

annotated
version

Machine Learning Course - CS-433

Linear Regression

Sept 10, 2024

Martin Jaggi

Last updated on: September 10, 2024

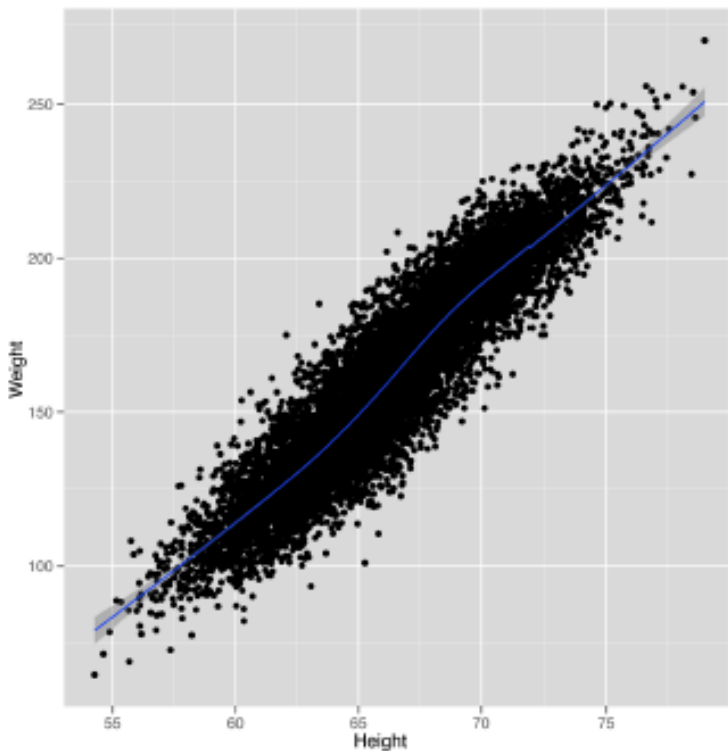
credits to Mohammad Emtiyaz Khan

EPFL

1 Model: Linear Regression

What is it?

Linear regression is a [model](#) that assumes a linear relationship between inputs and the output.



Why learn about linear regression?

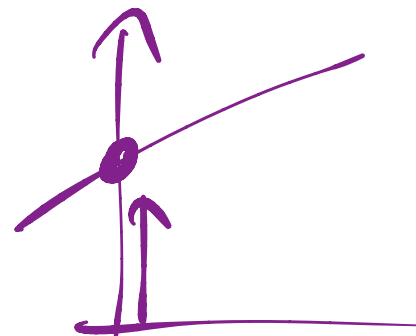
Plenty of reasons: simple, easy to understand, most widely used, easily generalized to non-linear models. Most importantly, you can learn almost all fundamental concepts of ML with regression alone.

Univariate linear regression

With only one input dimension, we get simple linear regression.

$$y_n \approx f(\mathbf{x}_n) := w_0 + w_1 x_{n1}$$

Here, $\mathbf{w} = (w_0, w_1)$ are the two parameters of the model. They describe f .

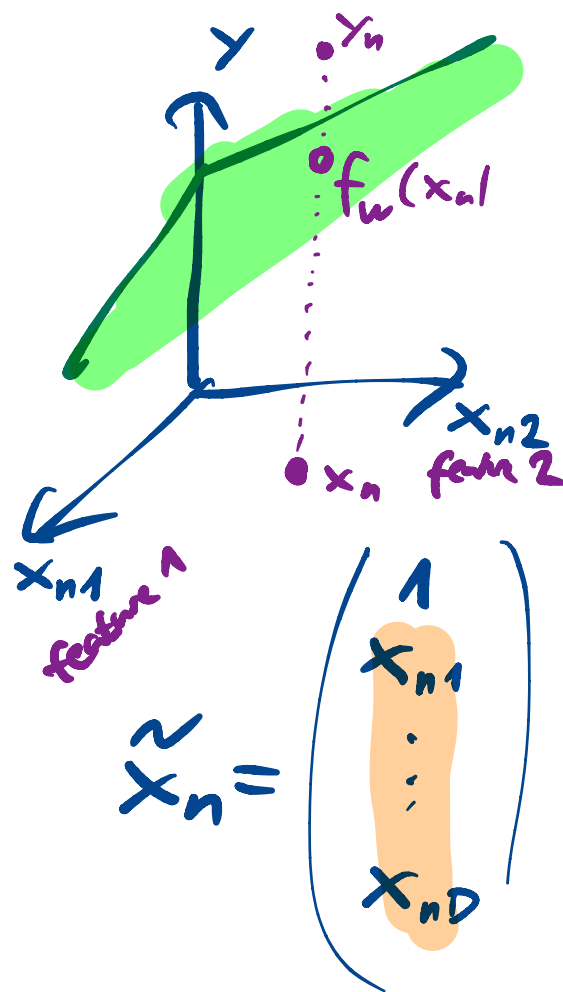


Multivariate linear regression

If our data has multiple input dimensions, we obtain multivariate linear regression.

$$\begin{aligned} y_n &\approx f(\mathbf{x}_n) = f_w(\mathbf{x}) \\ &:= w_0 + w_1 x_{n1} + \dots + w_D x_{nD} \\ &= w_0 + \mathbf{x}_n^\top \begin{pmatrix} w_1 \\ \vdots \\ w_D \end{pmatrix} \\ &=: \tilde{\mathbf{x}}_n^\top \tilde{\mathbf{w}} \end{aligned}$$

Note that we add a tilde over the input vector, and also the weights, to indicate they now contain the additional offset term (a.k.a. bias term).



$$\tilde{\mathbf{x}}_n = \begin{pmatrix} 1 \\ x_{n1} \\ \vdots \\ x_{nD} \end{pmatrix}$$

$$= \begin{pmatrix} 1 \\ \mathbf{x}_n \end{pmatrix} \in \mathbb{R}^{D+1}$$

$$f = \tilde{\mathbf{x}}_n^\top \tilde{\mathbf{w}}$$

Learning / Estimation / Fitting

Given data, we would like to find $\tilde{\mathbf{w}} = [w_0, w_1, \dots, w_D]$. This is called **learning** or **estimating** the parameters or **fitting** the model.

To do so, we need an optimization algorithm, which we will discuss in the chapter after the next.

Additional Notes

Alternative without using an 'offset' term

Above we have used $D + 1$ model parameters, to fit data of dimension D . An alternative also often used in practice, in particular for high-dimensional data, is to ignore the offset term w_0 .

$$\begin{aligned} y_n \approx f(\mathbf{x}_n) &:= w_1 x_{n1} + \dots + w_D x_{nD} \\ &= \mathbf{x}_n^\top \mathbf{w} \end{aligned}$$

in this case, we have just D parameters to learn, instead of $D + 1$.

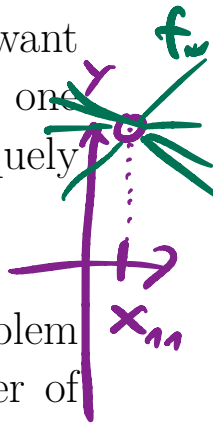
As a warning, you should be aware that for some machine learning models, the number of weight parameters (elements of \mathbf{w}) can in general be very different from D (being the dimension of the input data, the number of features). For an ML model where they do not coincide, think for example of a neural network (more details later).

Matrix multiplication

To go any further, one must revise matrix multiplication. Remember that multiplication of $M \times N$ matrix with a $N \times D$ matrix results in a $M \times D$ matrix. Also, two matrices of size $M \times N_1$ and $N_2 \times M$ can only be multiplied when $N_1 = N_2$.

Overparameterization — Or the $D > N$ issue

Consider the following simple situation: You have $N = 1$ and you want to fit $y_1 \approx w_0 + w_1 x_{11}$, i.e. you want to find $\mathbf{w} = (w_0, w_1)$ given one pair (y_1, x_{11}) . Is it possible to find such a line, and if so is it uniquely determined?



This identification issue is related to something called $D > N$ problem (in statistics typically named $p > n$). It means that the number of parameters exceeds number of data examples. In other words, we have more variables than we have data information. For many models, such as linear regression, this makes the task under-determined. We say that the model is over-parameterized for the task.

Using regularization is a way to avoid the identification issue described, which we will learn later. If we ignore the identification issue, overparameterization can also have advantages, such as on the training dynamics.