and this what the perceptron will try to learn and optimize for, without trying to model the exact underlying probability distribution given above. This is because perception is a simple linear model. Further there will be other advanced models that will try to capture the probabilistic aspect explicitly.

Let us begin by looking at the given data, which is available to us in the form of N input-output pairs :

$$(x_1, y_1), (x_2, y_2), (x_3, y_3) \dots (x_N, y_N)$$

As discussed, x_n is the application information of the n^{th} customer, which is represented by a column matrix as shown below :

$$x_n = \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ x_d \end{bmatrix}$$

Here x_1, x_2 and so on, upto x_d are the attributes of a customer like age, salary, years in residence, outstanding debt and so on. Whereas $y_n \in \{+1, -1\}$ is binary value representing the creditworthiness of a past customer in hindsight.

3.1 Perceptron Hypothesis Function

Now, what a perceptron basically does, is to assign weights to each of the customer attributes. These weights are multiplied to their respective attributes and then summed-up as follows :

$$w_1x_1 + w_2x_2 + \dots + w_dx_d$$

The above linear summation of weights and attributes can be called the credit score of an individual customer. The idea is pretty simple, If the credit score is greater than a certain threshold, then we approve the loan, else we deny it :

$$w_1x_1 + w_2x_2 + \dots + w_dx_d > \text{threshold (Approve Loan)}$$

$$w_1x_1 + w_2x_2 + \dots + w_dx_d < \text{threshold (Deny Loan)}$$

This is equivalent to :

$$w_1x_1 + w_2x_2 + \dots + w_dx_d - \text{threshold} > 0 \text{ (Approve Loan)}$$

$$w_1x_1 + w_2x_2 + \dots + w_dx_d - \text{threshold} < 0 \text{ (Deny Loan)}$$

We can simplify the above expression further, by taking the minus of threshold as the zeroth weight ($w_0 = -\text{threshold}$) and introducing an artificial attribute $x_0 = 1$ to go along with it, hence giving us :

$$w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_dx_d > 0 \text{ (Approve Loan)}$$

$$w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_dx_d < 0 \text{ (Deny Loan)}$$

Next, we re-write the above expression in vector (column matrix) notation to make it shorter and cleaner :

$$w_0x_0 + w_1x_1 + \dots + w_dx_d = \begin{bmatrix} w_0 & w_1 & \dots & w_d \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} = \mathbf{w}^T \mathbf{x}$$

The above formula outputs a real value that could be greater than, equal to or less than zero. In order to output a binary value, similar to y , we wrap the above formula inside a **sign** function.

$$\text{sign}(\mathbf{w}^T \mathbf{x})$$

The sign is just a simple function that takes in a real number value and outputs +1 if the value is greater than zero and -1 if it is less than or equal to zero, hence giving us the final form of our perceptron hypothesis function.

3.2 Perceptron Learning Algorithm

Examining the perceptron hypothesis function, that we have arrived at above :

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

We see that \mathbf{x} representing the customer information is already a given, along with y , which is the correct loan decision in hindsight. So, the underlying pattern that we are trying to learn will be reflected in our choice of weights, as represented by the weight vector \mathbf{w} . Hence, our task now is to arrive at a set of weights, that are able to classify all given data-points correctly.

So, how do we compute the correct weights? First, we initialize the weights randomly, let us say initialize them all to zero. The perceptron learning algorithm then iterates through the dataset, looking for misclassified points (\mathbf{x}_n, y_n) , where :

$$h(\mathbf{x}) \neq f(\mathbf{x})$$

that is;

$$\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$$

Whenever such a misclassified point is found, the weight vector is updated by applying the following **update rule** :

$$\mathbf{w}_{new} = \mathbf{w} + y_n \mathbf{x}_n$$

This process continues, until no misclassified points remain in the dataset. And that is it, we then have our required set of weight values.

3.3 Update Rule Justification

Let us examine how the above defined update rule actually nudges the perceptron classifier in the right direction at each misclassified point. If a point (\mathbf{x}_n, y_n) is misclassified, it means :

$$\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$$

Accordingly, the perceptron updates the weight vector as follows :

$$\mathbf{w}_{new} = \mathbf{w} + y_n \mathbf{x}_n$$

Substituting the updated weight vector from above and ignoring the sign function for now, we get the following :

$$\mathbf{w}_{new}^T \mathbf{x}_n = (\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n$$

Expanding the right-hand side, we get :

$$\mathbf{w}_{new}^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n + y_n \mathbf{x}_n^T \mathbf{x}_n$$

This is equivalent to :

$$\mathbf{w}_{new}^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n + y_n \|\mathbf{x}_n\|^2$$

Therefore, for a misclassified point, where $y_n = +1$ and $\text{sign}(\mathbf{w}^T \mathbf{x}_n) = -1$, it implies $\mathbf{w}^T \mathbf{x}_n \leq 0$ and after the update we get :

$$\mathbf{w}_{new}^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n + \|\mathbf{x}_n\|^2$$

The update is adding some positive value to $\mathbf{w}^T \mathbf{x}_n$ which might make $\mathbf{w}^T \mathbf{x}_n > 0$. Although, it does not guarantee an immediate sign flip, it clearly nudges the perceptron classifier in the right direction.

Similarly, when $y_n = -1$ and $\text{sign}(\mathbf{w}^T \mathbf{x}_n) = +1$, implying $\mathbf{w}^T \mathbf{x}_n > 0$, we get :

$$\mathbf{w}_{new}^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n - \|\mathbf{x}_n\|^2$$

Hence, showing that each update to the weight vector nudges the perceptron in the required direction, improving its performance on that particular data-point. The following image illustrates the above discussion :

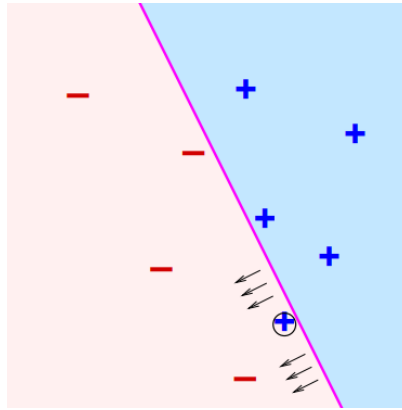


Figure 1: Perceptron

The pink line represents our linear perceptron classifier that is uniquely determined by the choice of weight vector \mathbf{w} . We have a misclassified **plus** point in the **minus** region and the update-rule tries to correct this by nudging the classifier in the required direction.

3.4 Convergence Theorem

The convergence theorem for perceptron, guarantees that if the data is **linearly separable**, repeated updates will eventually classify all points correctly.

Data being linearly separable simply means that, when the data-points are plotted on a two-dimensional graph, it is possible to draw a straight-line that separates them into two distinct categories. The data-points in the above image for example, are linearly separable.

4 Types of Learning

As discussed in the introduction, machine learning is about learning patterns from data. Alternatively, we can put it as **using a set of observations to uncover an underlying process**. In either case, we have established that having a set of observations or data-points is a non-negotiable. However, the form in which the

data is available to us can vary. Depending on this form and the nature of the learning task, machine learning is commonly divided into three main paradigms: supervised learning, unsupervised learning, and reinforcement learning.

4.1 Supervised Learning

In supervised learning, the model is trained on labeled data, where each example consists of : (**input, correct output**). The goal is to approximate a function from inputs to outputs that generalizes well to unseen data. For instance, the perceptron model above.

4.2 Unsupervised Learning

In unsupervised learning, the data is in the form of inputs only, without associated labels : (**input, no output**). The task is to discover hidden structures, patterns or clusters within the data. Examples include clustering, dimensionality reduction, and generative modeling.

4.3 Reinforcement Learning

In reinforcement learning, the model interacts with the environment by taking actions and receiving feedback in the form of rewards or penalties : (**action, feedback**). The model learns by trial and error to maximize its cumulative reward, gradually improving its decision-making. This setup is especially suited to sequential problems such as game playing, robotics and control systems.

5 References

1. CalTech Machine Learning Course - CS156, Lecture 1.
2. CalTech Machine Learning Course - CS156, Lecture 4.