

CSE 176 Introduction to Machine Learning Midterm Review

Exam Questions. 1. Matrix Inverse what is a norm? Three conditions 2 3. LI, Lz. Las norm 4. Eigen-decomposition. $A = \begin{bmatrix} 2 & 3 \\ 1 & 4 \end{bmatrix} \rightarrow find V, N.$ 5. Bayes' rule 6. Perception 7, KNN classifier



K-means/medians/modes objectives. 9. mean-shift update X < 10. GMM. how to update soft assignment and {U, S, PZ ZPi=1. M-step 11. adjacency matrix A. Degree matix D. Laplacian Matrix D-A=L 12. Normalized Cut objective 13. PCA (mean, covariance, eigendecom) 14. how to choose bearing rate. VS, Back Propagation

Different AI systems





Major Types of machine learning

□Supervised learning: Given pairs of input-output, learn to map the input to output

□Image classification

□Speech recognition

□Regression (continuous output)

Unsupervised learning: Given unlabeled data, uncover the underlying structure or distribution of the data

□Clustering

Dimensionality reduction

□Reinforcement learning: training an agent to make decisions within an environment to maximize a cumulative reward

□Game playing (e.g., AlphaGo) □Robot control





Linear Algebra, Probability and Statistics

Matrix inverse

□For a 2x2 matrix:

$$A=egin{pmatrix} a&b\c&d \end{pmatrix}$$

The inverse is :

Quiz: find the inv
$$A^{-1} = rac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

□Answer:

$$A = \begin{bmatrix} 3 & -1 & 1 \\ 2 & 0 & 2 \end{bmatrix}, B = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$
$$AB^{T} = \begin{bmatrix} 3 & -1 & 1 \\ 2 & 0 & 2 \end{bmatrix} \cdot \begin{bmatrix} 2 & 0 \\ 2 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0 \\ 6 & 2 \end{bmatrix}$$
$$(AB^{T})^{-1} = \frac{1}{5 \cdot 2 - 0 \cdot 6} \cdot \begin{bmatrix} 2 & 0 \\ -6 & 5 \end{bmatrix} = \frac{1}{10} \cdot \begin{bmatrix} 2 & 0 \\ -6 & 5 \end{bmatrix} = \begin{bmatrix} 0.2 & 0 \\ -0.6 & 0.5 \end{bmatrix}$$

Norms

- Used for measuring the size of a vector
- Norms map vectors to non-negative values
- Norm of vector $\mathbf{x} = [x_1, ..., x_n]^T$ is distance from origin to \mathbf{x}
 - It is any function f that satisfies:

$$f(\boldsymbol{x}) = \boldsymbol{0} \Rightarrow \boldsymbol{x} = \boldsymbol{0}$$

$$f(\boldsymbol{x} + \boldsymbol{y}) \le f(\boldsymbol{x}) + f(\boldsymbol{y})$$
 Triangle Inequality
$$\forall \boldsymbol{\alpha} \in R \quad f(\boldsymbol{\alpha} \boldsymbol{x}) = |\boldsymbol{\alpha}| f(\boldsymbol{x})$$



L^p Norm

• Definition:



- $-L^2$ Norm
 - Called Euclidean norm
 - Simply the Euclidean distance between the origin and the point *x*
 - written simply as ||x||
 - Squared Euclidean norm is same as $\mathbf{x}^{\mathsf{T}}\mathbf{x}$
 - $-L^1$ Norm
 - Sum of absolute value for each x_i



$$\left\| \left| \boldsymbol{x} \right| \right|_{\infty} = \max_{i} \left| x_{i} \right|$$

Called max norm



Slide from S

Eigendecomposition

- Suppose that matrix A has n linearly independent eigenvectors {v⁽¹⁾,...,v⁽ⁿ⁾} with eigenvalues {λ₁,...,λ_n}
- Concatenate eigenvectors to form matrix V
- Concatenate eigenvalues to form vector $\lambda = [\lambda_1, ..., \lambda_n]$
- Eigendecomposition of Ais given by

A=Vdiag(λ) V^{1}



Chain rule of conditional probability

 Any probability distribution over many variables can be decomposed into conditional distributions over only one variable

$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}) = P(\mathbf{x}^{(1)}) \prod_{i=2}^{n} P(\mathbf{x}^{(i)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i-1)})$$

An example with three variables

$$P(\mathbf{a}, \mathbf{b}, \mathbf{c}) = P(\mathbf{a} \mid \mathbf{b}, \mathbf{c})P(\mathbf{b}, \mathbf{c})$$

$$P(\mathbf{b}, \mathbf{c}) = P(\mathbf{b} \mid \mathbf{c})P(\mathbf{c})$$

$$P(\mathbf{a}, \mathbf{b}, \mathbf{c}) = P(\mathbf{a} \mid \mathbf{b}, \mathbf{c})P(\mathbf{b} \mid \mathbf{c})P(\mathbf{c})$$



Independence and conditional independence

- Independence: $x \perp y$
 - Two variables x and y are independent if their probability distribution can be expressed as a product of two factors, one involving only x and the other involving only y

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, \ p(\mathbf{x} = x, \mathbf{y} = y) = p(\mathbf{x} = x)p(\mathbf{y} = y)$$

- Conditional Independence: $x \perp y \mid z$
 - Two variables x and y are independent given variable z, if the conditional probability distribution over x. and y factorizes in this way for every z

 $\forall x \in \mathbf{x}, y \in \mathbf{y}, z \in \mathbf{z}, \ p(\mathbf{x} = x, \mathbf{y} = y \mid \mathbf{z} = z) = p(\mathbf{x} = x \mid \mathbf{z} = z)p(\mathbf{y} = y \mid \mathbf{z} = z)$



Bayes's rule

Bayes' theorem (alternatively Bayes' law or Bayes' rule), named after <u>Thomas Bayes</u>, describes the <u>probability</u> of an <u>event</u>, based on prior knowledge of conditions that might be related to the event.

For example, if the risk of health problems is known to increase with age, Bayes' theorem allows the risk to an individual of a known age to be assessed more accurately by conditioning it relative to their age.

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





Supervised Learning: Classification and Regression



Loss function

□Training dataset of *I* pairs of input/output examples

$$\{{\bf x}_n, {\bf y}_n\}_{n=1}^N$$

Loss function or cost function measures how bad model is:

$$\mathbf{W}^* = \operatorname{arg\,min}_{\mathbf{w}} \Sigma_n L(\mathbf{y}_n, \mathbf{h}(\mathbf{w}, \mathbf{x}_n))$$

 $\Box\,\Theta$ is also a common notation for weights



Supervised ML algorithm

- 1. A model $h(\mathbf{x}; \Theta)$ (hypothesis class) with parameters Θ . A particular value of Θ determines a particular hypothesis in the class. Ex: for linear models, $\Theta = \text{slope } w_1$ and intercept w_0 .
- 2. A loss function $L(\cdot, \cdot)$ to compute the difference between the desired output (label) y_n and our prediction to it $h(\mathbf{x}_n; \boldsymbol{\Theta})$. Approximation error (loss):

$$E(\mathbf{\Theta}; \mathcal{X}) = \sum_{n=1}^{N} L(y_n, h(\mathbf{x}_n; \mathbf{\Theta})) = \text{sum of errors over instances}$$

Ex: 0/1 loss for classification, squared error for regression.

3. An optimization procedure (learning algorithm) to find parameters Θ^* that minimize the error:

 $\Theta^* = \operatorname*{arg\,min}_{\Theta} E(\Theta; \mathcal{X})$



Linear classifier example: perceptron





Underfitting \rightarrow Overfitting

underfitting

"just right"



high training errorhigh test error



Iow training errorIow test error



overfitting

low training errorhigh test error



Model selection and generalization

- Machine learning problems (classification, regression and others) are typically ill-posed : the observed data is finite and does not uniquely determine the classification or regression function.
- How to choose the right inductive bias, in particular the right hypothesis class? This is the model selection problem.



Cross Validation

□Training set:

 \Box Used to train, i.e., to fit a hypothesis h \in H_i.

Optimize parameters of h given the model structure and

hyperparameters.

□Usually done with an optimization algorithm (the learning algorithm).

□ Validation set:

□Used to minimize the generalization error.

Optimize hyperparameters or model structure.

□Usually done with a "grid search". Ex: try all values of $H \in \{10, 50, 100\}$ and $\lambda \in \{10-5, 10-3, 10-1\}$.

Test set:

□Used to report the generalization error.

□We optimize nothing on it, we just evaluate the final model on it





KNN Classifier

K nearest neighbor algorithm

□Nearest neighbor often instable (noise)

□ For a test input *x*, assign the most common label amongst its k most similar training inputs





Effect of K

□Which partition do you prefer? Why?



□*K*controls the degree of smoothing.□What if K=N(number of data points)?



Choosing K

How should we choose K?

□Select K with highest test accuracy

□Can we simply split to training and testing set?

□Solution: split data into training, validation and test sets

- □Training set: compute nearest neighbour
- □Validation set: optimize hyperparameters such as K
- □Test set: measure performance



[4] The table below shows the test set for a **1-nearest-neighbor classifier** that uses Manhattan distance, i.e., the distance between two points at coordinates p and q is |p - q|. The only attribute, X, is real-valued, and the label, Y, has two classes, 0 and 1. Suppose a *subset* containing $n \le 8$ examples is selected from this set to train the classifier, and the accuracy of the classifier is 100 percent when tested on this set (with *all* 8 examples). What is the *smallest* possible value for n? In case of ties in distance, use the example with smallest X value as the neighbor.

X	-5	-4	-1	0	1	3	4	8
Y	0	1	0	0	0	0	0	1

- A. 2B. 3C. 4
- D. 5
- E. 6



KNN for high-dimensional data

Can we use KNN classifier for high-dimensional data?

□Assumption for KNN classifier to work:

□K nearest neighbors are *nearby*

□Are K nearest neighbors nearby for d>>0?





Clustering

K-means Clustering: Algorithm

- Initialization step •
 - pick *K* cluster centers randomly 1.
 - assign each sample to its closest center 2.



- Iteration steps
 - compute centers as cluster means $\mu_k = \frac{1}{|S^k|} \sum_{p \in S^k} f_p$ 1.
 - re-assign each sample to the closest mean 2.
- Iterate until clusters stop changing ٠

K-means objective



 μ_k : extra parameters (means)

$$E(S,\mu) = + + + +$$

$$= \sum_{k=1}^{K} \sum_{p \in S^k} \|f_p - \mu_k\|^2$$

(SSD)



Squared distance as log-likelihood

Assume K=2, $\Omega = S \cup \bar{S}$

single Gaussian of *fixed* covariance

$$\sum_{p \in S} \|f_p - \mu_S\|^2 + \sum_{p \in \overline{S}} \|f_p - \mu_{\overline{S}}\|^2$$

$$= -\sum_{p \in \mathbf{S}} \ln \mathcal{N}(f_p | \mu_{\mathbf{S}}) - \sum_{p \in \bar{\mathbf{S}}} \ln \mathcal{N}(f_p | \mu_{\bar{\mathbf{S}}})$$

single Gaussian



$$\theta_S = \{\mu_S\}$$

K-means as variance clustering criteria



[4] You want to cluster 7 points into 3 clusters using the *k*-Means Clustering algorithm. Suppose after the first iteration, clusters C_1 , C_2 and C_3 contain the following two-dimensional points:

 C_1 contains the 2 points: {(0,6), (6,0)}

C₂ contains the 3 points: {(2,2), (4,4), (6,6)}

 C_3 contains the 2 points: {(5,5), (7,7)}

What are the cluster centers computed for these 3 clusters?

A. C_1 : (3,3), C_2 : (4,4), C_3 : (6,6) B. C_1 : (3,3), C_2 : (6,6), C_3 : (12,12) C. C_1 : (6,6), C_2 : (12,12), C_3 : (12,12)

D. C₁: (0,0), C₂: (48,48), C₃: (35,35)

E. None of these



(generalization) Distortion Clustering can use different "distortion" measures

$$E(S,\mu) = \sum_{k=1}^{K} \sum_{p \in S_k} ||f_p - \mu_k||_d$$

interpretation of parameters μ_k	examples of distortion measure $\ \cdot\ _d$					
K-means	squared L_2 norm	$\ \cdot\ _d = \ \cdot\ ^2$				
e n K-medians	absolute L_2 norm	$\ \cdot\ _d = \ \cdot\ $				
K-modes	$\cdot \ ^2)$	$\ \cdot\ _d = 1 - \exp(-\ $				

NOTE: besides changing the distortion measure, there are different generalizations of K-means requiring **other interpretations of SSE objective**

Mean Shift

Iterative Mode Search

[Fukunaga and Hostetler 1975, Cheng 1995, Comaniciu & Meer 2002]



- 1. Initialize random seed, and fixed window
- 2. Calculate center of gravity 'x' of the window (the "mean")
- 3. Translate the search window to the mean
- 4. Repeat Step 2 until convergence

Mean Shift as K-modes

[Salah, Mitche, Ben-Ayed 2010]



Mean-shift segmentation relates to distortion clustering with a bounded loss (**K-modes**)



Mean shift trajectories

Kernel Density Estimation

Kernel density estimate with bandwidth σ : a mixture having one component for each data point:

$$p(\mathbf{x}) = \sum_{n=1}^{N} p(\mathbf{x}|n) p(n) = \frac{1}{N\sigma^{D}} \sum_{n=1}^{N} K\left(\frac{\mathbf{x} - \mathbf{x}_{n}}{\sigma}\right) \qquad \mathbf{x} \in \mathbb{R}^{D}.$$

Usually the kernel K is Gaussian: $K\left(\frac{\mathbf{x}-\mathbf{x}_n}{\sigma}\right) = (2\pi)^{-D/2} \exp\left(-\frac{1}{2} \|(\mathbf{x}-\mathbf{x}_n)/\sigma\|^2\right).$



Mean-shift

Mean-shift algorithm: starting from an initial value of \mathbf{x} , it iterates the following expression:

$$\mathbf{x} \leftarrow \sum_{n=1}^{N} p(n|\mathbf{x}) \mathbf{x}_{n} \quad \text{where} \quad p(n|\mathbf{x}) = \frac{p(\mathbf{x}|n)p(n)}{p(\mathbf{x})} = \frac{\exp\left(-\frac{1}{2}\|(\mathbf{x} - \mathbf{x}_{n})/\sigma\|^{2}\right)}{\sum_{n'=1}^{N} \exp\left(-\frac{1}{2}\|(\mathbf{x} - \mathbf{x}_{n'})/\sigma\|^{2}\right)}$$

 $\sum_{n=1}^{N} p(n|\mathbf{x})\mathbf{x}_n$ can be understood as the weighted average of the N data points using as weights the posterior probabilities $p(n|\mathbf{x})$. The mean-shift algorithm converges to a mode of $p(\mathbf{x})$. Which one it converges to depends on the initialization. By running mean-shift starting at a data point \mathbf{x}_n , we effectively assign \mathbf{x}_n to a mode. We repeat for all points $\mathbf{x}_1, \ldots, \mathbf{x}_N$.



K-means and MLE (maximum likelihood estimation)

"hard"
K-means
$$E(S,\mu) = -\sum_{k=1}^{K} \sum_{p \in S^k} \log P(f_p \mid \mu_k)$$

multi-variate (i.e. $x, \mu \in R^{T}$) Gaussian distribution (simple special case $\Sigma = \sigma^2 \mathbf{I}$)

$$\frac{1}{P(x|\mu)} = \frac{1}{\sqrt{(2\pi\sigma^2)^N}} \exp{-\frac{\|x-\mu\|^2}{2\sigma^2}}$$



Towards soft clustering... Gaussian Mixture Models (GMM)

- Soft clustering using Gaussian Mixture Model (GMM)
 - no "hard" assignments of points to K distinct (Gaussian) clusters S^k
 - all points are used to estimate parameters of one complex K-mode distribution (GMM)



Expectation-Maximization (EM)

GMM estimation - optimization of ML objective (sum of Negative Log Likelihoods, a.k.a. NLL loss)

$$E_{gmm}(\theta) := -\sum_{p} \log P_{gmm}(x_{p} \mid \theta) \equiv -\sum_{p} \log \left(\sum_{k} \rho_{k} P(x_{p} \mid \mu_{k}, \sigma_{k}) \right)$$

$$\downarrow L(\theta \mid \tilde{S}) \quad \text{for given } \tilde{\theta} = (\tilde{\mu}, \tilde{\sigma}, \tilde{\rho}) \quad \text{can find tight upper bound} \quad \tilde{S}_{p}^{k} = \frac{\tilde{\rho}_{k} P(x_{p} \mid \tilde{\mu}_{k}, \tilde{\sigma}_{k})}{\sum_{m} \tilde{\rho}_{m} P(x_{p} \mid \tilde{\mu}_{m}, \tilde{\sigma}_{m})} \quad \mathbf{E-step} \quad$$

Graph representation

- \Box A graph is defined as a tuple G = (V, E)
 - University where V is a set of nodes
 - \Box and E is a set of edges
- □ An example graph with 6 nodes and 7 edges



Minimum Cut

- □ In many applications it is desired to find the cut with minimum cost: *minimum cut*
- Well studied problem in graph theory, with many applications
- □ There exists efficient algorithms for finding minimum cuts



Normalized Cut

□ Normalize *cut cost* by volume of clusters

Compute the cut cost as a fraction of the total edge connections to all nodes in the graph





Summary of Normalized cut algorithm

- Given a set of features, construct a weighted graph by computing weight on each edge and then placing the data into W and D.
- Solve (D-W)x=λDx for eigen vectors with the smallest eigenvalues.
- Output Use the eigen vector corresponding to the second smallest eigenvalue to bipartition the graph into two groups.
- ${f 0}$ Recursively repartition the segmented parts if necessary.



Toward Kernel K-means



$$E(S,\mu) = \sum_{k=1}^{K} \sum_{p \in S^k} \|f_p - \mu_k\|^2$$

(Basic K-means)

$$E_k(S, \hat{\mu}) = \sum_{k=1}^K \sum_{p \in S^k} \|\phi(f_p) - \hat{\mu}_k\|^2$$

(Kernel K-means)

Explicit Kernel



kernel K-means or *average association*







Other kernel (graph) clustering objectives

Average Association

Average Cut









kernel (graph) clustering objectives

Average Association Average Cut





Normalized Cut

Normalized Association

$$\sum_{k=1}^{K} \frac{S^{k'} A \left(1 - S^{k}\right)}{d' S^{k}} \equiv K - \sum_{k=1}^{K} \frac{S^{k'} A S^{k}}{d' S^{k}}$$

normalization $d := A\mathbf{1}$



Dimensionality Reduction

Motivation for Dimensionality Reduction



Principal Component Analysis

Goal: Find *r*-dim projection that best preserves variance

- 1. Compute mean vector μ and covariance matrix Σ of original points
- 2. Compute eigenvectors and eigenvalues of Σ
- 3. Select top r eigenvectors
- 4. Project points onto subspace spanned by them:

$$y = A(x - \mu)$$

where y is the new point, x is the old one, and the rows of A are the eigenvectors

Covariance

- Variance and Covariance:
 - Measure of the "spread" of a set of points around their center of mass(mean)
- Variance:
 - Measure of the deviation from the mean for points in one dimension
- Covariance:
 - Measure of how much each of the dimensions vary from the mean with respect to each other
 - Covariance is measured between two dimensions
 - Covariance sees if there is a relation between two dimensions
 - Covariance between one dimension is the variance

Principal Component Analysis

Input: $\mathbf{x} \in \mathbb{R}^D$: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

Set of basis vectors: $\mathbf{u}_1, \ldots, \mathbf{u}_K$

Summarize a D dimensional vector X with K dimensional feature vector h(x)

$$h(\mathbf{x}) = \begin{bmatrix} \mathbf{u}_1 \cdot \mathbf{x} \\ \mathbf{u}_2 \cdot \mathbf{x} \\ \\ \cdots \\ \mathbf{u}_K \cdot \mathbf{x} \end{bmatrix}$$

You are given a design matrix $X = \begin{bmatrix} 6 & -4 \\ -3 & 5 \\ -2 & 6 \\ 7 & -3 \end{bmatrix}$. Let's use PCA to reduce the dimension from 2 to 1.

(1) [6 pts] Compute the covariance matrix for the sample points. (Warning: Observe that X is not centered.) Then compute the **unit** eigenvectors, and the corresponding eigenvalues, of the covariance matrix. *Hint:* If you graph the points, you can probably guess the eigenvectors (then verify that they really are eigenvectors).





Neural Network



 x_1

*x*₂

 x_3

Input layer



+1

Multi-layer perceptron



Example of Multi Layer Perceptron (MLP)



Activation function



Rectified Linear Unit



Multi-variate functions

Gradient Descent

Example: for a function of two variables



- direction of (negative) gradient at point $\mathbf{x} = (x_1, x_2)$ is direction of the steepest descent towards lower values of function $L_{5,62}$
- magnitude of gradient at $\mathbf{x} = (x_1, x_2)$ gives the value of the slope user

Multi-variate functions

Gradient Descent

Example: for a function of two variables



How to Set Learning Rate α ?

$$\mathbf{x}' = \mathbf{x} - \alpha \, \nabla L$$

If α too small, too many iterations to converge



 If α too large, may overshoot the local minimum and possibly never even converge



Variable Learning Rate

If desired, can change learning rate α at each iteration

 $\begin{aligned} \mathbf{k} &= 1 \\ \mathbf{x}^{(1)} &= \text{any initial guess} \\ \text{choose } \alpha, \varepsilon \\ \text{while } \alpha ||\nabla \mathbf{L}(\mathbf{x}^{(k)})|| &> \varepsilon \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha \nabla \mathbf{L}(\mathbf{x}^{(k)}) \\ \mathbf{k} &= \mathbf{k} + 1 \end{aligned}$

fixed α gradient descent variable α gradient descent



Learning Rate

 Monitor learning rate by looking at how fast the objective function decreases



Computing Derivatives: Small Example

- Small network f(x,y,z) = (x+y)z
- Rewrite using $\mathbf{q} = \mathbf{x} + \mathbf{y} \Rightarrow \mathbf{f}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \mathbf{q}\mathbf{z}$
- Want $\frac{\partial \mathbf{f}}{\partial \mathbf{x}^{'} \partial \mathbf{y}^{'} \partial \mathbf{z}}$

chain rule for
$$f(y(x))$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial x}$$

- Compute $\frac{\partial \mathbf{f}}{\partial}$ from the end backwards
 - for each edge, with respect to the main variable at edge origin
 - using chain rule with respect to the variable at edge end, if needed



Computing Derivatives: Vector Notation

- Let $\mathbf{f}(\mathbf{x}): \mathbf{R}^n \rightarrow \mathbf{R}^m$,
 - \mathbf{x} is \mathbf{n} -dimensional vector and output $\mathbf{f}(\mathbf{x})$ is \mathbf{m} -dimensional vector
- Jacobian matrix
 - has **m** rows and **n** columns
 - has $\frac{\partial \mathbf{f}_i}{\partial \mathbf{x}_j}$ in row **i**, column **j**



Computing Derivatives: Vector Notation

- $f(x): \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g(x): \mathbb{R}^k \rightarrow \mathbb{R}^n$
- $f(g(x)): \mathbb{R}^k \rightarrow \mathbb{R}^m$
- Chain rule for vector functions

 $\partial \mathbf{f} - \partial \mathbf{f} \partial \mathbf{g}$ $\partial \mathbf{x} \quad \partial \mathbf{g} \, \partial \mathbf{x}$ Jacobian matrices



(10 points) Consider boolean logical operators such as AND, OR, and XOR,

AND				OR			XOR		
x1	x2	у	x 1	x2	у	x 1	x2	у	
0	0	0	0	0	0	0	0	0	
0	1	0	0	1	1	0	1	1	
1	0	0	1	0	1	1	0	1	
1	1	1	1	1	1	1	1	0	

(a) (5 points) Can we implement OR using perception? If yes, show the corresponding perceptron network. If no, explain why.
 Solution:

(b) (5 points) Can we implement XOR using perception? If yes, show the corresponding perceptron network. If no, explain why.
 Solution:



(10 points) This problem explores computing derivatives on composite functions. Consider the function:

$$y = \exp[\exp[x] + \exp[x]^2] + \sin[\exp[x] + \exp[x]^2].$$

We can break this down into a series of intermediate computations so that:

$$\begin{array}{rcrcrcr} f_1 & = & \exp[x] \\ f_2 & = & f_1^2 \\ f_3 & = & f_1 + f_2 \\ f_4 & = & \exp[f_3] \\ f_5 & = & \sin[f_3] \\ y & = & f_4 + f_5 \end{array}$$

Compute the derivatives in order using the chain rule of gradients in each case to make use of the derivatives already computed.

$$rac{\partial y}{\partial f_5}, rac{\partial y}{\partial f_4}, rac{\partial y}{\partial f_3}, rac{\partial y}{\partial f_2}, rac{\partial y}{\partial f_1}, ext{ and } rac{\partial y}{\partial x}.$$

