Eric Zacharia

1a. The expected value of the Normal distribution is its mean, $\mu$, and its variance is its stdev squared, $\sigma^2$.

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

Entropy, $H = E\left(\lg_2 \frac{1}{P(x)}\right) = -\sum_x P(x) \lg_2 P(x)$

for a **discrete** r.v.

Plug

Entropy, $H(X) = -\int_{-\infty}^{\infty} f(x) \ln f(x)\, dx$

for a **continuous** r.v.

$$= -\int_{-\infty}^{\infty} f(x) \left(-\frac{1}{2}\ln \sigma\sqrt{2\pi} - \frac{(x-\mu)^2}{2\sigma^2}\right) dx$$

$$= \int_{-\infty}^{\infty} f(x)\,\frac{1}{2}\ln \sigma\sqrt{2\pi} + f(x)\,\frac{(x-\mu)^2}{2\sigma^2}\, dx$$

$$= \frac{1}{2}\ln \sigma\sqrt{2\pi} \int_{-\infty}^{\infty} f(x)\, dx + \frac{1}{2\sigma^2} \int_{-\infty}^{\infty} f(x)(x-\mu)^2\, dx$$

$$= \frac{1}{2}\ln \sigma\sqrt{2\pi} + \frac{1}{2} = 1.25\ln \sigma + \frac{1}{2} = \underline{\ln \sigma + C}$$

1b. Cross entropy loss for the test sequence of words:

$$-\frac{1}{n} \sum_{i=k}^{n} \lg_2 \hat{P}(w_i | c_i)$$

Perplexity of a model (w.r.t. corpus)

$$PP(W) = P(w_1 w_2 \ldots w_n)^{-\frac{1}{n}} = \sqrt[n]{\frac{1}{P(w_1 w_2 \ldots w_n)}}$$

by the chain rule,

$$\sqrt[n]{\prod_{i=1}^{n} \frac{1}{P(w_i | w_1 w_2 \ldots w_{i-1})}}$$

w.r.t. the test sequence, this becomes

$$\sqrt[n]{\prod_{i=k}^{n} \frac{1}{P(w_i | c_i)}} \quad , \text{where } c_i = w_{i-k+1:i-1} \quad \left( (k-1)\text{-gram context} \right)$$

Cross entropy loss is derived from the log perplexity.

$$\log PP(W) = \log \left( \sqrt[n]{\prod_{i=k}^{n} \frac{1}{P(w_i | c_i)}} \right)$$

$$= -\frac{1}{n} \sum_{i=k}^{n} \lg_2 P(w_i | c_i)$$

1c. $p(cw) =$ "The probability of the K-gram corresponding to <u>context of size K-1</u> followed by a <u>word</u>."
$\phantom{1c. p(cw) =}$ $c$
$\phantom{1c. p(cw) =}$ $w$

$p(w|c) =$ "The conditional probability of seeing a <u>word</u> after seeing a <u>context of size K-1</u>."
$\phantom{p(w|c) =}$ $w$ $\phantom{after seeing a}$ $c$

The entropy of English based on a K-gram language model:

$$H_K = - \sum_{c,w} p(cw) \lg_2 p(w|c)$$

Random process in which each word-context pair $(w, c)$ is generated independently with prob. $p(cw)$:

$$G_K = - \sum_{c,w} p(cw) \lg_2 p(cw)$$

To say "In general, $H_K < G_K$" implies

$\Rightarrow$ "$-\sum_{c,w} p(cw) \lg_2 p(w|c) < -\sum_{c,w} p(cw) \lg_2 p(cw)$"

$>$ flip

$\Rightarrow$ implies "$p(w|c) > p(cw)$"

<u>implies that seeing a word after its context is</u> <u>more probable than a K-gram that corresponds</u> <u>to a context followed by a word</u>, and that makes sense.

1c.

$$H_K = -\sum_{c,w} P(cw) \, lg_2 P(w/c) = -\sum_{c,w} P(cw) \, lg_2 \frac{P(cw)}{P(c)} = -P(c) \, lg_2 \frac{P(c)}{P(c)} = 0$$

$$G_K - G_{K-1} = -\sum_{c,w} P(cw) \, lg_2 P(w_{i-k+1:i-1} w_i) - \left( -\sum_{c,w} P(cw) \, lg_2 P(w_{i-k+2:i-1} w_i) \right)$$

$$= -P(c) \, lg_2 P(w_{i-k+1:i-1}) + P(c) \, lg_2 P(w_{i-k+2:i-1}) = 0$$

$$\therefore H_K = G_K - G_{K-1}$$