

## Continuous Joint Distributions, Central Limit Theorem

---

Before you leave lab, make sure you [click here](#) so that you're marked as having attended this week's section. The CA leading your discussion section can enter the password needed once you've submitted.

### 1 Warmups

#### 1.1 Sums of Random Variables

For each of the problems below, assume that  $X$  and  $Y$  are independent.

1. Let  $X \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $Y \sim \mathcal{N}(\mu_2, \sigma_2^2)$ . What is  $\mu$  and  $\sigma^2$  for  $X + Y \sim \mathcal{N}(\mu, \sigma^2)$ ?
2. Let  $X \sim \text{Uni}(0, 1)$  and  $Y \sim \text{Uni}(0, 1)$ . What is the PDF for  $X + Y$ ?
3. In general, two random variables  $X$  and  $Y$ , what is the PDF  $f$  of  $X + Y$ ?

#### 1.2 Food for Thought

Karel the dog eats an unpredictable amount of food. Every day, the dog is equally likely to eat between a continuous amount in the range 100 to 300g. How much Karel eats is independent of all other days. You only have 6.5kg of food for the next 30 days. What is the probability that 6.5kg will be enough for the next 30 days?

### 2 Problems

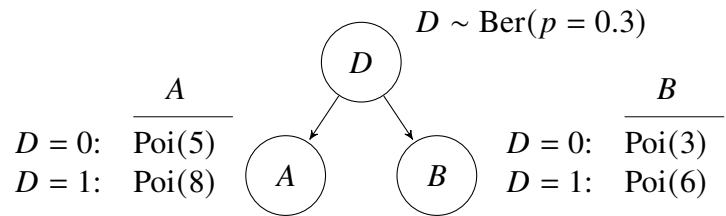
#### 2.1 Grading Exams

Jacob and Kathleen are planning to grade Problem 1 on your Week 7 exam, and they'll each grade their half independently of the other. Jacob takes  $X \sim \text{Exp}(\frac{1}{3})$  hours to finish his half while Kathleen takes  $Y \sim \text{Exp}(\frac{1}{4})$  hours to finish his half.

- a. Find the CDF of  $X/Y$ , which is the ratio of their grading completion times.
- b. What is the probability that Kathleen finishes before Jacob does?

## 2.2 Fish Sticks [courtesy of Lisa Yan]

Fish Sticks, the online platform designed to meet all of your fish stick needs, wants to model their hourly homepage traffic from Stanford. The company decides to model two different behaviors for homepage visits according to the Bayesian Network on the right:



$A$  and  $B$  are the numbers of Stanford students and faculty, respectively, who visit the Fish Sticks homepage in an hour. Since Fish Sticks does not know when Stanford people eat, the company models demand as a "hidden" Bernoulli random variable  $D$ , which determines the distribution of  $A$  and  $B$ . Recall that in a Bayesian Network, random variables are conditionally independent given their parents. For example, given  $D = 0$ ,  $A \sim \text{Poi}(5)$  and  $B \sim \text{Poi}(3)$ , two independent random variables.

- Given that 6 users from group  $A$  visit the homepage in the next hour, what is the probability that  $D = 0$ ?
- What is the probability that in the next hour, the *total* number of users who visit the homepage from groups  $A$  and  $B$  is equal to 12, i.e., what is  $P(A + B = 12)$ ?
- Now simulate  $P(A + B = \text{total})$ , where  $\text{total} = 12$ , by implementing the `infer_prob_total(total, ntrials)` function below using rejection sampling.
  - `total` is the total number of users from groups  $A$  and  $B$  in the event  $A + B = \text{total}$ .
  - `ntrials` is the number of observations to generate for rejection sampling.
  - `prob` is the return value to the function, where  $\text{prob} \approx P(A + B = \text{total})$ .
  - The function call is implemented for you at the bottom of the code block.

You can call the following functions from the `scipy` package:

- `stats.bernoulli.rvs(p)`, which randomly generates a 1 with probability  $p$ , and generates a 0 otherwise.
- `stats.poisson.rvs( $\lambda$ )`, which randomly generates a value according to a Poisson distribution with parameter  $\lambda$

You are not required to use lists or `numpy` arrays in this question (but you can if you want). **Pseudo-code is fine** as long as your code accurately conveys your approach.

```
import numpy as np
from scipy import stats

def infer_prob_total(total, ntrials):
    # here's where your implementation belongs
    return prob
```

```
total = 12
ntrials = 50000
print('Simulated% P(A + B)=', infer_prob_total(total, ntrials))
```

This is the full implementation right here:

```
import numpy as np
from scipy import stats

def infer_prob_total(total, ntrials):
    n_samples_event = 0
    for i in range(ntrials):
        d = stats.bernoulli.rvs(0.3)
        if d == 0:
            user_sum = stats.poisson.rvs(5) + stats.poisson.rvs(3)
        else:
            user_sum = stats.poisson.rvs(8) + stats.poisson.rvs(6)
        if user_sum == total: n_samples_event += 1
    prob = n_samples_event/ntrials
    return prob

ntrials = 50000
total = 12
print("Simulated P(A + B)=", infer_prob_total(total, ntrials))
```

### 3 ChatGPT, Watermarking, and Bayesian Inference

ChatGPT is a generative AI technology that can be coarsely summarized to be a chat bot with a seemingly boundless ability to discuss any topic—history, computer science, art, nuclear physics, probability, and even the ethics of using ChatGPT—in any one of several written languages, including English, French, Spanish, C++, JavaScript, Python, and some 100 others.

Unsurprisingly, we will soon be left to wonder whether a poem, a Tweet, a C++ function, or a college thesis is written by ChatGPT or a human being. Questions about authorship, accuracy, and attribution have prompted OpenAI, the company behind ChatGPT, to address these concerns by implementing ChatGPT to employ what’s termed **watermarking** and insert certain words more or less often than is customary in even the best of human-authored writing.

To illustrate, let’s assume that most humans use the word **the** an average of 4.8 times per 100 words, whereas ChatGPT might generate prose where **the** appears on average 6.5 times per 100 words. Similarly, humans use the word **of** on average about 3.9 times per 100 words, whereas ChatGPT might leverage **of** about 6.2 times per 100 words. Conversely, ChatGPT might generate the word **by** only 1.6 times per 100 words, whereas humans use the word **by** about 2.7 times per 100 words.

- a. You elect to model word frequency of all words using either a Poisson for paragraphs of 200 or so words—but as a Gaussian for larger documents—say 10000 words or more. Explain why the Gaussian might be the better choice for larger documents than the Poisson, whereas Poisson is more easily defended for smaller documents.
- b. A deeper statistical analysis of many human-written documents strongly suggests that  $H_{\text{the}} \sim \mathcal{N}(5, 1)$  and  $H_{\text{by}} \sim \mathcal{N}(4, 1)$ , whereas a separate but equally deep analysis strongly suggests that  $C_{\text{the}} \sim \mathcal{N}(3, 1)$  and  $C_{\text{by}} \sim \mathcal{N}(2, 1)$ . (The parameters are rounded values for simplicity and assumed to be per-100-words.)

Assuming a prior belief that a very long document was written by a human is 0.99, what is your posterior belief that the document was written by a human if the document contains an average of 5 **the**'s every 100 words but only 1 **by** every 100 words. You may assume all Gaussian distributions of interest are independent.

- c. In reality, the Gaussians here are not independent, since the presence of one word implies the absence of another. Assuming a correlation value  $\rho$  that is slightly negative, would you expect the observations stated in part b) to result in a larger posterior probability or a smaller one? Briefly explain why.