1. [6 points] Prove Bayes' Theorem. Briefly explain why it is useful for machine learning problems.

Bayes' Theorem is the probability of two events, A and B, occurring when there is nonzero probability of occurrence. It is useful for machine learning because as new data comes in, it will update probabilities to obtain more accurate results.

$$\frac{P(AB) = P(BA) \cdot P(A)}{P(B)}$$

Proof by Derivation

The rule for conditional probability is,

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

This can be written reversally as,

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

The numerators of both equations P(A|B) and P(B|A) are the same because the probability of A and B occurring is the same as the probability of B and A occurring. Therefore, we can rewrite our reverse equation as,

Solve for $P(A \cap B)$

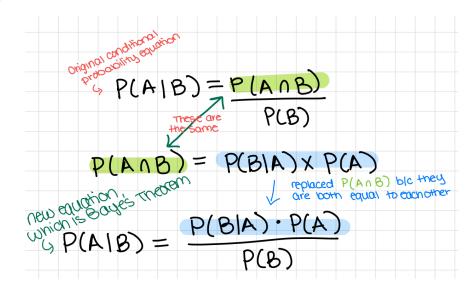
$$P(B|A) \times P(A) = \frac{P(A \cap B)}{P(A)} \times P(A)$$

$$P(A \cap B) = P(B|A) \times P(A)$$

Referring back to our rule for conditional probability

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

We see that $P(A \cap B)$ in the numerator can now be replaced by our solution $P(B|A) \times P(A)$. This will give,

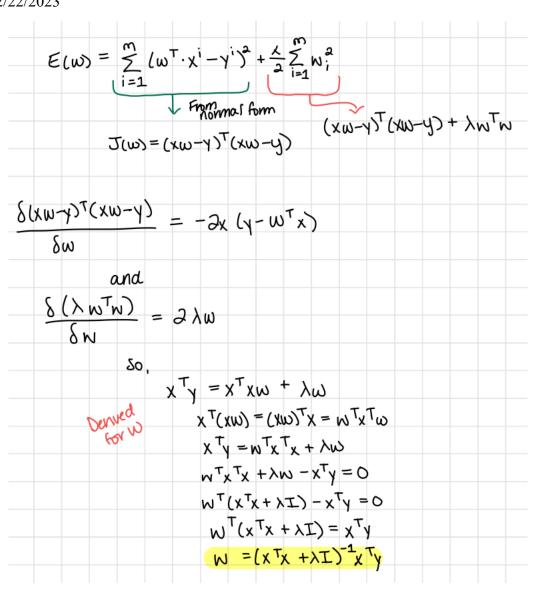


As shown above, the derivation gives us Bayes' Theorem,

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

$$P(B)$$

- 2. [10 points] In Module 2, we gave the normal equation (i.e., closed-form solution) for linear regression using MSE as the cost function. **Prove that the closed-form solution for Ridge Regression** is $\mathbf{w} = (\lambda I + X^T \cdot X)^{-1} \cdot X^T \cdot \mathbf{y}$, where I is the identity matrix, $X = (x^{(1)}, x^{(2)}, \dots, x^{(m)})^T$ is the input data matrix, $x^{(i)} = (1, x_1, x_2, \dots, x_n)$ is the i-th data sample, and $\mathbf{y} = (y^{(1)}, y^{(2)}, \dots, y^m)$. Assume the hypothesis function $h_w(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$, and $y^{(j)}$ is the measurement of
 - $h_w(x)$ for the *j*-th training sample. The cost function of the Ridge Regression is $E(\mathbf{w}) = \sum_{i=1}^m (\mathbf{w}^T \cdot \mathbf{x}^{(i)} \mathbf{y}^{(i)})^2 + \lambda \sum_{i=1}^m w_{i2}$.



- 3. [10 points] Assume we have K different classes in a multi-class Softmax Regression model. The posterior probability is $\hat{p}_k = \frac{\exp{(s_k(x))}}{\sum_{j=1}^K \exp{(s_j(x))}}$ for k = 1, 2, ..., K, where $s_k(x) = \theta_k^T \cdot x$, input x is an n-dimension vector, and K the total number of classes.
 - 1) To learn this Softmax Regression model, how many parameters we need to estimate? What are these parameters?

n by k, These parameters would be

2) Consider the cross-entropy cost function $J(\Theta)$ of m training samples $\{(x_i, y_i)\}_{i=1,2,...,m}$ as below. Derive the gradient of $J(\Theta)$ regarding to θ_k .

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log (p_k^{(i)})$$

where $y_k(i) = 1$ if the ith instance belongs to class k; 0 otherwise.

$$\nabla J(\theta) = \sqrt{m} \sum_{i=1}^{m} \left(\widehat{\rho}_{k}^{(i)} - y_{k}^{(i)} \right) x^{(i)}$$

$$\widehat{\rho}_{k} = \delta(S_{k}(x))_{k} = \frac{e^{x} \left(\Theta_{k}^{T} \cdot x^{(i)} \right)}{\sum_{j=1}^{k} e^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right)}$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{k} y_{k}^{(i)} \log \left(\frac{e^{x} \Theta_{k}^{T} \cdot x^{(i)}}{\sum_{j=1}^{k} e^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right)} \right)$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{k} y_{k}^{(i)} \left[\log \left(e^{x} \Theta_{j}^{T} \cdot x^{(i)} \right) \right] - \log \left(\sum_{j=1}^{k} e^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} \sum_{k=1}^{k} y_{k}^{(i)} \log_{j} \left(c^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right) \right) + \sum_{j=1}^{m} \sum_{k=1}^{k} y_{k}^{(j)} \log_{j} \left(\sum_{j=1}^{k} e^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right) \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} \sum_{k=1}^{k} y_{k}^{(i)} \cdot \Theta_{j}^{T} \cdot x^{(i)} \right) + \sum_{j=1}^{m} \sum_{k=1}^{k} y_{k}^{(j)} \log_{j} \left(\sum_{j=1}^{k} e^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right) \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} \sum_{k=1}^{k} y_{k}^{(i)} \cdot \Theta_{j}^{T} \cdot x^{(i)} \right) + \sum_{j=1}^{m} \sum_{k=1}^{k} y_{k}^{(i)} \log_{j} \left(\sum_{j=1}^{k} e^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right) \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} \sum_{k=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{m} \sum_{k=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{k=1}^{k} e^{x} \left(\Theta_{j}^{T} \cdot x^{(i)} \right) \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} \sum_{k=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{m} \sum_{k=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{k=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{k=1}^{m} y_{k}^{(i)} \right) \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} \sum_{k=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{m} \sum_{k=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{k=1}^{m} y_{k}^{(i)} \right) \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{m} \sum_{k=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{j=1}^{m} y_{k}^{(i)} \right) \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{j=1}^{m} y_{k}^{(i)} \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{j=1}^{m} y_{k}^{(i)} \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{m} y_{k}^{(i)} \log_{j} \left(\sum_{j=1}^{m} y_{k}^{(i)} \right) \right)$$

$$= -\frac{1}{m} \left(\sum_{j=1}^{m} y_{k}^{(i)} \cdot x_{k}^{(i)} + \sum_{j=1}^{$$