

Week 1 Lecture Notes

ML:Introduction

What is Machine Learning?

Two definitions of Machine Learning are offered. Arthur Samuel described it as: "the field of study that gives computers the ability to learn without being explicitly programmed." This is an older, informal definition.

Tom Mitchell provides a more modern definition: "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E ."

Example: playing checkers.

E = the experience of playing many games of checkers

T = the task of playing checkers.

P = the probability that the program will win the next game.

In general, any machine learning problem can be assigned to one of two broad classifications:

supervised learning, OR

unsupervised learning.

Supervised Learning

In supervised learning, we are given a data set and already know what our correct output should look like, having the idea that there is a relationship between the input and the output.

Supervised learning problems are categorized into "regression" and "classification" problems. In a regression problem, we are trying to predict results within a continuous output, meaning that we are trying to map input variables to some continuous function. In a classification problem, we are instead trying to predict results in a discrete output. In other words, we are trying to map input variables into discrete categories. Here is a description on Math is Fun on Continuous and Discrete Data.

Example 1:

Given data about the size of houses on the real estate market, try to predict their price. Price as a function of size is a continuous output, so this is a regression problem.

We could turn this example into a classification problem by instead making our output about whether the house "sells for more or less than the asking price." Here we are classifying the houses based on price into two discrete categories.

Example 2:

(a) Regression - Given a picture of Male/Female, We have to predict his/her age on the basis of given picture.

(b) Classification - Given a picture of Male/Female, we have to predict whether He/She is of High school, College, Graduate age. Another Example for Classification - Banks have to decide whether or not to give a loan to someone on the basis of his credit history.

Unsupervised Learning

Unsupervised learning, on the other hand, allows us to approach problems with little or no idea what our results should look like. We can derive structure from data where we don't necessarily know the effect of the variables.

We can derive this structure by clustering the data based on relationships among the variables in the data.

With unsupervised learning there is no feedback based on the prediction results, i.e., there is no teacher to correct you.

Example:

Clustering: Take a collection of 1800 essays written on the US Economy, and find a way to automatically group these essays into a small number that are somehow similar or related by different variables, such as word frequency, sentence length, page count, and so on.

Non-clustering: The "Cocktail Party Algorithm", which can find structure in messy data (such as the identification of individual voices and music from a mash of sounds at a cocktail party https://en.wikipedia.org/wiki/Cocktail_party_effect). Here is an answer on Quora to enhance your understanding. <https://www.quora.com/What-is-the-difference-between-supervised-and-unsupervised-learning-algorithms?>

ML:Linear Regression with One Variable

Model Representation

Recall that in regression problems, we are taking input variables and trying to fit the output onto a continuous expected result function.

Linear regression with one variable is also known as "univariate linear regression."

Univariate linear regression is used when you want to predict a **single output** value y from a **single input** value x . We're doing **supervised learning** here, so that means we already have an idea about what the input/output cause and effect should be.

The Hypothesis Function

Our hypothesis function has the general form:

$$\hat{y} = h_{\theta}(x) = \theta_0 + \theta_1 x$$

Note that this is like the equation of a straight line. We give to $\hat{h}_{\theta}(x)$ values for θ_0 and θ_1 to get our estimated output \hat{y} . In other words, we are trying to create a function called \hat{h}_{θ} that is trying to map our input data (the x 's) to our output data (the y 's).

Example:

Suppose we have the following set of training data:

input x	output y
0	4
1	7
2	7
3	6

Now we can make a random guess about our \hat{h}_{θ} function: $\theta_0 = 2$ and $\theta_1 = 2$. The hypothesis function becomes $\hat{h}_{\theta}(x) = 2 + 2x$.

So for input of 1 to our hypothesis, y will be 4. This is off by 3. Note that we will be trying out various values of θ_0 and θ_1 to try to find values which provide the best possible "fit" or the most representative "straight line" through the data points mapped on the x - y plane.

Cost Function

We can measure the accuracy of our hypothesis function by using a **cost function**. This takes an average (actually a fancier version of an average) of all the results of the hypothesis with inputs from x 's compared to the actual output y 's.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

To break it apart, it is $\frac{1}{2} J$ where J is the mean of the squares of $\hat{h}_{\theta}(x_i) - y_i$, or the difference between the predicted value and the actual value.

This function is otherwise called the "Squared-error function", or "Mean squared error". The mean is halved ($\frac{1}{2m}$) as a convenience for the computation of the gradient descent, as the derivative term of the square function will cancel out the $\frac{1}{2}$ term.

Now we are able to concretely measure the accuracy of our predictor function against the correct results we have so that we can predict new results we don't have.

If we try to think of it in visual terms, our training data set is scattered on the x - y plane. We are trying to make straight line (defined by $\hat{h}_{\theta}(x)$) which passes through this scattered set of data. Our objective is to get the best possible line. The best possible line will be such so that the average squared vertical distances of the scattered points from the line will be the least. In the best case, the line should pass through all the points of our training data set. In such a case the value of $J(\theta_0, \theta_1)$ will be 0.

ML:Gradient Descent

So we have our hypothesis function and we have a way of measuring how well it fits into the data. Now we need to estimate the parameters in hypothesis function. That's where gradient descent comes in.

Imagine that we graph our hypothesis function based on its fields θ_0 and θ_1 (actually we are graphing the cost function as a function of the parameter estimates). This can be kind of confusing; we are moving up to a higher level of abstraction. We are not graphing x and y itself, but the parameter range of our hypothesis function and the cost resulting from selecting particular set of parameters.

We put θ_0 on the x axis and θ_1 on the y axis, with the cost function on the vertical z axis. The points on our graph will be the result of the cost function using our hypothesis with these specific theta parameters.

We will know that we have succeeded when our cost function is at the very bottom of the pits in our graph, i.e. when its value is the minimum.

The way we do this is by taking the derivative (the tangential line to a function) of our cost function. The slope of the tangent is the derivative at that point and it will give us a direction to move towards. We make steps down the cost function in the direction with the steepest descent, and the size of each step is determined by the parameter α , which is called the learning rate.

The gradient descent algorithm is:

repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

where

$j \in \{0,1\}$ represents the feature index number.

Intuitively, this could be thought of as:

repeat until convergence:

$$\theta_j := \theta_j - \alpha [\text{slope of tangent aka derivative in } j \text{ dimension}] \text{ Slope of tangent aka derivative in } j \text{ dimension}$$

Gradient Descent for Linear Regression

When specifically applied to the case of linear regression, a new form of the gradient descent equation can be derived. We can substitute our actual cost function and our actual hypothesis function and modify the equation to (the derivation of the formulas are out of the scope of this course, but a really great one can be found here):

$$\begin{aligned} \text{repeat until convergence: } \\ \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_i(x_i) - y_i) \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_i(x_i) - y_i) x_i \end{aligned}$$

where m is the size of the training set, θ_0 a constant that will be changing simultaneously with θ_1 and x_i , y_i are values of the given training set (data).

Note that we have separated out the two cases for θ_j into separate equations for θ_0 and θ_1 ; and that for θ_1 we are multiplying x_i at the end due to the derivative.

The point of all this is that if we start with a guess for our hypothesis and then repeatedly apply these gradient descent equations, our hypothesis will become more and more accurate.

Gradient Descent for Linear Regression: visual worked example

Some may find the following video (<https://www.youtube.com/watch?v=Ng0xPM0350>) useful as it visualizes the improvement of the hypothesis as the error function reduces.

ML:Linear Algebra Review

Khan Academy has excellent Linear Algebra Tutorials (<https://www.khanacademy.org/linear-algebra/>)

Matrices and Vectors

Matrices are 2-dimensional arrays:

$$\begin{bmatrix} a & b & cd & e & fg & h & ij & k & l \end{bmatrix}$$

The above matrix has four rows and three columns, so it is a 4 x 3 matrix.

A vector is a matrix with one column and many rows:

$$\begin{bmatrix} wxyz \end{bmatrix}$$

So vectors are a subset of matrices. The above vector is a 4 x 1 matrix.

Notation and terms:

- A_{ij} refers to the element in the i th row and j th column of matrix A .
- A vector with ' n ' rows is referred to as an ' N -dimensional vector'
- θ_i refers to the element in the i th row of the vector.
- In general, all our vectors and matrices will be 1-indexed. Note that for some programming languages, the arrays are 0-indexed.
- Matrices are usually denoted by uppercase names while vectors are lowercase.
- "Scalar" means that an object is a single value, not a vector or matrix.
- \mathbb{R} refers to the set of scalar real numbers
- \mathbb{R}^n refers to the set of n -dimensional vectors of real numbers.

Addition and Scalar Multiplication

Addition and subtraction are **element-wise**, so you simply add or subtract each corresponding element:

$$\begin{bmatrix} a & bc & d \end{bmatrix} + \begin{bmatrix} w & xy & z \end{bmatrix} = \begin{bmatrix} a+w & b+x & c+y & d+z \end{bmatrix}$$

To add or subtract two matrices, their dimensions must be **the same**.

In scalar multiplication, we simply multiply every element by the scalar value:

$$\begin{bmatrix} a & bc & d \end{bmatrix} \times x = \begin{bmatrix} ax & bx & cx & dx \end{bmatrix}$$

Matrix-Vector Multiplication

We map the columns of the vector onto each row of the matrix, multiplying each element and summing the result.

$$\begin{bmatrix} a & bc & de & f \end{bmatrix} \times \begin{bmatrix} xyz \end{bmatrix} = \begin{bmatrix} ax+bx+cx+dx+fx \end{bmatrix}$$

The result is a **vector**. The vector must be the **second** term of the multiplication. The number of **columns** of the matrix must equal the number of **rows** of the vector.

An **$m \times n$ matrix** multiplied by an **$n \times 1$ vector** results in an **$m \times 1$ vector**.

Matrix-Matrix Multiplication

We multiply two matrices by breaking it into several vector multiplications and concatenating the result:

$$\begin{bmatrix} a & bc & de & f \end{bmatrix} \times \begin{bmatrix} w & xy & z \end{bmatrix} = \begin{bmatrix} aw+bx+cx+dw+fx & aw+bx+cx+dw+fx & aw+bx+cx+dw+fx & aw+bx+cx+dw+fx \end{bmatrix}$$

An **$m \times n$ matrix** multiplied by an **$n \times n$ matrix** results in an **$m \times n$ matrix**. In the above example, a 2 x 2 matrix times a 2 x 2 matrix resulted in a 2 x 2 matrix.

To multiply two matrices, the number of **columns** of the first matrix must equal the number of **rows** of the second matrix.

Matrix Multiplication Properties

- Not commutative. $A \times B \neq B \times A$
- Associative. $(A \times B) \times C = A \times (B \times C)$

The **Identity matrix**, when multiplied by any matrix of the same dimensions, results in the original matrix. It's just like multiplying numbers by 1. The identity matrix simply has 1's on the diagonal (upper left to lower right diagonal) and 0's elsewhere.

$$\begin{bmatrix} 1 & 0 & 00 & 1 & 00 & 0 & 1 \end{bmatrix}$$

When multiplying the identity matrix after some matrix $[A \times I]$, the square identity matrix should match the other matrix's **columns**. When multiplying the identity matrix before some other matrix $[I \times A]$, the square identity matrix should match the other matrix's **rows**.

Inverse and Transpose

The **inverse** of a matrix A is denoted A^{-1} . Multiplying by the inverse results in the identity matrix.

A non square matrix does not have an inverse matrix. We can compute inverses of matrices in octave with the `pinv(A)` function [1] and in matlab with the `inv(A)` function. Matrices that don't have an inverse are *singular* or *degenerate*.

The **transpose** of a matrix is like rotating the matrix 90° in clockwise direction and then reversing it. We can compute transposition of matrices in matlab with the `transpose(A)` function or `A'`:

$$A = \begin{bmatrix} a & bc & de & f \end{bmatrix}$$

$$A^T = \begin{bmatrix} a & c & eb & d & f \end{bmatrix}$$

In other words:

$$A_{ij} = A_{ji}^T$$