

CpE 520: HW#13

West Virginia University

Ali Zafari

Question 1

Part (a): Inverse of $\mathbf{y}_{n \times 1} = \mathbf{A}_{n \times m} \mathbf{x}_{m \times 1}$ where $n > m$ and $\text{rank} = m$

In this case there are more unknowns than equations so the system is underdetermined. Because the matrix A is full rank with $\text{rank} = n$, AA^T will be an $n \times n$ non-singular square matrix. Therefore, there will be an inverse for it like below:

$$(AA^T)(AA^T)^{-1} = I \quad (*)$$

As it is clear from the above equation, we may use $A^T(AA^T)^{-1}$ as a right inverse for matrix A . We first assume that the inverse we are seeking for has an equation like $\mathbf{x}_{m \times 1} = \mathbf{R}_{m \times n} \mathbf{y}_{n \times 1}$. Finding $R_{m \times n}$ will be the goal. We propose $R_{m \times n} = A^T(AA^T)^{-1}$ and through the equations below we will check if it will result in our main equation:

$$\boxed{\mathbf{x} = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{y}}$$

$$\mathbf{x} = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{y}$$

$$\xrightarrow{\times A} \mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{y}$$

$$\xrightarrow{\text{using } *} \mathbf{A}\mathbf{x} = \mathbf{y}$$

Part (b): Inverse of $\mathbf{y}_{n \times 1} = \mathbf{A}_{n \times m} \mathbf{x}_{m \times 1}$ where $n < m$ and $\text{rank} = m$

In this case there are more equations than unknowns, so the system is overdetermined. Because the matrix A is full rank with $\text{rank} = m$, $A^T A$ will be an $m \times m$ non-singular square matrix. Therefore, there will be an inverse for it like below:

$$(A^T A)^{-1}(A^T A) = I \quad (**)$$

As it is clear from the above equation, we may use $(A^T A)^{-1} A^T$ as a left inverse for matrix A .

$$\mathbf{A}\mathbf{x} = \mathbf{y}$$

$$\xrightarrow{\times A^T} \mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{y}$$

$$\xrightarrow[\text{using } *]{\times (A^T A)^{-1}} \mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

$$\boxed{\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}}$$

Question 2

Part (a): Difference Between Ridge and LASSO Regressions

Ridge regression uses the the sum of squared values of unknown parameters; however, LASSO regression uses sum of the absolute values of unknown parameters. Therefore, LASSO is capable to make some of the unknown coefficients zero. On the other hand, Ridge regression is able to make the coefficient really small but not absolutely zero.

Part (b): Expressions for Ridge and LASSO Regressions

Ridge Regression:

$$\tilde{\mathbf{x}} = \min_{\mathbf{x}} \{(\mathbf{y} - \mathbf{Ax})^2 + \lambda \sum_{i=1}^m x_i^2\}$$

LASSO Regression:

$$\tilde{\mathbf{x}} = \min_{\mathbf{x}} \{(\mathbf{y} - \mathbf{Ax})^2 + \lambda \sum_{i=1}^m |x_i|\}$$

Question 3

Convolution of two vectors:

$$\begin{bmatrix} 1 \\ 2 \\ 2 \\ 1 \end{bmatrix} * \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1+2 \\ 1+2+2 \\ 2+2+1 \\ 2+1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \\ 5 \\ 5 \\ 3 \\ 1 \end{bmatrix}$$

Question 4

By using 1x1 convolution in DNNs we can change the depth of channels of a layer so there will be much less computations for the DNN to be trained. We can mention an example here to clarify its purpose:

Example 1: input (256 channels) -> 1x1 convolution (64 channels) -> 4x4 convolution (256 channels)

Example 2: input (256 channels) -> 4x4 convolution (256 channels)

If we accept that the computation cost is proportional to the number of weights, we can find it for both the examples above:

number of weights for example 1:

$$(256 \times 64 \times 1 \times 1) + (64 \times 256 \times 4 \times 4) = 278,528$$

$$256 \times 256 \times 4 \times 4 = 1,048,576$$

Question 5

- **Average Pooling**

It computes the average of the values in its filter.

- **Global Average Pooling**

This layer will compute the average over any channel of the input individually.

- **Max Pooling**

It chooses a single maximum value in its filter.

Question 6

Softmax Regression is actually a generalization of logistic regression that could be used for multi-class classification, in contrast of what we have in logistic regression which classifies between two classes.

Question 7

Hinge loss expression:

$$L(d, y) = \max(0, 1 - dy)$$

in which d is the desired response and y is the computed output by the SVM.

This loss function is mainly used in support vector machines (SVM).

Question 8

L1 norm:

$$||\mathbf{x}||_1 = \sum_{i=1} |x_i|$$

L2 norm:

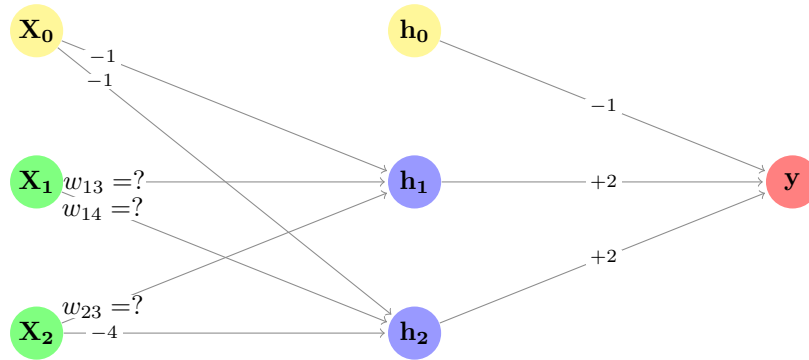
$$||\mathbf{x}||_2 = \sqrt{\sum_{i=1} |x_i|^2}$$

Question 9

Two main applications of autoencoders are **data denoising** and **dimensionality reduction** for data visualization. An advantage of autoencoders is that they do not require any engineering to reduce the dimensionality. So, by only providing enough amount of data they could be trained automatically.

Question 10

There are three unknown weights to be found on the below neural network, so it can perform XOR operation.



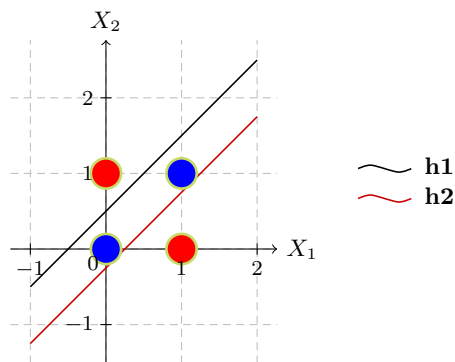
The equation for neurons of the network are written below (φ is the threshold function):

$$h_1 = \varphi(w_{13}X_1 + w_{23}X_2 - 1)$$

$$h_2 = \varphi(w_{14}X_1 - 4X_2 - 1)$$

$$y = \varphi(2h_1 + 2h_2 - 1)$$

With calculating the equations of the lines in the X_1 - X_2 coordinate system, all of the weights will be found.



$$-2X_1 + 2X_2 - 1 = 0 \rightarrow w_{13} = -2, w_{23} = 2$$

$$4X_1 - 4X_2 - 1 = 0 \rightarrow w_{14} = 4$$

Question 11

Derivative of sigmoid function:

$$\begin{aligned}
 \frac{\partial}{\partial x} f(x) &= \frac{\partial}{\partial x} \left(\frac{1}{1 + e^{-\lambda x}} \right) \\
 &= \frac{\lambda e^{-\lambda x}}{(1 + e^{-\lambda x})^2} \\
 &= \lambda \left(\frac{e^{-\lambda x}}{1 + e^{-\lambda x}} \right) \left(\frac{1}{1 + e^{-\lambda x}} \right) \\
 &= \lambda \left(\frac{1 + e^{-\lambda x} - 1}{1 + e^{-\lambda x}} \right) \left(\frac{1}{1 + e^{-\lambda x}} \right) \\
 &= \lambda \left(1 - \frac{1}{1 + e^{-\lambda x}} \right) \left(\frac{1}{1 + e^{-\lambda x}} \right) \\
 &= \lambda f(x)(1 - f(x))
 \end{aligned}$$

Question 12

Kernel function is a function which with the help of dot product could do the same operation as mapping from input space to a high dimensional output space.

Kernel Trick is a way by which we can simply compute the inner products of the input samples (with the help of kernel function) in their own input feature space without any need to compute the mapping explicitly. This operation will be less computationally expensive.

Question 13

we start with $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$:

$$\begin{aligned}
 \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) &= \begin{bmatrix} x_{i_1}^2 \\ \sqrt{2}x_{i_1}x_{i_2} \\ x_{i_2}^2 \end{bmatrix}^T \begin{bmatrix} x_{j_1}^2 \\ \sqrt{2}x_{j_1}x_{j_2} \\ x_{j_2}^2 \end{bmatrix} \\
 &= x_{i_1}^2 x_{j_1}^2 + 2x_{i_1}x_{j_1}x_{i_2}x_{j_2} + x_{i_2}^2 x_{j_2}^2 \\
 &= (x_{i_1}x_{j_1} + x_{i_2}x_{j_2})^2 \\
 &= (\mathbf{x}_i^T \mathbf{x}_j)^2 \\
 &= k(\mathbf{x}_i, \mathbf{x}_j)
 \end{aligned}$$