

# Linear Regression

Let's take a look at a probabilistic interpretation of linear regression. Why do we use least squares? What's the purpose of using squared error?

Let's assume that

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

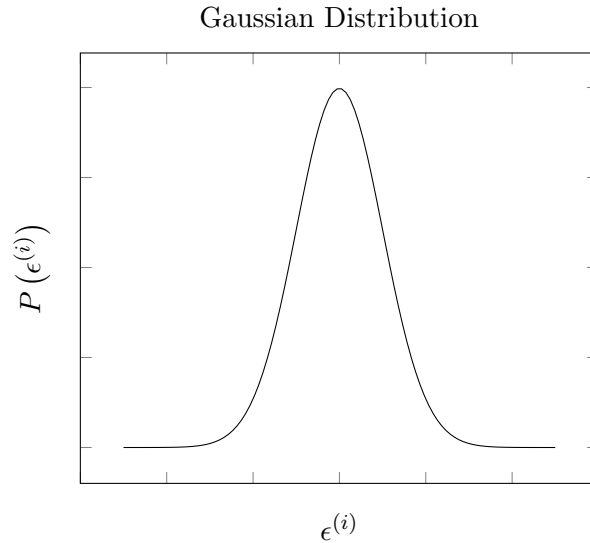
where  $\epsilon^{(i)}$  is the error term that includes unmodeled effects and random noise. We also assume that

$$\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

What this means is the probability density of the error  $P(\epsilon^{(i)})$  is equivalent to the Gaussian density equation with a mean of 0 and standard deviation  $\sigma$  ( $\sigma^2$  is the variance)

$$\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

which is a function that integrates to 1. We can see this distribution represented as a Gaussian as such



Under this set of assumptions, it is implied that

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

since  $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$ . Another way of representing this is with the distribution relation

$$y^{(i)} | x^{(i)}; \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2)$$

which tells us that our random variable is  $y^{(i)}$  and our mean is  $\theta^T x^{(i)}$  with a variance of  $\sigma^2$ . Therefore, we've concluded that, based on the previous assumptions, the probability density of  $y^{(i)}$  given  $x^{(i)}$  parameterized by  $\theta$  follows a Gaussian distribution.

Under the assumptions we just made, the likelihood of the parameters  $\mathcal{L}(\theta)$  is defined as the probability of the data  $P(\vec{y} | x; \theta)$  which is equivalent to

$$\prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

Substituting for the probability function in the product for what we determined previously,

$$\mathcal{L}(\theta) = \prod_{i=1}^m \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

where, as before,  $m$  is the number of training samples,  $\vec{y}$  is our target vector,  $x$  is our feature matrix, and  $\theta$  is our set of parameters.

We then define the log likelihood of the parameters  $\ell(\theta)$  to be  $\log \mathcal{L}(\theta)$ . Following the behavior of logarithm, we know that

$$\log \prod_{i=1}^m \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

is the same as

$$\sum_{i=1}^m \log \left[ \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \right]$$

because the log of a product is the same as the sum of the log of every term in the product. This can also be applied further to the expression, resulting in

$$\sum_{i=1}^m \log \frac{1}{\sigma\sqrt{2\pi}} + \sum_{i=1}^m \log \left[ \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \right]$$

We see that the sum term on the left does not include  $i$ , so it is simply just added to itself  $m$  times which we know as the definition of multiplication. Therefore this is equivalent to

$$m \log \frac{1}{\sigma\sqrt{2\pi}} + \sum_{i=1}^m \log \left[ \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \right]$$

which simplifies to

$$m \log \frac{1}{\sigma\sqrt{2\pi}} - \sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

When we want to get the maximum likelihood estimation, preferably we should use the log likelihood because it is a strictly monotonically increasing function and should also result in maximizing the likelihood. When we look back at our log likelihood equation,  $m \log \frac{1}{\sigma\sqrt{2\pi}}$  is a constant and our second term is negative but dependent on  $\theta$ . We can ignore  $\sigma$  because it is a constant. To maximize the likelihood, we want to minimize

$$\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

which we know is our cost function  $J(\theta)$ , which circles back to why we want to use squared error.