Introduction to Machine Learning

Information Theory KL for ML





Learning goals

- Understand why measuring distribution similarity is important in ML
- Understand the advantages of forward and reverse KL

MEASURING DISTRIBUTION SIMILARITY IN ML

 Information theory provides tools (e.g., divergence measures) to quantify the similarity between probability distributions



× 0 0 × × ×

- The most prominent divergence measure is the KL divergence
- In ML, measuring (and maximizing) the similarity between probability distributions is a ubiquitous concept, which will be shown in the following.

MEASURING DISTRIBUTION SIMILARITY IN ML / 2

• Probabilistic model fitting

Assume our learner is probabilistic, i.e., we model $p(y|\mathbf{x})$ (for example, logistic regression, Gaussian process, ...).



× × ×

We want to minimize the difference between $p(y|\mathbf{x})$ and the conditional data generating process $\mathbb{P}_{y|\mathbf{x}}$ based on the data stemming from $\mathbb{P}_{y,\mathbf{x}}$.

Many losses can be derived this way. (e.g., cross-entropy loss)

MEASURING DISTRIBUTION SIMILARITY IN ML / 3

• Feature selection In feature selection, we want to choose features the target strongly depends on.



× 0 0 × 0 × ×

We can measure dependency by measuring the similarity between $p(\mathbf{x}, y)$ and $p(\mathbf{x}) \cdot p(y)$.

We will later see that measuring this similarity with KL leads to the concept of mutual information.

MEASURING DISTRIBUTION SIMILARITY IN ML / 4

• Variational inference (VI) By Bayes' theorem it holds that the posterior density

$$p(\theta|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)}{\int p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)d\theta}$$

However, computing the normalization constant $c = \int p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)d\theta$ analytically is usually intractable.



In VI, we want to fit a density q_{ϕ} with parameters ϕ to $p(\theta | \mathbf{X}, \mathbf{y})$.

× 0 0 × 0 × ×

KL DIVERGENCE

Divergences can be used to measure the similarity of distributions.

For distributions p, q they are defined such that

- **2** D(p,q) = 0 iff p = q.

 \Rightarrow divergences can be (and often are) non-symmetrical.

If the same measure dominates the distributions p, q, we can use KL. For a target distribution p and parametrized distribution q_{ϕ} , we call

- $D_{KL}(p \| q_{\phi})$ forward KL,
- $D_{KL}(q_{\phi} \| p)$ reverse KL.

In the following, we highlight some properties of the KL that make it attractive from an ML perspective.

× 0 0 × × ×

KL DIVERGENCE / 2

• Forward KL for probabilistic model fitting

We have samples from the DGP $p(y, \mathbf{x})$ when we fit our ML model.

If we have a probabilistic ML model q_{ϕ} the expected forward KL

$$\mathbb{E}_{\mathbf{x}\sim
ho_{\mathbf{x}}} \mathcal{D}_{\mathcal{KL}}(
ho(\cdot|\mathbf{x}) \| q_{\phi}(\cdot|\mathbf{x})) = \mathbb{E}_{\mathbf{x}\sim
ho_{\mathbf{x}}} \mathbb{E}_{y\sim
ho_{y|\mathbf{x}}} \log\left(rac{
ho(y|\mathbf{x})}{q_{\phi}(y|\mathbf{x})}
ight).$$

We can directly minimize this objective since

$$egin{aligned}
abla_{\phi} & \mathbb{E}_{\mathbf{x} \sim
ho_{\mathbf{x}}} D_{\mathcal{K}L}(
ho(\cdot | \mathbf{x}) \| q_{\phi}(\cdot | \mathbf{x})) &= \mathbb{E}_{\mathbf{x} \sim
ho_{\mathbf{x}}} \mathbb{E}_{y \sim
ho_{y|\mathbf{x}}}
abla_{\phi} \log \left(
ho(y|\mathbf{x})
ight) \ &- \mathbb{E}_{\mathbf{x} \sim
ho_{\mathbf{x}}} \mathbb{E}_{y \sim
ho_{y|\mathbf{x}}}
abla_{\phi} \log \left(q_{\phi}(y|\mathbf{x})
ight) \ &= -
abla_{\phi} \mathbb{E}_{\mathbf{x} \sim
ho_{\mathbf{x}}} \mathbb{E}_{y \sim
ho_{y|\mathbf{x}}} \mathbb{E}_{y \sim
ho_{y|\mathbf{x}}} \log \left(q_{\phi}(y|\mathbf{x})
ight) \end{aligned}$$

 \Rightarrow We can estimate the gradient of the expected forward KL without bias, although we can not evaluate $p(y|\mathbf{x})$ in general.

× × 0 × × ×

KL DIVERGENCE / 3

• Reverse KL for VI

Here, we know our target density $p(\theta | \mathbf{X}, \mathbf{y})$ only up to the normalization constant, and we do not have samples from it.

We can directly apply the reverse KL since for any $c \in \mathbb{R}_+$

$$egin{aligned}
abla_{\phi} \mathcal{D}_{\mathcal{KL}}(q_{\phi} \| p) &=
abla_{\phi} \mathbb{E}_{ heta \sim q_{\phi}} \log\left(rac{q_{\phi}(heta)}{p(heta)}
ight) \ &=
abla_{\phi} \mathbb{E}_{ heta \sim q_{\phi}} \log\left(rac{q_{\phi}(heta)}{p(heta)}
ight) - \underbrace{
abla_{\phi} \mathbb{E}_{ heta \sim q_{\phi}} \log c}_{=0} \ &=
abla_{\phi} \mathbb{E}_{ heta \sim q_{\phi}} \log\left(rac{q_{\phi}(heta)}{c \cdot p(heta)}
ight). \end{aligned}$$

 \Rightarrow We can estimate the gradient of the reverse KL without bias (even if we only have an unnormalized target distribution)

× × ×

KL DIVERGENCE / 4

The asymmetry of the KL has the following implications

- Forward KL $D_{KL}(p || q_{\phi}) = \mathbb{E}_{\mathbf{x} \sim p} \log \left(\frac{p(\mathbf{x})}{q_{\phi}(\mathbf{x})} \right)$ is mass-covering since $p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{q_{\phi}(\mathbf{x})} \right) \approx 0$ if $p(\mathbf{x}) \approx 0$ and $q_{\phi}(\mathbf{x}) \gg p(\mathbf{x})$.
- Reverse KL $D_{KL}(q_{\phi} || p) = \mathbb{E}_{\mathbf{x} \sim q_{\phi}} \log \left(\frac{q_{\phi}(\mathbf{x})}{p(\mathbf{x})} \right)$ is mode-seeking (zero-avoiding) since $q_{\phi}(\mathbf{x}) \log \left(\frac{q_{\phi}(\mathbf{x})}{p(\mathbf{x})} \right) \gg 0$ if $p(\mathbf{x}) \approx 0$ and $q_{\phi}(\mathbf{x}) > 0$





Figure: Optimal q_{ϕ} when q_{ϕ} is restricted to be Gaussian.