

# CS 446 / ECE 449: Machine learning.

Exam #1: due at noon, CST, on March 18, 2021.



Name:

NetID:

M J T

- 
- **This is an individual exam:** you may not consult with anyone except course staff.
  - **Resources you may access:** course lecture slides (`Lxy_Slides_Split.pdf` recommended for fewest typos), your homework solutions, and two single-sided pages typed or hand-written at a reasonable size.
  - **Citing lecture:** if you skip steps via something from lecture, be explicit: cite the lecture and slide number.
  - **Where to go for clarifications (not “hints” or “help”!):**
    - Make a private (TA/instructor only) post on campuswire.
    - Attend one of the following campuswire office hours (includes extras over regular schedule):  
**3/16:** 15:30 – 17:30, 22:00 – 22:30; **3/17:** 11:00 – 11:30, 16:30 – 17:30; **3/18:** 10:00 – 11:00.
  - **Where not to go for help** includes (but is not limited to):
    - Discord.
    - Public posts on campuswire.
    - Any discussion with other students.
    - ...
  - **How to submit:**
    1. Prepare a PDF of your solutions; e.g., (a) print the exam, write on it, and take pictures, (b) write directly on the exam via a tablet, (c) write on blank sheets and take pictures.
    2. Upload to gradescope under `midterm`, similarly to homeworks. You may re-submit often.
  - **Double-ruled boxes:** if you write directly on the exam PDF, please limit your answers to these boxes. They were designed to give generous amounts of space.
  - **Point values:** problems have different point values (as indicated), but parts of problems are all equal! Do not lose too much time on parts that are much more difficult!
  - **Notation:**  $\|\cdot\|$  without a subscript means 2-norm (Euclidean norm), thus  $\|\mathbf{v}\|^2 = \|\mathbf{v}\|_2^2 = \sum_{i=1}^d v_i^2$ . “ $\mathbf{x}$  has unit norm” means  $\|\mathbf{x}\| = 1$ .  $\mathbf{e}_i$  is the  $i^{\text{th}}$  standard basis vector of appropriate dimension.  $k$ -nearest-neighbor uses 2-norm by default. The decision boundary of a linear SVM with no bias goes through the origin; with a bias, it need not pass through the origin.
-

$x =$  hard to ground?

Method 1  
pictures

(1)

convexity  
of  $\delta$ s

(2)

convolutional

$$\begin{matrix} & a \\ x & b \\ & c \\ x & d \\ & e \end{matrix}$$

(3)

boundaries  
+  
check  $\delta$ 's  
solution

$$\begin{matrix} & a \\ x & b \\ & c \\ & d \\ & e \end{matrix}$$

(4)

$$\begin{matrix} & a \\ x & b \\ & c \\ x & d \\ & e \\ & f \end{matrix}$$

second easiest

(I)

1, 2  $\{a, b, c\}$  Haocken

maybe hardest

2  $\{d, e\}$ , 4  $\{d, e\}$ ,  
pe: y.

easiest

(II)

3  $\{a, b, c, d, e\}$   
x: diagonal

(III)

MJT

maybe hardest

(IV)

4  $\{a, b, c, f\}$

Jiangyan

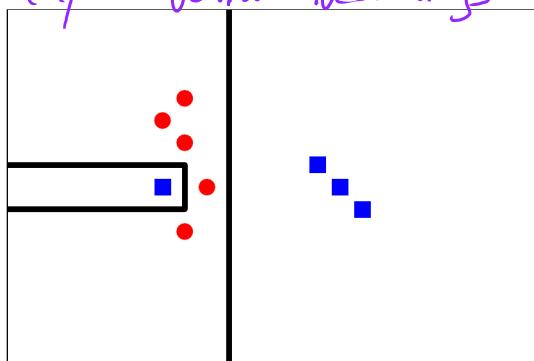
# 1. Decision boundaries (6 total points, 1 point each figure).

Match the following predictors to the following decision boundaries: (a) 1-nearest-neighbor, (b) axis-aligned decision tree, (c) linear SVM, (d) arbitrary linear separator, (e) RBF SVM, (f) ReLU network.

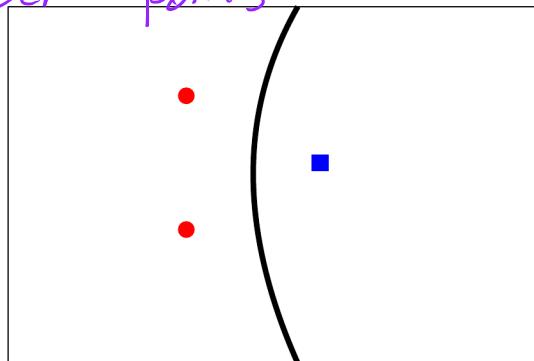
Note that matching means you use each answer exactly once; there is only one valid matching.

*Only valid markings receive points*

1.1



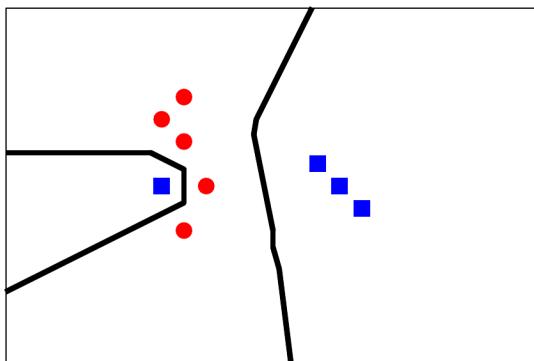
1.4



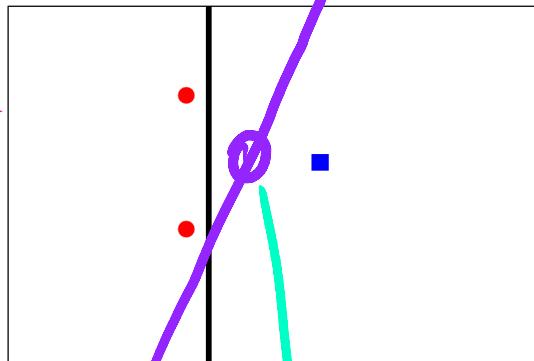
Circle one: (a) (b) (c) (d) (e) (f)

Circle one: (a) (b) (c) (d) (e) (f)

1.2



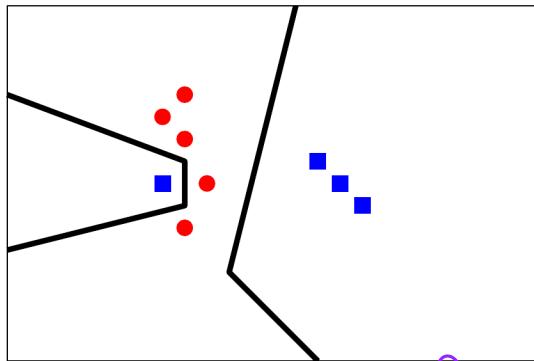
1.5



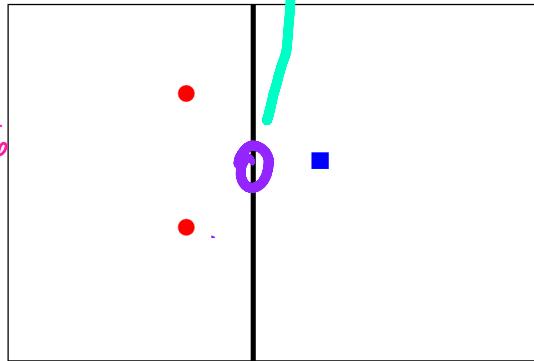
Circle one: (a) (b) (c) (d) (e) (f)

Circle one: (a) (b) (c) (d) (e) (f)

1.3



1.6



Circle one: (a) (b) (c) (d) (e) (f)

Circle one: (a) (b) (c) (d) (e) (f)

*partial credit delicate*

## 2. Short answer (15 total points, 3 per part).

- (a) Prove that  $g(x) := x^3$  is not convex (over  $\mathbb{R}$ ).

**Hint:** Produce a counterexample to the definition on slide 5/26 from lecture 5. (Your proof needs to be an explicit symbolic derivation, but consider first studying a plot on scratch paper.)

Collect  
counterexample  
but no  
proof:

total ≤ 1 points.

Consider  $x := -1$ ,  $x' := 0$ ,  $\alpha := \frac{1}{2}$ ,  
and  $z := \alpha x + (1-\alpha)x' = -\frac{1}{2}$ .

$$\begin{aligned} \text{Then } g(z) &= -\frac{1}{8} > -\frac{1}{2} = \alpha(-1) + (1-\alpha)(0) \\ &= \alpha g(x) + (1-\alpha)g(x'), \end{aligned}$$

Meaning  $g$  not convex.

Can also do other points, also  $g'$  &  $g''$ .

- (b) Consider the following least squares problem: data matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  has rank  $r < \min\{n, d\}$ , target vector  $\mathbf{y} \in \mathbb{R}^n$  is given, and let  $\hat{\mathbf{w}}_{\text{ols}}$  denote the OLS solution. Let  $\mathbf{v}$  denote any nonzero vector which is orthogonal to the right singular vectors  $(\mathbf{v}_1, \dots, \mathbf{v}_r)$  of  $\mathbf{X}$ , and define another solution  $\mathbf{w} := \hat{\mathbf{w}}_{\text{ols}} + \mathbf{v}$ . Show that  $\mathbf{w}$  is an optimal solution to this least squares problem.

As in lecture 2, it suffices to show that  $\mathbf{w}$  satisfies the normal equations. Let  $\mathbf{X} = \sum_{i=1}^r s_i u_i v_i^T$  be an SVD with the provided right singular vectors, whereby  $\mathbf{X}\mathbf{v} = \sum_{i=1}^r s_i u_i v_i^T \mathbf{v} = 0$ . Since  $\hat{\mathbf{w}} := \hat{\mathbf{w}}_{\text{ols}}$  satisfies the normal equations,

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T (\mathbf{X} \hat{\mathbf{w}} + \mathbf{X} \mathbf{v}) = \mathbf{X}^T \mathbf{X} (\hat{\mathbf{w}} + \mathbf{v}) = \mathbf{X}^T \mathbf{X} \mathbf{w}.$$

Scans there's a risk of cheating in here?

- (c) Consider a convolutional layer with one  $2 \times 2$  filter  $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$ , all entries nonzero, and the simplest settings: no bias, no padding, a stride of 1, no dilation. Construct a  $3 \times 3$  input image, not all entries equal to zero, which is in the *null space* of this filter, meaning the output image after applying the above filter is all zeros. For full points, include your input image and a brief verification that the output image is all zeros.

Either convolution convention may (filter reversed or not).

Consider input data  $X = \begin{bmatrix} c & 0 & 0 \\ -a & 0 & 0 \\ a^2/c & 0 & 0 \end{bmatrix}$ .

Then

$$\text{conv}\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}, X\right) = \begin{bmatrix} ac - ac & 0 \\ -a^2 + c\left(\frac{a^2}{c}\right) & 0 \end{bmatrix}$$

tough to grade. Many crazy solutions.

$$= \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

If fixed filter, then  $\leq 2$  total point

- (d) Consider the *three-dimensional XOR* problem in  $\mathbb{R}^3$ : there are 8 data points, consisting of all possible  $(\pm 1, \pm 1, \pm 1)$  triples, and the labels are the products of the three coordinates. Circle all of the classifiers/methods in the following list which can perfectly classify this problem:

- |  |   |
|--|---|
| i. 2-layer ReLU network,                 | iv. decision tree with 3 axis-aligned splits, |
| ii. fourth degree polynomial kernel SVM, | v. 1-nearest-neighbor,                        |
| iii. linear predictor,                   | vi. 4-nearest-neighbor,                       |

-1 points: 1-2 errors

-2 points: 3-4 errors

-3 points 5-6 errors.

- (e) Recall the primal and dual optimization problems corresponding to the hard-margin SVM for linearly separable data:

$$\min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 : \mathbf{w} \in \mathbb{R}^d, \forall i, \mathbf{w}^\top \mathbf{x}_i y_i \geq 1 \right\} = \max_{\alpha \in [0, \infty)^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^\top \mathbf{x}_j.$$

Hard  
prob  
is showing  
distinct  
dual  
choices  
are optimal;  
here I use  
weak  
duality!

Consider the following (linearly separable) data:  $d = 2, n = 2, \mathbf{x}_1 = \mathbf{e}_1, y_1 = +1, \mathbf{x}_2 = -\mathbf{x}_1 = -\mathbf{e}_1, y_2 = -y_1 = -1$ . Show that the dual optimum is not unique.

**Hint:** Here is one way to solve the problem (there are others). On scratch paper, draw the data and guess an optimal primal weight vector via the geometric view of the SVM (this primal weight vector will be unique). Then write out the dual, plug in the data, and solve for two distinct optimal values of  $\alpha$ . To write up your proof for your final solution, you can produce your primal solution and two distinct dual solutions, point out that the corresponding primal and dual objective values are equal (thus all three values are optimal), which implies the dual optimum is not unique.

Define  $\bar{\mathbf{w}} := \mathbf{e}_1, \quad \} \text{primal.}$

$\bar{\alpha} := \mathbf{e}_1, \quad \} \text{dual.}$

$\bar{\beta} := \mathbf{e}_2 \quad \}$

$\bar{\mathbf{w}}$  is feasible:  $\bar{\mathbf{w}}^\top \mathbf{e}_1 (+1) \geq 1, \bar{\mathbf{w}}^\top (-\mathbf{e}_1) (-1) \geq 1$ .

Primal value:  $\frac{1}{2} \|\bar{\mathbf{w}}\|^2 = \frac{1}{2}$ .

Dual values:  $D(\bar{\alpha}) = 1 - \frac{1}{2} y_1^2 \mathbf{x}_1^\top \mathbf{x}_1 = \frac{1}{2},$

$D(\bar{\beta}) = 1 - \frac{1}{2} y_2^2 \mathbf{x}_2^\top \mathbf{x}_2 = \frac{1}{2}.$

$\frac{1}{2} \|\bar{\mathbf{w}}\|^2 = D(\bar{\alpha}) = D(\bar{\beta}),$  thus respectively primal, dual, dual  
by weak duality  $\rightarrow$  optimal, and note  $\bar{\alpha} \neq \bar{\beta}$ .

Many little details. Any other approach  
quite messy & error prone.<sup>5</sup> Optimality conditions  
dictate since primal & dual are constrained.

two ways to solve

① Using duality

Start from distinct feasible dual variables  $\bar{\alpha} + \bar{\beta}$ ,

then argue they are optimal.

maybe like off points

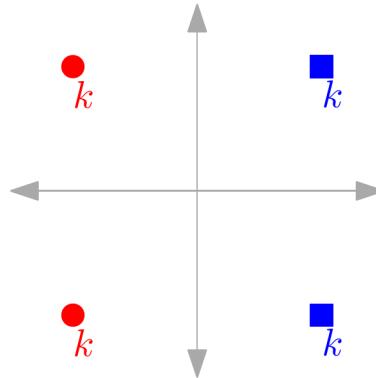
In reference solution above,  
I used  $\bar{w}$  & weak duality  
to assert their optimality. //

② Take gradient of dual objective, set to zero,  
pick two different solutions  
(I'm fine with this approach.)

### 3. Robustness (21 total points, 3 per part).

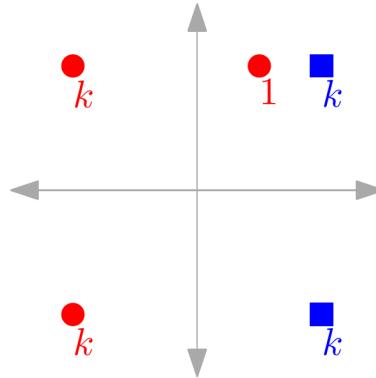
In this problem, we will study how different predictors respond to slight changes in the training data. We will use two sets of training data, one with  $4k$  points and another with  $4k + 1$  points, where  $k \geq 3$  is arbitrary.

- **Training set 1, used in parts (a) & (b).** This data consists of  $n_1 := 4k$  points in  $\mathbb{R}^2$ , forming the rows of input matrix  $\mathbf{X}_1 \in \mathbb{R}^{n_1 \times 2}$ , with corresponding label vector  $\mathbf{y}_1 \in \{\pm 1\}^{n_1}$ .



In detail, this data has  $k \geq 3$  copies of each of the four points of the form  $(\pm 2, \pm 2)$ . The labels are given by the colors/shapes: let red/circle denote  $-1$ , and blue/square denote  $+1$ . This data is linearly separable.

- **Training set 2, used in parts (c) – (g).** This data consists of  $n_2 := 4k + 1$  points in  $\mathbb{R}^2$ , forming the rows of input matrix  $\mathbf{X}_2 \in \mathbb{R}^{n_2 \times 2}$ , with corresponding label vector  $\mathbf{y}_2 \in \{\pm 1\}^{n_2}$ .

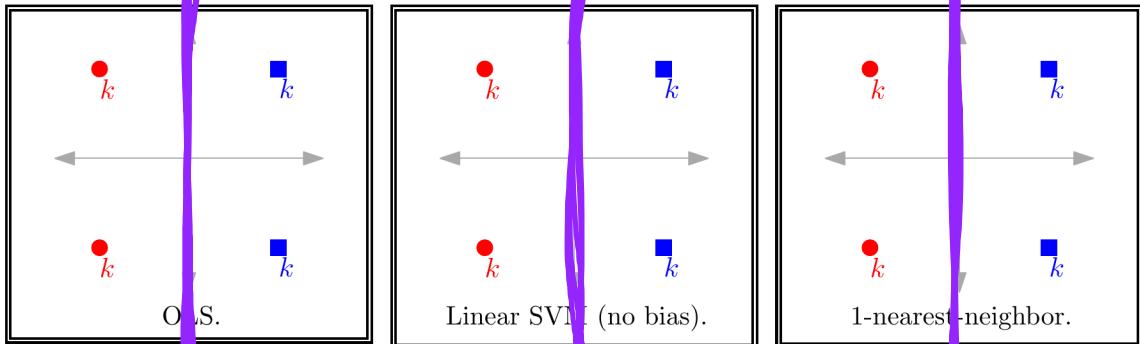


In detail,  $4k$  of the points are common with training set 1, however there is an additional red/circle point (with label  $-1$ ) at location  $(1, 2)$ . The data is still linearly separable.

These parts of the problem use training set 1, meaning inputs  $X_1$  and labels  $y_1$ .

- (a) Draw the decision boundary of OLS, linear SVM (with no bias), and 1-nearest-neighbor in the following three plots. You do not need to justify your answers.

**Note for anybody not using this PDF template:** Don't stress the details too much, just try to make your figures unambiguous and easy to grade.



- (b) Write down the OLS solution, and provide a verification that it satisfies the normal equations.

**Hint:** You don't need to do tedious calculations for this problem, here are two ways to do it. The first option is to compute all objects in the normal equations, and from there compute a solution and argue it is actually the OLS solution. A second option is to "guess" the SVD of  $X_1$ , which has a clean form, and from there get the OLS solution.

$$\text{Wlog } X = \begin{bmatrix} 2 & 2 \\ 1 & 1 \\ 2 & -2 \\ 1 & -1 \\ -2 & 2 \\ 1 & 1 \\ -2 & -2 \\ 1 & 1 \end{bmatrix}^n, \quad y = \begin{bmatrix} +1 \\ 1 \\ -1 \\ 1 \end{bmatrix}_{2n}^T.$$

$X$  has orthogonal columns, thus

$$X^T X = \begin{bmatrix} 16k & 0 \\ 0 & 16k \end{bmatrix}, \quad \text{which is invertible,}$$

therefore normal equations have unique  
solution  $(X^T X)^{-1} (X^T) y = \begin{bmatrix} 16k & 0 \\ 0 & 16k \end{bmatrix}^{-1} \begin{bmatrix} 8k \\ 4k \end{bmatrix} = \begin{bmatrix} k \\ k \end{bmatrix}$ , which must agree with  
ols solution.

(see next page for another.)

t from lecture:  
ols always satisfy  
normal.

Note  $X = S_1 u_1 v_1^T + S_2 u_2 v_2^T$  where

$$v_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad v_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$u_1 = \frac{1}{\sqrt{2k}} \begin{bmatrix} 1 \\ 1 \\ 0 \\ \vdots \\ 0 \\ -1 \end{bmatrix}^k \quad u_2 = \frac{1}{\sqrt{2k}} \begin{bmatrix} 0 \\ 1 \\ 1 \\ \vdots \\ -1 \\ 0 \end{bmatrix}^k$$

$$S_1 = 4\sqrt{k}, \quad S_2 = 4\sqrt{k}.$$

Note also

$$u_1^T y = \sqrt{2k}, \quad u_2^T y = \sqrt{2k},$$

and

$$\hat{w}_{OLS} := X^+ y = \frac{1}{S_1} v_1 u_1^T y + \frac{1}{S_2} v_2 u_2^T y$$

$$= \frac{1}{4\sqrt{k}} \cdot \frac{1}{\sqrt{2}} \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} \sqrt{2k} + \begin{bmatrix} 1 \\ -1 \end{bmatrix} \sqrt{2k} \right) = \frac{1}{4\sqrt{2k}} \begin{bmatrix} 2\sqrt{2k} \\ 0 \end{bmatrix}$$

(Quite annoying!)

$$= \begin{bmatrix} 1/2 \\ 0 \end{bmatrix}. //$$

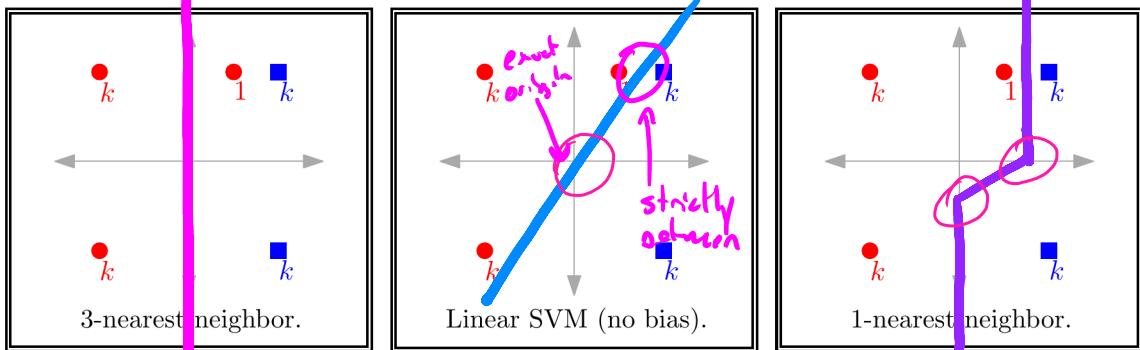
also, most say  
 uncs proud in lecture  
 that OLS satisfies  
 normal equations, or  
 check them.

These parts of the problem use training set 2, meaning inputs  $X_2$  and labels  $y_2$ .

Strict grading!

- (c) Draw the decision boundary of 3-nearest-neighbor, linear SVM (with no bias term), and 1-nearest-neighbor in the following three plots. You do not need to justify your answer.

Note for anybody *not* using this PDF template: Don't stress the details too much, just try to make your figures unambiguous and easy to grade.



- (d) State the number of points within a distance of 1 of the decision boundary for the SVM, and also the number of points within a distance of 1 of the decision boundary for the 3-nearest-neighbor classifier.

or they  
my  
SVM  
"3 points  
per SVM  
Tentative"

SVM:  $2k+1$  points.

3-nn: 1 point.

- (e) Limiting yourself to 1–5 sentences, state in intuitive terms whether you prefer the linear SVM solution or the 3-nearest-neighbor solution.

This problem will be graded leniently, but for full points, refer to parts (e) and (f), and how you feel they may relate to the performance on (unseen!) testing data.

MJ

If data (seen & unseen) has small perturbations & label flips, SVM is unstable (as above) & 3-nn is ideal (low test error).

If data dist non-3y, linear SVM could be good in some low-sample regimes.

Can discuss Soft-margin SVM?

#### 4. Vanishing/exploding gradients (18 total points, 3 per part).

Consider an  $L$ -layer ReLU network with input dimension, output dimension, and all widths equal to 2, and no biases, meaning it computes  $f(\mathbf{x}; \mathcal{W})$  given by

$$\mathbb{R}^2 \ni \mathbf{x} \mapsto f(\mathbf{x}; \mathcal{W}) := \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x}) \cdots) \in \mathbb{R}^2,$$

where each  $\mathbf{W}_i \in \mathbb{R}^{2 \times 2}$ , and each  $\sigma : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  applies a ReLU coordinate-wise (meaning  $\sigma(\mathbf{v}) = (\max\{0, v_1\}, \max\{0, v_2\})^\top \in \mathbb{R}^2$ ), and  $\mathcal{W} = (\mathbf{W}_L, \dots, \mathbf{W}_1)$  is notation for the tuple of parameters across layers.

Throughout this problem, you will be asked to choose parameters  $\mathcal{W} = (\mathbf{W}_L, \dots, \mathbf{W}_1)$  which satisfy certain conditions. In all parts of this problem:

- i. An input vector  $\mathbf{x} \in \mathbb{R}^2$  is given with unit norm, but is otherwise arbitrary;
- ii. You must choose matrices  $\mathcal{W} = (\mathbf{W}_L, \dots, \mathbf{W}_1)$ , where each matrix  $\mathbf{W}_i \in \mathbb{R}^{2 \times 2}$  has rank 2, with two distinct singular values exactly equal to 2 and  $1/2$  respectively;
- iii. You must provide both an answer and a brief derivation/verification for full credit.

The parts of the problem are as follows.

- (a) Let unit norm input  $\mathbf{x}$  be given, and let  $\mathbf{z} \in \mathbb{R}^2$  be another unit norm vector which is orthogonal to  $\mathbf{x}$ . Define a matrix

$$\mathbb{R}^{2 \times 2} \ni \mathbf{M} := 2\mathbf{e}_1 \mathbf{x}^\top + \frac{1}{2}\mathbf{e}_2 \mathbf{z}^\top.$$

Provide the singular values of  $\mathbf{M}$ , and compute  $\mathbf{M}\mathbf{x}$  and  $\mathbf{M}\mathbf{z}$ .

**Hint:** This is not a trick question; compare the definition of  $\mathbf{M}$  to the definition of the SVD from slide 9/26 in lecture 2.

e<sub>1</sub> & e<sub>2</sub> are orthonormal, as are x & z,  
 therefore 2e<sub>1</sub>x<sup>T</sup> + 1/2e<sub>2</sub>z<sup>T</sup> is an SVD  
 for M, and singular values are 2 & 1/2.  
 Moreover

$$Mx = 2e_1 x^T x + \frac{1}{2} e_2 z^T x = 2 + 0 = 2,$$

$$Mz = 2e_1 x^T z + \frac{1}{2} e_2 z^T z = 0 + \frac{1}{2} = \frac{1}{2}.$$

- (b) Let unit norm input  $\mathbf{x}$  be given. Provide a choice of  $\mathcal{W}$  so that  $\|f(\mathbf{x}; \mathcal{W})\| = 2^L$ .

**Hint:** The previous part can help you choose  $\mathbf{W}_1$  (though, if you use it, you must still provide verification); how should you choose the weight matrices in other layers?

No  
need  
for  
induction.

at most  
1 point  
if no  
verification

Choose  $\mathbf{W}_1 = \mathbf{M} = 2\mathbf{e}_1 \mathbf{x}^T + \frac{1}{2}\mathbf{e}_2 \mathbf{z}^T$ ,  
where  $\mathbf{z}$  unit &  $\mathbf{x}^T \mathbf{z} = 0$ .

Choose  $\mathbf{W}_i := 2\mathbf{e}_1 \mathbf{e}_1^T + \frac{1}{2}\mathbf{e}_2 \mathbf{e}_2^T, i \geq 2$ .

By construction, singular values  $(2, 1)$ .

Moreover  $\mathbf{W}_{i+1}(\mathbf{W}_{i-1} \cdots \mathbf{W}_3(\mathbf{W}_2(\mathbf{W}_1 \mathbf{x}))) \cdots = \begin{bmatrix} 2^L \\ 0 \end{bmatrix}$ ,  
which has norm  $2^L$ .

- (c) Let unit norm input  $\mathbf{x}$  be given. Provide a choice of  $\mathcal{W}$  so that  $\|f(\mathbf{x}; \mathcal{W})\| = 2^{-L}$ .

Pick  $\mathbf{z}$  as before, define

$$\mathbf{W}_1 := 2\mathbf{e}_1 \mathbf{z}^T + \frac{1}{2}\mathbf{e}_2 \mathbf{z}^T$$

$$\text{and } \mathbf{W}_i := 2\mathbf{e}_1 \mathbf{e}_1^T + \frac{1}{2}\mathbf{e}_2 \mathbf{e}_2^T,$$

Satisfying conditions by construction.

By induction, output of  $\mathbf{W}_i$  is  $\begin{bmatrix} 0 \\ 2^{-i} \end{bmatrix}$ :

Base:  $\mathbf{W}_{i=1} = 0 + \frac{1}{2}\mathbf{e}_2$ ; IS:  $\mathbf{W}_i \circ \left( \begin{bmatrix} 0 \\ 2^{-i} \end{bmatrix} \right) = \mathbf{W}_i 2^{-i} \mathbf{e}_2 \Rightarrow 2^{-i} \mathbf{e}_2$

- (d) Let unit norm input  $\mathbf{x}$  be given, and let  $\mathbf{W}_1 \in \mathbb{R}^{2 \times 2}$  be any  $2 \times 2$  matrix satisfying the conditions on the first layer which hold throughout this problem. Show that it is impossible to have  $\mathbf{W}_1\mathbf{x} = 0$ .

**Hint:** What is the rank of  $\mathbf{W}_1$ ? What is the largest possible rank of a  $2 \times 2$  matrix?

lenient

any  
proof

/  
argument  
is  
fine

$\mathbf{W}_1$  is full rank, thus has  
no non-trivial right null space.

OR:

Implies  $\mathbf{x} = \mathbf{W}^{-1}\mathbf{W}_1\mathbf{x} = \mathbf{W}^{-1}\mathbf{0} = \mathbf{0}$ ,  
but  $\|\mathbf{x}\| = 1$ , contradiction.

- (e) Let unit norm input  $\mathbf{x}$  be given. Provide a choice of  $\mathcal{W}$  so that  $\|f(\mathbf{x}; \mathcal{W})\| = 0$ .

**Hint:** You can't enforce  $\mathbf{W}_1\mathbf{x} = 0$ , since that would contradict the previous part. Keep in mind what you have extra in this part: you have  $L$  layers, not just one, and you have ReLUs.

$\{ L \geq 1 \}$

No  
need  
for  
explicit  
induction

Choose  $z$  as before and  $\mathbf{W}_1 = -2\mathbf{e}_1\mathbf{x}^T + \frac{1}{2}\mathbf{e}_2\mathbf{z}^T$   
and  $\mathbf{W}_2 := 2\mathbf{e}_1\mathbf{e}_1^T + \frac{1}{2}\mathbf{e}_2\mathbf{e}_2^T$  as before;  
all conditions met, and  
 $\sigma(\mathbf{W}_1\mathbf{x}) = \sigma\left(\begin{bmatrix} -2 \\ 0 \end{bmatrix}\right) = \{0\}$ , and  

$$\left\| \mathbf{W}_2 \sigma(\dots \underbrace{\mathbf{W}_2 \sigma(\mathbf{W}_1\mathbf{x})}_{\{0\}} \dots) \right\| = 0.$$

- (f) In lecture we only wrote out an abstract form of a deep network gradient; here is a clean explicit way for the simple networks in this problem. First, define *diagonal activation matrices*  $(\mathbf{D}_L, \dots, \mathbf{D}_1)$  as

$$\begin{aligned}\mathbf{D}_1 &:= \text{diag}(\mathbb{1}[\mathbf{W}_1 \mathbf{x} \geq 0]) \in \mathbb{R}^{2 \times 2}, \\ \mathbf{D}_2 &:= \text{diag}(\mathbb{1}[\mathbf{W}_2 \mathbf{D}_1 \mathbf{W}_1 \mathbf{x} \geq 0]) \in \mathbb{R}^{2 \times 2}, \\ &\vdots \\ \mathbf{D}_i &:= \text{diag}(\mathbb{1}[\mathbf{W}_i \mathbf{D}_{i-1} \cdots \mathbf{D}_1 \mathbf{W}_1 \mathbf{x} \geq 0]) \in \mathbb{R}^{2 \times 2}, \\ &\vdots \\ \mathbf{D}_{L-1} &:= \text{diag}(\mathbb{1}[\mathbf{W}_{L-1} \mathbf{D}_{L-2} \cdots \mathbf{D}_1 \mathbf{W}_1 \mathbf{x} \geq 0]) \in \mathbb{R}^{2 \times 2}.\end{aligned}$$

(These matrices depend on both  $\mathbf{x}$  and  $\mathcal{W}$ , but we hide this notation.) The gradient with respect to layer  $i$  of the first output coordinate of the whole network, can be written as

$$\frac{df(\mathbf{x}; \mathcal{W})_1}{d\mathbf{W}_i} = (\mathbf{e}_1^\top \mathbf{W}_L \mathbf{D}_{L-1} \cdots \mathbf{W}_{i+1} \mathbf{D}_{i+1})^\top (\mathbf{D}_{i-1} \mathbf{W}_{i-1} \cdots \mathbf{D}_1 \mathbf{W}_1 \mathbf{x})^\top.$$

(You don't need to verify this formula in this problem, you just need to apply it.)

Let unit norm input  $\mathbf{x}$  be given, suppose  $L > 2$ , and consider  $\frac{df(\mathbf{x}; \mathcal{W})_1}{d\mathbf{W}_2}$ , the gradient with respect to  $\mathbf{W}_2$  (of the first output coordinate), the parameters in layer 2, as given by the earlier formula. Provide a choice of  $\mathcal{W}$  so that

$$\left\| \frac{df(\mathbf{x}; \mathcal{W})_1}{d\mathbf{W}_2} \right\|_F = 2^{L-1}.$$

**Hint:** Try your solution from part (b).

No verification  
give  
at most  
1 point.

Choose solution in part (b)  
 By induction,  $\mathbf{z}_{:,i} = \sigma(\mathbf{w}_{:,i} \mathbf{z}_{:-i}) = \begin{bmatrix} 2^i \\ 0 \end{bmatrix}$ ,

$$\mathbf{D}_i := \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}.$$

Thus

$$\left\| \underbrace{\mathbf{e}_1^\top \mathbf{w}_L \cdots \mathbf{D}_{i+1}}_{\begin{bmatrix} 2^{L-i} \\ 0 \end{bmatrix}^\top} \left( \mathbf{D}_{i-1}^\top \mathbf{w}_{i-1} \cdots \mathbf{D}_1^\top \mathbf{w}_1 \mathbf{x} \right)^\top \right\|_F = \begin{bmatrix} 2^{i-1} \\ 0 \end{bmatrix} \rightarrow 2^{L-1}.$$