

Question 1:

We start with

$$P(A|B) = \frac{P(A) P(B|A)}{P(B)}$$

$$P(B) P(A|B) = P(A \text{ and } B) = P(A) P(B|A)$$

We can make the two equal to each other:

$$P(B) P(A|B) = P(A) P(B|A)$$

thus

$$P(A|B) = \frac{P(A) P(B|A)}{P(B)}$$

Bayes theorem is important for machine learning because it helps us make predictions and decisions based on available data while considering uncertainty. It allows us to start with an initial guess and update our prediction as we get new data. You can update the prediction using Bayes theorem. You can essentially use bayes theorem to build models that learn from data and improve predictions over time. These models can not only make predictions but also tell you how uncertain those predictions are.

Question 2:

$$E(w) = ((w^T X - y)^T (w^T X - y)) + \lambda w^T w$$

$$= (((Xw)^T - y^T) (Xw - y)) + \lambda w^T w$$

$$= (w^T X^T X w - y^T X w - w^T X^T y - y^T y) + \lambda w^T w$$

We take the derivative with respect to w:

$$\frac{d}{dw} ((w^T X^T X w - y^T X w - w^T X^T y - y^T y) + \lambda w^T w)$$

The derivative of $w^T X^T X w$ with respect to w = $2X^T X w$

The derivative of $-y^T X w$ with respect to w is $-X^T y$

The derivative of $-w^T X^T y$ with respect to w is $-X^T y$

The derivative of $-y^T y$ with respect to w is 0 because it doesn't depend on w

The derivative of $\lambda w^T w$ is $2\lambda w$

Summing all that up, we get

$$2X^T X w - X^T y - X^T y + 2\lambda w$$

After combining like terms, we get:

$$2X^T X w - 2X^T y + 2\lambda w = 0$$

Now we can factor out 2

$$2(X^T X w - X^T y + \lambda w) = 0$$

Now we can divide both sides by 2

To get

$$X^T X w - X^T y + \lambda w = 0$$

Moving terms involving w to the left

$$X^T X w + \lambda w - X^T y = -0$$

We can factor out w

$$(X^T X + \lambda I) w = X^T y$$

We can multiply both sides by the inverse of the matrix to get:

$$w = (X^T X + \lambda I)^{-1} X^T y$$

Question 3:

We are trying to derive the gradient of $J\theta$ with respect to θ_k .

Using the multivariable chain rule we can deduce that

$$\frac{\partial J\theta}{\partial \theta_k} = \frac{\partial J\theta}{\partial p_{ki}} * \frac{\partial \hat{p}_{ki}}{\partial sk(xi)} * \frac{\partial sk(xi)}{\partial \theta_k}$$

$$\frac{\partial J\theta}{\partial \theta} = \sum_{i=1}^m (\hat{p}_{ki} - y_{ki}) x(i)$$

First we look at $\frac{\partial J\theta}{\partial \hat{p}_{ki}}$

$$\frac{\partial J\theta}{\partial \hat{p}_{ki}} = \frac{\partial}{\partial \hat{p}_{ki}} - \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_{ki} \log(\hat{p}_{ki})$$

$$= - \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \frac{y_{ki}}{\hat{p}_{ki}}$$

Next we look at $\frac{\partial \hat{p}_{ki}}{\partial sk(xi)}$

Here, we're calculating the derivative of the softmax with respect to its inputs where m is the number of inputs.

$$\frac{\partial \hat{p}_{ki}}{\partial sk(xm)} = \frac{\partial}{\partial sk(xm)} \frac{\exp(sk(xi))}{\sum_{j=1}^k \exp(sj(x))}$$

We can use the quotient rule which states that:

$$\frac{dy}{dx} = \frac{v * \frac{du}{dx} - u * \frac{dv}{dx}}{v^2}$$

$$f(x) = \exp(sk(xi)) \text{ and } g(x) = \sum_{j=1}^k \exp(sj(x))$$

$$\frac{\sum_{j=1}^k \exp(sj(x)) * \frac{\partial}{\partial sk} \exp(sk(xi)) - \exp(sk(xi)) * \frac{\partial}{\partial sk} \sum_{j=1}^k \exp(sj(x))}{(\sum_{j=1}^k \exp(sj(x)))^2}$$

The derivative of

$$\sum_{j=1}^k \exp(sj(x)) = \exp(sk(xm))$$

$$= \frac{\sum_{j=1}^k \exp(sj(x)) * \exp(sk(xi)) - \exp(sk(xi)) * \exp(sj(xm))}{(\sum_{j=1}^k \exp(sj(x)))^2}$$

* In the above, we can get rid of the sum because all other terms other than sk are 0s. Only sk is 1.

The result can be factored

So we get

$$\begin{aligned}
 &= \frac{\exp(sk(xi)) \left(\sum_{j=1}^k \exp(sj(x)) - \exp(sk(xm)) \right)}{\left(\sum_{j=1}^k \exp(sj(x)) \right)^2} \\
 &= \frac{\exp(sk(xi))}{\sum_{j=1}^k \exp(sj(x))} * \frac{\sum_{j=1}^k \exp(sj(x) - \exp(sk(xm)))}{\sum_{j=1}^k \exp(sj(x))} \\
 &= \frac{\exp(sk(xi))}{\sum_{j=1}^k \exp(sj(x))} * \left(1 - \frac{\exp(sk(xm))}{\sum_{j=1}^k \exp(sj(x))} \right) \\
 &= \hat{p}_{ki} * (1 - \hat{p}_{km})
 \end{aligned}$$

Since i = m

The above equals $\hat{p}_{ki} * (1 - \hat{p}_{ki})$

When i is not equal to m:

$$\begin{aligned}
 \frac{\partial \hat{p}_{ki}}{\partial sk(xm)} &= \frac{\partial}{\partial sk(xm)} \frac{\exp(sk(xi))}{\sum_{j=1}^k \exp(sj(x))} \\
 &= \frac{\sum_{j=1}^k \exp(sj(x)) * \frac{\partial}{\partial sk} \exp(sk(xi)) - \exp(sk(xi)) * \frac{\partial}{\partial sk} \sum_{j=1}^k \exp(sj(x))}{\left(\sum_{j=1}^k \exp(sj(x)) \right)^2}
 \end{aligned}$$

$$\frac{\partial}{\partial sk(xi)} \exp(sk(xi)) = 0 \text{ now since it doesn't depend on } sj(xm)$$

$$= \frac{\sum_{j=1}^k \exp(sj(x)) * \frac{\partial}{\partial sk} \exp(sk(xi)) - \exp(sk(xi)) * \frac{\partial}{\partial sk} \sum_{j=1}^k \exp(sj(x))}{(\sum_{j=1}^k \exp(sj(x)))^2}$$

$$= \frac{0 - \exp(sk(xi)) * \frac{\partial}{\partial sk} \sum_{j=1}^k \exp(sj(x))}{(\sum_{j=1}^k \exp(sj(x)))^2}$$

$$\frac{\partial}{\partial sk} \sum_{j=1}^k \exp(sj(x)) = \exp(sk(xm))$$

So

$$= \frac{-\exp(sk(xi))}{\sum_{j=1}^k \exp(sj(x))} * \frac{\exp(sk(xm))}{\sum_{j=1}^k \exp(sj(x))}$$

$$= -\hat{p}_{ki} * \hat{p}_{km}$$

Lastly

$$\frac{\partial}{\partial \theta k^T x^i} = \theta k^T x^i = x^i$$

If we put it all together, we get:

$$\begin{aligned} \frac{\partial J\theta}{\partial \theta k} &= \frac{1}{m} \sum_{i=1}^m -\frac{y_{ki}}{\hat{p}_{ki}} * \frac{\partial \hat{p}_{ki}}{\partial sk(xm)} \\ &= \frac{1}{m} \sum_{i=1}^m (-\frac{y_{ki}}{\hat{p}_{ki}} * \hat{p}_{ki} (1 - \hat{p}_{ki}) + \sum_{m \neq i}^m -y_{km} * \frac{1}{\hat{p}_{km}} (-\hat{p}_{km} * \hat{p}_{ki})) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{m} \sum_{i=1}^m (-y_{ki} (1 - \hat{p}_{ki}) + \sum_{m \neq i}^m y_{km} * \hat{p}_{ki}) \\
&= \frac{1}{m} \sum_{i=1}^m (-y_{ki} + (y_{ki} * p_{ki})) + \sum_{m \neq i}^m y_{km} * \hat{p}_{ki} \\
&= \frac{1}{m} \sum_{i=1}^m (-y_{ki} + \sum_m y_{km} * \hat{p}_{ki}) \\
&= \frac{1}{m} \sum_{i=1}^m (\hat{p}_{ki} \sum_m y_{km} - y_{ki}
\end{aligned}$$

In the above, the summation of y_{km} is equal to 1 so we get:

$$\frac{1}{m} \sum_{i=1}^m (\hat{p}_{ki} - y_{ki})$$

Lastly we must multiply this by $\frac{\partial}{\partial \theta k^T x^i} = \theta k^T x^i = x^i$

$$\text{To get } \frac{1}{m} \sum_{i=1}^m (\hat{p}_{ki} - y_{ki}) * x^i$$