

# Supervised Learning: Linear Regression

April 7, 2021

## 1 Motivation

Imagine you are trying to predict the prices of houses. The only data you have available to you are the prices of a bunch of houses, and the square footage of all the houses. If given another house, alongside the square footage of said house, can you predict the price of the house?

If we assume that increasing square footage is approximately proportional to the price by a constant (the data is approximately linear), we can plot each point  $P_i = (x, y) = (sqft, \$\$)$  in a scatter plot and find a line of best fit to the data.

## 2 Fitting a Line

For any set of data *assuming the data is approximately linear*, we can use a linear model to predict the output of any new data that comes in. As a reminder, the linear model we are working with (for now) is:

$$\hat{y} = w_1x + w_2$$

where

$\hat{y} :=$  Predicted output of the linear model.

Note that if we increase  $w_1$ , we rotate the line CCW, if we decrease  $w_1$ , we rotate the line CW, and increasing/decreasing  $w_2$  pulls up/down the line.

## 3 Moving a Line Towards a Point

Let's assume the line we have to work with is:  $y(x) = w_1x + w_2$  ( $y(x)$  will be treated as  $y$  unless explicitly used as a function), where  $w_1, w_2$  are arbitrary constants, like a really bad guess as to what the actual line would be. Place a point reasonably far away from the line,  $P = (p, q)$ . How do we, using the data that we have, move the line such that it gets closer to the point?

### 3.1 Absolute Trick

The Absolute Trick is as such: Every iteration, we change  $w_1$  by the value  $p$  (the x coordinate), and change the value  $w_2$  by the value 1 (we will get to cases where we increase/decrease in a second). So our model transforms to:

$$y = w_1x + w_2 \rightarrow y = (w_1 \pm p)x + (w_2 \pm 1), \text{ where each } \pm \text{ is independent of one another}$$

If we continue under the assumption that the point is far away, we might get a good approximation. In all likelihood, however, we are likely to overshoot. So instead we iterate over this multiple times, taking steps towards the right solution. We can do this by multiplying  $p$  and 1 by  $\alpha$ , the learning rate. So our formula becomes:

$$y = w_1x + w_2 \rightarrow y = (w_1 \pm p\alpha)x + (w_2 \pm \alpha)$$

To determine whether to use  $+$  or  $-$ , we need to look at how we want the line to move. For  $w_1$ , we see which direction (CCW or CW) we should rotate the line to reach the point (smallest angle to sweep across so the line meets the point). For  $w_2$ , do the same thing, but this time your distance is up/down.

## 3.2 Square Trick

The above trick works to get near to the point, but we never factor in the y distance from the line (other than to determine the sign). What if we adjust  $w_1$  and  $w_2$  the same way as above, but also account for the vertical distance from the point to the line? We can do this by multiplying our  $p\alpha$  and  $\alpha$  term by  $q - \hat{q}$ , where  $\hat{q} = y(p)$ . So our formula for the square trick becomes:

$$y = w_1x + w_2 \rightarrow y = (w_1 + p\alpha(q - \hat{q}))x + (w_2 + \alpha(q - \hat{q}))$$

Note that we don't need to use  $\pm$  anymore because multiplying by  $q - \hat{q}$  means that if our point is below our line, this difference will yield a negative number, and if the point is above our line, the difference will yield a positive number.

## 3.3 Why not just do $\frac{q}{p}x$ as our line?

Obviously these tricks are pretty bad standalone operations, but their power is in the fact that they aren't exact. When we do linear regression, we're modifying our line of best fit so it coincides with all the data, and, unless all of our data points are co-linear (HIGHLY unlikely!), we want our line of best fit to be close to all of our data points. Think of each point pulling towards a line, the farther away the point from the line, the harder it pulls the line towards it.

## 4 What is the Line of Best Fit

Now that we have an iterative approach to pulling a line closer to a point, we need to define what a line of best fit is. If we define error (roughly) as the difference between our predicted y for each x value in our data,  $\hat{y} = y(x)$ , and the actual y value from each data point in our value, then our line of best fit is such that  $\sum f(|y - \hat{y}|)$ , such that  $f$  is a monotonic function, is at a minimum.

## 5 Linear Least Squares

Turns out, finding the exact line of best fit for a set of points  $(x_i, y_i)$  is an already solved problem, with a nice, closed form solution.

## 6 Gradient Descent

### 6.1 Types of Gradient Descent

## 7 Generalizing to Higher Dimensions

## 8 Warnings

## 9 Polynomial Regression

## 10 Regularization

## 11 Feature Scaling

## 12 External Resources