

Statistical Data Analysis

Prof. Dr. Nizamettin AYDIN

naydin@yildiz.edu.tr

<http://www3.yildiz.edu.tr/~naydin>

Random Variables and Probability Distributions

Random variables

- We are interested in calculating the probabilities associated with both **quantitative** and **qualitative** events.
- For example,
 - we can determine the probability that a machinist selected at random from the workers in a large automotive plant would suffer an accident during an 8-hour shift.
 - We can also find the probability that a machinist selected at random would work more than 80 hours without suffering an accident.
- These qualitative and quantitative events can be classified as **events** (or **outcomes**) associated with **qualitative** and **quantitative variables**.

Qualitative Random variables

- For example,
 - in the automotive plant accident study, the randomly selected machinist's accident report would consist of checking one of the following:
 - No Accident, Minor Accident, or Major Accident.
 - Thus, the data on 100 machinists in the study would be **observations on a qualitative variable** because the possible responses are the different categories of accident and are not different in any measurable, numerical amount.
- Because we cannot predict with certainty what type of accident a particular machinist will suffer, the variable is classified as a **qualitative random variable**.

Qualitative Random variables

- Other examples of qualitative random variables that are commonly measured are
 - political party affiliation,
 - socioeconomic status,
 - the species of insect discovered on an apple leaf,
 - the brand preferences of customers.
 - ...
- There are a finite (and typically quite small) number of possible outcomes associated with any qualitative variable.

Quantitative Random variables

- Many times the events of interest in an experiment are quantitative outcomes associated with a **quantitative random variable**, since the possible responses vary in numerical magnitude.
 - For example, in the automotive plant accident study, the number of consecutive 8-hour shifts between accidents for a randomly selected machinist is an observation on a **quantitative random variable**.
 - Events of interest, such as the number of 8-hour shifts between accidents for a randomly selected machinist, are observations on a **quantitative random variable**.

Quantitative Random variables

- Other examples of quantitative random variables are:
 - the change in earnings per share of a stock over the next quarter,
 - the length of time a patient is in remission after a cancer treatment,
 - the yield per acre of a new variety of wheat,
 - the number of persons voting for the incumbent in an upcoming election.
 - ...

7

Random variables

- Formally, a **random variable** X assigns a numerical value to each possible outcome (and event) of a **random phenomenon**.
- For instance, we can define X based on possible genotypes of a bi-allelic gene A as follows:

$$X = \begin{cases} 0 & \text{for genotype } AA, \\ 1 & \text{for genotype } Aa, \\ 2 & \text{for genotype } aa. \end{cases}$$

- In this case, the random variable assigns **0** to the outcome AA , **1** to the outcome Aa , and **2** to the outcome aa .

8

Random variables

- The way we specify random variables based on a specific random phenomenon is not unique.
- Alternatively, we can define a random variable Y as:

$$Y = \begin{cases} 0 & \text{for genotypes } AA \text{ and } aa, \\ 1 & \text{for genotype } Aa. \end{cases}$$

- In this case, Y assigns **0** to the homozygous event and assigns **1** to the heterozygous event.

9

Random variables

- When the underlying outcomes are numerical, the values the random variable assigns to each outcome can be the same as the outcome itself.
 - For the die Rolling example, we can define a random variable Z to be equal to $1, 2, \dots, 6$ for outcomes $1, 2, \dots, 6$, respectively.
 - Alternatively, we can define a random variable W and set W to **1** when the outcome is an odd number and to **2** when the outcome is an even number.
- The set of values that a random variable can assume is called its **range**.
 - For the above examples, the range of X is $\{0, 1, 2\}$, and the range of Z is $\{1, 2, \dots, 6\}$.

10

Random variables

- After we define a random variable, we can find the probabilities for its possible values based on the probabilities for its underlying random phenomenon.
- This way, instead of talking about the probabilities for different outcomes and events,
 - we can talk about the probability of different values for a random variable.
- {For example,
 - suppose $P(AA) = 0.49$, $P(Aa) = 0.42$, and $P(aa) = 0.09$.
 - Then, we can say that $P(X = 0) = 0.49$,
 - i.e., X is equal to 0 with probability of 0.49. }
 - Note that the total probability for the random variable is still 1.

11

Random variables

- The probability distribution of a random variable specifies its possible values (i.e., its range) and their corresponding probabilities.
 - For the random variable X defined based on genotypes, the probability distribution can be simply specified as follows:

$$P(X = x) = \begin{cases} 0.49 & \text{for } x = 0, \\ 0.42 & \text{for } x = 1, \\ 0.09 & \text{for } x = 2. \end{cases}$$

- Here, x denotes a specific value (i.e., 0, 1, or 2) of the random variable.

12

Discrete vs. continuous random variables

- We divide the random variables into two major groups:
 - discrete and continuous.
- When observations on a quantitative random variable can assume only a countable number of values, the variable is called a **discrete random variable**.
 - These variables can be categorical (nominal or ordinal), such as genotype, or counts, such as the number of patients visiting an emergency room per day

13

Discrete vs. continuous random variables

- When observations on a quantitative random variable can assume any one of the uncountable number of values in a line interval, the variable is called a **continuous random variable**.
 - Typical continuous random variables are temperature, pressure, height, weight, and distance.
- The distinction between discrete and continuous random variables is pertinent when we are seeking the probabilities associated with specific values of a random variable.

14

Probability distribution

- The probability distribution of a random variable provides the required information to find the probability of its possible values.
- We need to know the probability of observing a particular sample outcome in order to make an inference about the population from which the sample was drawn.
- To do this, we need to know the probability associated with each value of the variable.
- Viewed as relative frequencies, these probabilities generate a distribution of theoretical relative frequencies called the **probability distribution** of the variable.

15

Probability distribution

- Probability distributions differ for discrete and continuous random variables.
 - For discrete random variables, we will compute the probability of specific individual values occurring.
 - For continuous random variables, the probability of an interval of values is the event of interest.
- The probability distributions discussed here are characterized by one or more **parameters**.
- The parameters of probability distributions we assume for random variables are usually unknown.

16

Probability distribution

- Typically, we use Greek alphabets such as μ and σ to denote these parameters and distinguish them from known values.
 - We usually use μ to denote the mean of a random variable and use σ^2 to denote its variance.
- For a population of size N , the **mean** and **variance** are calculated as follows:

$$\mu = \frac{\sum_{i=1}^N x_i}{N} \quad \sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

17

Discrete probability distributions

- The probability distribution for a discrete random variable displays the probability $P(y)$ associated with each value of y .
 - This display can be presented as a table, a graph, or a formula.
- The probability distribution of a discrete random variable is fully defined by the **probability mass function (pmf)**.
 - This is a function that specifies the probability of each possible value within range of random variable.

18

Discrete probability distributions

- For the genotype example, the pmf of the random variable X is

$$P(X = x) = \begin{cases} 0.49 & \text{for } x = 0, \\ 0.42 & \text{for } x = 1, \\ 0.09 & \text{for } x = 2. \end{cases}$$

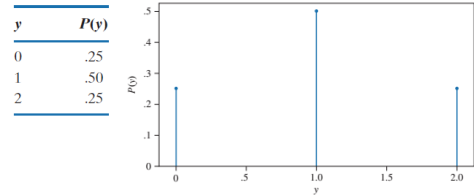
- As another example, suppose Y is a random variable that is equal to 1 when a newborn baby has low birthweight, and is equal to 0 otherwise.
 - We say Y is a **binary random variable**.
- Further, assume that the probability of having a low birthweight for babies is 0.3.
 - Then the pmf for the random variable Y is

$$P(Y = y) = \begin{cases} 0.7 & \text{for } y = 0, \\ 0.3 & \text{for } y = 1. \end{cases}$$

19

Discrete probability distributions

- Example:
 - Probability distribution for the number of heads when two coins are tossed



20

Properties of Discrete Random Variables

- The probability distribution for the discrete random variable given in previous slide illustrates three important properties of discrete random variables.
 - The probability associated with every value of y lies between 0 and 1.
 - The sum of the probabilities for all values of y is equal to 1.
 - The probabilities for a discrete random variable are additive.
 - Hence, the probability that $y = 1$ or 2 is equal to $P(1) + P(2)$

21

Bernoulli Distribution

- Binary random variables are abundant in scientific studies.
 - Examples include disease status (healthy and diseased), gender (male and female), survival status (dead, survived), and a gene with two possible alleles (A and a).
- The binary random variable X with possible values 0 and 1 has a **Bernoulli distribution** with parameter θ ,
 - where, $P(X = 1) = \theta$ and $P(X = 0) = 1 - \theta$.
- We denote this as $X \sim \text{Bernoulli}(\theta)$, where $0 \leq \theta \leq 1$.
 - Here θ is unknown parameter.
- If θ were known, we could fully specify the probability mass function:

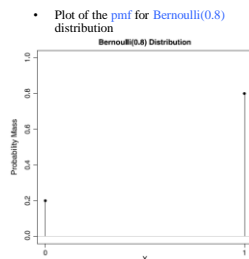
$$P(X = x) = \begin{cases} 1 - \theta & \text{for } x = 0 \\ \theta & \text{for } x = 1 \end{cases}$$

22

Bernoulli Distribution

- For example, let X be a random variable representing the five-year survival status of breast cancer patient,
 - where $X = 1$ if the patient survived and $X = 0$ otherwise.
- Suppose that the probability of survival is $\theta = 0.8$: $P(X = 1) = 0.8$
- Therefore, the probability of not surviving is
 - $P(X = 0) = 1 - \theta = 0.2$
- Then X has a Bernoulli distribution with parameter $\theta = 0.8$, and we denote this as $X \sim \text{Bernoulli}(0.8)$.
- The pmf for this distribution is

$$P(X = x) = \begin{cases} 0.2 & \text{for } x = 0 \\ 0.8 & \text{for } x = 1 \end{cases}$$



23

Bernoulli Distribution

- The mean of a binary random variable, X , with **Bernoulli(θ)** distribution is θ .
 - We show this as $\mu = \theta$.
 - In this case, the mean can be interpreted as the proportion of the population who have the outcome of interest.
- The variance of a random variable with **Bernoulli(θ)** distribution is
 - $\sigma^2 = \theta(1 - \theta) = \mu(1 - \mu)$
- The standard deviation is obtained by taking the square root of variance
 - $\sigma = \sqrt{\theta(1 - \theta)} = \sqrt{\mu(1 - \mu)}$

24

Bernoulli Distribution

- In the above example, $\mu = 0.8$.
 - 80% of patients survive.
- The variance of the random variable is
$$\sigma^2 = 0.8 \times 0.2 = 0.16,$$
- Its standard deviation is $\sigma = 0.4$.
- This reflects the extent of variability in survival status from one person to another.
 - For this example, the amount of variation is rather small.
 - Therefore, we expect to see many survivals ($X = 1$) with occasional death ($X = 0$).

25

Bernoulli Distribution

- For comparison, suppose that the probability of survival for bladder cancer is $\theta = 0.6$.
- Then, the variance becomes
$$\sigma^2 = 0.6 \times (1 - 0.6) = 0.24.$$
- This reflects a higher variability in the survival status for bladder cancer patients compared to that of breast cancer patients.

26

Binomial Distribution

- A sequence of binary random variables X_1, X_2, \dots, X_n is called **Bernoulli trials**
 - if they all have the same Bernoulli distribution and are independent.
- The random variable representing the number of times the outcome of interest occurs in n Bernoulli trials (i.e., the sum of Bernoulli trials) has a **Binomial(n, θ) distribution**,
 - where θ is the probability of the outcome of interest (a.k.a. the probability of success).

27

Binomial Distribution

- A **binomial distribution** is defined by the number of Bernoulli trials n and the probability of the outcome of interest θ for the underlying Bernoulli trials.
- The **pmf** of a **Binomial(n, θ)** specifies the probability of each possible value (integers from 0 through n) of the random variable.
- The theoretical (population) mean of a random variable Y with **Binomial(n, θ)** distribution is $\mu = n\theta$.
- The theoretical (population) variance of Y is $\sigma^2 = n\theta(1 - \theta)$.

28

Binomial Distribution

- A **binomial experiment** is one that has the following properties:
 - The experiment consists of n identical trials.
 - Each trial results in one of two outcomes.
 - We will label one outcome a **success** and the other a **failure**.
 - The probability of success on a single trial is equal to p , and p remains the same from trial to trial.
 - The trials are independent;
 - that is, the outcome of one trial does not influence the outcome of any other trial.
 - The random variable y is the number of successes observed during the n trials.

29

Binomial Distribution -Example

- A large power utility company uses gas turbines to generate electricity.

The engineers employed at the company monitor the reliability of each turbine

 - that is, the probability that the turbine will perform properly under standard operating conditions over a specified period of time.

The engineers wanted to estimate the probability a turbine will operate successfully for 30 days after being put into service.

The engineers randomly selected 75 of the 100 turbines currently in use and examined the maintenance records. They recorded the number of turbines that did not need repairs during the 30-day time period.
- Is this a binomial experiment?

30

Binomial Distribution -Example

- For solution, we check this experiment against the five characteristics of a binomial experiment.
 - Are there identical trials?
 - The 75 trials could be assumed identical only if the 100 turbines are the same type of turbine, are the same age, and are operated under the same conditions.
 - Does each trial result in one of two outcomes?
 - Yes. Each turbine either does or does not need repairs in the 30-day time period.
 - Is the probability of success the same from trial to trial?
 - No. If we let success denote a turbine "did not need repairs," then the probability of success can change considerably from trial to trial.
 - For example, suppose that 15 of the 100 turbines needed repairs during the 30-day inspection period.
 - Then p , the probability of success for the first turbine examined, would be $85/100=0.85$.
 - If the first trial is a failure (turbine needed repairs), the probability that the second turbine examined did not need repairs is $85/99=0.859$.
 - Suppose that after 60 turbines have been examined, 50 did not need repairs and 10 needed repairs.
 - The probability of success of the next (61st) turbine would be $35/40=0.875$.

31

Binomial Distribution -Example

- Were the trials independent?
 - Yes, provided that the failure of one turbine does not affect the performance of any other turbine.
 - However, the trials may be dependent in certain situations. For example,
 - suppose that a major storm occurs that results in several turbines being damaged.
 - Then the common event, a storm, may result in a common result, the simultaneous failure of several turbines.
 - Was the random variable of interest to the engineers the number of successes in the 75 trials?
 - Yes. The number of turbines not needing repairs during the 30-day period was the random variable of interest.
- This example shows how the probability of success can change substantially from trial to trial in situations in which the sample size is a relatively large portion of the total population size.
- This experiment does not satisfy the properties of a binomial experiment.

32

Binomial Distribution

- Although it is possible to approximate $P(y)$, the probability associated with a value of y in a binomial experiment, by using a relative frequency approach, it is easier to use a general formula for binomial probabilities.
- The probability of observing y successes in n trials of a binomial experiment is

$$P(y) = \frac{n!}{y!(n-y)!} \theta^y (1-\theta)^{n-y}$$

where

- n = number of trials
- θ = probability of success on a single trial
- $1 - \theta$ = probability of failure on a single trial
- y = number of successes in n trials
- $n!$ = $n(n-1)(n-2) \dots (3)(2)(1)$

33

Binomial Distribution - Example

- A new variety of turf grass has been developed for use on golf courses, with the goal of obtaining a germination rate of 85%.
- To evaluate the grass, 20 seeds are planted in a greenhouse so that each seed will be exposed to identical conditions.
- If the 85% germination rate is correct,
 - what is the probability that 18 or more of the 20 seeds will germinate?
 - what is the average number of seeds that will germinate in the sample of 20 seeds?
 - what is the variance of seeds that will germinate in the sample of 20 seeds?
 - what is the standard deviation of seeds that will germinate in the sample of 20 seeds?

34

Binomial Distribution - Example

- Solution:**
- $P(y) = \frac{n!}{y!(n-y)!} \theta^y (1-\theta)^{n-y}$
- $n = 20$, $\theta = 0.85$, $y = 18, 19$, and 20
- $P(y = 18) = \frac{20!}{18!(20-18)!} (0.85)^{18} (1-0.85)^{20-18} = 0.229$
- $P(y = 19) = \frac{20!}{19!(20-19)!} (0.85)^{19} (1-0.85)^{20-19} = 0.137$
- $P(y = 20) = \frac{20!}{20!(20-20)!} (0.85)^{20} (1-0.85)^{20-20} = 0.038$
- $P(y \geq 18) = P(y = 18) + P(y = 19) + P(y = 20) = 0.405$
- The following commands in R will compute the binomial probabilities:
 - To calculate $P(X = 18)$, use the command `dbinom(18, 20, 0.85)`
 - To calculate $P(X \leq 17)$, use the command `pbinom(17, 20, 0.85)`
 - To calculate $P(X \geq 18)$, use the command `1 - pbinom(17, 20, 0.85)`

35

Binomial Distribution - Example

- The average number of seeds that will germinate in the sample of 20 seeds is

$$\mu = n\theta = 20 \times 0.85 = 17$$
- The variance of seeds that will germinate in the sample of 20 seeds is

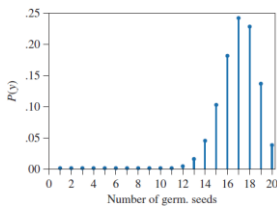
$$\sigma^2 = n\theta(1 - \theta) = 20 \times 0.85 (1 - 0.85) = 2.55$$
- The standard deviation of seeds that will germinate in the sample of 20 seeds is

$$\sigma = \sqrt{n\theta(1 - \theta)} = \sqrt{\sigma^2} = \sqrt{2.55} = 1.60$$

36

Binomial Distribution - Example

- Suppose we examine the germination records of a large number of samples of 20 seeds each.
- If the germination rate has remained constant at 85%, then the average number of seeds that germinate should be close to 17 per sample.



- If in a particular sample of 20 seeds we determine that only 12 had germinated, would the germination rate of 85% seem consistent with our results?
- Using a computer software program, we can generate the probability distribution for the number of seeds that germinate in the sample of 20 seeds, as shown in the Figure

37

Binomial Distribution - Example

- Suppose that a sample of households is randomly selected from all the households in the city in order to estimate the percentage in which the head of the household is unemployed.
- To illustrate the computation of a binomial probability, suppose that the unknown percentage is actually 10% and that a sample of $n = 5$ (we select a small sample to make the calculation manageable) is selected from the population.
- What is the probability that all five heads of households are employed?

38

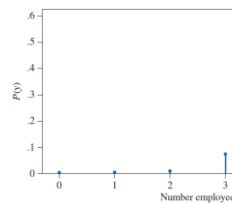
Binomial Distribution - Example

- Solution:**
- We must carefully define which outcome we wish to call a success.
 - For this example, we define a success as being employed.
- Then the probability of success when one person is selected from the population is $\theta = 0.9$ (because the proportion unemployed is 0.1).
- We wish to find the probability that $y = 5$ (all five are employed) in five trials.
- $$P(y) = \frac{n!}{y!(n-y)!} \theta^y (1-\theta)^{n-y}$$
- $$P(y = 5) = \frac{5!}{5!(5-5)!} (0.9)^5 (1-0.9)^{5-5} = 0.59$$

39

Binomial Distribution - Example

- The binomial probability distribution for $n = 5$, $\theta = 0.9$ is shown in the figure.



– Here, the probability of observing five employed in a sample of five is shown to be 0.59.

40

Poisson Distribution

- In 1837, S. D. Poisson developed a discrete probability distribution, suitably called the **Poisson distribution**, which has as one of its important applications the modeling of events of a particular time over a unit of time or space
- For example, the number of automobiles arriving at a toll booth during a given 5-minute period of time.
 - The event of interest would be an arriving automobile, and the unit of time would be 5 minutes.

41

Poisson Distribution

- A second example would be the situation in which an environmentalist measures the number of PCB particles discovered in a liter of water sampled from a stream contaminated by an electronics production plant.
 - The event would be a PCB particle discovered.
 - The unit of space would be 1 liter of sampled water.

42

Poisson Distribution

- Let y be the number of events occurring during a fixed time interval of length t or a fixed region R of area or volume $m(R)$.
- Then the probability distribution of y is Poisson, provided certain conditions are satisfied:
 - Events occur one at a time; two or more events do not occur precisely at the same time or in the same space.
 - The occurrence of an event in a given period of time or region of space is independent of the occurrence of the event in a nonoverlapping time period or region of space;
 - that is, the occurrence (or nonoccurrence) of an event during one period or in one region does not affect the probability of an event occurring at some other time or in some other region.
 - The expected number of events during one period or in one region, μ , is the same as the expected number of events in any other period or region.

43

Poisson Distribution - Example

- A large industrial plant is being planned in a rural area. As a part of the environmental impact statement, a team of wildlife scientists is surveying the number and types of small mammals in the region.
- Let y denote the number of field mice captured in a trap over a 24-hour period.
- Suppose that y has a Poisson distribution with $\mu = 2.3$; that is, the average number of field mice captured per trap is 2.3.
 - What is the probability of finding exactly four field mice in a randomly selected trap?
 - What is the probability of finding at most four field mice in a randomly selected trap?
 - What is the probability of finding more than four field mice in a randomly selected trap?

45

Poisson Distribution - Example

- The Poisson probabilities can be computed using the following R commands.


```
> dpois(4, 2.3)
[1] 0.1169022
> ppois(4, 2.3)
[1] 0.7993471
> 1 - ppois(4, 2.3)
[1] 0.08375072
```
- When n is large and θ is small in a binomial experiment, $n \geq 100$, $\theta \leq 0.01$, and $n\theta \leq 20$,
 - the Poisson distribution provides a reasonable approximation to the binomial distribution.
- In applying the Poisson approximation to the binomial distribution, use $\mu = n\theta$.

47

Poisson Distribution

- Assuming that the above conditions hold, the Poisson probability of observing y events in a unit of time or space is given by the formula

$$P(y) = \frac{\mu^y e^{-\mu}}{y!}$$

where e is a naturally occurring constant approximately equal to 2.71828 and μ is the average value of y .

44

Poisson Distribution - Example

- The probability that a trap contains exactly four field mice is computed to be

$$P(y = 4) = \frac{\mu^y e^{-\mu}}{y!} = \frac{(2.3)^4 e^{-2.3}}{4!} = \frac{(27.9841)(0.10002588)}{24} = 0.1169$$
- The probability of finding at most four field mice in a randomly selected trap is,

$$P(y \leq 4) = P(y = 0) + P(y = 1) + P(y = 2) + P(y = 3) + P(y = 4)$$

$$P(y \leq 4) = 0.1003 + 0.2306 + 0.2652 + 0.2033 + 0.1169 = 0.9163$$
- The probability of finding more than four field mice in a randomly selected trap, using the idea of complementary events, is

$$P(y > 4) = 1 - P(y \leq 4) = 1 - 0.9163 = 0.0837$$

Thus, it is a very unlikely event to find five or more field mice in a trap.

46

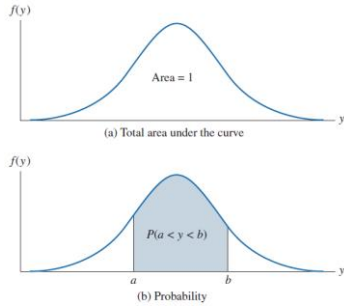
Continuous probability distributions

- For discrete random variables, the pmf provides the probability of each possible value.
- For continuous random variables, the number of possible values is uncountable, and the probability of any specific value is zero.
- For these variables, we are interested in the probability that the value of the random variable is within a specific interval from x_1 to x_2 ;
 - we show this probability as $P(x_1 < X \leq x_2)$.

48

Continuous probability distributions

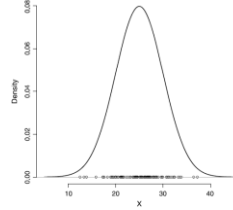
- Probability distribution for a continuous random variable



49

Continuous probability distributions

- For continuous random variables, we use **probability density functions (pdf)** to specify the distribution.
- Using the **pdf**, we can obtain the probability of any interval.

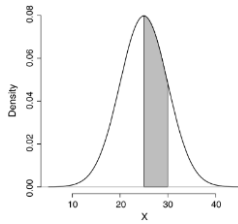


- The assumed probability distribution for BMI (Body Mass Index), which is denoted as X , along with random sample of 100 values, which are shown as circles along the horizontal axis

50

Continuous probability distributions

- The total area under the probability density curve is 1.
- The curve (and its corresponding function) gives the probability of the random variable falling within an interval.

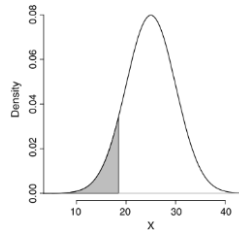


- This probability is equal to the area under the probability density curve over the interval.
- The shaded area is the probability that a person's BMI is between 25 and 30.
- People whose BMI is in this range are considered as overweight.
- Therefore, the shaded area gives the probability of being overweight

51

Lower tail probability

- The probability of observing values less than or equal to a specific value x , is called the lower tail probability and is denoted as $P(X \leq x)$

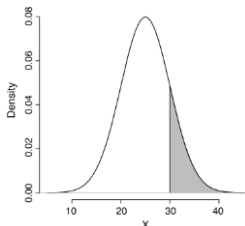


- This probability is found by measuring the area under the curve to the left of x .
- For example, the shaded area in the left panel of the figure is the lower tail probability of having a BMI less than or equal to 18.5 (i.e., being underweight), $P(X \leq 18.5)$.

52

Upper tail probability

- The probability of observing values greater than x , is called the upper tail probability and is denoted as $P(X > x)$



- This probability is found by measuring the area under the curve to the right of x .
- For example, the shaded area in the right panel of the figure is the upper tail probability of having a BMI greater than 30 (i.e., being obese), $P(X > 30)$.

53

Probability of intervals

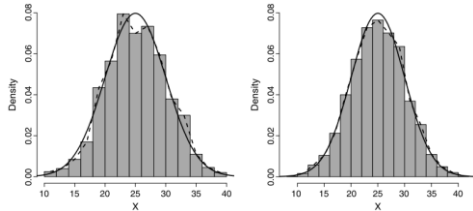
- The probability of any interval from x_1 to x_2 , where $x_1 < x_2$, can be obtained using the corresponding lower tail probabilities for these two points as follows:

$$P(x_1 < X \leq x_2) = P(X \leq x_2) - P(X \leq x_1)$$
- For example, suppose that we wanted to know the probability of a BMI between 25 and 30.
- This probability $P(25 < X \leq 30)$ is obtained by subtracting the lower tail probability of 25 from the lower tail probability of 30:

$$P(25 < X \leq 30) = P(X \leq 30) - P(X \leq 25)$$

54

Probability Density Curves and Density Histograms



- **Left panel:** Histogram of BMI for 1000 observations.
 - The dashed line connects the height of each bar at the midpoint of the corresponding interval
 - The smooth solid curve is the density curve for the probability distribution of BMI
- **Right panel:** Histogram of BMI for 5000 observations.
 - The histogram and its corresponding dashed line provide better approximations to the density curve
- Recall that the height of each bar is the density for the corresponding interval, and the area of each bar is the relative frequency for that interval.
- The density histogram and the dashed line, which shows the density for each interval based on the observed data, provide reasonable approximations to the density curve.
- Also, the area of each bar, which is equal to the relative frequency for the corresponding interval, is approximately equal to the area under the curve over that interval.

55

The 68-95-99.7% rule

- The **68–95–99.7%** rule for normal distributions specifies that
 - **68% of values fall within 1 standard deviation of the mean:**
 $P(\mu - \sigma < X \leq \mu + \sigma) = 0.68$
 - **95% of values fall within 2 standard deviations of the mean:**
 $P(\mu - 2\sigma < X \leq \mu + 2\sigma) = 0.95$
 - **99.7% of values fall within 3 standard deviations of the mean:**
 $P(\mu - 3\sigma < X \leq \mu + 3\sigma) = 0.997$

57

Example

- For example, suppose we know that the population mean and standard deviation for SBP are $\mu = 125$ and $\sigma = 15$, respectively.
 - That is, $X \sim N(125, 15^2)$,
 - where X is the random variable representing SBP.
- Therefore, the probability of observing an SBP in the range $\mu \pm \sigma$ is 0.68:
 $P(125 - 15 < X \leq 125 + 15) = P(110 < X \leq 140) = 0.68$.
- This probability corresponds to the central area shown in the Fig. b in the previous slide.

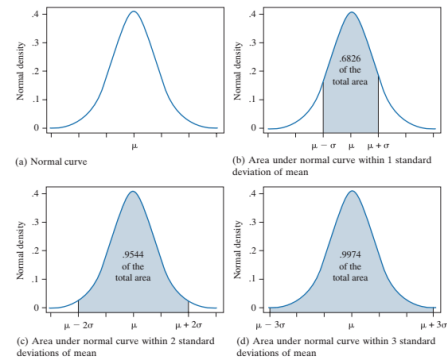
59

Normal distribution

- A **normal distribution** and its corresponding pdf are fully specified by the mean μ and variance σ^2 .
- A random variable X with normal distribution is denoted $X \sim N(\mu, \sigma^2)$,
 - where μ is a real number, but σ^2 can take positive values only.
- The normal density curve is always symmetric about its mean μ , and its spread is determined by the variance σ^2 .
- A normal distribution with a mean of **0** and a standard deviation (or variance) of **1** is called the **standard normal distribution** and denoted $N(0, 1)$.

56

The 68-95-99.7% rule



58

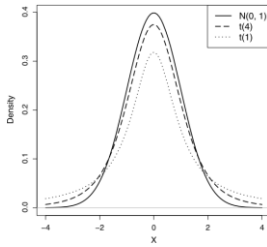
Example

- The probability of observing an SBP in the range $\mu \pm 2\sigma$ is **0.95**:
 $P(125 - 2 \times 15 < X \leq 125 + 2 \times 15) = P(95 < X \leq 145) = 0.95$.
- This probability is shown in the Fig. c in the previous slide.
- Lastly, the probability of observing an SBP in the range $\mu \pm 3\sigma$ is **0.997**:
 $P(125 - 3 \times 15 < X \leq 125 + 3 \times 15) = P(80 < X \leq 170) = 0.997$.
- Therefore, we rarely (probability of **0.003**) expect to see SBP values less than **80** or greater than **170**.

60

Student's t-distribution

- Another continuous probability distribution that is used very often in statistics is the Student's t -distribution or simply the t -distribution.



- Comparing the pdf of a standard normal distribution to t -distributions with 1 degree of freedom and then with 4 degrees of freedom.
- The t -distribution has heavier tails than the standard normal;
 - however, as the degrees of freedom increase, the t -distribution approaches the standard normal

61

Student's t-distribution

- A t -distribution is specified by only one parameter called the degrees of freedom, df .
- The t -distribution with df degrees of freedom is usually denoted as $t(df)$ or tdf , where df is a positive real number ($df > 0$).
- The mean of this distribution is $\mu = 0$,
- The variance is determined by the degrees of freedom parameter, $\sigma^2 = df/(df - 2)$,
 - which is of course defined when $df > 2$.

62

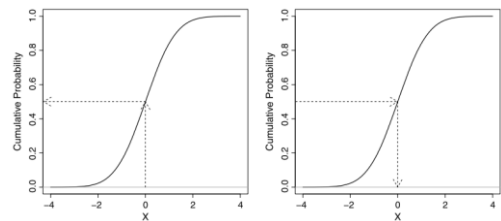
Cumulative distribution function

- We saw that by using lower tail probabilities, we can find the probability of any given interval.
- Indeed, all we need to find the probabilities of any interval is a function that returns the lower tail probability at any given value of the random variable: $P(X \leq x)$.
- This function is called the cumulative distribution function (cdf) or simply the distribution function.

63

Quantiles

- We can use the cdf plot in the reverse direction to find the value of the random variable for a given lower tail probability.



64

Quantiles

- In previous slide:
 - Left panel:**
 - Plot of the cdf for the standard normal distribution, $N(0, 1)$.
 - The cdf plot of the cdf can be used to find the lower tail probability.
 - For instance, following the arrow from $x = 0$ (on the horizontal axis) to the cumulative probability (on the vertical axis) gives us the probability $P(X \leq 0) = 0.5$.
 - Right panel:**
 - Given the lower tail probability of 0.5 on the vertical axis, we obtain the corresponding quantile $x = 0$ on the horizontal axis

65

Scaling and shifting random variables

- If $Y = aX + b$, then

$$\mu_Y = a\mu_X + b$$

$$\sigma_Y^2 = a^2\sigma_X^2$$

$$\sigma_Y = |a|\sigma_X$$
- The process of shifting and scaling a random variable to create a new random variable with mean zero and variance one is called standardization.
 - For this, we first subtract the mean μ and then divide the result by the standard deviation σ .

$$Z = (X - \mu)/\sigma$$
 - If $X \sim N(\mu, \sigma^2)$, then $Z \sim N(0, 1)$.

66

Adding/subtracting random variables

- If $W = X + Y$, then

$$\mu_W = \mu_X + \mu_Y$$

- If the random variables X and Y are independent, then we can find the variance of W as follows:

$$\sigma_W^2 = \sigma_X^2 + \sigma_Y^2$$

- If $X \sim N(\mu_X, \sigma_X^2)$, and $Y \sim N(\mu_Y, \sigma_Y^2)$, then assuming that the two random variables are independent, we have

$$W = X + Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

67

Adding/subtracting random variables

- If we subtract Y from X , then

$$\mu_W = \mu_X - \mu_Y$$

- If the random variables X and Y are independent, then we can find the variance of W as follows:

$$\sigma_W^2 = \sigma_X^2 + \sigma_Y^2$$

- If $X \sim N(\mu_X, \sigma_X^2)$, and $Y \sim N(\mu_Y, \sigma_Y^2)$, then assuming that the two random variables are independent, we have

$$W = X - Y \sim N(\mu_X - \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

68