

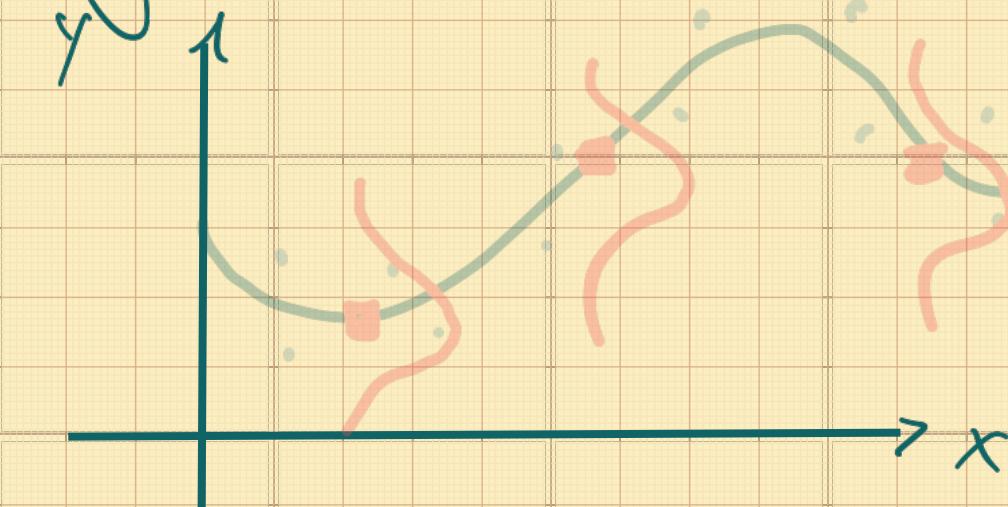
Linear Least Squares

Probabilistic P.D.V.: Signal/data contains noise, e.g. measurement error

→ regression with likelihood:

$$p(y|x) = \mathcal{N}(y | f(x), \sigma^2)$$

here, x is again the input, y is the output/target and we have $y = f(x) + \epsilon$, where ϵ is an i.i.d. Gaussian noise with zero μ and variance σ^2



- Note: $f(x)$ is the hidden generating function, that we can only observe through sampling
- Note: we assume that σ^2 is known and constant
- Linear regression means linear in the parameters Θ : (features can be combined non-linearly):

$$f(x) \approx w_1 x + w_0 \cdot 1 = x^T \Theta$$

$$= (x^T)^T \begin{pmatrix} w_1 \\ w_0 \end{pmatrix}$$

↙
features

↘
parameters

$$\text{so: } p(y|x, \theta) = \mathcal{N}(y | x^T \theta, \sigma^2), *$$

$$y = f(x) + \epsilon = x^T \theta + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

↑
linear in θ

Remark : iff $\sigma^2 \rightarrow 0$: $p(y|x, \theta) \rightarrow \delta$

How is (*) useful in finding a model?

⇒ classic estimator: MLE (see lecture statistics II)

maximizing the likelihood / prediction

of the training data D given the unknown, but
fixed parameters θ :

$$\theta_{\text{MLE}} := \arg \max_{\theta} p(D, \theta)$$

so more specifically:

$$\theta_{MLE} = \arg \max_{\theta} P(\underline{Y} | \underline{\bar{X}}, \theta)$$

feature
1
↑

set of N
training samples

$$\left(\begin{array}{c} (y_1, y_2, \dots, y_N)^T \\ \vdots \\ (x_1^0, x_1^1, x_2^0, x_2^1, \dots, x_N^0, x_N^1, \dots) \end{array} \right)$$

Since we assume the samples to be i.i.d., we can factorize the likelihood as:

$$P(\underline{Y} | \underline{\bar{X}}, \theta) = \prod_{m=1}^N P(y_m | x_m, \theta)$$

$$= \prod_{m=1}^N \mathcal{N}(y_m | x_m^T \theta, \sigma^2)$$

**

Remark: $P(\underline{Y} | \underline{X}, \Theta)$ is not a pdf in Θ , but it is one in \underline{Y}

Instead of maximizing (**), it is numerically easier to minimize the log-likelihood.

Since $\log(x)$ is strictly monotone, the extrema of x are identical to those of $\log(x)$!

[MMLB, §.2]

$$\begin{aligned} NLL: NLL(\Theta) &:= -\log P(\underline{Y} | \underline{X}, \Theta) \\ &= -\log \prod_{n=1}^N p(y_n | x_n, \Theta) \\ &= -\sum_{n=1}^N \log p(y_n | x_n, \Theta) \end{aligned}$$

Since we have $p(y_n | x_n, \theta) = \mathcal{N}(y_n | x_n^T \theta, \sigma^2)$:

$$\begin{aligned} NLL(\theta) &= -\sum_{n=1}^N \log \mathcal{N}(y_n | x_n^T \theta, \sigma^2) \\ &= -\sum_{n=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(y_n - x_n^T \theta)^2}{2\sigma^2} \end{aligned}$$

$$= -\sum_{n=1}^N \log \exp \left(\frac{-(y_n - x_n^T \theta)^2}{2\sigma^2} \right) + \text{const.}$$

$$= \frac{1}{2\sigma^2} \sum_{n=1}^N \frac{(y_n - x_n^T \theta)^2}{2\sigma^2} + \text{const'}$$

\downarrow
no influence
on minimum
location!

so, we need to find the minimum of

$$NLL^*(\theta) = \sum_{n=1}^N (y_n - \mathbf{x}_n^T \theta)^2$$

This is simply an L_2 -error / quadratic form minimization as we have seen before!

To make the connection clearer, rewrite NLL^* in vector form:

$$NLL^*(\theta) = (\mathbf{Y} - \mathbf{X}\theta)^T (\mathbf{Y} - \mathbf{X}\theta)$$

We have seen before that this can be seen as a quadratic form, for which the minimum is found by solving the associated linear system, which leads again to the

normal equation of LLS:

$$\hat{\theta}_{MLE} = \frac{\underline{\underline{X}}^T \underline{\underline{Y}}}{\underline{\underline{X}}^T \underline{\underline{X}}} \quad \rightarrow$$

see talk on
normal form.

→ More on MLE and MAP estimator +
regularized least squares : Bishop Sec 3.1.