

# CS 446/ECE 449: Machine Learning

## Problem Sets

## 1. [13 points] Regression

Suppose we are given a set of observations  $\{(x^{(i)}, y^{(i)})\}$  where  $x, y \in \mathbb{R}$  and  $i \in \{1, 2, \dots, N\}$ . Consider the following program:

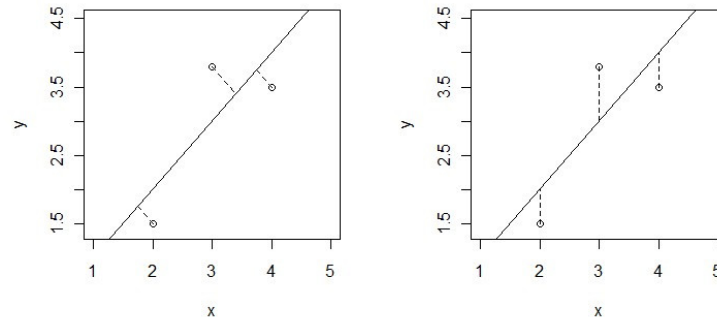
$$\operatorname{argmin}_{w_1, w_2} \sum_{i=1}^N \left( y^{(i)} - w_1 \cdot x^{(i)} - w_2 \right)^2. \quad (1)$$

- (a) What is the minimum number of observations required for a unique solution?

Your answer:

- (b) Which of the following plots shows the correct residual that is minimized with the above program? Circle the correct answer.

Your answer:



- (c) Suppose we now want to fit a quadratic model to the observed data. Modify the program given in Eq. 1 accordingly. Derive the closed form solution for this case. You may assume that you have a sufficient number of data samples. Use matrix vector notation, i.e.,  $\mathbf{w}$ ,  $\mathbf{X}$ , and  $\mathbf{Y}$  and define them carefully.

Your answer:

- (d) Briefly describe the problem(s) that we encounter if we were to fit a high degree polynomial to a data that is known to be linear.

Your answer:

- (e) The program above (Eq. 1) assumes  $x \in \mathbb{R}$ . State the program for  $\mathbf{x} \in \mathbb{R}^D$  and specify the dimensions of  $\mathbf{w}$ .

Your answer:

- (f) If  $\dim(\mathbf{x}) = D > N$ , how could the program be modified such that a unique solution can be obtained?

Your answer:

- (g) Is there a closed form solution to the new program? If so derive it.

Your answer:

2. [7 points] Regression

Suppose we are given a set of observations  $\{(x^{(i)}, y^{(i)})\}$  where  $x, y \in \mathbb{R}$  and  $i \in \{1, 2, \dots, N\}$ . Consider the following program:

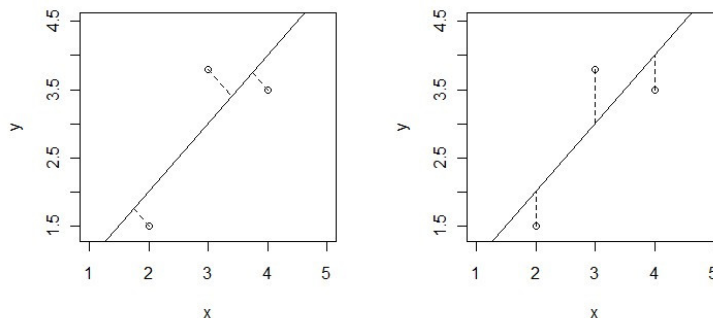
$$\underset{w_1, w_2}{\operatorname{argmin}} \sum_{i=1}^N \left( y^{(i)} - w_1 \cdot x^{(i)} - w_2 \right)^2. \quad (2)$$

- (a) (1 point) What is the minimum number of distinct observations required for a unique solution?

Your answer:

- (b) (1 point) Which of the following two plots shows the correct residual that is minimized with the above program? Circle the correct answer.

Your answer:



- (c) (2 points) The program above (Eq. 2) assumes  $x \in \mathbb{R}$ . State the program for  $\mathbf{x} \in \mathbb{R}^d$  and specify the dimensions of  $\mathbf{w}$ . Assume the feature vector  $\mathbf{x} \in \mathbb{R}^d$  includes the 1 for the bias.

Your answer:

- (d) (1 point) If  $\dim(\mathbf{x}) = d > N$ , how could the program be regularized such that a unique solution can be obtained? (Hint: use  $\ell_2$  regularization)

Your answer:

- (e) (2 points) **Derive** a closed form solution  $w^*$  for the new regularized program. Use the following notation:  $\mathbf{X} \in \mathbb{R}^{d \times N}$  is the input data matrix,  $\mathbf{Y} \in \mathbb{R}^N$  is the vector of labels,  $\mathbf{I}$  is the identity matrix. Your final expression should be in terms of these provided quantities.

Your answer:

3. [15 points] Softmax Regression

We are given a dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{|\mathcal{D}|}$  with feature vectors  $\mathbf{x}^{(i)} \in \mathbb{R}^d$  and their corresponding labels  $y^{(i)} \in \{1, \dots, K\}$ . Here,  $K$  denotes the number of classes. The distribution over  $y^{(i)}$  is given via

$$p(y^{(i)} = k | \mathbf{x}^{(i)}) = \mu_k(\mathbf{x}^{(i)}) = \frac{e^{\mathbf{w}_k^T \mathbf{x}^{(i)}}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}^{(i)}}. \quad (3)$$

Our goal is to learn the weight parameters  $\{\mathbf{w}_j\}_{j=1}^K$ , with  $\mathbf{w}_j \in \mathbb{R}^d$  ( $\forall j$ ).

- (a) Show that the negative conditional log-likelihood  $\ell(\mathbf{w}_1, \dots, \mathbf{w}_K)$  is given by the expression,

$$-\log p(\mathcal{D}) = -\log p(\{y^{(i)}\}_{i=1}^{|\mathcal{D}|} | \{\mathbf{x}^{(i)}\}_{i=1}^{|\mathcal{D}|}) = -\sum_{i=1}^{|\mathcal{D}|} \sum_{k=1}^K \mathbb{1}_{\{y^{(i)}=k\}} \mathbf{w}_k^T \mathbf{x}^{(i)} + \sum_{i=1}^{|\mathcal{D}|} \log \left( \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}^{(i)}} \right).$$

Here,  $\mathbb{1}_{\{y^{(i)}=k\}}$  is an indicator variable, *i.e.*,  $\mathbb{1}_{\{y^{(i)}=k\}}$  is equal to 1 if  $(y^{(i)} = k)$  and equal to 0 otherwise. *Show intermediate steps and state any used assumptions.*

Your answer:

- (b) We want to minimize the negative log-likelihood. To combat overfitting, we add a regularizer to the objective function. The regularized objective is  $\ell_r(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{1}{|\mathcal{D}|} \ell(\mathbf{w}_1, \dots, \mathbf{w}_K) + \lambda \sum_{k=1}^K \|\mathbf{w}_k\|^2$ . Justify that  $\lambda$  should be a strictly positive scalar, *i.e.*,  $\lambda > 0$ .

Your answer:

- (c) Show that the gradient of the regularized loss  $\ell_r$  is

$$\nabla_{\mathbf{w}_k} \ell_r(\mathbf{w}_1, \dots, \mathbf{w}_K) = 2\lambda \mathbf{w}_k + \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} (\mu_k(\mathbf{x}^{(i)}) - \mathbb{1}_{\{y^{(i)}=k\}}) \mathbf{x}^{(i)}.$$

Your answer:

- (d) State the gradient update (formula) for **gradient descent** for the regularized loss  $\ell_r$ . Use a learning rate of  $\alpha$ .

Your answer:

- (e) State the gradient update for **stochastic gradient descent** for the regularized loss  $\ell_r$ . Use a batch size equal to 1 and a learning rate of  $\alpha$ .

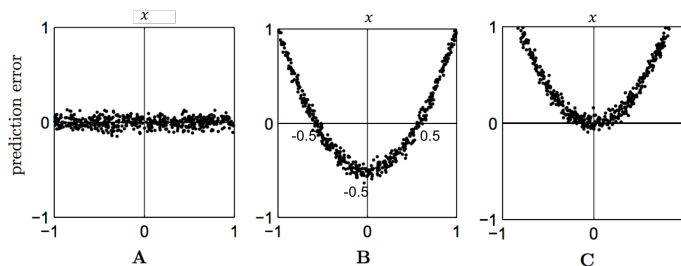
Your answer:

- (f) Consider the case where every instance can be assigned more than one label, for example a movie can be assigned the labels *action* and *comedy* simultaneously. Would you choose to use the classifier from Eq. (3) or  $K$  1-vs.-all logistic regression classifiers?

Your answer:

4. [15 points] Regression

- (a) [3 points] Each plot below claims to represent prediction errors as a function of  $x$  for a trained regression model based on some dataset. Some of these plots could potentially be prediction errors for linear ( $y = w_1x + w_0$ ) or quadratic ( $y = w_2x^2 + w_1x + w_0$ ) regression models, while others could not. The regression models are trained with the least squares loss (L2 Loss). Please indicate compatible models and plots.



Your answer:

	A	B	C
linear regression	( )	( )	( )
quadratic regression	( )	( )	( )

- (b) [4 points] Consider a simple one dimensional logistic regression model:

$$P(y = 1|x, w_1, w_2) = g(w_2 + w_1x), \quad (4)$$

where  $g(z) = (1 + \exp(-z))^{-1}$  is the logistic function. Fig. 1 shows two possible conditional distributions  $P(y = 1|x, w_1, w_2)$ , viewed as a function of  $x$ , that we can get by changing the parameters  $w_1$  and  $w_2$ . Indicate the number of classification errors for each conditional on the samples A, B and C with coordinates  $x_A = -1, x_B = 0$  and  $x_C = 1$  and labels  $y_A = 0, y_B = 1$  and  $y_C = 0$ . Consider a threshold of 0.5, *i.e.*,  $y$  is assigned the label 1 if  $P(y = 1|x, w_1, w_2) \geq 0.5$ .

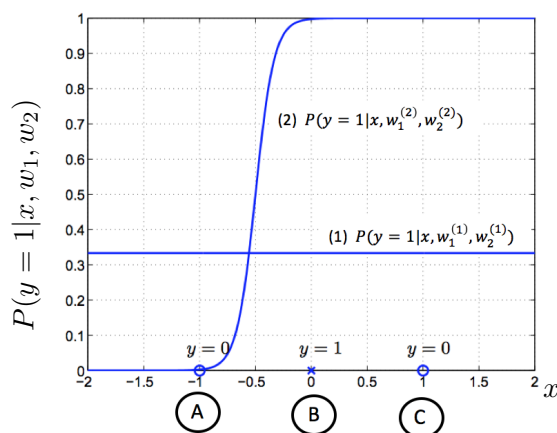


Figure 1: Two possible logistic regression solutions for the three labeled points.

Your answer:

Conditional (1) makes ( ) classification error(s) on sample(s) ( ).  
 Conditional (2) makes ( ) classification error(s) on sample(s) ( ).

- (c) [3 points] One of the conditionals in Fig. 1 corresponds to the **maximum likelihood** setting of the parameters  $w_1$  and  $w_2$  based on the labeled data in the figure. Which one is the maximum likelihood solution (1 or 2)?

Hint: Justify by providing  $P(y_A, y_B, y_C | x, w_1, w_2)$  for the two conditionals.

Your answer:

- (d) [2 points] If we additionally tell you that we added a regularization penalty  $\frac{1}{2}w_1^2$  to the log-likelihood estimation criterion, would this new information affect your choice of the solution in part (c)? Answer by Yes/No. Justify.

Hint: Check the values of  $w_1^{(1)}$  and  $w_1^{(2)}$ . See Fig. 1 for their definition.

Your answer:

- (e) [3 points] In the following, consider an 8-dimensional logistic regression with parameters  $w = \{w_i\}_{i=0}^7$  fitted on a dataset  $\{(x_k, y_k)\}_{k=1}^N$  and the three objective functions:

$$L_1 : \min_w - \sum_{k=1}^N \log P(y_k | x_k, w), \quad (5)$$

$$L_2 : \min_w - \sum_{k=1}^N \log P(y_k | x_k, w) + \lambda \sum_{j=0}^7 w_j^2, \quad (6)$$

$$L_3 : \min_w - \sum_{k=1}^N \log P(y_k | x_k, w) + \lambda \sum_{j=0}^7 |w_j|, \quad (7)$$

where  $\lambda \in \mathbb{R}$  and  $\lambda > 0$ . The following table contains the weights learned for all three objective functions (not in any particular order).

	Column A	Column B	Column C
$w_0$	0.61	0.23	0.01
$w_1$	0.84	0.12	0.00
$w_2$	1.4	0.12	0.00
$w_3$	0.32	0.09	0.28
$w_4$	0.74	0.34	0.01
$w_5$	0.93	0.52	0.00
$w_6$	0.58	0.16	0.13
$w_7$	0.24	0.14	0.00

Beside each objective write the appropriate column label (A, B, or C).

Your answer:

Objective  $L_1 \rightarrow$  Column ( )

Objective  $L_2 \rightarrow$  Column ( )

Objective  $L_3 \rightarrow$  Column ( )



5. [9 points] Regression

- (a) Consider the following dataset  $\mathcal{D}$  in the one-dimensional space.

$i$	$x^{(i)}$	$y^{(i)}$
1	0	-1
2	1	2
3	1	0

Table 1: Data for  $\mathcal{D}$

For a set of observations  $\mathcal{D} = \{(y^{(i)}, x^{(i)})\}$ , where  $y^{(i)}, x^{(i)} \in \mathbb{R}$  and  $i \in \{1, 2, \dots, |\mathcal{D}|\}$ , we optimize the following program.

$$\underset{w_1, w_2}{\operatorname{argmin}} \sum_{(y^{(i)}, x^{(i)}) \in \mathcal{D}} (y^{(i)} - w_1 \cdot x^{(i)} - w_2)^2 \quad (8)$$

Find the optimal  $w_1^*, w_2^*$  given the aforementioned dataset  $\mathcal{D}$  and justify your answer.

**Compute the scalars  $w_1^*$  and  $w_2^*$ .**

Your answer:

$$w_1^* =$$

$$w_2^* =$$

- (b) What is the minimum number of observations that are required to obtain a unique solution for the program in Eq. (8)?

Your answer:

- (c) Consider another dataset  $\mathcal{D}_1$ , where  $x^{(i)}, y^{(i)} \in \mathbb{R}$

$i$	$x^{(i)}$	$y^{(i)}$
1	0	0
2	1	1
3	2	4
4	3	9
5	4	16

Table 2: Data for  $\mathcal{D}_1$

Clearly  $\mathcal{D}_1$  can not be fit exactly with a linear model. In class, we discussed a simple approach of building a nonlinear model while still using our linear regression tools. How would you use the linear regression tools to obtain a nonlinear model which better fits  $\mathcal{D}_1$ , *i.e.*, what feature transform would you use? Provide your reasons and write down the resulting program that you would optimize using a notation which follows Eq. (8), *i.e.*, make all the trainable parameters explicit. **Do NOT plug the datapoints from  $\mathcal{D}_1$  into your program and solve for its parameters. Just provide the program.**

Your answer:

- (d) Write down a program equivalent to the one derived in part (c) using matrix-vector notation. Carefully define the matrices and vectors which you use, their dimensions and their entries. Show how you fill the matrices and vectors with the data. Derive the closed form solution for this program using the symbols which you introduced. **Do NOT compute the solution numerically.**

Your answer:

- (e) Briefly describe the problem(s) that we will encounter if we were to fit a very high degree polynomial to the dataset  $\mathcal{D}_1$ ?

Your answer:

6. [14 points] Linear and Logistic Regression

- (a) Using vector  $y \in \mathbb{R}^N$ , matrix  $X \in \mathbb{R}^{N \times d}$ , and vector  $w \in \mathbb{R}^d$ , linear regression can be formulated as

$$\min_w \frac{1}{2} \|y - Xw\|_2^2 \quad (9)$$

What does  $N$  and  $d$  refer to and what is the optimal argument  $w^*$  to the program given in Eq. 9.

Your answer:

- (b) What is an issue when using the program given in Eq. 9 for classification rather than regression?

Your answer:

- (c) Assume we are given a dataset  $\mathcal{D} = \{(x_i, y_i)\}$  containing sample pairs composed of datapoints  $x_i \in \mathbb{R}^d$  and class labels  $y_i \in \{0, 1\}$ . Further, assume we are given a probabilistic model  $p(y_i|x_i)$ . Under the assumption that all samples are independent and identically distributed (i.i.d.), what is the probability/likelihood of the dataset under the model  $p(y_i|x_i)$ .

Your answer:

- (d) Assume our probabilistic model depends on some parameters  $w \in \mathbb{R}^d$  and is given by

$$p(y_i|x_i) = \frac{1}{1 + \exp(-y_i w^\top x_i)}$$

What is the negative log-likelihood of the i.i.d. dataset  $\mathcal{D}$  under this model, and how do we want to choose the parameters  $w$  of the model?

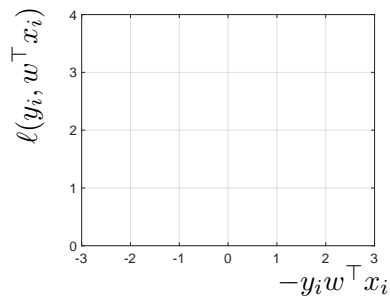
Your answer:

- (e) Without further assumptions and restrictions on the program you derived, is it possible to analytically compute the solution to the program? Yes or No? Justify your answer.

Your answer:

- (f) Determine  $R$  and  $\ell$  by comparing your program to the following formulation and complete the given illustration:

$$\min_{w \in \mathbb{R}^d} R(w) + \sum_{(x_i, y_i) \in \mathcal{D}} \ell(y_i, w^\top x_i)$$



Your answer:

7. [12 points] Binary Classifiers

- (a) Is it possible to use a linear regression model for binary classification? If so, how do we map the regression output  $\mathbf{w}^\top \mathbf{x}$  to the class labels  $y \in \{-1, 1\}$ ?

Your answer:

- (b) Assume  $y \in \{-1, 1\}$ . Consider the following program for logistic regression:

$$\min_{\mathbf{w}} f(\mathbf{w}) := \sum_i \log \left( 1 + \exp(-y^{(i)} \mathbf{w}^T \phi(x^{(i)})) \right).$$

Assume the data is not linearly separable. Is there an analytical solution? If so, derive it. If not, sketch an algorithm that solves the program iteratively. (Make sure to include any important mathematical expressions.)

Your answer:

- (c) The above program for binary classification, makes assumptions on the samples/data points. What are those assumptions.

Your answer:

- (d) What is the probability model assumed for logistic regression and linear regression? Give names rather than equations.

Your answer:

8. [8 points] Binary Classifiers

Based on a data set,  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ , where  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{0, 1\}$ , and samples are i.i.d., we want to train a logistic regression model. We define our probabilistic model to have the form:

$$\hat{y}_i = g(\mathbf{w}^T \mathbf{x}_i) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}},$$

where  $\hat{y}_i$  is the probability given data  $\mathbf{x}_i$  and model parameters  $\mathbf{w} \in \mathbb{R}^d$ . Given this notation we define the probability of predicting  $y_i$  via

$$P[Y = y_i | X = \mathbf{x}_i] = (\hat{y}_i)^{y_i} \cdot (1 - \hat{y}_i)^{(1-y_i)}.$$

We want to find the model parameters  $\mathbf{w}$ , such that the likelihood of the data set  $D$  is maximized, which is formulated as

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( - \sum_{i=1}^N (y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)) \right).$$

(a) (6 points) Let the program above be referred to as:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} L(y, \mathbf{x}, \mathbf{w})$$

Is  $L(y, \mathbf{x}, \mathbf{w})$  convex with respect to  $\mathbf{w}$ ? Prove it is convex or non-convex without using knowledge of convexity for any function. (Hint: use the Hessian.)

Your answer:

- (b) (2 points) Can we find a closed form analytic solution for  $\mathbf{w}$ ? How to train the model  $\mathbf{w}$  based on the data set  $D$ ? State your approach and write down the equation.

Your answer:



9. [10 points] Binary Classifiers

- (a) Assume  $y \in \{-1, 1\}$ . Consider the following program for linear regression:

$$\min_{\mathbf{w}} f(\mathbf{w}) := \frac{1}{2} \sum_i \left(1 - y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}\right)^2$$

Is the objective function  $f(\mathbf{w})$  convex in  $\mathbf{w}$  assuming everything else given and fixed? (Yes or No)

Your answer:

- (b) Assume  $y \in \{-1, 1\}$ . Consider the following program for logistic regression:

$$\min_{\mathbf{w}} f(\mathbf{w}) := \sum_i \log \left(1 + \exp(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)})\right)$$

Is the objective function  $f(\mathbf{w})$  convex in  $\mathbf{w}$  assuming everything else given and fixed? (Yes or No)

Your answer:

- (c) We want to use gradient descent to address the above **logistic regression** program. What is the gradient  $\nabla_{\mathbf{w}} f(\mathbf{w})$ ? Use the symbols and notation which was used in the cost function.

Your answer:

- (d) What is the probability model assumed for logistic regression and linear regression? Give names rather than equations.

Your answer:

10. [17 points] Optimization

- (a) [2 points] Fig. 2 shows the cost curves of a deepnet training with two different optimization algorithms, *i.e.*, **gradient descent** and **stochastic gradient descent**. Which of the graphs corresponds to which optimization algorithm? Explain.

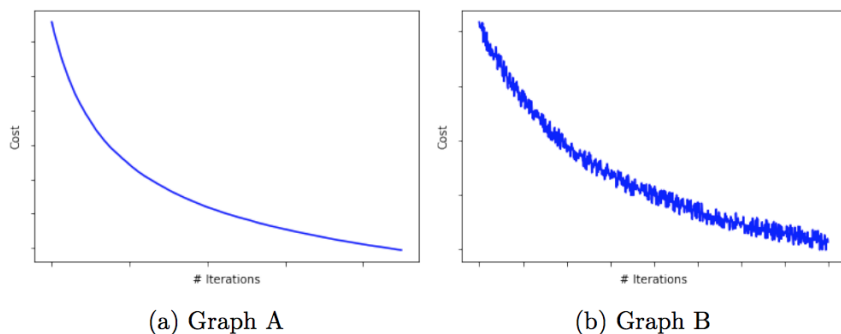


Figure 2: Training cost as a function of the number of iterations.

Your answer:

- (b) [1 points] Fig. 3 below shows the cost curve during a deepnet training. Which of the following options could have caused the sudden drop in the cost: (1) learning rate-decay (decreasing the learning rate), (2) increasing the batch-size. Explain.

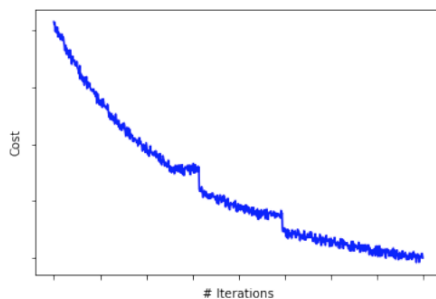


Figure 3: Training cost as a function of the number of iterations.

Your answer:

- (c) [3 points] Consider the following program:

$$\min_w \frac{1}{2} \max_{i \in \{1, \dots, k\}} (a_i^T w + b_i) + \frac{1}{2} \|w\|^2, \quad (10)$$

where  $w, a_i \in \mathbb{R}^d$  and  $b_i \in \mathbb{R}$  for  $1 \leq i \leq k$ . To rewrite the optimization problem as a quadratic programming program, we substitute the maximum term with an auxiliary slack variable  $\xi \in \mathbb{R}$ . Fill in the missing constraint in the program below:

$$\begin{aligned} \min_{w, \xi} \quad & \xi + \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & \dots\dots\dots \\ & \text{(missing constraint)} \end{aligned} \quad (11)$$

Your answer:

- (d) [2 points] Consider the following program:

$$\begin{aligned} \min_{w, \xi} \quad & \xi + \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & A^T w + b \leq \xi \mathbf{1} \end{aligned} \quad (12)$$

Write down the Lagrangian  $L(w, \xi, u)$  of the program in Eq. (12), where  $u \in \mathbb{R}^k$  denotes the Lagrange multiplier,  $u \geq 0$ ,  $A \in \mathbb{R}^{d \times k}$ ,  $b \in \mathbb{R}^k$  and  $\mathbf{1} \in \mathbb{R}^k$  is a vector of all ones.

Your answer:

- (e) [3 points] Write down the dual program of the program in Eq. (12) as a function of  $L(w, \xi, u)$ , the **min** and **max** operators, as well as  $w$ ,  $\xi$  and  $u$ . Make sure to specify all the constraints.

Your answer:

- (f) [2 points] Show that the solution to the program  $\min_w L(w, \xi, u)$  satisfies the constraint  $w = -Au$ .

Hint: Compute  $\frac{\partial L(w, \xi, u)}{\partial w}$ .

Your answer:

- (g) [2 points] Show that solving  $\min_{\xi} L(\xi, u)$  results in the constraint  $\|u\|_1 = 1$ .

Your answer:

- (h) [2 points] Using parts (e)-(g), show that the dual to Eq. (12) is:

$$\begin{aligned} \max_u \quad & -\frac{1}{2}u^T A^T A u + u^T b \\ \text{s.t.} \quad & u \geq 0, \|u\|_1 = 1 \end{aligned} \tag{13}$$

Your answer:

11. [13 points] Support Vector Machines

- (a) Give a high level explanation of the binary support vector machine in a few words. Use a diagram to help illustrate your explanation.

Your answer:

- (b) Consider the following unconstrained program for a binary SVM

$$\min_{\mathbf{w}} \|\mathbf{w}\|_2^2 + \sum_i \max\{0, 1 - y^{(i)} \mathbf{w}^T \phi(x^{(i)})\} \quad (14)$$

One possible approach to tackle this problem is via gradient descent. However, a term in the cost function is not differentiable. Which is the offending term and how can we overcome this?

Your answer:

- (c) Another approach to addressing the program given in Eq. 14 is the use of the dual problem. Rewrite the program given in Eq. 14 as an equivalent constrained program, next obtain the dual objective. Carefully define the Lagrangian and the KKT conditions as well as their solution. Note that we are looking for the dual to the program in Eq. 14 and not a more general form.

Your answer:

- (d) What observation can you make about the dual and why is this useful?

Your answer:

12. [15 points] Binary and Multi-class Support Vector Machine for Object Detection

- (a) Explain the intuition behind a binary support vector machine in a few words and draw an illustration that underlines your explanation.

Your answer:

- (b) The program for a binary SVM is as follows when given a dataset  $\mathcal{D} = \{(x_i, y_i)\}$  containing pairs of input data  $x_i \in \mathbb{R}^d$  and corresponding label  $y_i \in \{0, 1\}$ :

$$\min_{w \in \mathbb{R}^d, \xi_i \geq 0} \frac{C}{2} \|w\|_2^2 + \sum_{(x_i, y_i) \in \mathcal{D}} \xi_i \quad \text{s.t.} \quad y_i w^\top x_i \geq 1 - \xi_i \quad \forall (x_i, y_i) \in \mathcal{D} \quad (15)$$

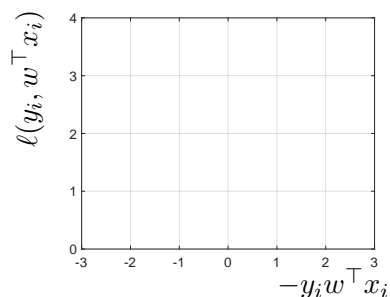
Transform this program into an unconstrained program (not using Lagrangian).

Your answer:

- (c) Relate your program to the task of empirical risk minimization which is given in the following:

$$\min_{w \in \mathbb{R}^d} R(w) + \sum_{(x_i, y_i) \in \mathcal{D}} \ell(y_i, w^\top x_i)$$

and complete the enclosed plot while paying attention to the given axis labels.



Your answer:

- (d) Provide intuition and the program formulation for the extension of the task given in Eq. 15 to multiple classes, i.e., for the case of  $y_i \in \{0, 1, \dots, K - 1\}$ . Keep your formulation similar to the form given in Eq. 15.

Your answer:



13. [11 points] L2 SVM

We are given a dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{|\mathcal{D}|}$  with feature vectors  $\mathbf{x}^{(i)} \in \mathbb{R}^d$  and their corresponding labels  $y^{(i)} \in \{-1, 1\}$ . The following is the primal formulation of L2 SVM, a variant of the standard SVM obtained by **squaring** the hinge loss,

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\xi}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{2} \sum_{i=1}^{|\mathcal{D}|} (\xi^{(i)})^2 \\ \text{s.t.} \quad & y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)}, \quad i \in \{1, \dots, |\mathcal{D}|\} \\ & \xi^{(i)} \geq 0, \quad i \in \{1, \dots, |\mathcal{D}|\} \end{aligned} \tag{16}$$

Here,  $\mathbf{w} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  are learnable weights and  $\boldsymbol{\xi} = (\xi^{(1)}, \dots, \xi^{(|\mathcal{D}|)}) \in \mathbb{R}^{|\mathcal{D}|}$  are slack variables. We will first show that removing the last set of constraints  $\boldsymbol{\xi} \geq 0$  does not change the optimal solution of the problem, *i.e.*, we will show that for the optimal solution  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*)$ , the inequality  $\boldsymbol{\xi}^* \geq 0$  always holds.

- (a) Assume that the dataset consists of 3 samples, that  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*)$  is the optimal solution to the problem without the last set of constraints, *i.e.*,  $\boldsymbol{\xi} \geq 0$ , and that  $(\xi^{(3)})^* = 0$ . Write down the resulting expression of the loss as a function of  $\mathbf{w}^*$ ,  $(\xi^{(1)})^*$  and  $(\xi^{(2)})^*$ .

Your answer:

- (b) Assume that the dataset consists of 3 samples, that  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*)$  is the optimal solution to the problem without the last set of constraints, *i.e.*,  $\boldsymbol{\xi} \geq 0$ , and that  $(\xi^{(3)})^* < 0$ . Write-down the expression of the resulting loss and compare it with the one from **part a**. Which of the losses has a larger value?

Your answer:

- (c) Consider the general case where  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*)$  is the optimal solution to the problem without the constraints  $\boldsymbol{\xi} \geq 0$ . Suppose, there exists some  $(\xi^{(j)})^* < 0$ . Show that  $(\hat{\mathbf{w}}, \hat{b}, \hat{\boldsymbol{\xi}})$  with  $\hat{\mathbf{w}} = \mathbf{w}^*$ ,  $\hat{b} = b^*$ ,  $\hat{\xi}^{(i)} = (\xi^{(i)})^*$  ( $\forall i \neq j$ ) and  $\hat{\xi}^{(j)} = 0$ , is a feasible solution.

Your answer:

- (d) Compare the losses obtained for  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*)$  and  $(\hat{\mathbf{w}}, \hat{b}, \hat{\boldsymbol{\xi}})$ . Conclude that the optimal solution does not change when removing the constraints  $\boldsymbol{\xi} \geq 0$ .

Your answer:

- (e) After removing the constraints  $\boldsymbol{\xi} \geq 0$ , we get a simpler program

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\xi}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{2} \sum_{i=1}^{|\mathcal{D}|} (\xi^{(i)})^2 \\ \text{s.t.} \quad & y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)}, \quad i \in \{1, \dots, |\mathcal{D}|\}. \end{aligned} \tag{17}$$

Give the Lagrangian of the program above as a function of  $\mathbf{w}$ ,  $b$ ,  $\xi^{(i)}$ ,  $y^{(i)}$ ,  $\mathbf{x}^{(i)}$ ,  $C$ ,  $\mathcal{D}$  and Lagrange multipliers  $\alpha^{(i)}$ . What's the range of the Lagrange multipliers  $\alpha^{(i)}$ ?

Your answer:

(f) Show that the dual of the program in Eq. (17) is

$$\begin{aligned} \max_{\boldsymbol{\alpha} \geq 0} \quad & -\frac{1}{2} \boldsymbol{\alpha}^T (\mathbf{Q} + \mathbf{I}_{\frac{1}{C}}) \boldsymbol{\alpha} + \mathbf{1}^T \boldsymbol{\alpha} \\ \text{s.t.} \quad & \mathbf{y}^T \boldsymbol{\alpha} = 0 \end{aligned}$$

with  $Q_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ ,  $\boldsymbol{\alpha} \in \mathbb{R}^{|\mathcal{D}|}$  is a vector of Lagrange multipliers and  $\mathbf{y} \in \mathbb{R}^{|\mathcal{D}|}$  a vector of labels.

Your answer:

14. [15 points] Support Vector Machines for Function Estimation.

Suppose we are given a dataset  $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots\}$ , where  $\mathcal{D} \subset \mathbb{R}^d \times \mathbb{R}$ . We want to find a model  $f(\mathbf{x}^{(i)}) = \mathbf{w}^\top \mathbf{x}^{(i)}$  such that the deviation from the corresponding  $y^{(i)}$  is at most  $\epsilon$ , where  $\epsilon > 0$ , while maintaining a small  $\|\mathbf{w}\|_2^2$ .

We can write this goal as the following optimization problem:

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \quad (18)$$

$$\text{subject to} \quad \mathbf{w}^\top \mathbf{x}^{(i)} - y^{(i)} \leq \epsilon, \quad i = 1, \dots, |\mathcal{D}| \quad (19)$$

$$y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} \leq \epsilon, \quad i = 1, \dots, |\mathcal{D}|. \quad (20)$$

- (a) (2 points) Consider the following toy-dataset to build intuition, where both  $x^{(i)}$  and  $y^{(i)}$  are in  $\mathbb{R}$ . Find the optimal  $w^*$  for the optimization problem given in Eq. (18)-(20), where  $\epsilon = 2$ . Explain your reasoning. **(Only use the toy-dataset for this subproblem.)**

$i$	$x^{(i)}$	$y^{(i)}$
1	0	2
2	2	0

Table 3: Toy-dataset  $\mathcal{D}$

Your answer:

- (b) (1 point) The above optimization problem given in Eq. (18) – Eq. (20) may not always be feasible. Consider the following optimization problem:

$$\underset{\mathbf{w}, \xi, \hat{\xi}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i^{|\mathcal{D}|} (\xi^{(i)} + \hat{\xi}^{(i)})$$

$$\text{subject to} \quad \mathbf{w}^\top \mathbf{x}^{(i)} - y^{(i)} \leq \epsilon + \xi^{(i)}, \quad i = 1, \dots, |\mathcal{D}| \quad (21)$$

$$y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} \leq \epsilon + \hat{\xi}^{(i)}, \quad i = 1, \dots, |\mathcal{D}|$$

$$\xi^{(i)}, \hat{\xi}^{(i)} \geq 0, \quad i = 1, \dots, |\mathcal{D}|.$$

What are  $\xi^{(i)}$  and  $\hat{\xi}^{(i)}$  called?

Your answer:

- (c) (2 points) Explain why the optimization problem given in Eq. (21) handles the infeasible constraints, and why  $\xi^{(i)}$  and  $\hat{\xi}^{(i)}$  has to be greater than 0.

Your answer:

- (d) (1 point) Let the optimization problem given in Eq. (21) be the primal problem. What is the Lagrange function  $\mathcal{L}$  for this primal problem? For the Lagrange multipliers use the symbols  $\alpha^{(i)}, \hat{\alpha}^{(i)}, \mu^{(i)}, \hat{\mu}^{(i)}$ , specifically (don't leave out any Lagrange multipliers, i.e., treat all constraints explicitly):

- Use  $\alpha^{(i)}$  for constraints involving  $\mathbf{w}$  and  $\xi^{(i)}$ .
- Use  $\hat{\alpha}^{(i)}$  for constraints involving  $\mathbf{w}$  and  $\hat{\xi}^{(i)}$ .
- Use  $\mu^{(i)}$  for constraints involving only  $\xi^{(i)}$ .
- Use  $\hat{\mu}^{(i)}$  for constraints involving only  $\hat{\xi}^{(i)}$ .

Your answer:

- (e) (7 points) Derive the dual optimization problem, and simplify (*i.e.*, the final optimization problem should only include  $\alpha^{(i)}$ ,  $\hat{\alpha}^{(i)}$ ,  $\mathbf{x}^{(i)}$ ,  $y^{(i)}$ ,  $\epsilon$ ). (You must show all your work.)

Your answer:

- (f) (2 points) Explain how this dual optimization problem can be extended to non-linear features. Rewrite the optimization problem to involve kernel  $\mathcal{K}$ .

Your answer:

15. [10 points] Support Vector Machine

- (a) Recall, a hard-margin support vector machine in the primal form optimizes the following program

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1, \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \quad (22)$$

What is the Lagrangian,  $L(\mathbf{w}, b, \alpha)$ , of the constrained optimization problem in Eq. (22)?

Your answer:

- (b) Consider the Lagrangian

$$L(\mathbf{w}, \alpha) := \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \alpha^{(i)} (1 - y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)}) \quad (23)$$

where  $\alpha^{(i)}$  are elements of  $\alpha$ . *Note: This Lagrangian is not the same as the solution in the previous part.*

**Derive** the dual program for the Lagrangian given in Eq. (23). Provide all its constraints if any.

Your answer:

(c) Recall that a kernel SVM optimizes the following program

$$\max_{\alpha} \sum_{i=1}^{|\mathcal{D}|} \alpha^{(i)} - \frac{1}{2} \sum_{i,j=1}^{|\mathcal{D}|} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \quad (24)$$

s.t.  $\alpha^{(i)} \geq 0$  and  $\sum_i \alpha^{(i)} y^{(i)} = 0$

We have chosen the kernel to be

$$\kappa(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z})^2 + 1$$

Consider the following dataset  $\mathcal{D}_2$  in the one-dimensional space;  $x^{(i)}, y^{(i)} \in \mathbb{R}$ .

$i$	$x^{(i)}$	$y^{(i)}$
1	$\frac{1}{2}$	+1
2	-1	+1
3	$\sqrt{3}$	-1
4	4	-1

What are the optimal primal parameters,  $\mathbf{w}^*$ ,  $b^*$  when optimizing the program in Eq. (24) on the dataset  $\mathcal{D}_2$ . Note:  $b$  is NOT included in the margin or the features (treat it explicitly).

**Hint:** First, construct a feature vector  $\phi(\mathbf{x}) \in \mathbb{R}^2$  such that  $\kappa(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z})$  for the given one dimensional dataset. Then use this feature vector to transform the data  $\mathcal{D}_2$  into feature space and plot the result. Read of the bias term  $b$  and the optimal weight vector  $\mathbf{w}^*$ .

Your answer:

(d) (Continuing from previous part) Which of the points in  $\mathcal{D}_2$  are support vectors? What are  $\alpha^{(1)}$  and  $\alpha^{(2)}$ ?

**Hint:** To find  $\alpha^{(2)}$  make use of the relationship between the primal solution and the dual variables, i.e.,  $\mathbf{w}^* = \sum_{i=1}^N \alpha^{(i)} y^{(i)} \phi(\mathbf{x}^{(i)})$ . Assume  $\mathbf{w}^* = [-4 \ 0]^T$  if you couldn't solve part (c).

Your answer:



16. [7 points] Multiclass Classification

Consider the objective function of a multiclass SVM given by

$$\min_{w, \xi^{(i)} \geq 0} \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \xi^{(i)}$$

$$\text{s.t. } w_{y^{(i)}}^\top \phi(x^{(i)}) - w_{\hat{y}}^\top \phi(x^{(i)}) \geq 1 - \xi^{(i)} \quad \forall i \in \{1, \dots, n\}; \hat{y} \in \{0, \dots, K-1\}$$

where  $w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{K-1} \end{bmatrix}$ .

- (a) What's the optimal value of  $\xi^{(i)}$ , given  $\phi(x^{(i)})$ ,  $y^{(i)}$ , and  $w$ ?

Your answer:

- (b) Rewrite the objective function in unconstrained form, using the optimal value of  $\xi^{(i)}$ .

Your answer:

- (c) Briefly explain using English language the reason for using  $w_{y^{(i)}}^\top \phi(x^{(i)}) - w_{\hat{y}}^\top \phi(x^{(i)}) \geq 1 - \xi^{(i)}$  in a multiclass SVM formulation, *i.e.*, what does this constraint encourage?

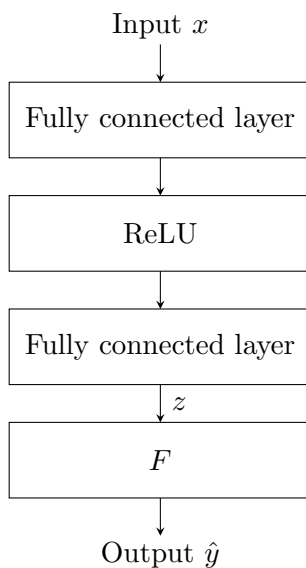
Your answer:

- (d) Suppose we want to train a set of one-vs-rest classifiers and a set of one-vs-one classifiers on a dataset of 5,000 samples and 10 classes, each class having 500 samples. Suppose the running time of the underlying binary classifier we use is  $n^2$  in nanoseconds, where  $n$  is the size of the training dataset. Which one is faster, training of the one-vs-rest classifiers or training of the one-vs-one classifiers? Explain your reason.

Your answer:

17. [6 points] Multiclass Classification via Neural Networks

Suppose we use a multi-layer neural network to classify any input image into one of the following three classes: *apple*, *pear* & *orange*. The neural network architecture is summarized in the following figure:



The output  $\hat{y} = F(z)$  is a three-dimensional vector  $(\Pr(\textit{apple}|x), \Pr(\textit{pear}|x), \Pr(\textit{orange}|x))$ , where  $\Pr(c|x)$  denotes the probability of  $x$  being in class  $c$ .

- (a) (1 point) Which function should be used as activation function  $F$ , for multiclass classification? (a) logistic (b) softmax (c) ReLU (d) sigmoid. Suppose the input to  $F$  is  $z = (z_1, z_2, z_3)$ , write down the expression of  $F(z)$ . (Hint:  $F(z)$  should sum to 1.)

Your answer:

- (b) (2 points) Suppose we have an alternative activation function  $G$ :

$$G(z) = \left( \frac{z_1}{z_1 + z_2 + z_3}, \frac{z_2}{z_1 + z_2 + z_3}, \frac{z_3}{z_1 + z_2 + z_3} \right)$$

which normalizes vector  $z$  to sum to 1. Consider the following two inputs  $z^{(1)} = (0.01, 0.01, 0.02)$  and  $z^{(2)} = (1.01, 1.01, 1.02)$ . Use the given inputs to answer whether  $F$  and  $G$  are *translation invariant*, i.e. the value of the function does **not** change when we add a constant to all its inputs  $z_i$ . Use this fact to give an advantage of using  $F$  over  $G$ . (Hint: you do not need to exactly evaluate the expressions.)

Your answer:

- (c) (2 points) Suppose for an input image, the second fully-connected layer outputs  $z = (1, 10^{-5}, 10^{-5})$ , which means it is very confident that the image is *apple*, while the true label  $y = \textit{orange}$ . Considering this input, give another advantage of using  $F$  over  $G$ , by evaluating (1) the cross entropy between the true label and classifier prediction  $\text{CE}(y, F(z))$ ,  $\text{CE}(y, G(z))$  and (2) their derivatives w.r.t.  $z_3$ , where  $z = (z_1, z_2, z_3)$ .

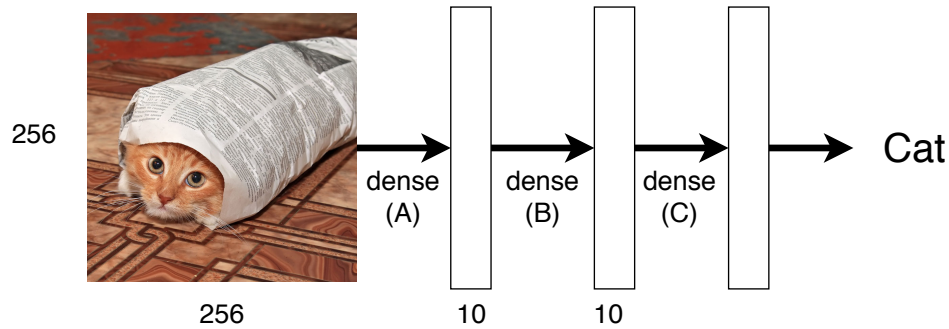
Your answer:

- (d) (1 point) As an alternative to a multiclass neural network, we can use one-vs-rest multiclass classification, which fits a single-output neural network for each class. Suppose we use the same number of hidden units in the two approaches. Which approach will have faster prediction (output  $\hat{y}$  given  $x$ ) speed? Explain your reason.

Your answer:

18. [5 points] Modeling of Deep Neural Networks

We are interested in building a deep net for image classification. The input is a  $256 \times 256$  RGB image, and there are 7 different possible output classes. The model architecture is shown below:



- (a) How many trainable parameters are in layer “dense (A)”?

Your answer:

- (b) How many trainable parameters are in layer “dense (B)”?

Your answer:

- (c) How many trainable parameters are in layer “dense (C)”?

Your answer:

- (d) We redesigned the network by replacing dense (A) and dense (B) layers with a convolution layer of four  $3 \times 3$  filters, stride 1 and no padding. What are the dimensions for the output of this convolution layer?

Your answer:

19. [9 points] Representation of Deep Neural Networks

We will use the XOR dataset,  $\mathcal{D}$ , shown in Table 4. Each example  $\mathbf{x}^{(i)} \in \mathbb{R}^2$  and the label  $y \in \{0, 1\}$ . We are interested in designing classification models for this dataset.

$i$	$\mathbf{x}_0^{(i)}$	$\mathbf{x}_1^{(i)}$	$y$
0	0	0	0
1	0	1	1
2	1	0	1
3	1	1	0

Table 4: The XOR Dataset  $\mathcal{D}$ .

(a) Consider a model with the following parameterization:

$$p(y^{(i)}|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x}^{(i)} - b)}, \quad (25)$$

where  $\mathbf{w} \in \mathbb{R}^2$  and  $b \in \mathbb{R}$ .

What is the highest accuracy for this model on the XOR dataset? **Note:** To compute accuracy, we use a threshold of 0.5, *i.e.*, the final prediction of the model is  $\delta[p(y^{(i)}|\mathbf{x}) > 0.5]$ , where  $\delta$  denotes the indicator function.

Your answer:

(b) Consider another model with the parametrization shown below:

$$\tilde{y}^{(i)} = \frac{1}{1 + \exp(-a_2^{(i)})} \quad (26)$$

$$a_2^{(i)} = \theta^\top \max(\mathbf{a}_1^{(i)}, 0) + b \quad (27)$$

$$\mathbf{a}_1^{(i)} = \mathbf{W}\mathbf{x}^{(i)} + \mathbf{c} \quad (28)$$

where  $\theta \in \mathbb{R}^2$ ,  $b \in \mathbb{R}$ ,  $\mathbf{W} \in \mathbb{R}^{2 \times 2}$  and  $\mathbf{c} \in \mathbb{R}^2$ .

Find a  $\theta$  and  $b$  that achieve 100 % accuracy on the XOR dataset, given  $\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $\mathbf{c} = [0, -1]^\top$ . To show your work, write out  $\mathbf{a}_1^{(i)}$  and  $a_2^{(i)}$  for the four datapoints in the XOR dataset and your choice of  $\theta$ , and  $b$ .

**Note:** To compute accuracy, we use a threshold of 0.5, *i.e.*, the final prediction of the model is  $\delta[\tilde{y}^{(i)} > 0.5]$ , where  $\delta$  denotes the indicator function.

Your answer:

- (c) To learn the parameters of the model in (b), the learning problem is formulated as the following program:

$$\min_{\theta, b, \mathbf{W}, \mathbf{c}} \mathcal{L} := \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \|y^{(i)} - \tilde{y}^{(i)}\|_2^2. \quad (29)$$

Write out  $\frac{\partial \mathcal{L}}{\partial \tilde{y}^{(i)}}$ . You may use the terms  $y^{(i)}$  and  $\tilde{y}^{(i)}$ .

Your answer:

- (d) Write out  $\frac{\partial \tilde{y}^{(i)}}{\partial a_2^{(i)}}$ . You may use the terms  $a_2^{(i)}$ .

Your answer:

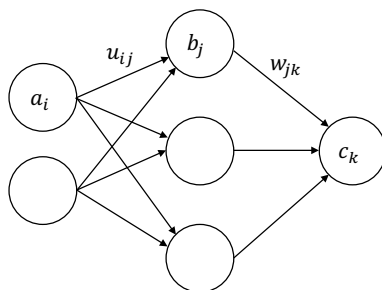
- (e) Write out  $\frac{\partial \mathcal{L}}{\partial \mathbf{c}}$ . You may use the terms  $\frac{\partial \mathcal{L}}{\partial \tilde{y}^{(i)}}$ ,  $\frac{\partial \tilde{y}^{(i)}}{\partial a_2}$ ,  $\mathbf{W}$ ,  $\mathbf{x}^{(i)}$ ,  $\mathbf{c}$ ,  $\theta$ ,  $b$ ,  $\mathbf{a}_1^{(i)}$  and  $\delta[\cdot]$  as the indicator function.

Your answer:



20. [8 points] Backpropagation

Consider the deep net in the figure below consisting of an input layer, an output layer, and a hidden layer. The feed-forward computations performed by the deep net are as follows: every input  $a_i$  is multiplied by a set of fully-connected weights  $u_{ij}$  connecting the input layer to the hidden layer. The resulting weighted signals are then summed and combined with a bias  $e_j$ . This results in the activation signal  $z_j = e_j + \sum_i a_i u_{ij}$ . The hidden layer applies activation function  $g$  on  $z_j$  resulting in the signal  $b_j$ . In a similar fashion, the hidden layer activation signals  $b_j$  are multiplied by the weights connecting the hidden layer to the output layer  $w_{jk}$ , a bias  $f_k$  is added and the resulting signal is transformed by the output activation function  $g$  to form the network output  $c_k$ . The loss between the desired target  $t_k$  and the output  $c_k$  is given by the MSE:  $E = \frac{1}{2} \sum_k (c_k - t_k)^2$ , where  $t_k$  denotes the ground truth signal corresponding to  $c_k$ . Training a neural network involves determining the set of parameters  $\theta = \{U, W, e, f\}$  that minimize  $E$ . This problem can be solved using gradient descent, which requires determining  $\frac{\partial E}{\partial \theta}$  for all  $\theta$  in the model.



- (a) For  $g(x) = \sigma(x) = \frac{1}{1+e^{-x}}$ , compute the derivative  $g'(x)$  of  $g(x)$  as a function of  $\sigma(x)$ .

Your answer:

- (b) Compute  $\frac{\partial E}{\partial w_{jk}}$ . Use  $c_k, t_k, g', f_k, b_j, w_{jk}$ . Note the use of  $g'$  to simplify the expression.

Your answer:

- (c) Compute  $\frac{\partial E}{\partial f_k}$ . Use  $c_k, t_k, g', f_k, b_j, w_{jk}$ . Note the use of  $g'$  to simplify the expression.

Your answer:

- (d) Compute  $\frac{\partial E}{\partial u_{ij}}$ . Use  $c_k, t_k, g', f_k, b_j, w_{jk}, a_i$ . Note the use of  $g'$  to simplify the expression.

Your answer:

- (e) Compute  $\frac{\partial E}{\partial e_j}$ . Use  $c_k, t_k, g', f_k, b_j, w_{jk}, a_i$ . Note the use of  $g'$  to simplify the expression.

Your answer:

21. [13 points] Deep Nets

- (a) (2 points) You are working on an automated check-out kiosk for a supermarket, and are building a classifier for apples, bananas and oranges. Suppose your classifier obtains a training set error of 0.5%, and a test set error of 7%. Name two promising things to try to improve your classifier?

Your answer:

- (b) (3 points) Suppose gradient descent in a deep network is taking excessively long to find a value of the parameters that achieves a small value for the cost function. Name three techniques that could help improve the convergence.

Your answer:

- (c) (1 point) Does it make sense to initialize all weights in a deep network to 0?

Your answer:

- (d) (1 point) State one advantage of linear rectified activation compared to logistic sigmoid activation.

Your answer:

- (e) (2 points) State two roles of pooling layers in a CNN?

Your answer:

- (f) (1 point) A convolutional neural network has 4 consecutive  $3 \times 3$  convolutional layers with stride 1 and no pooling. How large is the support of (the set of image pixels which activate/influence) a neuron in the 3rd non-image layer of this network?

Your answer:

- (g) (1 point) What is the output image size resulting from applying three  $5 \times 5 \times 3$  filters to a  $32 \times 32 \times 3$  input image (no padding, stride is 1)?

Your answer:

- (h) (2 points) What are the padding and stride sizes that produce an output size of  $32 \times 32 \times 3$  given a filter size of  $5 \times 5 \times 3$  and input dimensions of  $32 \times 32 \times 3$ ?

Your answer:

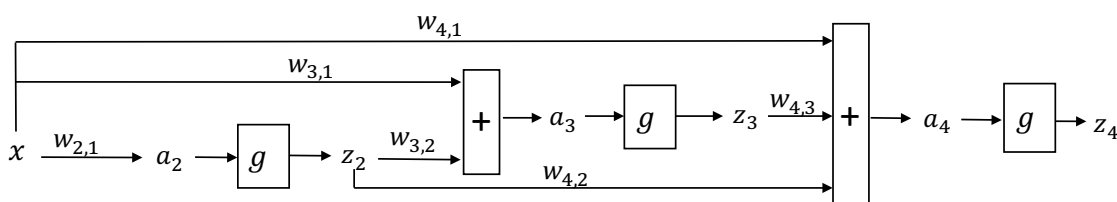
22. [13 points] Backpropagation

Consider the neural network given in the figure below. The network has a scalar input variable  $x \in \mathbb{R}$  and a scalar target  $t \in \mathbb{R}$  and is defined as follows:

$$z_j = \begin{cases} x, & \text{if } j = 1 \\ g(a_j) & \text{if } j \in \{2, 3, 4\} \text{ with } a_j = \sum_{i=1}^{j-1} w_{j,i} z_i \end{cases} \quad (30)$$

Suppose that the network is trained to minimize the L2 loss per sample, *i.e.*,  $E = \frac{1}{2}(z_4 - t)^2$ . The error gradient can be written as:

$$\frac{\partial E}{\partial w_{j,i}} = \delta_j z_i \quad (31)$$



- (a) [2 pts] For  $g(x) = \sigma(x) = \frac{1}{1+e^{-x}}$ , compute the derivative  $g'(x)$  of  $g(x)$  as a function of  $\sigma(x)$ .

Your answer:

- (b) [2 pts] Compute  $\delta_4$  as a function of  $z_4$ ,  $t$  and  $g'(a_4)$ .

Your answer:

- (c) [2 pts] Compute  $\delta_3$  as a function of  $\delta_4$ ,  $w_{4,3}$  and  $g'(a_3)$ .

Your answer:

- (d) [3 pts] Compute  $\delta_2$  as a function of  $\delta_3$ ,  $\delta_4$ ,  $w_{3,2}$ ,  $w_{4,2}$  and  $g'(a_2)$ .

Your answer:

- (e) [4 pts] Write down a recursive formula for computing  $\delta_j$  for  $j \in \{2, \dots, M-1\}$ , as a function of  $\delta_k$ ,  $w_{k,j}$  and  $g'(a_j)$  for  $k \in \{j+1, \dots, M\}$ .

Your answer:

23. [15 points] Inference in Discrete Markov Random Fields

- (a) Inference in Markov random fields amounts to finding the highest scoring configuration for a set of variables. Suppose we have two variables  $x_1 \in \{0, 1\}$  and  $x_2 \in \{0, 1\}$  and their local evidence functions  $\theta_1(x_1)$  and  $\theta_2(x_2)$  as well as a pairwise function  $\theta_{1,2}(x_1, x_2)$ . Using this setup, inference solves  $\arg \max_{x_1, x_2} \theta_1(x_1) + \theta_2(x_2) + \theta_{1,2}(x_1, x_2)$ . Using

$$\theta_1(x_1) = \begin{cases} 1 & \text{if } x_1 = 0 \\ 2 & \text{otherwise} \end{cases} \quad \theta_2(x_2) = \begin{cases} 1 & \text{if } x_2 = 0 \\ 2 & \text{otherwise} \end{cases} \quad \theta_{1,2}(x_1, x_2) = \begin{cases} -1 & \text{otherwise} \\ 2 & \text{if } x_1 = 0 \text{ \& } x_2 = 1 \end{cases}$$

what is the integer linear programming formulation of the inference task. Make cost function and constraints explicit for the given problem, i.e., no general formulation.

Your answer:

- (b) What is the solution (value and argument) to the program in part (a).

Your answer:

- (c) Why do we typically not use the integer linear program for reasonably sized MRFs and what other methods do you know to approximately solve MRFs. Name at least three other methods.

Your answer:

24. [10 points] Inference in Discrete Markov Random Fields

- (a) Inference in Markov random fields amounts to finding the highest scoring configuration for a set of variables. Suppose we have two variables  $x_1 \in \{0, 1\}$  and  $x_2 \in \{0, 1\}$  and their local evidence functions  $\theta_1(x_1)$  and  $\theta_2(x_2)$  as well as pairwise function  $\theta_{1,2}(x_1, x_2)$ . Using this setup, inference solves  $\arg \max_{x_1, x_2} \theta_1(x_1) + \theta_2(x_2) + \theta_{1,2}(x_1, x_2)$ . Using

$$\theta_1(x_1) = \begin{cases} -1 & \text{if } x_1 = 0 \\ 1 & \text{otherwise} \end{cases} \quad \theta_2(x_2) = \begin{cases} -1 & \text{if } x_2 = 0 \\ 1 & \text{otherwise} \end{cases}$$

$$\theta_{1,2}(x_1, x_2) = \begin{cases} 2 & \text{if } x_1 = 1 \text{ \& } x_2 = 0 \\ 1 & \text{if } x_1 = 0 \text{ \& } x_2 = 1 \\ -1 & \text{otherwise} \end{cases}$$

what is the integer linear programming (ILP) formulation of the inference task? Make cost function and constraints explicit for the given problem, i.e., do not use a general formulation.

Your answer:

- (b) If the two variables instead took on values  $x_1, x_2 \in \{0, 1, 2, 3\}$ , how many constraints would the integer linear program have?

Your answer:

- (c) Let's say we wanted to use a different method to solve this inference problem. Can we use a dynamic programming method? Why or why not?

Your answer:



- (d) Name two other inference methods that may be more efficient than ILP, and name one advantage and one disadvantage for each.

Your answer:

25. **[20 points]** Structured Prediction

We define a distribution  $p(x_1, x_2)$  over two discrete random variables  $x_1 \in \{0, 1\}$  and  $x_2 \in \{0, 1\}$  as:

$$p(x_1, x_2) = \frac{1}{Z} \exp(-E(x_1, x_2)), \quad (32)$$

where  $E : \{0, 1\}^2 \rightarrow \mathbb{R}$ .

- (a) [3 points] What is  $Z$  called? Give the expression of  $Z$  as a function of  $E(x_1, x_2)$ .

Your answer:

- (b) [1 point] We are interested in maximizing the likelihood  $p(x_1, x_2)$  with respect to  $x_1$  and  $x_2$ . Justify that this is equivalent to solving the program:

$$\min_{x_1 \in \{0, 1\}, x_2 \in \{0, 1\}} E(x_1, x_2). \quad (33)$$

Your answer:

(c) [8 points]  $E(x_1, x_2)$  is an energy function modeled as:

$$E(x_1, x_2) = -(\theta_1(x_1) + \theta_2(x_2) + \theta_{1,2}(x_1, x_2)). \quad (34)$$

The potentials  $\theta_1(x_1)$ ,  $\theta_2(x_2)$  and  $\theta_{1,2}(x_1, x_2)$  are given as follows:

$$\theta_1(x_1) = \begin{cases} 1 & \text{if } x_1 = 0 \\ 2 & \text{otherwise.} \end{cases} \quad (35)$$

$$\theta_2(x_2) = \begin{cases} 1 & \text{if } x_2 = 0 \\ 2 & \text{otherwise.} \end{cases} \quad (36)$$

$$\theta_{1,2}(x_1, x_2) = \begin{cases} -2 & \text{otherwise} \\ 1 & \text{if } x_1 = 0 \text{ and } x_2 = 1. \end{cases} \quad (37)$$

What is the Integer Linear Programming formulation of the inference task. Make cost function and constraints explicit for the given problem, *i.e.*, no general formulation.

Your answer:

- (d) [3 points] What is the solution (value and argument) to the program in Eq. (33). Give the values of  $x_1^*$ ,  $x_2^*$  and  $E(x_1^*, x_2^*)$ .

Your answer:

- (e) [1 point] Name a disadvantage of using an Integer Linear Programming solver.

Your answer:

- (f) [2 points] Is  $E(x_1, x_2)$  sub-modular? Justify.

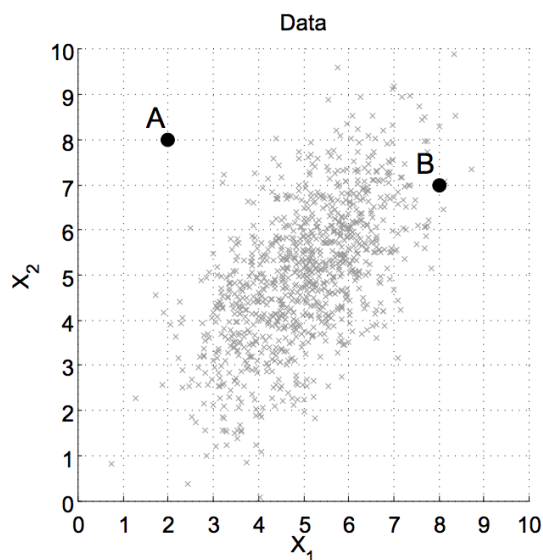
Your answer:

- (g) [2 points] Draw the Markov Random Field associated with the energy  $E(x_1, x_2)$ . Can we use a dynamic programming solver for the program above? Justify.

Your answer:

26. [7 points] Principle Component Analysis (PCA)

Plotted in the figure below is a two dimensional data set drawn from a multivariate Normal distribution.



- (a) [2 points] What is the mean of this distribution? Estimate the answer visually and round to the nearest integer.

Your answer:

$$\mathbb{E}[X_1] = \mu_1 = \dots$$

$$\mathbb{E}[X_2] = \mu_2 = \dots$$

- (b) [1 points] Circle the right answer. What would the off-diagonal co-variance  $cov(X_1, X_2)$  be?

Your answer:

- (a) negative
- (b) positive
- (c) approximately zero

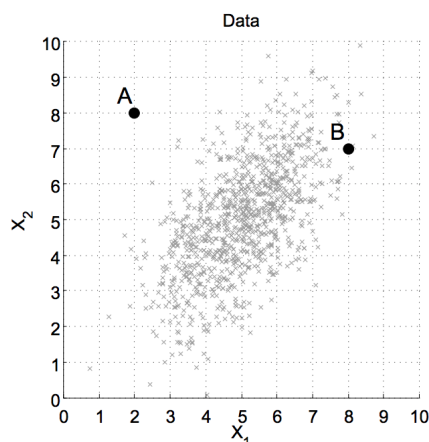
- (c) [2 points] Let  $w_1$  and  $w_2$  be the directions of the first and second principal component. These directions define a change of basis,

$$Z_1 = (X - \mu)^T w_1, \quad (38)$$

$$Z_2 = (X - \mu)^T w_2. \quad (39)$$

Sketch and label  $w_1$  and  $w_2$  on the following figure. The arrows should originate from the mean of the distribution. You do not need to solve the SVD, instead visually estimate the directions. Present  $w_1$  and  $w_2$  as unit norm vectors.

Your answer:



- (d) [1 points] Circle the right answer. The co-variance  $cov(Z_1, Z_2)$  is:

Your answer:

- (a) negative
- (b) positive
- (c) approximately zero

- (e) [1 points] Which point (A or B) would have the higher reconstruction error after projecting onto the first principal component direction  $w_1$ ?

Your answer:

- (1) Point A
- (2) Point B

27. [11 points] k-means

- (a) (3 points) Fill in the blanks in the following statement using the options provided in the boxes.

k-means is a/an **(a.)** *supervised, unsupervised, semi – supervised* learning algorithm used for **(b.)** *classification, clustering, regression*. It aims to partition  $n$  observations into **(c.)**  $n, d, k$  clusters in which each observation belongs to the cluster with the nearest center that is the **(d.)** *median, mode, mean*. The k-means objective function is related to the **(e.)** *Hamming, Euclidean, Manhattan* distance and the goal is to **(f.)** *maximize, minimize, orthogonalize* the objective function.

Choose the correct option for each blank from the box following it. Write down your answers below.

Your answer:

- (b) (4 points) Briefly describe the k-means algorithm. We want you to just write down the 4-5 lines. What is the run time, for  $n$  samples in  $d$  dimensions?

Your answer:

(c) (4 points) In class, you saw that the k-means objective was defined using the squared L2-norm, i.e.  $\sum_{i=1}^n \min_j \|x_i - \mu_j\|_2^2$ . What are optimal cluster centers with this objective function?

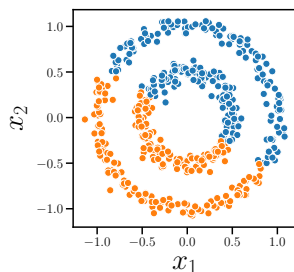
Now if we define a new objective function using the L1-norm instead, i.e.  $\sum_{i=1}^n \min_j \|x_i - \mu_j\|_1$ , what would be the expression for the optimal cluster centers? **Derive** your answer assuming  $x_i - \mu_j \neq 0, \forall i, j$ , i.e. the cluster centers don't overlap with the given input datapoints.

Your answer:

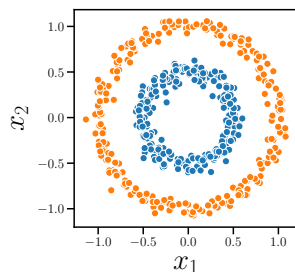


28. [15 points] K-means and mixture models

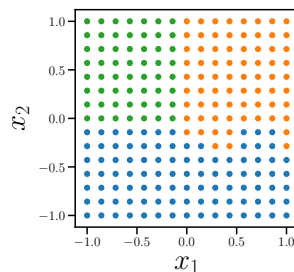
(a) [4 points] Below are four cluster visualizations:



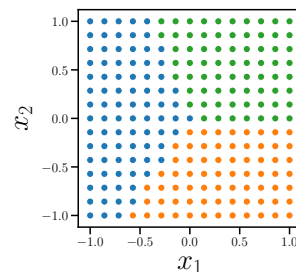
(a)



(b)



(c)



(d)

Please indicate which visualizations are generated from the K-means algorithm in the visualized feature space,  $x = (x_1, x_2)$ , and why.

**Hint:** Our K-means algorithm optimizes the following:

$$\min_{\mu} \min_r \sum_i \sum_k \frac{1}{2} r_{ik} \|x^{(i)} - \mu_k\|_2^2 \quad \text{s.t.} \quad \begin{cases} r_{ik} \in \{0, 1\} & \forall i, k \\ \sum_{k=1}^K r_{ik} = 1 & \forall i \end{cases}.$$

Your answer:

(b) [8 points] Recall, a mixture model is given by

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \cdot p(\mathbf{x}|z = k), \quad (40)$$

where  $z \in \{1, 2, \dots, K\}$ .

Suppose the vector  $\mathbf{x}$  is partitioned into two parts  $\mathbf{x} := [\mathbf{x}_a, \mathbf{x}_b]$ .

Show that  $p(\mathbf{x}_b|\mathbf{x}_a)$  is also a mixture distribution, *i.e.*, the conditional density has the form

$$p(\mathbf{x}_b|\mathbf{x}_a) = \sum_{k=1}^K \hat{\pi}_k \cdot p(\mathbf{x}_b|\mathbf{x}_a, z = k).$$

To this end, provide an expression for  $\hat{\pi}_k$  in terms of  $p(\mathbf{x})$ ,  $p(\mathbf{x}_a|z)$  and  $p(z = k)$ .

**Hint:**  $p(\mathbf{x}_b, z|\mathbf{x}_a)$  may be useful.

Your answer:

(c) [3 points] How are K-means and Gaussian Mixture Model related? (There are three conditions.) Explain them shortly.

Your answer:

29. [16 points] Gaussian Mixture Models & K-Means

Consider a Gaussian mixture model with  $K$  components ( $k \in \{1, \dots, K\}$ ), each having mean  $\mu_k$ , variance  $\sigma_k^2$ , and mixture weight  $\pi_k$ . Further, we are given a dataset  $\mathcal{D} = \{x_i\}$ , where  $x_i \in \mathbb{R}$ . We also use  $z_{ik}$  to denote the latent variables.

- (a) What is the log-likelihood of the data according to the Gaussian Mixture Model? (use  $\mu_k$ ,  $\sigma_k$ ,  $\pi_k$ ,  $K$ ,  $x_i$  and  $\mathcal{D}$ ). Don't use any abbreviations.

Your answer:

- (b) Assume  $K = 1$ , find the maximum likelihood estimate for the parameters  $(\mu_1, \sigma_1^2, \pi_1)$ .

Your answer:

- (c) What is the probability distribution on the latent variables, *i.e.*, what is the distribution  $p(z_{i,1}, z_{i,2}, \dots, z_{i,K})$  underlying Gaussian mixture models. Also give its name.

Your answer:

- (d) For general  $K$ , what is the posterior probability  $p(z_{ik} = 1 | x^{(i)})$ ? To simplify, wherever possible, use  $\mathcal{N}(x_i | \mu_k, \sigma_k)$ , a Gaussian distribution over  $x_i \in \mathbb{R}$  having mean  $\mu_k$  and variance  $\sigma_k^2$ .

Your answer:

- (e) How are kMeans and Gaussian Mixture Model related? (There are three conditions)

Your answer:

- (f) Show that the objective for kMeans and Gaussian Mixture Model are equivalent under the conditions you provided in the previous part (e).

Your answer:

30. [10 points] Generative models

In class, you have studied various generative models including GMMs, VAEs, and GANs. In this problem we will analyze the links between them.

- (a) (1 point) GMMs parameters are learned by maximizing the log-likelihood of the data. Suppose we are given  $N$  samples  $x_i$  ( $i = 1, 2, \dots, N$ ) from a random variable  $X \sim P$  and we want to fit a Gaussian mixture model with  $K$  components. Write the log-likelihood of the data in terms of the mixture parameters  $\theta = (\mu_k, \sigma_k, \pi_k$  for  $k = 1, 2, \dots, K$ ).

Your answer:

- (b) (1 point) Other generative models minimize a distance. E.g. GANs may minimize the KL divergence between a ground truth distribution  $P$  and a generated distribution  $G$ . Write the KL divergence between  $P$  and  $G$ .

Your answer:

- (c) (5 points) Let the distribution represented by the GMM be  $M$ . We find a solution  $\hat{\theta}$  for the parameters as

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^N \log M(x_i) \quad (41)$$

Show that as  $N \rightarrow \infty$ , this formulation is equivalent to minimizing the divergence  $KL(P||M)$ , where  $P$  is the ground truth distribution.

Your answer:

- (d) (3 points) For GMMs, we know the model of the distribution. Therefore, we can sample from it easily. GANs take a different approach where they transform a known distribution into the data distribution. Suppose we have a uniform distribution  $U[0, 1]$ . We want to map this distribution using a function  $f$  to another distribution whose pdf is  $p$  and cdf is  $C$ . Find  $f$  in terms of  $p$  and  $C$ .

Your answer:

31. [7 points] Variational Autoencoder (VAE)

- (a) Show that the KL divergence between two discrete distributions  $p(x)$  and  $q(x)$  defined on the domain  $x \in X$  is non-negative.

Your answer:

- (b) What is the reparameterization trick? Why do we need it in VAEs?

Your answer:

- (c) Given a code snippet of the `_sample_z` function from mp3, which builds the graph to perform the reparameterization trick. Circle the line that is incorrect and explain why it is wrong. You may assume all the shapes and syntax are correct.

```
def _sample_z(self, z_mean, z_log_var):
    """
    Sample z using reparameterization trick.

    Args:
        z_mean (tf.Tensor): The latent mean,
            tensor of dimension (None, 2)
        z_log_var (tf.Tensor): The latent log variance,
            tensor of dimension (None, 2)
    Returns:
        z (tf.Tensor): Random z sampled of dimension (None, 2)
    """
    eps = np.random.randn(self.z_shape[0], self.z_shape[1])
    z = z_mean + (tf.exp(z_log_var / 2) * eps)
    return z
```

Your answer:

32. [15 points] Variational Auto-encoders for Image Generation

- (a) What is the expression for the Kullback-Leibler divergence  $D_{KL}(p, q)$  between two 1-dimensional distributions  $p$  and  $q$  defined over the domain of real numbers  $x$ ?

Your answer:

- (b) What is the expression for a 1-dimensional Gaussian distribution defined over the domain of real numbers  $x$  which has mean  $\mu$  and standard deviation  $\sigma$ , often abbreviated  $\mathcal{N}(x | \mu, \sigma)$ .

Your answer:

- (c) Explain the manifold assumption underlying variational auto-encoders.

Your answer:

- (d) Show the following identity:

$$\int_z q(z|x) \log \frac{p(x, z)}{q(z|x)} dz + \int_z q(z|x) \log \frac{q(z|x)}{p(z|x)} dz = \log p(x)$$

Your answer:

- (e) Given your derivation in the previous part, provide a lower bound on the log-probability  $\log p(x)$  and explain.

Your answer:



33. [14 points] Variational Autoencoders (VAEs)

Suppose we are given a dataset with  $N$  data points  $\{x_i\}_{i=1}^N$  where each of the data points is a  $D$ -dimensional vector. We use VAEs to learn the distribution of the data. Let  $z$  denote the unobserved latent variable. We refer to the approximated posterior  $q_\phi(z|x)$  as the encoder and to the conditional distribution  $p_\theta(x|z)$  as the decoder. Use the above notations to answer the following questions.

- (a) (2 points) The empirical lower bound (ELBO) of  $p_\theta(x_i)$  is

$$\mathcal{L}(\theta, \phi, x_i) = -D_{KL}(q_\phi(z|x_i)||p(z)) + E_{q_\phi(z|x_i)} \left[ \log p_\theta(x_i|z) \right]$$

Write down the minimization program that VAEs solve in terms of  $p_\theta(\cdot)$  and  $q_\phi(\cdot)$  given the dataset  $\mathbf{x} = \{x_i\}_{i=1}^N$ .

Your answer:

- (1 point) Write down the reconstruction error in the loss function of VAEs in your previous answer, again given the dataset  $\mathbf{x} = \{x_i\}_{i=1}^N$ .

Your answer:

- (1 point) Write down the formula to compute the reconstruction error empirically by drawing  $M$  samples from the distribution  $q_\phi(z|x_i)$ . Denote these  $M$  samples as  $z_{i,m}$ , where  $m = 1, 2, \dots, M$ .

Your answer:

- (b) (3 points) Let  $f(z_{i,m}) \in \mathbf{R}^D$  be the reconstructed sample with respect to  $z_{i,m}$ , which is the output of the decoder. What is the empirical reconstruction error if we assume  $p_\theta(x_i|z)$  to be a Gaussian distribution  $\mathcal{N}(f(z), \sigma^2 \mathbf{I})$ , where  $\sigma$  is a constant and  $\mathbf{I}$  is the  $D$ -by- $D$  identity matrix (simplify as much as possible)?

Your answer:

- (c) (4 points) Now consider all the data points  $x_i$  to be binary, i.e.,  $\forall i, x_i \in \{0, 1\}^D$ . If we want to have the empirical reconstruction error to be the cross entropy loss, what should we assume  $p_\theta(x_i|z)$  to be? What is the name of the distribution? Let the output of the decoder be  $g = f(z) \in [0, 1]^D$ , where the values are all between 0 and 1. If you need, use  $x_i^{(d)}$  to denote the  $d$ -th element in the vector  $x_i$ .

Your answer:

- (d) (3 points) The following shows a code snippet of the `_sample_z` function from `mp10`, which builds the graph to perform the reparametrization trick. Is this implementation correct? If not, circle each of the place that you think is incorrect and explain your reason. You may assume all the shapes and syntax are correct.

```
def _sample_z(self, z_mean, z_log_var):
    """
    Sample z using the reparametrization trick.
    Args:
        z_mean (tf.Tensor): The latent mean,
            tensor of dimension (None, 2)
        z_log_var (tf.Tensor): The latent log variance,
            tensor of dimension (None, 2)
    Returns:
        z (tf.Tensor): Random z sampled of dimension (None, 2)
    """
    eps = np.random.randn(self.z_shape[0], self.z_shape[1])
    z = z_mean * tf.sqrt(z_log_var) * eps
    return z
```

Your answer:

34. [20 points] Variational Autoencoders (VAEs)

We use VAEs to learn the distribution of the data,  $x$ . Let  $z$  denote the unobserved latent variable. We refer to the approximated posterior  $q_\phi(z|x)$  as the encoder and to the conditional distribution  $p_\theta(x|z)$  as the decoder. Use these names to answer the following questions.

- (a) [4 points] We are interested in modeling data,  $x \in \{0, 1\}^G$ . Hence, we choose  $p_\theta(x|z)$  to follow  $G$  independent Bernoulli distributions. Recall, a Bernoulli distribution has a probability density function of

$$P(k) = \begin{cases} 1 - p & \text{if } k = 0 \\ p & \text{if } k = 1 \end{cases}.$$

Write down the explicit form for  $p_\theta(x|z)$ . Use  $\hat{y}_j$  to denote the  $j^{\text{th}} \in [1, G]$  dimension of the decoder's output.

Your answer:

- (b) [4 points] We further assume that  $z \in \mathbb{R}^2$  and that  $q_\phi(z|x)$  follows a multi-variate Gaussian distribution with an identity covariance matrix. What is the output dimension of the encoder and why?

Your answer:

- (c) [8 points] Recall, the evidence lower bound (ELBO) of the log likelihood,  $\log p_\theta(x)$ , is

$$\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x)||p(z)). \quad (42)$$

We can also write the ELBO as

$$\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z) + \log p(z) - \log(q_\phi(z|x))]. \quad (43)$$

**Practically**, will training a VAE using the formulation in Eq. (42) be the same as the one in Eq. (43)? If not, why use one formulation over another?

Your answer:

- (d) [4 points] Observe that the ELBO in Eq. (42) works for any  $q_\phi$  distribution. Is it a good idea to choose  $q_\phi(z|x) := \mathcal{N}(\mathbf{0}, \mathbf{I})$ ? In other words, why is an encoder necessary?

Your answer:

35. [8 points] Generative Adversarial Network (GAN)

- (a) What is the key difference between VAE and GAN?

Your answer:

- (b) What is the cost function for classical GANs? Use  $D_w(x)$  as the discriminator and  $G_\theta(x)$  as the generator.

Your answer:

- (c) Assume arbitrary capacity for both discriminator and generator. In this case we refer to the discriminator using  $D(x)$ , and denote the distribution on the data domain induced by the generator via  $p_G(x)$ . State an equivalent problem to the one asked for in part (a), by using  $p_G(x)$ .

Your answer:

- (d) Assume arbitrary capacity, derive the optimal discriminator  $D^*(x)$  in terms of  $p_{data}(x)$  and  $p_G(x)$ .

You may need the Euler-Lagrange equation:

$$\frac{\partial L(x, D, \dot{D})}{\partial D} - \frac{d}{dx} \frac{\partial L(x, D, \dot{D})}{\partial \dot{D}} = 0$$

where  $\dot{D} = \partial D / \partial x$ .

Your answer:

- (e) Assume arbitrary capacity and an optimal discriminator  $D^*(x)$ , show that the optimal generator,  $G^*(x)$ , generates the distribution  $p_G^* = p_{data}$ , where  $p_{data}(x)$  is the data distribution

You may need the Jensen-Shannon divergence:

$$\text{JSD}(p_{\text{data}}, p_G) = \frac{1}{2} D_{KL}(p_{\text{data}}, M) + \frac{1}{2} D_{KL}(p_G, M) \quad \text{with} \quad M = \frac{1}{2}(p_{\text{data}} + p_G)$$

Your answer:

- (f) What is the optimal discriminator  $D^*(x)$ , assuming arbitrary capacity and optimal generator?

Your answer:

36. [9 points] Generative Adversarial Nets (GANs) for Image Generation

- (a) Explain the intuition underlying generative adversarial nets (GANs).

Your answer:

- (b) What is the relation between the cost function for GANs, i.e.,

$$\max_{\theta} \min_w - \sum_x \log D_w(x) - \sum_z \log(1 - D_w(G_{\theta}(z)))$$

and the following program:

$$\max_{p_G} \min_D - \int_x p_{\text{data}}(x) \log D(x) dx - \int_x p_G(x) \log(1 - D(x)) dx \quad (44)$$

Your answer:

- (c) Use the Euler-Lagrange formalism which says that the stationary point of  $S(D) = \int_x L(x, D, \dot{D}) dx$  can be obtained from

$$\frac{\partial L(x, D, \dot{D})}{\partial D} - \frac{d}{dx} \frac{\partial L(x, D, \dot{D})}{\partial \dot{D}} = 0$$

to demonstrate that a stationary point of GANs (use Eq. 44) is obtained for  $p_D = p_G$ . Note that  $\dot{D} = \partial D / \partial x$ . Hints: solve for the optimal discriminator first. You may need the Jensen-Shannon divergence:

$$\text{JSD}(p_{\text{data}}, p_G) = \frac{1}{2} D_{KL}(p_{\text{data}}, M) + \frac{1}{2} D_{KL}(p_G, M) \quad \text{with} \quad M = \frac{1}{2}(p_{\text{data}} + p_G)$$

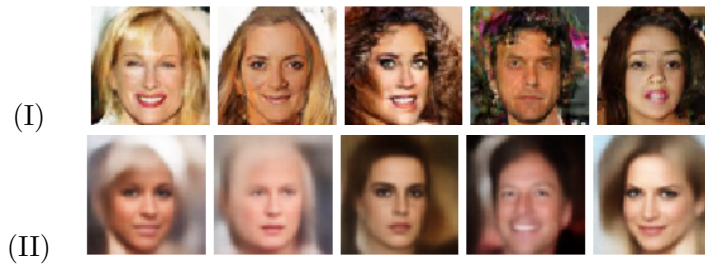
Your answer:

37. [20 points] Generative Adversarial Nets (GANs)

- (a) [2 points] What is the key difference between a VAE and a GAN?

Your answer:

- (b) [2 points] Below are images sampled from a VAE and a GAN trained on the CelebA dataset.



Which row of samples are generated from a GAN and why?

Your answer:



- (c) [4 points] The loss function for a classical GAN is

$$\max_{\theta} \min_w - \sum_{x \in \mathcal{D}} \log D_w(x) - \sum_{z \sim \mathcal{N}(0,1)} \log(1 - D_w(G_{\theta}(z))), \quad (45)$$

where  $D_w$  denotes the discriminator,  $G_{\theta}$  denotes the generator, and  $\mathcal{D} = \{(x)\}$  denotes the dataset.

What is an issue when there is a poor generator and a good discriminator? What is a common heuristic to counter this issue?

Your answer:

- (d) [8 points] Please write down the mini-batch gradient descent training algorithm for optimizing the program given in Eq. (45).

Your answer:

---

**Algorithm 1:** Minibatch gradient descent training

---

Initialize \_\_\_\_\_

**for** *number of training iterations* **do**

• Sample minibatch, \_\_\_\_\_

• Sample minibatch, \_\_\_\_\_

• Update \_\_\_\_\_

$w +=$  \_\_\_\_\_

• Update \_\_\_\_\_

$\theta +=$  \_\_\_\_\_

**end**

---

- (e) [4 points] To theoretically analyze the cost function in Eq. (45), we consider the following program:

$$\max_{p_G} \min_D - \int_x p_{\text{data}}(x) \log D(x) dx - \int_x p_G(x) \log(1 - D(x)) dx. \quad (46)$$

Write down the meaning and relation to Eq. (45) for each of the following symbols:  $p_{\text{data}}$ ,  $p_G$ ,  $D$ , and  $G$ .

Your answer:

38. [10 points] Markov Decision Processes (blue color: rewards; red color: decision probabilities)

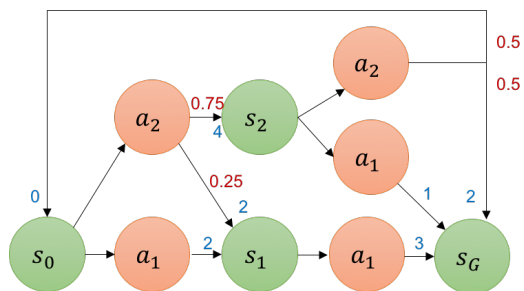


Figure 4: MDP for this problem. Start state is  $s_0$ .

- (a) Explain the difference between a deterministic and a stochastic Markov decision process. Is the given MDP stochastic or deterministic. Explain your answer.

Your answer:

- (b) What are three mechanisms to find the optimal policy  $\pi^*$  for a given MDP?

Your answer:

- (c) For the policy  $\pi(s_0) = a_1$ ,  $\pi(s_1) = a_1$ , what is the policy graph and the resulting value function  $V^\pi(s_1)$  and  $V^\pi(s_0)$ ?

Your answer:

- (d) For the policy  $\pi(s_0) = a_2$ ,  $\pi(s_2) = a_1$ , what is the policy graph and the resulting value function  $V^\pi(s_2)$ ,  $V^\pi(s_1)$  and  $V^\pi(s_0)$ ?

Your answer:

- (e) For the policy  $\pi(s_0) = a_2$ ,  $\pi(s_2) = a_2$ , what is the policy graph and the resulting value function  $V^\pi(s_2)$ ,  $V^\pi(s_1)$  and  $V^\pi(s_0)$ ?

Your answer:

- (f) What is the optimal policy for the MDP given in Fig. 5? Briefly explain your answer.

Your answer:

39. [13.5 points] Markov Decision Processes (transition probabilities are given in boxes)

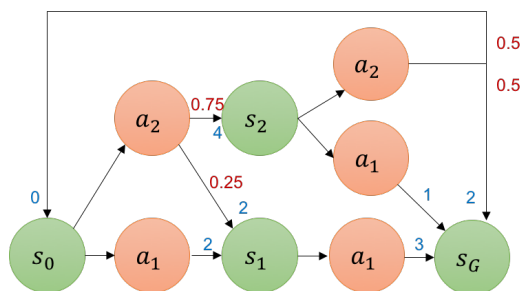


Figure 5: MDP for this problem. Start state is  $s_0$ .

- (a) (1 point) Explain the difference between a deterministic and a stochastic Markov decision process. Is the given MDP stochastic or deterministic. Explain your answer.

Your answer:

- (b) (1.5 points) What are three mechanisms to find the optimal policy  $\pi^*$  for a given MDP?

Your answer:

- (c) (3 points) For the policy  $\pi(s_0) = a_1$ ,  $\pi(s_1) = a_1$ , what is the policy graph and the resulting value function  $V^\pi(s_1)$  and  $V^\pi(s_0)$ ?

Your answer:

- (d) (3 points) For the policy  $\pi(s_0) = a_2$ ,  $\pi(s_2) = a_1$ , what is the policy graph and the resulting value function  $V^\pi(s_2)$ ,  $V^\pi(s_1)$  and  $V^\pi(s_0)$ ?

Your answer:

- (e) (3 points) For the policy  $\pi(s_0) = a_2$ ,  $\pi(s_2) = a_2$ , what is the policy graph and the resulting value function  $V^\pi(s_2)$ ,  $V^\pi(s_1)$  and  $V^\pi(s_0)$ ?

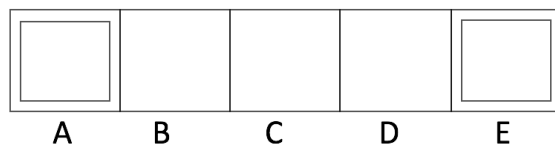
Your answer:

- (f) (2 points) What is the optimal policy for the MDP given in Fig. 5? Briefly explain your answer.

Your answer:

40. [19 points] Markov Decision Process (MDP)

Consider the MDP defined on the grid below. In all cases double-rectangle states are exit states. From an exit state (A, E), the only action available is **Exit**, which results in the listed immediate reward and ends the game (by moving into a terminal state  $T$ ). From non-exit states, the agent can choose either **Left** or **Right** actions with equal probability, which move the agent in the corresponding direction. The only non-zero rewards come from exiting the grid ( $R(A, \text{Exit}, T) = 1$  and  $R(E, \text{Exit}, T) = 2$ ). Throughout this problem, assume that value iteration begins with initial values  $V_0(s) = 0$  for all states  $S = \{A, B, C, D, E, T\}$ .



- (a) [3.5 points] Fill in the blank in the MDP diagram below with the missing actions (on the edges), rewards (in squares) and states (in circles).

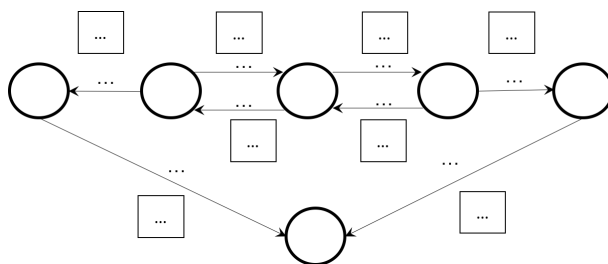


Figure 6: MDP representation

Your answer:

- (b) [1 points] Given the Bellman equation

$$V^*(s) = \max_{a \in A_s} \sum_{s' \in S} P(s'|s, a)(\gamma V^*(s') + R(s, a, s')) \quad (47)$$

with  $S$  being the set of states,  $A_s$  the set of available actions from the state  $s$ ,  $V$  the value function,  $R$  the reward,  $P$  the transition probability and  $\gamma$  a discount factor. Write down the update-rule for the **value iteration** algorithm at time  $t$ .

Your answer:

- (c) [1 points] What is the value of  $P(s'|s, a)$  for the different states of the grid?

Your answer:

- (d) [3.5 points] Let the discount factor be  $\gamma = 0.5$ . Fill in the missing values for each state following the value iteration algorithm in the following table:

Your answer:

t	A	B	C	D	E
0	0	0	0	0	0
1	...	...	...	...	...
2	...	...	...	...	...
3	...	...	...	...	...

- (e) [2 points] Starting from state  $C$ , what is the optimal sequence of actions?

Your answer:



- (f) [1 points] Instead of finding the optimal policy, assume we want to do **policy evaluation** for a policy  $\pi$ . Write down the iterative refinement formula for the Bellman equation.

Your answer:

- (g) [3.5 points] Evaluate the following policy:  $\pi(B) = \text{Left}$ ,  $\pi(C) = \text{Left}$  and  $\pi(D) = \text{Right}$  for  $\gamma = 0.5$ .

Your answer:

t	A	B	C	D	E
0	0	0	0	0	0
1	...	...	...	...	...
2	...	...	...	...	...
3	...	...	...	...	...

- (h) [2 points] Recalling that the discount factor must be in range  $0 \leq \gamma \leq 1$ , for what range of values for  $\gamma$  is the optimal action  $\pi^*(B) = \text{Right}$ ?

Your answer:

- (i) [1.5 point] How is policy iteration different from value iteration?

Your answer:

41. [13 points] Q-Learning

- (a) State the Bellman optimality principle as a function of the optimal Q-function  $Q^*(s, a)$ , the expected reward function  $R(s, a, s')$  and the transition probability  $P(s'|s, a)$ , where  $s$  is the current state,  $s'$  is the next state and  $a$  is the action taken in state  $s$ .

Your answer:

- (b) In case the transition probability  $P(s'|s, a)$  and the expected reward  $R(s, a, s')$  are unknown, a stochastic approach is used to approximate the optimal Q-function. After observing a transition of the form  $(s, a, r, s')$ , write down the update of the Q-function at the observed state-action pair  $(s, a)$  as a function of the learning rate  $\alpha$ , the discount factor  $\gamma$ ,  $Q(s, a)$  and  $Q(s', a')$ .

Your answer:

- (c) What is the advantage of an epsilon-greedy strategy?

Your answer:

- (d) What is the advantage of using a replay-memory?

Your answer:

- (e) Consider a system with two states  $S_1$  and  $S_2$  and two actions  $a_1$  and  $a_2$ . You perform actions and observe the rewards and transitions listed below. Each step lists the current state, reward, action and resulting transition as:  $S_i; R = r; a_k : S_i \rightarrow S_j$ . Perform Q-learning using a learning rate of  $\alpha = 0.5$  and a discount factor of  $\gamma = 0.5$  for each step by applying the formula from part (b). The Q-table entries are initialized to zero. Fill in the tables below corresponding to the following four transitions. What is the optimal policy after having observed the four transitions?

- i.  $S_1; R = -10; a_1 : S_1 \rightarrow S_1$
- ii.  $S_1; R = -10; a_2 : S_1 \rightarrow S_2$
- iii.  $S_2; R = 18.5; a_1 : S_2 \rightarrow S_1$
- iv.  $S_1; R = -10; a_2 : S_1 \rightarrow S_2$

$Q$	$S_1$	$S_2$
$a_1$	.	.
$a_2$	.	.

$Q$	$S_1$	$S_2$
$a_1$	.	.
$a_2$	.	.

$Q$	$S_1$	$S_2$
$a_1$	.	.
$a_2$	.	.

$Q$	$S_1$	$S_2$
$a_1$	.	.
$a_2$	.	.

Your answer:

42. [10 points] Q-Learning

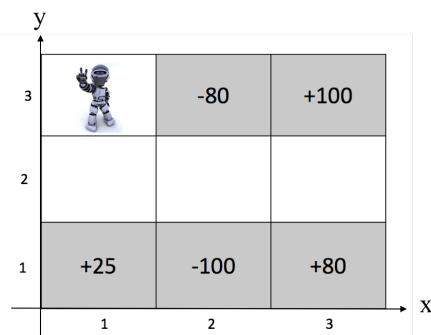


Figure 7: Find me the optimal policy. May the force be with us!

Consider the grid-world given in Figure 7 and our artificial agent who is trying to learn the optimal policy. If an action results in landing into one of the shaded states the corresponding reward is awarded during that transition. All shaded states are terminal states. The other states have the North, East, South, West actions available, which deterministically move the agent to the corresponding neighboring state (or have him stay in place if the action tries to move out of the grid). Assume a discount factor  $\gamma = 0.5$  and the Q-learning rate of  $\alpha = 0.5$  for all calculations. The agent starts in state  $(x, y) = (1, 3)$ .

- (a) (2 points) Write down the expression of the optimum value  $V^*(s)$  at state  $s$  as a function of the reward  $r_{ss'}^a$ ,  $\gamma$  and the next state's optimum value  $V^*(s')$  in this grid-world.

Your answer:

- (b) (1 point) What is the optimal value  $V^*$  at state  $(3, 2)$ ?

Your answer:

- (c) (1 point) What is the optimal value  $V^*$  at state  $(2, 2)$ ?

Your answer:

- (d) (1 point) What is the optimal value  $V^*$  at state  $(1, 3)$ ?

Your answer:

The agent starts from the top left corner and you are given the following three episodes from runs of the agent through this grid-world. Each line in an Episode is a tuple containing  $(s, a, s', r)$ .

Episode1	Episode2	Episode3
$(1, 3), S, (1, 2), 0$	$(1, 3), S, (1, 2), 0$	$(1, 3), S, (1, 2), 0$
$(1, 2), E, (2, 2), 0$	$(1, 2), E, (2, 2), 0$	$(1, 2), E, (2, 2), 0$
$(2, 2), S, (2, 1), -100$	$(2, 2), E, (3, 2), 0$	$(2, 2), E, (3, 2), 0$
	$(3, 2), N, (3, 3), +100$	$(3, 2), S, (3, 1), +80$

- (e) (2 points) Write down the Q-learning gradient update rule for  $Q(s, a)$  as a function of  $\alpha, \gamma, Q(s', a')$  and  $R(s, a, s')$ .

Your answer:

- (f) (3 points) Assume the entire aforementioned experience replay to be given and all the Q values are initialized to zero. Perform a single step of gradient updates for the following Q-values:  $Q((3, 2), N)$ ,  $Q((1, 2), S)$  and  $Q((2, 2), E)$ . The aforementioned Q-values are updated in the given order, i.e., when updating  $Q((2, 2), E)$ ,  $Q((3, 2), N)$  and  $Q((1, 2), S)$  are already updated and their updated values should be used.

Your answer: