

18 mins gone

Apparently 15 mins spent
by Sir in writing

If f is a function & p is a probability distribution
for a r.v. X

$$E(f) = \sum f(x) p(x) \\ = \int f(x) p(x) dx$$

$$H(X) = E(-\log(p(X)))$$

$$H(Y/X) = E(-\log(p(Y/X)))$$

$$= \iint p(x,y) (\log p(Y/X)) dy dx$$

(just note this for now)

$$f(x_1) \quad f(x_2) \quad \dots \quad f(x_n)$$

$$E(f) \approx \frac{1}{N} \sum_i f(x_i)$$

Let $x \sim N(\mu, \sigma^2)$

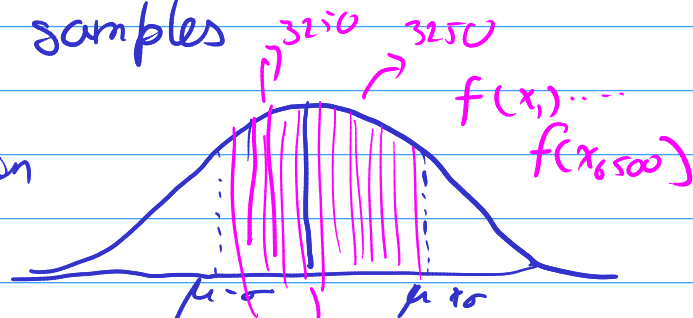
Collect 10,000 samples

When f is a function

$$E(f) = \int f(x) p(x) dx$$

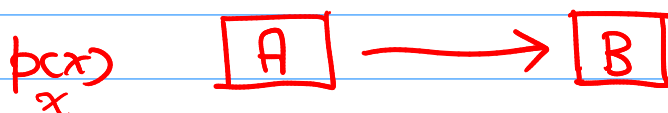
if you don't know $p(x)$, how do we approximate $E(f)$ with $\frac{1}{N} \sum_{i=1}^N f(x_i)$

This is how:



Averaging out, we can get $f(\mu)$
LMAO

$$-\int p(x) \log(q(x)) dx = (-\int p(x) \log(p(x)) dx)$$



a b c d e f g h
000 001 010 011 100 101 110 111



Let all of them be in a basket

Let them all follow a distribution p .

if the 8 chars are uniform distribution

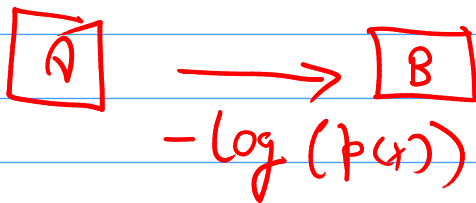
$$\text{Avg code length} = \frac{3(8)(1000)}{8000} = 3$$

What's the encoding mechanism for least number of bits: Ans: $-\log(p(x))$

Shannon's:

$$\sum_x -p(x) \log(p(x)) = \text{Entropy}$$

↙ min no. of bits



$q(x) : -\log(q(x))$

$-\int p(x) \log(q(x)) dx$

Average information for $p(x)$ $-\sum p(x) \log(p(x))$
 if $-\log(p(x)) = f(x)$

\downarrow
 Expectation of code length
 then

$= -\int p(x) f(x) dx = E(f)$

if $p(x)$ dictates the universe
 and we encode in $q(x)$
 then avg code length

code length: $-\log(q(x))$
 probability
 should be dictated
 by universe : $p(x)$

$\Rightarrow -\int p(x) \log(q(x)) dx \rightarrow \text{more}$

vs $-\int p(x) \log(p(x)) dx \rightarrow \text{lessen (by Shannon's encoding theorem)}$

increase of code length:

$(-\int p(x) \log(q(x)) dx) - (-\int p(x) \log(p(x)) dx)$

$$= - \int p(x) \log \left(\frac{q(x)}{p(x)} \right) dx \Rightarrow \text{This is the KL divergence of } p \text{ and } q$$

$KL(p||q)$

We don't have $p(x)$

$q(x|\theta)$, and we need to find the probability distributions for this!

so, $\min_{\theta} KL(p(x)||q(x|\theta))$

HOW TO FIND KL divergence without p ?
(That's for another lecture)

$$E(f) \simeq \frac{1}{N} \sum_{n=1}^N f(x_n) \quad \text{JK it is @ bottom}$$

From here we can estimate

$$E(f) \simeq - \int f(x) p(x) dx \simeq - \frac{1}{N} \sum_{n=1}^N f(x_n)$$

$\min_{\theta} (KL \text{ divergence}) =$

$$\min_{\theta} \quad - \frac{1}{N} \sum_{n=1}^N \underbrace{f(x_n)}_{\log(q(x_n|\theta))} + \frac{1}{N} \sum_{n=1}^N \underbrace{g(x_n)}_{\log(p(x_n))}$$

$$= \min_{\theta} \quad - \frac{1}{N} \sum_{n=1}^N \log q(x_n|\theta) + \frac{1}{N} (\log p(x_n))$$

$$= \min_{\theta} \quad \frac{1}{N} \sum_{n=1}^N (-\log q(x_n|\theta) + \log(p(x_n)))$$

Now if we have θ_1 & θ_2 , we can entirely disregard $\log(p(x_n))$

to find out which is better, since $p(x_n)$ is independent of θ , (But we need to sample points from p)

$$\Rightarrow \min_{\theta} \frac{1}{N} \sum_{n=1}^N (-\log(q(x_n/\theta)))$$

likelihood lol!

$$\min_{\theta} -\frac{1}{N} (\log q(x_1/\theta) + \log(q(x_2/\theta)) + \dots + \log(q(x_n/\theta)))$$

$\min_{\theta} (\text{negative log likelihood})!$

\Rightarrow KL divergence is small

we see another face of log likelihood.

We however didn't actually use KL divergence from the START of the course, it was rather intuitive.

\Rightarrow Average code length is closest to minimum code length

$$p(y/x)$$

given

x , find average code length to specify y

will be a function

$$\int -p(y/x) \log(p(y/x)) dy$$

weighted average

now

$$H(y/x) = \int \left(\int -p(y/x) \log(p(y/x)) dy \right) p(x) dx$$

$f(x)$

\rightarrow expected value of average length

$$= \int \int p(y/x) (-\log(p(y/x))) dy dx$$

$$H(y/x) = - \iint p(x,y) \log(p(y/x)) dy dx$$

$$H(x/y) = - \iint p(x,y) \log(p(x/y)) dx dy$$

$$p(x,y) = p(x)p(y) \Rightarrow x, y \text{ are independent}$$

$$H_x(x,y) = - \int p(x,y) \log(p(x,y)) dx$$

↪ for a given x

$$H(x,y) = - \iint p(x,y) \log(p(x,y)) dy dx$$

$$= - \iint p(x,y) \log(p(y/x)) dx dy$$

$$- \iint p(x,y) \log(p(x)) dx dy$$

$$= H(y/x) - \iint p(x,y) \log(p(x)) dx dy$$

$$= H(y/x) - \int \log(p(x)) \left(\int p(x,y) dy \right) dx$$

$$= H(y/x) - \int p(x) \log(p(x)) dx$$

$$\boxed{H(x,y) = H(y/x) + H(x)} \rightarrow \text{Conditional Entropy}$$

Independence

(We can throw away a feature similar to another)

$p(x,y)$ $p(x)p(y)$ [assumed marginalized]
is the KL divergence less

or a feature that doesn't contribute

To find the KL divergence in simple terms
as a measure of entropy

$$KL(p(x,y) \parallel p(x)p(y))$$

$$= - \iint p(x,y) \log \left(\frac{p(x)p(y)}{p(x,y)} \right) dx dy$$

expand

$$= - \iint p(x,y) \left[\log(p(x)) + \log(p(y)) - \log(p(x)p(y/x)) \right] dx dy$$

$$= - \iint p(x,y) \left[\log(p(x)) + \log(p(y)) - \log(p(x)) - \log(p(y/x)) \right] dx dy$$

$$= - \iint p(x,y) \left[\log(p(y)) - \log(p(y/x)) \right] dx dy$$

$$= - \iint p(x,y) \left[\log(p(y)) \right] dx dy + \iint p(x,y) \log(p(y/x)) dx dy$$

$$= - \int \log(p(y)) \left(\int p(x,y) dx \right) dy - H(y/x)$$

$$= - \int \log(p(y)) p(y) dy - H(y/x)$$

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

If we know the entropies of these
we know the KL divergence

The Independent mutual information is hence given by

$$\boxed{I(x,y) = H(y) - H(y/x) = H(x) - H(x/y)} \rightarrow \text{Similar elimination}$$

if x & y are really independent

$$H(y|x) = H(y) \quad I(x, y) = 0$$

$I(x, y)$ is the reduction of entropy of y , when x is involved

if y is dependent on x , then $I(x, y)$ can have many values, and this can give insights.

After this we look at a simple algorithm for classification & feature engineering, data preprocessing etc.