

# HOMework 3

## LINEAR ALGEBRA \*

10-606 MATHEMATICAL FOUNDATIONS FOR MACHINE LEARNING

### START HERE: Instructions

- **Collaboration Policy:** Please read the collaboration policy in the syllabus.
- **Late Submission Policy:** See the late submission policy in the syllabus.
- **Submitting your work:** You will use Gradescope to submit answers to all questions.
  - **Written:** For written problems such as short answer, multiple choice, derivations, proofs, or plots, please use the provided template. Submissions can be handwritten onto the template, but should be labeled and clearly legible. If your writing is not legible, you will not be awarded marks. Alternatively, submissions can be written in  $\text{\LaTeX}$ . Each derivation/proof should be completed in the boxes provided. To receive full credit, you are responsible for ensuring that your submission contains exactly the same number of pages and the same alignment as our PDF template.
  - **Latex Template:** <https://www.overleaf.com/read/nsqfhfwxyrtf>

Question	Points
Linear Systems and Linear Algebra	19
Plotting Linear (Affine) Functions	7
Projection Matrices	10
Linear Regression	4
Matrix Memories	5
Neural Networks	5
Total:	50

---

\*Compiled on Wednesday 2<sup>nd</sup> August, 2023 at 17:26

**Instructions for Specific Problem Types**

For “Select One” questions, please fill in the appropriate bubble completely:

**Select One:** Who taught this course?

- ☒ Matt Gormley
- ☐ Marie Curie
- ☐ Noam Chomsky

If you need to change your answer, you may cross out the previous answer and bubble in the new answer:

**Select One:** Who taught this course?

- ☒ Henry Chai
- ☐ Marie Curie
- ☒ Noam Chomsky

For “Select all that apply” questions, please fill in all appropriate squares completely:

**Select all that apply:** Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☐ I don't know

Again, if you need to change your answer, you may cross out the previous answer(s) and bubble in the new answer(s):

**Select all that apply:** Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☐ I don't know

For questions where you must fill in a blank, please make sure your final answer is fully included in the given space. You may cross out answers or parts of answers, but the final answer must still be within the given space.

**Fill in the blank:** What is the course number?

10-606

10-606~~7~~

## 1 Linear Systems and Linear Algebra (19 points)

1. For each statement, indicate whether it is true or false.

- (a) (1 point) A system of  $n$  linear equations with  $n$  unknowns has at least one solution.  
☐ True ☐ False
- (b) (1 point) A system of  $n$  linear equations with  $n$  unknowns has at most one solution.  
☐ True ☐ False
- (c) (1 point) If a square matrix is full rank, it cannot be inverted.  
☐ True ☐ False
- (d) (2 points) Rank of matrix product  $\mathbf{AB}$  can be greater than the rank of either  $\mathbf{A}$  or  $\mathbf{B}$ .  
☐ True ☐ False

2. (5 points) Consider the following matrix:

$$\mathbf{A} = \begin{bmatrix} -1 & 0 & 4 \\ a & b & c \\ d & e & f \end{bmatrix}.$$

Complete the matrix such that the rank of  $A$  equals 1. Restrictions:

- $a, b, c, d, e, f \in \mathbb{Z}$  (integers)
- $a, b, c, d, e, f \notin -3, 2, 5$
- $a, b, c, d, e, f$  are all different

3. (3 points) Consider the following matrix:

$$\mathbf{B} = \begin{bmatrix} 1 & 3 & 1 & 3 \\ 4 & 4 & -1 & 3 \\ 0 & 0 & -2 & -1 \\ -2 & -1 & 0 & 0 \end{bmatrix}$$

Compute the matrix rank of  $\mathbf{B}$ . You may do it by hand or by using Python and numpy. You must show your work. If you use Python and numpy, Just paste your working code in the Work box.

Rank

Work

4. (6 points) Rewrite the following system in matrix form and solve it by Gaussian Elimination.

$$\begin{aligned} 3x_1 - 7x_2 - 2x_3 &= 5 \\ -3x_1 + 5x_2 + x_3 &= 0 \\ 6x_1 - 4x_2 &= 1 \end{aligned}$$

For each row operation you use to put the matrix in upper triangular form, state the row operation and show the resulting matrix and RHS. For example, you can use the following format to indicate that you are replacing row2 with the result of row1 + row2:

$$r_2 \leftarrow r_1 + r_2 : \begin{pmatrix} 3 & -7 & -2 & 5 \\ 0 & -2 & -1 & 5 \\ 6 & -4 & 0 & 1 \end{pmatrix}.$$

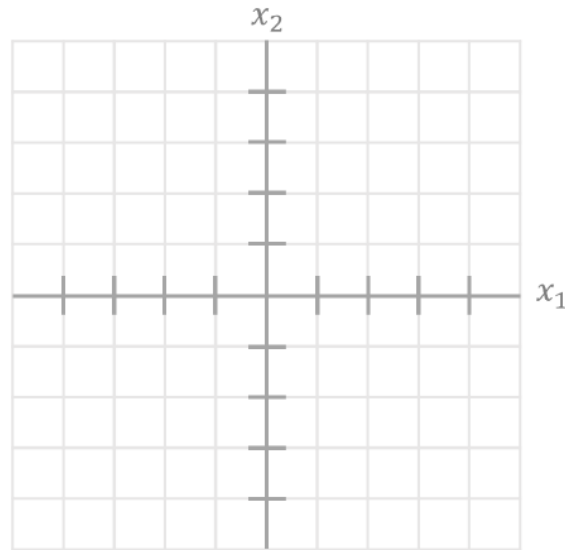
Then, as you back-substitute to find the solutions, state the value you assign to each variable in turn, and what equation and previously-assigned variables you are using. For example,  $x_2 + \frac{1}{12}x_3 = -\frac{5}{2}$ , and  $x_3 = -16$ , so  $x_2 = \frac{11}{2}$ . The solution to this equation system is  $x_1 = \frac{23}{6}$ ,  $x_2 = \frac{11}{2}$ ,  $x_3 = -16$ .

Solution

Work

## 2 Plotting Linear (Affine) Functions (7 points)

1. (3 points) Given  $\mathbf{w} = [2, -4]^T$  and  $b = -8$ , draw the line  $\{\mathbf{x} | \mathbf{w}^T \mathbf{x} + b = 0, \mathbf{x} \in \mathbb{R}^2\}$  in the 2-D plane. Clearly indicate the axes and origin of the plane.



2. Suppose we want to move the line closer to the origin.
- (a) (2 points) For a fixed  $\mathbf{w}$ , how should we change  $|b|$ ?
- ☐ Increase   ☐ Decrease   ☐ Not Possible
- (b) (2 points) For a fixed  $b$ , how should we change  $\|\mathbf{w}\|_2$ ?
- ☐ Increase   ☐ Decrease   ☐ Not Possible

### 3 Projection Matrices (10 points)

A projection matrix  $\mathbf{P}$  maps the vector of response values (dependent variable values) to the vector of fitted values (or predicted values). It describes the influence each response value has on each fitted value. Suppose we are estimating a linear model using least squares:

$$\mathbf{y} = \beta\mathbf{X} + \epsilon$$

Then in this case our projection matrix is given by,

$$\mathbf{P} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$$

1. (5 points) Prove that the projection matrix is necessarily symmetric. A matrix  $\mathbf{A}$  is said to be symmetric if  $\mathbf{A} = \mathbf{A}^T$ .

2. (5 points) Prove that the projection matrix is necessarily idempotent. A matrix  $\mathbf{A}$  is said to be idempotent if  $\mathbf{A}\mathbf{A} = \mathbf{A}$ .



## 4 Linear Regression (4 points)

In statistics, linear regression is a linear approach to modeling the relationship between a scalar response (or dependent variable) and one or more explanatory variables (or independent variables). Relationships are modeled using linear predictor functions whose unknown model parameters are estimated from the data. In this problem, we will try to fit a regression model to a training dataset  $(\mathbf{X}, \mathbf{Y})$  of  $n$  datapoints where  $\mathbf{X} \in \mathbb{R}^{n \times m}$  and  $\mathbf{Y} \in \mathbb{R}^{n \times 1}$ . Now, given a training set, we would like to learn a parameter  $\boldsymbol{\theta}$ ,  $h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta} \cdot \mathbf{x}$ , such that  $h_{\boldsymbol{\theta}}(\mathbf{x})$  is close to  $y$ , at least for the training examples we have. To formalize this, we will define a function that measures, for each value of  $\boldsymbol{\theta}$ , how close the  $h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$ 's are to the corresponding  $y^{(i)}$ 's. We define the objective function:

$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n ((h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)}))^2$$

To minimize  $J$ , we find its gradient with respect to  $\boldsymbol{\theta}$ .

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \frac{1}{2} (\mathbf{X}\boldsymbol{\theta} - \mathbf{Y})^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{Y}) \\ &= \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{Y} \end{aligned}$$

1. (2 points) To minimize  $J$ , we set its gradient to zero, and obtain the normal equations

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{Y} = 0$$

where  $\mathbf{X} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{Y} \in \mathbb{R}^{n \times 1}$ ,  $\boldsymbol{\theta} \in \mathbb{R}^{m \times 1}$ . **Solve for  $\boldsymbol{\theta}$ .**

Now suppose we used an  $l_2$  regularizer. Now, our objective function is,

$$J'(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n ((h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - (y^{(i)}))^2 + \frac{1}{2} \lambda \|\boldsymbol{\theta}\|^2$$

and its gradient is,

$$\nabla_{\boldsymbol{\theta}} J'(\boldsymbol{\theta}) = \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{Y} + \lambda \boldsymbol{\theta}$$

1. (2 points) To minimize  $J'$ , we set this new gradient of regularized linear regression to zero and obtain the equation:

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{Y} + \lambda \boldsymbol{\theta} = 0$$

where  $\mathbf{X} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{Y} \in \mathbb{R}^{n \times 1}$ ,  $\boldsymbol{\theta} \in \mathbb{R}^{m \times 1}$ . **Solve for  $\boldsymbol{\theta}$ .**

## 5 Matrix Memories (5 points)

Matrix Memories store a single pattern pair  $\mathbf{s} \rightarrow \mathbf{t}$  by encoding the outer product of a target vector  $\mathbf{t}$  with the input stimulus  $\mathbf{s}$ . In this problem, we will examine exactly when it is possible to store multiple pattern pairs. Suppose we wish to store the  $K$  pattern pairs below. Note that the  $k$ th stimulus  $\mathbf{s}^{(k)} = \begin{bmatrix} s_1^{(k)} & s_2^{(k)} & s_3^{(k)} & \dots & s_m^{(k)} \end{bmatrix}^T$  is paired with the  $k$ th target  $\mathbf{t}^{(k)} = \begin{bmatrix} t_1^{(k)} & t_2^{(k)} & t_3^{(k)} & \dots & t_m^{(k)} \end{bmatrix}^T$ —the superscript  $(k)$  is simply a label indicating which pair we are referring to.

$$\begin{bmatrix} s_1^{(1)} & s_2^{(1)} & s_3^{(1)} & \dots & s_m^{(1)} \end{bmatrix}^T \rightarrow \begin{bmatrix} t_1^{(1)} & t_2^{(1)} & t_3^{(1)} & \dots & t_m^{(1)} \end{bmatrix}^T$$

$$\begin{bmatrix} s_1^{(2)} & s_2^{(2)} & s_3^{(2)} & \dots & s_m^{(2)} \end{bmatrix}^T \rightarrow \begin{bmatrix} t_1^{(2)} & t_2^{(2)} & t_3^{(2)} & \dots & t_m^{(2)} \end{bmatrix}^T$$

$$\begin{bmatrix} s_1^{(K)} & s_2^{(K)} & s_3^{(K)} & \dots & s_m^{(K)} \end{bmatrix}^T \rightarrow \begin{bmatrix} t_1^{(K)} & t_2^{(K)} & t_3^{(K)} & \dots & t_m^{(K)} \end{bmatrix}^T$$

A Matrix Memory for multiple pattern pairs encodes the weight matrix as the sum of the outer products of the target/stimulus pairs:  $\mathbf{W} = \sum_{k=1}^K \mathbf{t}^{(k)} \otimes \mathbf{s}^{(k)} = \sum_{k=1}^K \mathbf{t}^{(k)} (\mathbf{s}^{(k)})^T$ , where  $(\mathbf{s}^{(k)})^T$  is the transpose of  $\mathbf{s}^{(k)}$ . Under this definition we have that the  $i, j$ th entry in  $\mathbf{W}$  is given by:

$$W_{ij} = \sum_{k=1}^K t_i^{(k)} (s_j^{(k)})^T$$

The Matrix Memory takes a new stimulus vector  $\mathbf{s}$  as input and computes the output response as  $\mathbf{r} = \mathbf{W}\mathbf{s}$ . If we wish to compute the response vector corresponding to the  $(l)$ th original pattern pair, we do the same  $\mathbf{r}^{(l)} = \mathbf{W}\mathbf{s}^{(l)}$ .

- (5 points) **Prove that** for the  $l$ th response vector  $\mathbf{r}^{(l)}$  for the  $l$ th stimulus  $\mathbf{s}^{(l)}$  will equal the target  $\mathbf{t}^{(l)}$  if all pairs of stimulus vectors  $(\mathbf{s}^{(k)}, \mathbf{s}^{(l)}) \forall k \neq l$  are orthonormal to each other.

## 6 Neural Networks (5 points)

A neural network by its definition always includes a nonlinear function like the sigmoid function. However, we could create one that has no nonlinear function. Suppose we have a function  $f$  defined as  $\mathbf{y} = f(\mathbf{x}) = \mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$ , where  $\mathbf{W}^{(1)}$ ,  $\mathbf{b}^{(1)}$ ,  $\mathbf{W}^{(2)}$ ,  $\mathbf{b}^{(2)}$  are the parameters of the function  $f$  for  $x \in \mathbb{R}^{d \times 1}$ ,  $\mathbf{W}^{(1)} \in \mathbb{R}^{k \times d}$ ,  $\mathbf{W}^{(2)} \in \mathbb{R}^{m \times k}$ ,  $\mathbf{b}^{(1)} \in \mathbb{R}^{k \times 1}$  and  $\mathbf{b}^{(2)} \in \mathbb{R}^{m \times 1}$ . We have a second function  $g$  defined as  $\mathbf{y} = g(\mathbf{x}) = \mathbf{U}\mathbf{x} + \mathbf{c}$  with parameters  $\mathbf{U}$  and  $\mathbf{c}$ .

1. (5 points) Now suppose we know the parameters of  $f$ . How can we define the parameters of  $g$  to ensure that  $g(\mathbf{x}) = f(\mathbf{x}), \forall \mathbf{x}$ ?

## 7 Collaboration Questions

After you have completed all other components of this assignment, report your answers to these questions regarding the collaboration policy. Details of the policy can be found in the syllabus.

1. Did you receive any help whatsoever from anyone in solving this assignment? If so, include full details.
2. Did you give any help whatsoever to anyone in solving this assignment? If so, include full details.
3. Did you find or come across code that implements any part of this assignment? If so, include full details.

Your Answer