

Lecture 10 Worksheet

Chrysafis Vogiatzis

Every worksheet will work as follows.

1. You will be asked to form a group with other students in the class: you can make this as big or as small as you'd like, but groups of 4-5 work best.
2. Read through the worksheet, discussing any questions with the other members of your group.
 - You can call me at any time for help!
 - I will also be interrupting you for general guidance and announcements at random points during the class time.
3. Answer each question (preferably in the order provided) to the best of your knowledge.
4. While collaboration between students is highly encouraged and expected, each student has to submit their own version.
5. You will have 24 hours (see gradescope) to submit your work.

Activity 1: Normally distributed random variables

We begin this section by mentioning something very important:

Adding two (or more) independent normal distributed random variables together results in a normal distribution.

In mathematical terms: if $X_i, i = 1, \dots, n$ are independent random variables distributed normally, then $Z = \sum_{i=1}^n X_i$ is also normally distributed. The same is true for any linear combination, i.e., $Z = \sum_{i=1}^n a_i \cdot X_i$ is normally distributed. Recall that to fully describe a normally distributed random variable, we simply need its expectation and variance.

Problem 1: Expectation and variance of the sum of normally distributed random variables

Consider two normally distributed random variables $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$, which are independent. What is the expectation and variance of $X + Y$?

Answer to Problem 1.

In general, $Z = \sum_{i=1}^n a_i \cdot X_i$ is normally distributed with $\mu = E[Z] = \sum_{i=1}^n a_i \cdot E[X_i]$ and $\sigma^2 = \text{Var}[Z] = \sum_{i=1}^n a_i^2 \cdot \text{Var}[X_i]$.

Problem 2: Application

You own a portfolio of stocks that consists of 5 stocks S_1 and 10 stocks S_2 . Let X_1 be the price of stock S_1 and X_2 the price of stock S_2 one year from now. X_1 is normally distributed with $\mathcal{N}(50, 100)$ ¹ and X_2 is normally distributed with $\mathcal{N}(60, 100)$. Last, assume that the two stocks are totally unrelated (i.e., they are independent of one another).

¹ That is, its mean is 50 and its variance is 100.

What is the probability that your portfolio is worth more than \$1000 one year from now? What is the probability that your portfolio is worth less than \$800 one year from now?

Answer to Problem 2.

Problem 3: Comparing random variables

Consider two normally distributed random variables X and Y , which are independent. Can you make the claim that $X - Y$ is normally distributed? What is the expectation and variance of $X - Y$?²

Answer to Problem 3.

² Sure we can! Remember that if X_1, X_2, \dots are normally distributed then any linear combination of them is also normally distributed!

Problem 4: Comparing stock prices

Again, consider that X_1 is the price of stock S_1 and X_2 is the price of stock S_2 one year from now. Both are normally distributed: $X_1 \sim \mathcal{N}(50, 100)$ and $X_2 \sim \mathcal{N}(60, 100)$. Finally, assume that X_1, X_2 are independent.

What is the probability that S_1 (a single stock of that) is worth more than S_2 (again, consider only a single stock of that one, too) a year from now? ³

Answer to Problem 4.

³ Think of $Z = X_1 - X_2$... If Z is a normally distributed random variable, then we can calculate the probability that $P(Z > 0)$, no?

What about the case where two random variables (distributed normally) are **not independent**? One could make the argument that two stocks affect each other, and knowing that one has gone up may affect our perspective of the other one going up, too. In this case, where independence is hard to assume and use, can we still compare two stocks (or calculate the probability of our portfolio being above a desired value a year from now)?

The answer is yes! However, we need to introduce ideas like **dependence** and measure it with **covariance** and **correlation**. More on that in Lecture 12...

Activity 2: Using the central limit theorem

In the online course world, certain classes have thousands of students: these classes are sometimes referred to as *Massive Open Online Courses* (MOOC). On average there are 8000 students in a course. For the purposes of this exercise, we assume that a student successfully finishes a MOOC 7.5% of the time and that all students behave independently. Assume that probability is the same regardless of the course. Finally, assume that every class has **exactly** 8000 students.

Answer the following questions.

Problem 5: Setting up the distribution

What is the best distribution to model the number of students to successfully finish a single class? What is the mean and the variance of the number of students who successfully finish a single class? ⁴

Answer to Problem 5.

⁴ Recall: you have 8000 students, each of whom may succeed or fail.. Which distribution is this? And what is its mean?

Problem 6: Setting up the central limit theorem

Now, consider the case of an online course provider, such as Coursera. Assume the total number of classes the provider offers is equal to 1000. For each individual class, the number of students who successfully finish that specific class follows the distribution you found in Problem 5.

What distribution does the average number of students graduating from all class of the provider follow? The average number of students graduating in all 1000 classes can be found by summing the number of graduates in each class and dividing by 1000. ⁵

Answer to Problem 6.

⁵ Does the central limit theorem apply? If so, then the distribution is...

Activity 3: Setting up hypothesis testing

Recall that we have made the hypothesis that a student successfully finishes a class 7.5% of the time. We have also found that the number of students successfully finishing one class follows a binomial distribution (see Problem 5) and the average number of students across 1000 courses follows a normal distribution (based on Problem 6).

The online course provider has decided to check whether this “7.5%” success rate is true or not. They have decided to survey 15 of the 1000 courses they offer and they found the following: ⁶

Course 1	588	Course 2	645	Course 3	632
Course 4	623	Course 5	635	Course 6	641
Course 7	644	Course 8	611	Course 9	630
Course 10	628	Course 11	569	Course 12	637
Course 13	635	Course 14	677	Course 15	610

⁶ You will not need the numbers until Problem 8.

Problem 7: Checking the numbers

What should the average number of graduating students in these 15 classes be distributed as? Once again, to find the average number of graduates, simply add up the numbers for the 15 classes and divide by $n = 15$. ⁷ What is the mean and the variance of the distribution?

⁷ Treat 15 as a large enough number for the central limit theorem to apply.

Answer to Problem 7.

Problem 8: Using the central limit theorem

From the numbers they found (see the tabulated data from Problem 7) that on average a pretty big 627 students in each class successfully finishes it ⁸.

⁸ Recall they were expecting 600 on average, or $8000 \cdot 0.075$.

Based on the distribution you have identified in Problem 7, what is the probability that there are more than 627 students on average finishing each class? ⁹

⁹ Remember that z values that do not appear in the table (i.e., are larger than 3.9) can be treated as corresponding to 1 (100%)!

Answer to Problem 8.

Problem 9: Huh?

Hopefully, you have gotten a very, very small ¹⁰ probability for Problem 8. This implies that the numbers they got appear to be *highly* improbable. While we are at it, let's calculate one more probability. What is the probability that the number of students that (on average) finish successfully each course is more than 610?

¹⁰ Even zero!

Answer to Problem 9.

Problem 10: Rejecting a hypothesis

Considering the small probability you got for the average to be as high as 610 or more, what can you deduce for the success rate of 7.5%? Should they believe it? Do *you* think it is valid? Or is the true success rate higher/lower? ¹¹

¹¹ Think: if the success rate is higher, and the mean of the normal distribution (for the average of 15 classes) is also higher, then does that mean the probabilities in Problem 8 and 9 go up or down?

Answer to Problem 10.

Congratulations, you just performed a fully-fledged data analysis experiment! We will focus a lot on this part after our second midterm.

