

Problem Set 6

Due: Wednesday, May 18, by 11:59pm

Instructions

Solutions format, collaboration policy, and late policy. See PSet 1 for further details. The same requirements and policies still apply. Also follow the typesetting instructions from the prior PSets.

Solutions submission. You must submit your solution via Gradescope. In particular:

- For the solutions to Task 1-5, submit under “PSet 6 [Written]” a *single* PDF file containing the solution to all tasks in the homework. Each numbered task should be solved on its own page (or pages). Follow the prompt on Gradescope to link tasks to your pages. Do not write your name on the individual pages – Gradescope will handle that.
- For the programming part (Task 6), submit your code under “PSet 6 [Coding]” as one file, `min_hash.py`.

Task 1 – Normal, normal, normal

[15 pts]

Your answers should be correct to four decimal places.

- Apparently IQ is roughly normally distributed, with a mean of 100 and a standard deviation of about 15. What fraction of people would be classified as “genius” if that means IQ of 140 or above? (These numbers are made up.)
- The height of a group of people is approximately normally distributed with a mean of 69.2 inches (5' 9.2") and a standard deviation of 2.67 inches. Approximately what fraction of these people are 5'2" tall or less? If we form a basketball team by picking 5 people from this group at random, calculate the probability that at least one of them is over 6' 3" (75 inches) tall.

Task 2 – CLT in “real life”

[10 pts]

Error-correcting codes¹ are used in order to compensate for errors in transmission of messages (and in recovery of stored data from unreliable hardware). You are on a mission to Mars and need to send regular updates to mission control. Most of the packets actually don't get through, but you are using an error-correcting code that can let mission control recover the original message you send so as long as at least 1024 packets are received (not erased). Suppose that each packet gets erased independently with probability 0.7. How many packets should you send such that you can recover the message with probability at least 99%

Use the Central Limit Theorem to approximate the answer, using the continuity correction. Your answers should be correct to 4 decimal places.

¹From the Wikipedia page: “In computing, telecommunication, information theory, and coding theory, an *error correcting code* (ECC) is used for controlling errors in data over unreliable or noisy communication channels. The central idea is that the sender encodes the message with redundant information in the form of an ECC. The redundancy allows the receiver to detect a limited number of errors that may occur anywhere in the message, and often to correct these errors without retransmission.”

Task 3 – Statistics Books

[25 pts]

Alice is going shopping for statistics books for H hours, where H is a random variable, equally likely to be 1, 2 or 3. The number of books B she buys is random and depends on how long she is in the store for. We are told that

$$\mathbb{P}(B = b \mid H = h) = \frac{1}{h}, \quad \text{for } b = 1, \dots, h.$$

- Find the joint distribution of B and H using the chain rule.
- Find the marginal distribution of B .
- Find the conditional distribution of H given that $B = 1$ (i.e., $\mathbb{P}(H = h \mid B = 1)$ for each possible h in 1,2,3). Use the definition of conditional probability and the results from previous parts.
- Suppose that we are told that Alice bought either 1 or 2 books. Find the expected number of hours she shopped conditioned on this event. Use the definition of conditional expectation and Bayes Theorem.
- The bookstore has a discounting policy that gives an extra 10% off the total purchase price if Alice buys two books and 20% off the total purchase price if she buys three books. Suppose that Alice's decision about what books to buy does not depend on their price and that, in an ironic twist, the bookstore owner also prices each statistics book randomly with a mean price of \$40 per book. What is the expected amount of money Alice spends (assuming that book purchases are tax-free)?

Warning: Be sure to use a formal derivation. Your work should involve the law of total expectation conditioning on the number of books bought, and make use of random variables X_i , where X_i is the amount of money she spends on the i th book she purchases.

Task 4 – Duck Hunt

[15 pts]

Ten hunters are waiting for ducks to fly by. When a flock of ducks flies overhead, hunters fire at the same time, but each chooses their target uniformly at random, independently of the others. If each hunter independently hits their target with probability 0.6, use the law of total expectation to compute the expected number of ducks that are hit. Assume that the number of ducks in a flock is a Poisson random variable with mean 6.

Task 5 – Joint Densities

[15 pts]

Suppose that X, Y are jointly continuous rv's with joint density

$$f_{X,Y}(x, y) = \begin{cases} 4xy & 0 \leq x \leq 1, 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

(Observe that this is a probability density function since it is non-negative and

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(xy) dy dx = \int_0^1 \int_0^1 4xy dy dx = \int_0^1 \left(2xy^2 \Big|_{y=0}^{y=1} \right) dx = \int_0^1 2x dx = x^2 \Big|_{x=0}^{x=1} = 1.)$$

Your answers below should **not** be evaluated unless otherwise specified. Your answers should usually be in terms of integrals or double integrals.

- Write an expression which we can evaluate to find $\mathbb{P}(Y \geq X)$. Hint: draw the region of the joint density, and the desired region.
- Write an expression which we can evaluate to find the marginal density $f_X(x)$. Specify the value of $f_X(x)$ for all $x \in \mathbb{R}$. Do the same for $f_Y(y)$.
- Are X and Y independent? Justify your answer. (You may need to evaluate an integral or two to do this.)

Task 6 – Distinct Elements

[20 pts]

Recall the setup for the MinHash algorithm presented in class. The universe of is the set \mathcal{U} (think of this as the set of all 8-byte integers), and we have a single **uniform** hash function $h : \mathcal{U} \rightarrow [0, 1]$. That is, for an integer y , pretend $h(y)$ is a **continuous** $\text{Unif}(0, 1)$ random variable. That is, $h(x_1), h(x_2), \dots, h(x_N)$ for any N **distinct** elements are iid continuous $\text{Unif}(0, 1)$ random variables, but since the hash function always gives the same output for some given input, if, for example, the i -th user ID x_i and the j -th user ID x_j are the same, then $h(x_i) = h(x_j)$ (i.e., they are the “same” $\text{Unif}(0, 1)$ random variable).

Then, the MinHash algorithm is realized by the following pseudocode, which explains its two key functions:

1. `UPDATE(x)`: How to update your variable when you see a new stream element.
2. `ESTIMATE()`: At any given time, how to estimate the number of distinct elements you’ve seen so far.

Note that this differs from the syntax used on the slides, but captures the same algorithm.

MinHash Operations

```
function INITIALIZE()
    val ← ∞
function UPDATE(x)
    val ← min {val, h(x)}
function ESTIMATE() return round( $\frac{1}{\text{val}} - 1$ )
for  $i = 1, \dots, N$ : do
    UPDATE( $x_i$ )
return ESTIMATE()
```

▷ Loop through all stream elements
▷ Update our single float variable
▷ An estimate for n , the number of distinct elements.

To help you out with the following questions, we have set up an [edstem lesson](#). However, you are required to upload your final solution to Gradescope (see instructions above).

- a) Implement the functions `UPDATE` and `ESTIMATE` in the MinHash class of [min.hash.py](#).
- b) The estimator we used in a) has high variance, and therefore it may not always give good answer. As outlined in class, we improve this by considering k variables

$$\text{val}_1, \text{val}_2, \dots, \text{val}_k$$

where each of val_i , $1 \leq i \leq k$ is an i.i.d. random variable with the distribution of the minimum of $m \leq N$ independent $\text{Unif}(0, 1)$ variables, obtained by hashing the N elements in the stream with independent hash functions h^1, \dots, h^k . Our final estimate will then be

$$\hat{n} = \frac{1}{\widehat{\text{val}}} - 1 \quad \text{where} \quad \widehat{\text{val}} = \frac{1}{k} \sum_{i=1}^k \text{val}_i.$$

Implement the functions `UPDATE` and `ESTIMATE` in the `MultMinHash` class of [min.hash.py](#) using the improved estimator.

Refer to [Section 9.5](#) of the book for more details on the distinct elements algorithm.