CMPT-413 Computational Linguistics

Anoop Sarkar http://www.cs.sfu.ca/~anoop

February 13, 2008

Quick Guide to Probability Theory

Log Probability

Basics of Information Theory

Quick guide to probability theory



- P(baby is a girl) = 0.5 percentage of total number of babies that are girls
- P(baby girl is named Kiki) = 0.001 percentage of total number of babies that are named Kiki



Joint probability

► P(X,Y) means probability that X and Y are both true

 P(baby girl, blue eyes) percentage of total number of babies that are girls and have blue eyes



Conditional probability

P(X | Y) means probability that X is true when we already know that Y is true

- P(baby is named Kiki | baby is a girl) = 0.002
- P(baby is a girl | baby is named Kiki) = 1



Conditional probability

Conditional and joint probabilities are related:

$$P(X \mid Y) = \frac{P(X, Y)}{P(Y)}$$



Bayes rule

Conditional probability re-written as likelihood times prior:

$$P(X \mid Y) = \frac{P(Y \mid X) \times P(X)}{P(Y)}$$



Bayes Rule

$$P(X \mid Y) = \frac{P(X, Y)}{P(Y)}$$
(1)

$$P(Y \mid X) = \frac{P(Y, X)}{P(X)}$$
(2)

$$P(X, Y) = P(Y, X)$$
(3)

$$P(X \mid Y) \times P(Y) = P(Y \mid X) \times P(X)$$
(4)

$$P(X \mid Y) = \frac{P(Y \mid X) \times P(X)}{P(Y)}$$
(5)

$$P(X \mid Y) = P(Y \mid X) \times P(X)$$
(6)

Basic Terms

- ▶ P(e) a priori probability or just prior
- ▶ P(f | e) conditional probability. The chance of f given e
- ► P(e, f) joint probability. The chance of e and f both happening.
- ▶ If e and f are *independent* then we can write $P(e, f) = P(e) \times P(f)$
- If e and f are not independent then we can write P(e, f) = P(e) × P(f | e) P(e, f) = P(f) × ?

Basic Terms

Addition of integers:

$$\sum_{i=1}^{n} i = 1 + 2 + 3 + \ldots + n$$

Product of integers:

$$\prod_{i=1}^n i = 1 \times 2 \times 3 \times \ldots \times n$$

► Factoring:

$$\sum_{i=1}^{n} i \times k = k + 2k + 3k + \ldots + nk = k \sum_{i=1}^{n} i$$

Product with constant:

$$\prod_{i=1}^{n} i \times k = 1k \times 2k \dots \times nk = k^{n} \times \prod_{i=1}^{n} i$$

Probability: Axioms

- P measures total probability of a set of events
- ► $P(\emptyset) = 0$
- P(all events) = 1
- $P(X) \leq P(Y)$ for any $X \subseteq Y$
- $P(X) + P(Y) = P(X \cup Y)$ provided that $X \cap Y = \emptyset$
- P(GC drives drunk & GC is in Hawaii) + P(GC drives drunk & GC is not in Hawaii) = P(GC drives drunk)

Probability Axioms

All events sum to 1:

$$\sum_{e} P(e) = 1$$

Marginal probability P(f):

$$P(f) = \sum_{e} P(e, f)$$

Conditional probability:

$$\sum_{e} P(e \mid f) = \sum_{e} \frac{P(e, f)}{P(f)} = \frac{1}{P(f)} \sum_{e} P(e, f) = 1$$

▶ Computing *P*(*f*) from axioms:

$$P(f) = \sum_{e} P(e) \times P(f \mid e)$$

Probability: Bias and Variance

- P(GC drives drunk | GC is in Hawaii, GC is alone, GC is low in polls, ...)
- As we add more material to the right of | :
 - probability could increase or decrease
 - probability usually gets more relevant (less bias)
 - probability usually gets less reliable (more variance)
 - removing items from the right of | makes it easier to get an estimate (more bias but less variance)

Probability: The Chain Rule

- P(GC is in Hawaii,GC is alone,GC is low in polls | GC drives drunk)
- We cannot remove items from the left of | (verify that it violates the definitions we have given based on sets)
- In this case we can use the chain rule of probability to rescue us
- P(GC in Hawaii,GC alone,GC low in polls | GC drives drunk) = P(GC in Hawaii | GC alone,GC low in polls,GC drives drunk) × P(GC alone | GC low in polls, GC drives drunk) × P(GC low in polls | GC drives drunk)

Probability: The Chain Rule

 P(GC in Hawaii,GC alone,GC low in polls | GC drives drunk) = P(GC in Hawaii | GC alone,GC low in polls,GC drives drunk) × P(GC alone | GC low in polls, GC drives drunk) × P(GC low in polls | GC drives drunk)

• Remember:
$$P(X \mid Y) = \frac{P(X,Y)}{P(Y)}$$

►
$$\frac{HALD}{D} = \frac{HALD}{ALD} \times \frac{ALD}{LD} \times \frac{LD}{D}$$

(simply cancel out the matching terms)

Probability: The Chain Rule

$$P(e_1, e_2, \dots, e_n) = P(e_1) \times P(e_2 \mid e_1) \times P(e_3 \mid e_1, e_2) \dots$$

$$P(e_1, e_2, \dots, e_n) = \prod_{i=1}^n P(e_i \mid e_{i-1}, e_{i-2}, \dots, e_1)$$

Probability: Random Variables and Events

- What is y in P(y) ?
- Shorthand for value assigned to a random variable Y, e.g. Y = y
- y is an element of some implicit event space: \mathcal{E}

Probability: Random Variables and Events

The marginal probability P(y) can be computed from P(x, y) as follows:

$$P(y) = \sum_{x \in \mathcal{E}} P(x, y)$$

Finding the value that maximizes the probability value:

$$\hat{x} = rac{rg \max}{x \in \mathcal{E}} P(x)$$

Quick Guide to Probability Theory

Log Probability

Basics of Information Theory

Practical problem with tiny P(e) numbers: underflow

One solution is to use log probabilities:

$$log(P(e)) = log(p_1 \times p_2 \times \ldots \times p_n)$$

= log(p_1) + log(p_2) + \dots + log(p_n)

Note that:

 $x = \exp(\log(x))$

Also more efficient: addition instead of multiplication

р	$\log(p)$		
0.0	$-\infty$		
0.1	-3.32		
0.2	-2.32		
0.3	-1.74		
0.4	-1.32		
0.5	-1.00		
0.6	-0.74		
0.7	-0.51		
0.8	-0.32		
0.9	-0.15		
1.0	0.00		

- ► So: $(0.5 \times 0.5 \times ... 0.5) = (0.5)^n$ might get too small but (-1 1 1 1) = -n is manageable
- Another useful fact when writing code (log₂ is *log to the base 2*):

$$\log_2(x) = \frac{\log_{10}(x)}{\log_{10}(2)}$$

- Adding probabilities is expensive to compute: logadd(x, y) = log(exp(x) + exp(y))
- ▶ A more efficient soln, let *big* be a large constant e.g. 10³⁰:

 $\begin{array}{l} \text{function } logadd(x, y) : \# \text{ returns } log(\exp(x) + \exp(y)) \\ \text{if } (y - x) > log(big) \text{ return } y \\ \text{elsif } (x - y) > log(big) \text{ return } x \\ \text{else } \text{ return } \\ min(x, y) + log(\exp(x - min(x, y)) + \exp(y - min(x, y))) \\ \text{endif} \end{array}$

► There is a more efficient way of computing log(exp(x - min(x, y)) + exp(y - min(x, y)))

function
$$logadd(x, y)$$
:
if $(y - x) > log(big)$ return y
elsif $(x - y) > log(big)$ return x
elsif $(x \ge y)$ return $x + log(1 + exp(y - x))$
note that $max(x, y) = x$ and $y - x \le 0$
else return $y + log(exp(x - y) + 1)$
note that $max(x, y) = y$ and $x - y \le 0$
endif

Also, in ANSI C, log1p efficiently computes log(1 + x)http://www.ling.ohio-state.edu/~jansche/src/logadd.c Quick Guide to Probability Theory

Log Probability

Basics of Information Theory

Information Theory

- Information theory is the use of probability theory to quantify and measure "information".
- Consider the task of efficiently sending a message. Sender Alice wants to send several messages to Receiver Bob. Alice wants to do this as efficiently as possible.
- Let's say that Alice is sending a message where the entire message is just one character a, e.g. aaaa.... In this case we can save space by simply sending the length of the message and the single character.

Information Theory

- Now let's say that Alice is sending a completely random signal to Bob. If it is random then we cannot exploit anything in the message to compress it any further.
- The upper bound on the number of bits it takes to transmit some infinite set of messages is what is called entropy.
- This formulation of entropy by Claude Shannon was adapted from thermodynamics, converting information into a quantity that can be measured.
- Information theory is built around this notion of message compression as a way to evaluate the amount of information.

- Consider a probability distribution p
- Entropy of p is:

$$H(p) = -\sum_{x \in \mathcal{E}} p(x) \log_2 p(x)$$

- Any base can be used for the log, but base 2 means that entropy is measured in bits.
- Entropy answers the question: What is the upper bound on the number of bits needed to transmit messages from event space *E*, where *p*(*x*) defines the probability of observing *x*.

- Alice wants to bet on a horse race. She has to send a message to her bookie Bob to tell him which horse to bet on.
- There are 8 horses. One encoding scheme for the messages is to use a number for each horse. So in bits this would be 001,010,...

(lower bound on message length = 3 bits in this encoding scheme)

Can we do better?

Horse 1	$\frac{1}{2}$	Horse 5	$\frac{1}{64}$
Horse 2	$\frac{1}{4}$	Horse 6	$\frac{1}{64}$
Horse 3	$\frac{1}{8}$	Horse 7	$\frac{1}{64}$
Horse 4	$\frac{1}{16}$	Horse 8	$\frac{1}{64}$

- If we know how likely we are to bet on each horse, say based on the horse's probability of winning, then we can do better.
- Let p be the probability distribution given in the table above. The entropy of p is H(p)

$$\begin{aligned} H(p) &= \\ &= -\sum_{i=1}^{8} p(i) \log_2 p(i) \\ &= -\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{4} \log_2 \frac{1}{4} + \frac{1}{8} \log_2 \frac{1}{8} + \frac{1}{16} \log_2 \frac{1}{16} + 4(\frac{1}{64} \log_2 \frac{1}{64})\right) \\ &= -\left(\frac{1}{2} \times -1 + \frac{1}{4} \times -2 + \frac{1}{8} \times -3 + \frac{1}{16} \times -4 + 4(\frac{1}{64} \times -6)\right) \\ &= -\left(-\frac{1}{2} - \frac{1}{2} - \frac{3}{8} - \frac{1}{4} - \frac{3}{8}\right) \\ &= 2 \text{ bits} \end{aligned}$$

▶ What is the entropy when the horses are equally likely to win? $H(uniform \ distribution) = -8(\frac{1}{8} \times -3) = 3 \ bits$

e.g., most likely horse gets code 0, next most likely gets 10, and then 110, 1110, ...

many possible coding schemes, this is a simple code to illustrate number of bits needed for a large number of messages ...

- Assume there are 320 messages (one for each race): code 0 occurs 160 times, code 10 occurs 80 times, code 110 occurs 40 times, code 1110 occurs 20 times, code 11110 occurs 5 times.
- ► Total number of bits for all messages: 160*len(0) + 80*len(10) + 40*len(110) + 20*len(1110) + 5*len(11110)
- ▶ Number of bits: 160*1 + 80*2 + 40*3 + 20*4 + 5*5 = 545
- ► Total number of bits per message (per race): ⁵⁴⁵/₃₂₀ ≈ 1.7 bits (always less than 2 bits)

Perplexity

- The value $2^{H(p)}$ is called the **perplexity** of a distribution p
- Perplexity is the weighted average number of choices a random variable has to make.
- Choosing between 8 equally likely horses (H=3) is $2^3 = 8$.
- Choosing between the biased horses from before (H=2) is $2^2 = 4$.

Relative Entropy

- In real life, we cannot know for sure the exact winning probability for each horse.
- Let's say pt is the true probability and pe is our estimate of the true probability (say we got pe by observing previous races with these horses)
- ► We define the *distance* between p_t and p_e as the **relative** entropy: written as D(p_t || p_e)

$$D(p_t \| p_e) = -\sum_{x \in \mathcal{E}} p_t(x) \log_2 \frac{p_e(x)}{p_t(x)}$$

The relative entropy is also called the Kullback-Leibler divergence.

Cross Entropy and Relative Entropy

The relative entropy can be written as the sum of two terms:

$$D(p_t || p_e) = -\sum_{x \in \mathcal{E}} p_t(x) \log_2 \frac{p_e(x)}{p_t(x)} \\ = -\sum_x p_t(x) \log_2 p_e(x) - \sum_x p_t(x) \log_2 p_t(x)$$

We know that H(pt) = -∑x pt(x) log₂ pt(x)
 Let us define Hpt(pe) = -∑x pe(x) log₂ pt(x)

$$D(p_t || p_e) = H_{p_t}(p_e) + H(p_t)$$

• The term $H_{p_t}(p_e)$ is called the **cross entropy**.

Cross Entropy and Relative Entropy

The relative entropy between p_e and p_t can be written as the sum of two terms:

 $\begin{array}{l} \textbf{relative entropy}(p_t, p_e) = \textbf{cross entropy}(p_t, p_e) + \textbf{entropy}(p_t) \\ D(p_t \| p_e) = H_{p_t}(p_e) + H(p_t) \end{array}$

•
$$H_{p_t}(p_e) \ge H(p_t)$$
 always.

- $D(p_t \| p_e) \ge 0$ always, and $D(p_t \| p_e) = 0$ iff $p_t = p_e$
- ▶ D(p_t || p_e) is not a true distance:
 - It is asymmetric: $D(p_t || p_e) \neq D(p_e || p_t)$,
 - ► It does not obey the triangle inequality: $D(p||r) \leq D(p||q) + D(q||r)$

Conditional Entropy and Mutual Information

• *Entropy* of a random variable X:

$$H(X) = -\sum_{x \in \mathcal{E}} p(x) \log_2 p(x)$$

Conditional Entropy between two random variables X and Y:

$$H(X \mid Y) = -\sum_{x,y \in \mathcal{E}} p(x,y) \log_2 p(x \mid y)$$

Mutual Information between two random variables X and Y:

$$I(X;Y) = D(p(x,y) || p(x)p(y)) = \sum_{x} \sum_{y} p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)}$$