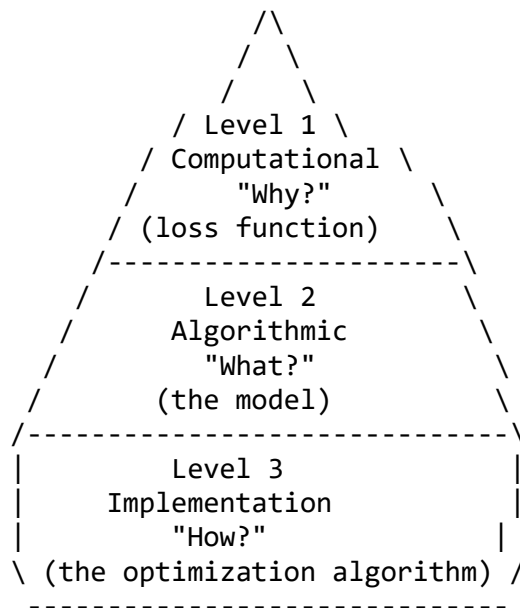


Aside: Marss's Levels of analysis

1. Computations "why?" eg loss function
2. Algorithmic "what?" eg the model
3. Implementation "how?" eg the optimization algorithm



How is dataset generated?

Let's say $((x, y) p(x, y))$, where (x) and (y) is some data.

Training: $(= \{(x, y) \}$

Assuming they are i.i.d.,

$$p(\mathcal{D}) = \prod p(x, y) = \prod p(x) p(y | x)$$

We are trying to learn $(p(y | x))$, basically it's a model of the true $(p(y | x))$.

A good model should make the data look probable.

Therefore, we must choose (θ) such that:

$$p(\mathcal{D}) = \prod p(x) p_{\theta}(y | x)$$

is maximized.

But the problem is that $(p(x))$, multiplying many such terms would make the result closer to zero.

Therefore, we take logarithm:

$$\begin{aligned}\log p(\mathcal{D}) &= \sum \log p(x_i) + \log p_{\theta}(y_i | x_i) \\ &= \sum \log p_{\theta}(y_i | x_i) + \text{const} \quad (p(x_i) \text{ doesn't depend on } \theta)\end{aligned}$$

Then,

$$\begin{aligned}\theta &\leftarrow \operatorname{argmax}_{\theta} \sum \log p_{\theta}(y_i | x_i) \quad (\text{MLE}) \\ \theta &\leftarrow \operatorname{argmin}_{\theta} - \sum \log p_{\theta}(y_i | x_i) \quad (\text{NLL})\end{aligned}$$

Therefore, what is loss function?

The loss function quantifies how bad θ is. We want the least bad θ .