

Negative Log Likelihood

Likelihood: $L(\text{model} \mid \text{data}) = P(\text{data} \mid \text{model})$

- How likely is a model given this data? \Rightarrow Used to compare diff. models
- Assume data is fixed \rightarrow find probability of getting this data based on model.

In ML, we want to maximize $L(\theta \mid x_1, x_2, \dots)$

- Given data is i.i.d., $L(\theta \mid x_1, x_2, \dots) = L(\theta \mid x_1) \cdot L(\theta \mid x_2) \cdot \dots \cdot L(\theta \mid x_n)$
 \Rightarrow Problem of underflow! keep multiplying probs = v. small value
- Log Likelihood turns products into sums:
 $\Rightarrow L(\theta \mid x_1, x_2, \dots) = \prod_{i=1}^n L(\theta \mid x_i)$
 $\Rightarrow \log L(\theta \mid x_1, x_2, \dots) = \sum_{i=1}^n \log L(\theta \mid x_i)$
- Hence we minimize "negative log likelihood" \Leftrightarrow maximize $\log L(\theta \mid x_1, x_2, \dots)$

NLL is the same as Cross-Entropy:

- Let actual vals = $y^{(n)} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \rightarrow$ one-hot vector
- Let predicted vals = $\hat{y}^{(n)} = \text{softmax}(\hat{z}) = \frac{1}{\sum e^{z_i}} \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ e^{z_3} \end{bmatrix} \rightarrow$ softmax is great as these are interpreted as probabilities!
- $\text{NLL} = -\log P_\theta(y \mid x_1, x_2, \dots)$
 $= -\log \prod_{n=1}^N P_\theta(y^{(n)} \mid x^{(n)})$
 $= -\sum_{n=1}^N \log P_\theta(y^{(n)} \mid x^{(n)}) \rightarrow$ prob my model assigns to $y^{(n)}$ given $x^{(n)}$
 $= -\sum_{n=1}^N \sum_{k=1}^K \underline{y_k^{(n)} \cdot \log \hat{y}_k^{(n)}} \rightarrow$ equal to prob. assigned to corr. class in
 \downarrow softmax!! ($y_k^{(n)}$ = one-hot; extracts prob of corr. class)
This is CROSS-ENTROPY!!
 $H(p, q) = -\sum_{k=1}^K p_k \log q_k$