

CM146, Fall 2017
Problem Set 0: Math prerequisites
Due Oct 09, 2017 at 11:59 pm

Submission instructions

- Submit your solutions electronically on the course Gradescope site as PDF files.
- If you plan to typeset your solutions, please use the LaTeX solution template. If you must submit scanned handwritten solutions, please use a black pen on blank white paper and a high-quality scanner app.

Although many students find a machine learning class to be rewarding, we do assume that you have a basic familiarity with several types of math. Before taking the class, you should evaluate whether you have the mathematical background the class depends upon.

- **Multivariate Calculus** (at the level of a first undergraduate course, *e.g.*, Math 32A and 32B at UCLA). For example, we rely on you being able to take derivatives and integrals. During the class you might be asked, for example, to derive gradients of multivariate functions.
- **Linear Algebra** (at the level of a first undergraduate course, *e.g.*, Math 33A). For example, we assume you know how to multiply vectors and matrices, and that you understand matrix inversion, eigenvectors and eigenvalues. During the class, you might also be asked to also learn about methods for matrix factorization.
- **Probability and Statistics** (at the level of a first undergraduate course, *e.g.*, Statistics 100A). For example, we assume you know how to find the mean and variance of a set of data, that you are familiar with common probability distributions such as the Gaussian and Uniform distributions, and that you understand basic notions such as conditional probabilities and Bayes rule. During the class, you might be asked to calculate the likelihood (probability) of a data set with respect to some given probability distribution, and to then derive the parameters of the distribution that maximize this likelihood.

This assignment helps you self-evaluate whether you have the background to succeed in the class. For each of these mathematical topics, we provide below (1) a minimum background test and (2) a moderate background test. If you pass the moderate background test, you are in excellent shape to take the class. If you pass the minimum background but not the moderate background test, then you can still take the class, but you should expect to devote extra time to fill in necessary math background. If you cannot pass the minimum background test, we suggest you fill in your math background before taking the class.

You may find the following resources helpful:

This assignment is adapted from course material by William Cohen, Ziv Bar-Joseph (CMU) and Jessica Wu (Harvey Mudd).

- Andrew Ng's CS229 Course (Stanford)
 - Linear Algebra Review (<http://cs229.stanford.edu/section/cs229-linalg.pdf>)
 - Probability Theory Review (<http://cs229.stanford.edu/section/cs229-prob.pdf>)

Additional resources are available on the course syllabus.

Necessary Minimum Background Test [45 pts]

While you are welcome to use online resources, such as Wolfram-Alpha, you should be able to solve these problems by hand.

1 Multivariate Calculus [2 pts]

Consider $y = x \sin(z)e^{-x}$. What is the partial derivative of y with respect to x ?

2 Linear Algebra [8 pts]

Consider the matrix \mathbf{X} and the vectors \mathbf{y} and \mathbf{z} below:

$$\mathbf{X} = \begin{pmatrix} 2 & 4 \\ 1 & 3 \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 1 \\ 3 \end{pmatrix} \quad \mathbf{z} = \begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

- (a) What is the inner product $\mathbf{y}^T \mathbf{z}$?
- (b) What is the product $\mathbf{X}\mathbf{y}$?
- (c) Is \mathbf{X} invertible? If so, give the inverse; if not, explain why not.
- (d) What is the rank of \mathbf{X} ?

3 Probability and Statistics [10 pts]

Consider a sample of data S obtained by flipping a coin five times. $X_i, i \in \{1, \dots, 5\}$ is a random variable that takes a value 0 when the outcome of coin flip i turned up heads, and 1 when it turned up tails. Assume that the outcome of each of the flips does not depend on the outcomes of any of the other flips. The sample obtained $S = (X_1, X_2, X_3, X_4, X_5) = (1, 1, 0, 1, 0)$.

- (a) What is the sample mean for this data?
- (b) What is the unbiased sample variance ?
- (c) What is the probability of observing this data assuming that a coin with an equal probability of heads and tails was used? (*i.e.*, The probability distribution of X_i is $P(X_i = 1) = 0.5$, $P(X_i = 0) = 0.5$.)
- (d) Note the probability of this data sample would be greater if the value of the probability of heads $P(X_i = 1)$ was not 0.5 but some other value. What is the value that maximizes the probability of the sample S ? [Optional: Can you prove your answer is correct?]
- (e) Given the following joint distribution between X and Y , what is $P(X = T|Y = b)$?

$P(X, Y)$		Y		
		a	b	c
X	T	0.2	0.1	0.2
	F	0.05	0.15	0.3

4 Probability axioms [5 pts]

Let A and B be two discrete random variables. In general, are the following true or false? (Here A^c denotes complement of the event A .)

(a) $P(A \cup B) = P(A \cap (B \cap A^c))$

(b) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

(c) $P(A) = P(A \cap B) + P(A^c \cap B)$

(d) $P(A|B) = P(B|A)$

(e) $P(A_1 \cap A_2 \cap A_3) = P(A_3|(A_2 \cap A_1))P(A_2|A_1)P(A_1)$

5 Discrete and Continuous Distributions[5 pts]

Match the distribution name to its formula.

- | | |
|-----------------|--|
| (a) Gaussian | (i) $p^x(1-p)^{1-x}$, when $x \in \{0, 1\}$; 0 otherwise |
| (b) Exponential | (ii) $\frac{1}{b-a}$ when $a \leq x \leq b$; 0 otherwise |
| (c) Uniform | (iii) $\binom{n}{x} p^x (1-p)^{n-x}$ |
| (d) Bernoulli | (iv) $\lambda e^{-\lambda x}$ when $x \geq 0$; 0 otherwise |
| (e) Binomial | (v) $\frac{1}{\sqrt{(2\pi)\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$ |

6 Mean and Variance[5 pts]

(a) What is the mean and variance of a *Bernoulli*(p) random variable?

(b) If the variance of a zero-mean random variable X is σ^2 , what is the variance of $2X$? What about the variance of $X + 2$?

7 Algorithms [10 pts]

(a) **Big-O notation**

For each pair (f, g) of functions below, list which of the following are true: $f(n) = O(g(n))$, $g(n) = O(f(n))$, or both. Briefly justify your answers.

i. $f(n) = \ln(n), g(n) = \lg(n)$. Note that \ln denotes log to the base e and \lg denotes log to the base 2.

ii. $f(n) = 3^n, g(n) = n^{10}$

iii. $f(n) = 3^n, g(n) = 2^n$

(b) **Divide and Conquer**

Assume that you are given an array with n elements all entries equal either to 0 or +1 such that all 0 entries appear before +1 entries. You need to find the index where the transition happens, *i.e.*, you need to report the index with the last occurrence of 0. Give an algorithm that runs in time $O(\log n)$. Explain your algorithm in words, describe why the algorithm is correct, and justify its running time.

Moderate Background Test [35 pts]

8 Probability and Random Variables [5 pts]

- (a) **Mutual and Conditional Independence** If X and Y are independent random variables, show that $\mathbb{E}[XY] = \mathbb{E}[X]E[Y]$.

- (b) **Law of Large Numbers and Central Limit Theorem**

Provide one line justifications.

- i. If a fair die is rolled 6000 times, the number of times 3 shows up is close to 1000.

- ii. If a fair coin is tossed n times and \bar{X} denotes the average number of heads, then the distribution of \bar{X} satisfies

$$\sqrt{n}(\bar{X} - \tfrac{1}{2}) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, \tfrac{1}{4})$$

9 Linear Algebra [20 pts]

(a) **Vector Norms** [4 pts]

Draw the regions corresponding to vectors $\mathbf{x} \in \mathbb{R}^2$ with following norms (you can hand draw or use software for this question):

- i. $\|\mathbf{x}\|_2 \leq 1$ (Recall $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$.)
- ii. $\|\mathbf{x}\|_0 \leq 1$ (Recall $\|\mathbf{x}\|_0 = \sum_{i:x_i \neq 0} 1$.)
- iii. $\|\mathbf{x}\|_1 \leq 1$ (Recall $\|\mathbf{x}\|_1 = \sum_i |x_i|$.)
- iv. $\|\mathbf{x}\|_\infty \leq 1$ (Recall $\|\mathbf{x}\|_\infty = \max_i |x_i|$.)

(b) **Matrix Decompositions [6 pts]**

i. Give the definition of the eigenvalues and the eigenvectors of a square matrix.

ii. Find the eigenvalues and eigenvectors of

$$\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}.$$

- iii. For any positive integer k , show that the eigenvalues of \mathbf{A}^k are $\lambda_1^k, \lambda_2^k, \dots, \lambda_n^k$, the k^{th} powers of the eigenvalues of matrix \mathbf{A} , and that each eigenvector of \mathbf{A} is still an eigenvector of \mathbf{A}^k .

(c) **Vector and Matrix Calculus [5 pts]**

Consider the vectors \mathbf{x} and \mathbf{a} and the symmetric matrix \mathbf{A} .

i. What is the first derivative of $\mathbf{a}^T \mathbf{x}$ with respect to \mathbf{x} ?

ii. What is the first derivative of $\mathbf{x}^T \mathbf{A} \mathbf{x}$ with respect to \mathbf{x} ? What is the second derivative?

(d) **Geometry [5 pts]**

- i. Show that the vector \mathbf{w} is orthogonal to the line $\mathbf{w}^T \mathbf{x} + b = 0$. (Hint: Consider two points $\mathbf{x}_1, \mathbf{x}_2$ that lie on the line. What is the inner product $\mathbf{w}^T(\mathbf{x}_1 - \mathbf{x}_2)$?)

- ii. Argue that the distance from the origin to the line $\mathbf{w}^T \mathbf{x} + b = 0$ is $\frac{b}{\|\mathbf{w}\|_2}$.

Programming Skills

Start familiarizing yourself with the Python libraries `numpy` and `matplotlib` by completing the following exercises. (You do not have to submit your code.)

You may find the following references helpful:

- http://docs.scipy.org/doc/numpy/reference/generated/numpy.random.multivariate_normal.html
- <http://docs.scipy.org/doc/numpy/reference/generated/numpy.linalg.eig.html>

10 Sampling from a Distribution [2.5 pts]

For questions (a-e), only submit your plots. You do not need to submit code.

- (a) Draw 1000 samples $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ from a 2-dimensional Gaussian distribution with mean $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$ and identity covariance matrix, i.e. $p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^2}} \exp\left(-\frac{\|\mathbf{x}\|^2}{2}\right)$, and make a scatter plot (x_1 vs x_2).
- (b) How does the scatter plot change if the mean is $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$?
- (c) How does the (original) scatter plot change if you double the variance of each component?
- (d) How does the (original) scatter plot change if the covariance matrix is changed to $\begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$?
- (e) How does the (original) scatter plot change if the covariance matrix is changed to $\begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix}$?

11 Eigendecomposition [2.5 pts]

Write a python program to compute the eigenvector corresponding to the largest eigenvalue of the following matrix and submit the computed eigenvector.

$$\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 1 & 3 \end{pmatrix}$$

.

12 Data [5 pts]

There are now lots of really interesting data sets publicly available to play with. They range in size, quality and the type of features and have resulted in many new machine learning techniques being developed.

Find a public, free, supervised (i.e. it must have features *and* labels), machine learning dataset. You may NOT list a data set from 1) The UCI Machine Learning Repository or 2) from Kaggle.com. Once you have found the data set, provide the following information:

- (a) The name of the data set.
- (b) Where the data can be obtained.
- (c) A brief (i.e. 1-2 sentences) description of the data set including what the features are and what is being predicted.
- (d) The number of examples in the data set.
- (e) The number of features for each example. If this is not concrete (i.e. it is text), then a short description of the features.

For this question, do not just copy and paste the description from the website or the paper; reference it, but use your own words. Your goal here is to convince the staff that you have taken the time to understand the data set, where it came from, and potential issues involved.