

Bi/BE/CS 183 2021-2022
Instructor: Lior Pachter
TAs: Tara Chari, Meichen Fang, Zitong (Jerry) Wang

Problem Set 5 - Midterm

You are not allowed to collaborate with others for this Midterm, though you may ask clarification questions on Piazza or at Office Hours. This Midterm is open note, which includes Lecture slides and your personal notes.

Submit your solutions as a single PDF file via Canvas by **10am Friday February 11th**.

- If writing up problems by hand, please use a pen and not a pencil, as it is difficult to read scanned submission of pencil work. Typed solutions are preferred.
- For problems that require coding, Colab notebooks will be provided. Please copy and save the shared notebook and edit your own copy, which you should then submit by including a clickable link in your submitted homework. Prior to submission make sure that you code runs from beginning to end without any error reports.

Problem 1 (10 points)

Given a Poisson random variable X , with probability mass function $P(X = x) = e^{-\lambda} \lambda^x / x!$, show that $E(X) = \lambda$. (Hint: $e^x = \sum_{n=0}^{\infty} x^n / n!$)

Problem 2 (16 points)

Suppose $x = (x_1, x_2, \dots, x_n)$ are i.i.d. observations from a Poisson random variable with unknown parameter λ .

- (a) (8 points) Write down the log-likelihood function $L(\lambda | x)$ for this set of observations.
- (b) (8 points) Find the maximum-likelihood estimator of λ by taking the appropriate derivative of $L(\lambda | x)$. (Hint: the estimator you obtain should be a function of the set of observations x)

Problem 3 (24 points)

Consider two independent geometric random variables, X and Y , with identical probability mass function $P(X = k) = (1 - p)^{k-1} p$, where $k = 1, 2, 3, \dots$ and p is a parameter of the distribution.

- (a) (8 points) Show that $P(X + Y = k) = (k - 1)(1 - p)^{k-2} p^2$ for $k = 2, 3, \dots$
- (b) (8 points) We say that Z has a negative binomial distribution with parameters r, p if,

$$P(Z = z) = \binom{z-1}{r-1} p^r (1-p)^{z-r}, z = r, r+1, \dots$$

Show that $X + Y$ has a negative binomial distribution with parameters $r = 2$ and p (note that this parametrization of the negative binomial distribution is different from the one given in problem 5).

- (c) (8 points) see Colab notebook here

Problem 4 (10 points)

Construct the suffix tree for the word “Yellowwooddoor”.

Problem 5 (40 points)

For this problem you will be exploring various models which can be used to describe count data i.e. the gene-count matrices we use in single-cell genomics. The Problem 5 notebook is here.

Single-cell gene counts, which describe stochastically sampled, discrete measurements of UMI counts, are often modeled as being generated from a negative binomial (or Gamma-Poisson) distribution. However, there is a common assumption that droplet-based methods for single-cell RNA seq incur an overabundance of zeros (more zero counts) than would be predicted by random sampling. Thus it is also common to see single-cell data modeled with zero-inflated negative binomials (the ZINB distribution, with an extra parameter for the probability of zero counts). Here you will analyze these zero-inflation assumptions, following work done in [1].

Your edited version of the notebook *must be submitted* for this problem. Reminder to check that your notebook runs all the way through with the the `Runtime` → `Restart` and `Runtime` → `Run All` commands.

References

1. Svensson, V. Droplet scRNA-seq is not zero-inflated. en. *Nat. Biotechnol.* **38**, 147–150 (Feb. 2020).