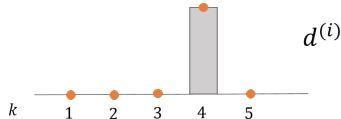


# Introduction to Machine Learning

## Information Theory

## Information Theory for Machine Learning



### Learning goals

- Minimizing KL = maximizing log-likelihood
- Minimizing KL = minimizing cross-entropy
- Minimizing CE between modeled and observed probabilities = log-loss minimization



# KL VS CROSS-ENTROPY

From this here we can see much more:

$$\arg \min_{\theta} D_{KL}(p||q_{\theta}) = \arg \min_{\theta} -\mathbb{E}_{x \sim p} \log q(x|\theta) = \arg \min_{\theta} H(p||q_{\theta})$$

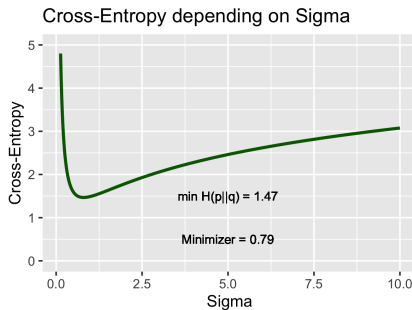
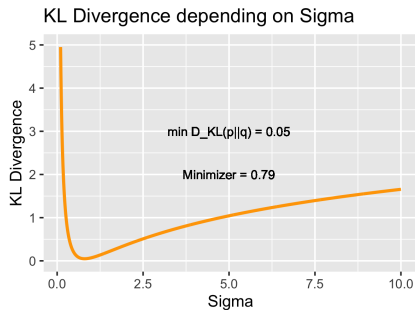
- So minimizing KL is the same as minimizing CE, is the same as maximum likelihood!
- We could now motivate CE as the "relevant" term that you have to minimize when you minimize KL - after you drop  $\mathbb{E}_p \log p(x)$ , which is simply the neg. entropy  $H(p)$ !
- Or we could say: CE between  $p$  and  $q$  is simply the expected negative log-likelihood of  $q$ , when our data comes from  $p$ !



# KL VS CROSS-ENTROPY EXAMPLE

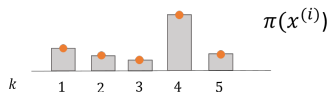
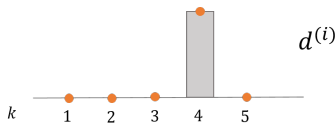
Let  $p(x) = N(0, 1)$  and  $q(x) = LP(0, \sigma)$  and consider again

$$\arg \min_{\theta} D_{KL}(p||q_{\theta}) = \arg \min_{\theta} -\mathbb{E}_{x \sim p} \log q(x|\theta) = \arg \min_{\theta} H(p||q_{\theta})$$



# CROSS-ENTROPY VS. LOG-LOSS

- Consider a multi-class classification task with dataset  $\mathcal{D} = ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$ .
- For  $g$  classes, each  $y^{(i)}$  can be one-hot-encoded as a vector  $d^{(i)}$  of length  $g$ .  $d^{(i)}$  can be interpreted as a categorical distribution which puts all its probability mass on the true label  $y^{(i)}$  of  $\mathbf{x}^{(i)}$ .
- $\pi(\mathbf{x}^{(i)}|\theta)$  is the probability output vector of the model, and also a categorical distribution over the classes.

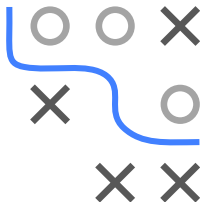


# CROSS-ENTROPY VS. LOG-LOSS / 2

To train the model, we minimize KL between  $d^{(i)}$  and  $\pi(\mathbf{x}^{(i)}|\theta)$  :

$$\arg \min_{\theta} \sum_{i=1}^n D_{KL}(d^{(i)} \parallel \pi(\mathbf{x}^{(i)}|\theta)) = \arg \min_{\theta} \sum_{i=1}^n H(d^{(i)} \parallel \pi(\mathbf{x}^{(i)}|\theta))$$

We see that this is equivalent to log-loss risk minimization!



$$\begin{aligned} R &= \sum_{i=1}^n H(d^{(i)} \parallel \pi_k(\mathbf{x}^{(i)}|\theta)) \\ &= \sum_{i=1}^n \left( - \sum_k d_k^{(i)} \log \pi_k(\mathbf{x}^{(i)}|\theta) \right) \\ &= \sum_{i=1}^n \underbrace{\left( - \sum_{k=1}^g [y^{(i)} = k] \log \pi_k(\mathbf{x}^{(i)}|\theta) \right)}_{\text{log loss}} \\ &= \sum_{i=1}^n (-\log \pi_{y^{(i)}}(\mathbf{x}^{(i)}|\theta)) \end{aligned}$$

# CROSS-ENTROPY VS. BERNOULLI LOSS

For completeness sake:

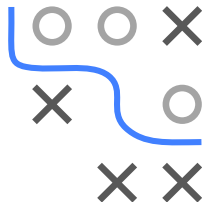
Let us use the Bernoulli loss for binary classification:

$$L(y, \pi(\mathbf{x})) = -y \log(\pi(\mathbf{x})) - (1 - y) \log(1 - \pi(\mathbf{x}))$$

If  $p$  represents a  $\text{Ber}(y)$  distribution (so deterministic, where the true label receives probability mass 1) and we also interpret  $\pi(\mathbf{x})$  as a Bernoulli distribution  $\text{Ber}(\pi(\mathbf{x}))$ , the Bernoulli loss  $L(y, \pi(\mathbf{x}))$  is the cross-entropy  $H(p \parallel \pi(\mathbf{x}))$ .



## ENTROPY AS PREDICTION LOSS



$$\begin{aligned}\mathcal{R} &= \frac{1}{n} \sum_{i=1}^n \left( - \sum_{k=1}^g [y^{(i)} = k] \log \pi_k \right) = - \frac{1}{n} \sum_{k=1}^g \sum_{i=1}^n [y^{(i)} = k] \log \pi_k \\ &= - \sum_{k=1}^g \frac{n_k}{n} \log \pi_k = - \sum_{k=1}^g \pi_k \log \pi_k = H(\pi)\end{aligned}$$