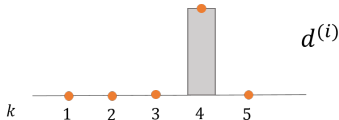


Introduction to Machine Learning

Information Theory for Machine Learning



Learning goals

- Minimizing KL is equivalent to maximizing the log-likelihood
- Minimizing KL is equivalent to minimizing cross-entropy
- Minimizing cross-entropy between modeled and observed probabilities is equivalent to log-loss minimization

KL VS MAXIMUM LIKELIHOOD

Minimizing KL between the true distribution $p(x)$ and approximating model $q(x|\theta)$ is equivalent to maximizing the log-likelihood.

$$\begin{aligned} D_{KL}(p\|q_\theta) &= \mathbb{E}_{x \sim p} \left[\log \frac{p(x)}{q(x|\theta)} \right] \\ &= \mathbb{E}_{x \sim p} \log p(x) - \mathbb{E}_{x \sim p} \log q(x|\theta) \end{aligned}$$

The first term above does not depend on θ . Therefore,

$$\begin{aligned} \arg \min_{\theta} D_{KL}(p\|q_\theta) &= \arg \min_{\theta} -\mathbb{E}_{x \sim p} \log q(x|\theta) \\ &= \arg \max_{\theta} \mathbb{E}_{x \sim p} \log q(x|\theta) \end{aligned}$$

For a finite dataset of n samples from p , this is approximated as

$$\arg \max_{\theta} \mathbb{E}_{x \sim p} \log q(x|\theta) \approx \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}^{(i)}|\theta).$$

KL VS CROSS-ENTROPY

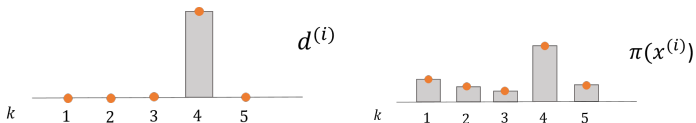
From this here we can actually see much more:

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

- So minimizing with respect to KL is the same as minimizing with respect to cross-entropy!
- That implies minimizing with respect to cross-entropy is the same as maximum likelihood!
- Remember, how we only characterized cross-entropy through source coding / bits? We could now motivate cross-entropy 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 entropy $H(p)$!
- Or we could say: Cross-entropy between p and q is simply the expected negative log-likelihood of q , when our data comes from p !

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

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

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

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

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

CROSS-ENTROPY VS. BERNOULLI LOSS

For completeness sake:

Let us use the Bernoulli loss for binary classification:

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

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 $H_{\pi(\mathbf{x})}(p)$.

ENTROPY AS PREDICTION LOSS

Assume log-loss for a situation where you only model with a constant probability vector π . We know the optimal model under that loss:

$$\pi_k = \frac{n_k}{n} = \frac{\sum_{i=1}^n [y^{(i)} = 1]}{n}$$

What is the (average) risk of that minimal constant model?

ENTROPY AS PREDICTION LOSS

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

So entropy is the (average) risk of the optimal "observed class frequency" model under log-loss!