

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2018

Information Theory

15th Century Voynich Manuscript



Appears to be written in Hebrew.

Announced this past week by Greg Kondrak of the University of Alberta.

Logistic Regression as a Major Advance

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Train}} (y - f_{\Phi}(x))^2$$

Switched to

$$Q_{\Phi}(y|x) = \operatorname{softmax}_{\hat{y}} f_{\Phi}(\hat{y}|x)$$

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Train}} -\log Q_{\Phi}(y|x)$$

Binary Classification

We have a population distribution over (x, y) with $y \in \{-1, 1\}$.

We compute a single number $f_\Phi(x)$ where

for $f_\Phi(x) \geq 0$ predict $y = 1$

for $f_\Phi(x) < 0$ predict $y = -1$

Softmax for Binary Classification

$$\begin{aligned} Q_{\Phi}(y|x) &= \frac{1}{Z} e^{yf(x)} \\ &= \frac{e^{yf(x)}}{e^{yf(x)} + e^{-yf(x)}} \\ &= \frac{1}{1 + e^{-2yf(x)}} \\ &= \frac{1}{1 + e^{-m(y)}} \end{aligned} \quad m(y|x) = 2yf(x) \text{ is the margin}$$

Logistic Regression for Binary Classification

$$\begin{aligned}\Phi^* &= \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Train}} \ln 1/Q_{\Phi}(y|x) \\ &= \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Train}} \ln \left(1 + e^{-m(y|x)} \right)\end{aligned}$$

$$\ln \left(1 + e^{-m(y|x)} \right) \approx 0 \quad \text{for } m(y|x) \gg 1$$

$$\ln \left(1 + e^{-m(y|x)} \right) \approx -m(y|x) \quad \text{for } m(y|x) \ll -1$$

Log Loss vs. Hinge Loss

Log loss:

$$\begin{aligned}\Phi^* &= \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Train}} - \ln Q_{\Phi}(y|x) \\ &= \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Train}} \ln \left(1 + e^{-m(y|x)} \right) \quad \text{binary case}\end{aligned}$$

Hinge Loss:

$$\begin{aligned}\Phi^* &= \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Train}} \max(0, 1 - m(y|x)) \\ m(y|x) &= \min_{\hat{y}} f(y|x) - f(\hat{y}|x)\end{aligned}$$

Log Loss vs. Hinge Loss

We will show that log loss is a consistent probability estimator.

For log loss, and for infinite training data, $Q^*(y|x)$ is the true population conditional probability.

Hinge loss is a consistent classifier but not a consistent probability estimator.

$$\begin{aligned} f^*(y^*|x) &= 1 \\ f^*(\hat{y}|x) &= 0 \quad \text{for } \hat{y} \neq y^* \end{aligned}$$

Entropy

Consider a probability distribution Pop on a finite set S .

Consider a code C assigning a bit string code word $C(y_1, \dots, y_B)$ to each possible batch of B elements with $y_i \sim \text{Pop}$.

Source coding theorem: As $B \rightarrow \infty$ the optimal coding uses exactly

$$H(\text{Pop}) = E_{y \sim \text{Pop}} - \log_2 \text{Pop}(y)$$

bits per batch element.

Prefix Free Codes

Let S be a finite set.

Let C be assignment of a bit string $C(y)$ to each $y \in S$.

C is called *prefix-free* if for $x \neq y$ we have that $C(x)$ is not a prefix of $C(y)$.

A concatenation of sequence of prefix-free code words can be uniquely segmented (parsed) back into a sequence of code words.

Prefix-Free Codes as Trees and as Probabilities

A prefix-free code defines a binary branching tree — branch on the first code bit, then the second, and so on.

The leaves of this tree are labeled with the elements of S .

The code defines a probability distribution on S by randomly selecting branches.

We have $Q_C(y) = 2^{-|C(y)|}$.

The Source Coding Theorem

(1) There exists a prefix-free code C such that

$$|C(y)| \leq (-\log_2 \text{Pop}(y)) + 1$$

and hence

$$E_{y \sim \text{Pop}} |C(y)| \leq H(\text{Pop}) + 1$$

(2) For any prefix-free code C

$$E_{y \sim \text{Pop}} |C(y)| \geq H(\text{Pop})$$

Code Construction

We construct a code by iterating over $y \in S$ in order of decreasing probability (most likely first).

For each y select a code word $C(y)$ (a tree leaf) with length (depth)

$$|C(y)| = \lceil -\log_2 \text{Pop}(y) \rceil$$

and where $C(y)$ is not an extension of (under) any previously selected code word.

Code Existence Proof

At any point before coding all elements of S we have

$$\sum_{y \in \text{Defined}} 2^{-|C(y)|} \leq \sum_{y \in \text{Defined}} \text{Pop}(y) < 1$$

Therefore there exists an infinite descent into the tree that misses all previous code words.

Hence there exists a code word $C(x)$ not under any previous code word with $|C(x)| = \lceil -\log_2 \text{Pop}(y) \rceil$.

Furthermore $C(x)$ is at least as long as all previous code words and hence $C(x)$ is not a prefix of any previously selected code word.

Huffman Coding

Maintain a list of trees T_1, \dots, T_N .

Initially each tree is just one root node labeled with an element of S .

Each tree T_i has a weight equal to the sum of the probabilities of the nodes on the leaves of that tree.

Repeatedly merge the two trees of lowest weight into a single tree until all trees are merged.

Optimality of Huffman Coding

Theorem: The Huffman code T for Pop is optimal — for any other tree T' we have $d(T; \text{Pop}) \leq d(T'; \text{Pop})$.

Proof: The algorithm maintains the invariant that there exists an optimal tree including all the subtrees on the list.

To prove that a merge operation maintains this invariant we consider any tree containing the given subtrees.

Consider the two subtrees T_i and T_j of minimal weight. Without loss of generality we can assume that T_i is at least as deep as T_j .

Swapping the sibling of T_i for T_j brings T_i and T_j together and can only improve the average depth.

Back to Log Loss

Log loss has both a conditional and an unconditional version.

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{(x,y) \sim P_{\text{op}}} - \log Q_{\Phi}(y|x)$$

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim P_{\text{op}}} - \log Q_{\Phi}(y)$$

Conditional test loss can often be made small (MNIST, speech recognition) but can also be inherently large (image colorization).

Unconditional log loss is typically inherently large (language modeling).

Log Loss is Cross Entropy

$$H(\text{Pop}) \doteq E_{y \sim \text{Pop}} - \log \text{Pop}(y)$$

$$H(\text{Pop}, Q) \doteq E_{y \sim \text{Pop}} - \log Q(y)$$

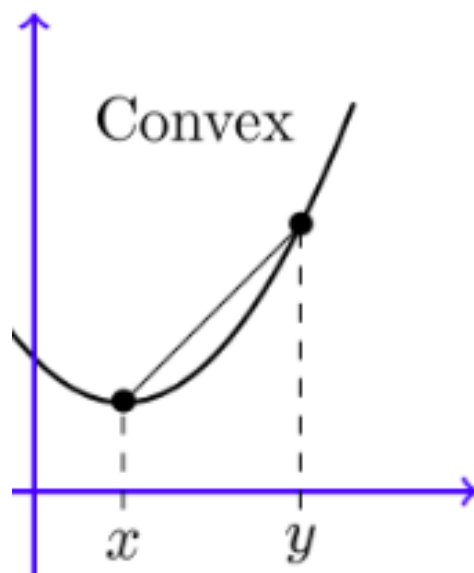
$$E_{y \sim \text{Pop}} - \log Q_\Phi(y) = H(\text{Pop}, Q_\Phi)$$

$$E_{(x,y) \sim \text{Pop}} - \log Q_\Phi(y|x) = E_{x \sim \text{Pop}} H(\text{Pop}(y|x), Q_\Phi(y|x))$$

KL Divergence

$$\begin{aligned} KL(\text{Pop}, Q) &\doteq E_{y \sim \text{Pop}} \log_2 \frac{\text{Pop}(y)}{Q(y)} \\ &= E_{y \sim \text{Pop}} \log \text{Pop}(y) - \log Q(y) \\ &= (E_{y \sim \text{Pop}} - \log Q(y)) - (E_{y \sim \text{Pop}} - \log \text{Pop}(y)) \\ &= H(\text{Pop}, Q) - H(P) \end{aligned}$$

Jensen's Inequality



For f convex (upward curving) we have

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

KL Divergence

$$\begin{aligned} KL(P, Q) &= \mathbb{E}_{y \sim P} \left[-\log \frac{Q(y)}{P(y)} \right] \\ &\geq -\log \mathbb{E}_{x \sim P} \left[\frac{Q(y)}{P(y)} \right] \\ &= -\log \sum_y P(y) \frac{Q(y)}{P(y)} \\ &= -\log \sum_y Q(y) \\ &= 0 \end{aligned}$$

Fundamentals

$$KL(\text{Pop}, Q) \geq 0$$

$$H(\text{Pop}, Q) = H(\text{Pop}) + KL(\text{Pop}, Q)$$

$$\geq H(\text{Pop})$$

$$\underset{Q}{\operatorname{argmin}} H(\text{Pop}, Q) = \text{Pop}$$

$$\underset{Q(y|x)}{\operatorname{argmin}} E_{x \sim \text{Pop}} H(\text{Pop}(y|x), Q(y|x)) = \text{Pop}(y|x)$$

Variational Methods

In deep learning the term “variational” often just means making the following approximation.

$$H(P) \approx \inf_{\Phi} H(P, Q_{\Phi})$$

Asymmetry of Cross Entropy

Consider

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

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

For (1) Q_{Φ} must cover all of the support of P .

For (2) Q_{Φ} concentrates all mass on the point maximizing P .

Unsupervised Learning

Unsupervised learning is sometimes equated with unconditional log loss (density estimation).

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim P_{\text{op}}} - \log Q_{\Phi}(y)$$

Unsupervised Learning

■ “Pure” Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

Unsupervised Learning

By “unsupervised learning” we will mean learning from **massively available** data. This is not a mathematical definition.

Massive: images, audio, text, video, click-through data.

Less Massive: car control data, stereo image pairs, closed captioned video, captioned images.

Big: Manually annotated images or audio.

Small: manually annotated text — parse trees, named entities, semantic roles, coreference, entailment.

Smallest: Manually annotated text in an obscure language.

Colorization

$$\Phi^* = \operatorname{argmin}_{\Phi} \mathbb{E}_{(x,y) \sim \text{Pop}} [-\log Q_{\Phi}(y|x)]$$



We have massive data for colorization.

But any colorization is inevitably a guess.

Differential Entropy

Consider a continuous density $p(x)$. For example

$$p(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{\frac{-x^2}{2\sigma^2}}$$

Differential entropy is often defined as

$$H(p) \doteq \int \left(\ln \frac{1}{p(x)} \right) p(x) dx$$

Finite Differential Entropy is Not Meaningful

$$\begin{aligned} H(\mathcal{N}(0, \sigma)) &= + \int \left(\ln(\sqrt{2\pi}\sigma) + \frac{x^2}{2\sigma^2} \right) p(x) dx \\ &= \ln(\sigma) + \ln(\sqrt{2\pi}) + \frac{1}{2} \end{aligned}$$

But if we take $y \doteq x/2$ we get $H(y) = H(x) - \ln 2$.

Also for $\sigma \ll 1$, we get $H(p) < 0$

Hence differential entropy then depends on the choice of units — a distributions on lengths will have a different entropy when measuring in inches than when measuring in feet.

Differential Entropy is Always Infinite

Consider quantizing the the real numbers into bins.

A continuous probability density p assigns a probability $p(B)$ to each bin.

As the bin size decreases toward zero the entropy of the bin distribution increases toward ∞ .

A meaningful convention is that $H(p) = +\infty$ for any continuous density p .

Differential KL-divergence is Meaningful

$$KL(p, q) = \int \left(\ln \frac{p(x)}{q(x)} \right) p(x) dx$$

This integral can be computed by dividing the real numbers into bins and computing the KL divergence between the distributions on bins.

The KL divergence between the bin distribution typically approaches a finite limit as the bin size goes to zero.

KL-Divergence can also be Infinite

$$KL(p, q) = E_{x \sim p} \log \frac{p(x)}{q(x)}$$

In either the discrete or continuous case, if a set is assigned nonzero probability by p but zero probability by q then $KL(p, q) = +\infty$.

If every set assigned nonzero probability by p is also assigned nonzero probability by q then we say that p is absolutely continuous with respect to q .

Random Variables

We consider variables where a single draw from the population determines a value for each variable.

This is the formal definition of a “random variable”.

Each random variable has a probability distribution defined by the distribution on the population.

We write $H(x)$ for the entropy of the distribution on x .

Mutual Information

For two random variables x and y there is a distribution on pairs (x, y) determined by the population distribution.

Mutual information concerns the relationship between the distribution on (x, y) and the marginal distributions on x and y .

For the discrete case we can write.

$$I(x, y) \doteq H(x) + H(y) - H(x, y)$$

This can be viewed as a quantity of non-independence — independent variables have zero mutual information.

Conditional Entropy

For the discrete case conditional entropy $H(y|x)$ is defined by

$$\begin{aligned} H(y|x) &\doteq \sum_x \text{Pop}(x) \sum_y \text{Pop}(y|x) - \log \text{Pop}(y|x) \\ &= E_{x \sim \text{Pop}} E_{y \sim \text{Pop}|x} - \log \text{Pop}(y|x) \\ &= E_{x \sim \text{Pop}} H(\text{Pop}(y|x)) \end{aligned}$$

More Identities

For the discrete case we have.

$$I(x, y) = H(x) - H(x|y)$$

$$= H(y) - H(y|x)$$

$$= KL(\text{Pop}(x, y), \text{Pop}(x) \times \text{Pop}(y))$$

The last identity can be taken as a definition of $I(x, y)$ in the continuous case.

END