# TTIC 31230, Fundamentals of Deep Learning

David McAllester, Autumn 2020

Some Information Theory

1

# Why Information Theory?

The fundamental equation involves cross-entropy.

Cross-entropy is an information-theoretic concept.

Information theory arises in many places and many forms in deep learning.

# Entropy of a Distribution

The entropy of a distribution $P$ is defined by

$$H(P) = E_{y \sim P}\left[-\ln P(y)\right] \text{ in units of ``nats''}$$

$$H_2(P) = E_{y \sim P}\left[-\log_2 P(y)\right] \text{ in units of bits}$$

# Why Bits?

Why is $-\log_2 P(y)$ a number of bits?

Example: Let $P$ be a uniform distribution on 256 values.

$$E_{y \sim P} \left[ -\log_2 P(y) \right] = -\log_2 \frac{1}{256} = \log_2 256 = 8 \text{ bits} = 1 \text{ byte}$$

$$1 \text{ nat} = \frac{1}{\ln 2} \text{ bits} \approx 1.44 \text{ bits}$$

# Shannon's Source Coding Theorem

Why is $-\log_2 P(y)$ a number of bits?

A prefix-free code for $\mathcal{Y}$ assigns a bit string $c(y)$ to each $y \in \mathcal{Y}$ such that no code string is prefix of any other code string.

For a probability distribution $P$ on $\mathcal{Y}$ we consider the average code length $E_{y \sim P} [\,|c(y)|\,]$.

Theorem: For any $c$ we have $E_{y \sim P} |c(y)| \geq H_2(P)$.

Theorem: There exists $c$ with $E_{y \sim P} |c(y)| \leq H_2(P) + 1$.

# Cross Entropy

Let $P$ and $Q$ be two distribution on the same set.

$$H(P, Q) = E_{y \sim P} \left[ -\ln Q(y) \right]$$

$$\Phi^* = \operatorname*{argmin}_{\Phi} H(\mathrm{Pop}, P_\Phi)$$

H(P,Q) also has a data compression interpretation.

$H(P, Q)$ can be interpreted as 1.44 times the number of bits used to code draws from $P$ when using the imperfect code defined by $Q$.

# Entropy, Cross Entropy and KL Divergence

Let $P$ and $Q$ be two distribution on the same set.

Entropy : $\qquad\qquad H(P) = E_{y \sim P}\ [-\ln\ P(y)]$

CrossEntropy : $\qquad H(P,Q) = E_{y \sim P}\ [-\ln\ Q(y)]$

KL Divergence : $KL(P,Q) = H(P,Q) - H(P)$

$$= E_{y \sim P}\quad \ln\ \frac{P(y)}{Q(y)}$$

We have $H(P,Q) \geq H(P)$ or equivalently $KL(P,Q) \geq 0$.

# The Universality Assumption

$$\Phi^* = \underset{\Phi}{\mathrm{argmin}} \ H(\mathrm{Pop}, P_\Phi) = \underset{\Phi}{\mathrm{argmin}} \ H(\mathrm{Pop}) + KL(\mathrm{Pop}, P_\Phi)$$

Universality assumption: $P_\Phi$ can represent any distribution and $\Phi$ can be fully optimized.

This is clearly false for deep networks. But it gives important insights like:

$$P_{\Phi^*} = \mathrm{Pop}$$

This is the motivatation for the fundamental equation.

# Asymmetry of Cross Entropy

Consider

$$\Phi^* = \operatorname*{argmin}_{\Phi} \ H(P, Q_\Phi) \qquad (1)$$

$$\Phi^* = \operatorname*{argmin}_{\Phi} \ H(Q_\Phi, P) \qquad (2)$$

For (1) $Q_\Phi$ must cover all of the support of $P$.

For (2) $Q_\Phi$ concentrates all mass on the point maximizing $P$.
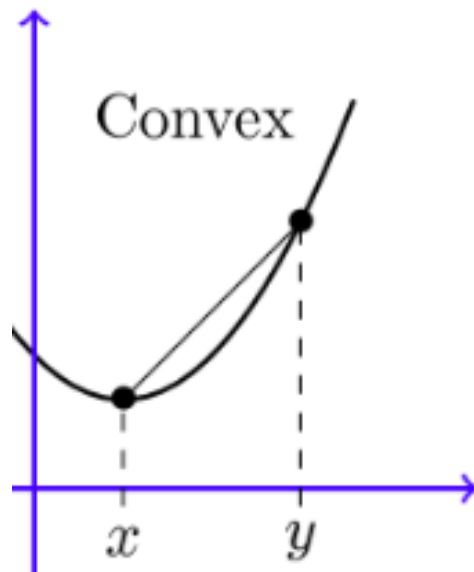
# Asymmetry of KL Divergence

Consider

$$\Phi^* = \operatorname*{argmin}_{\Phi}\ KL(P, Q_\Phi)$$

$$= \operatorname*{argmin}_{\Phi}\ H(P, Q_\Phi) \qquad\qquad (1)$$

$$\Phi^* = \operatorname*{argmin}_{\Phi}\ KL(Q_\Phi, P)$$

$$= \operatorname*{argmin}_{\Phi} H(Q_\Phi, P) - H(Q_\Phi) \quad (2)$$

If $Q_\Phi$ is not universally expressive we have that (1) still forces $Q_\Phi$ to cover all of $P$ (or else the KL divergence is infinite) while (2) allows $Q_\Phi$ to be restricted to a single mode of $P$ (a common outcome).

# Proving $KL(P, Q) \geq 0$: Jensen's Inequality



For $f$ convex (upward curving) we have

$$E[f(x)] \geq f(E[x])$$

# Proving $KL(P, Q) \geq 0$

$$KL(P, Q) = E_{y \sim P} \left[ -\ln \frac{Q(y)}{P(y)} \right]$$

$$\geq -\ln E_{y \sim P} \frac{Q(y)}{P(y)}$$

$$= -\ln \sum_y P(y) \frac{Q(y)}{P(y)}$$

$$= -\ln \sum_y Q(y)$$

$$= 0$$

# Appendix: The Rearrangement Trick

$$KL(P, Q) = E_{x \sim P} \left[ \ln \frac{P(x)}{Q(x)} \right]$$

$$= E_{x \sim P} \left[ (-\ln Q(x)) - (-\ln P(x)) \right]$$

$$= \left( E_{x \sim P} \left[ -\ln Q(x) \right] \right) - \left( E_{x \sim P} \left[ -\ln P(x) \right] \right)$$

$$= H(P, Q) - H(P)$$

In general $E_{x \sim P} \ln \left( \prod_i A_i \right) = E_{x \sim P} \sum_i \ln A_i$

# Summary

$$\Phi^* = \text{argmin}_\Phi \; H(\text{Pop}, P_\Phi) \text{ unconditional}$$

$$\Phi^* = \text{argmin}_\Phi \; E_{x \sim \text{Pop}} \; H(\text{Pop}(y|x), P_\Phi(y|x)) \text{ conditional}$$

$$\text{Entropy}: \qquad H(P) = E_{y \sim P} \; [-\ln \; P(y)]$$

$$\text{CrossEntropy}: \quad H(P, Q) = E_{y \sim P} \; [-\ln \; Q(y)]$$

$$\text{KL Divergence}: \; KL(P, Q) = H(P, Q) - H(P)$$

$$= E_{y \sim P} \quad \ln \; \frac{P(y)}{Q(y)}$$

$$H(P, Q) \geq H(P), \quad KL(P, Q) \geq 0, \quad \text{argmin}_Q \; H(P, Q) = P$$

END