

Chapter 5: Sampling Theory & Chapter 2: Probability Distribution Functions

Class Notes

Monday 9.17.2018



Dr. Basilio

* * *

Chapter 2: Random Variables

Discrete vs Continuous Variables

Definition 1: Discrete-vs-Continuous-Variable

- **Variable:** a function defined on the sample space. That is, given any event A from a sample space S , a random variable assigns a number to even A . Using function notation, we write this as $X(A)$.
- **Discrete variable:** a variable that can attain only specific values. Example: think the values of the roll of a dice.
- **Continuous variable:** a variable can attain infinitely many values over a certain span or range. Example: the height of a person.
- **RANDOM variable:** a variable defined on a sample space that is comprised of random process or experiment—that is, experiment where you don't know what the outcome is until it is completed. Example: flipping a coin is a random experiment.

Chapter 4: Probability Distribution Functions

Binomial Distribution

Definition 2: Binomial-Distribution

For a Binomial distribution it is important that we have a random process or experiment and we will run the experiment many times but each trail must be an **independent trail** which means previous trails do not have any influence on future trails.

- Let n be the total number of trails run in the experiment
- Let X be a random variable of a single “successful” trail
- Let p be the probability of the successful trail X
- Let q be the probability of trail X failing. (NOTE: $p + q = 1$, or $q = 1 - p$)
- Let x be the number of successful trials of X . So notice that x can take values from 0 up to n , i.e. $x = 0, 1, 2, 3, \dots, n$.
- Let $P(X = x)$ denote the **probability of exactly x successful trails out n in a random experiment with independent trails**, then

$$P(X = x) = \binom{n}{x} p^x q^{n-x} \quad (1)$$

Recall that: $\binom{n}{x} = {}_n C_x = \frac{n!}{x!(n-x)!}$.

Definition 3: Binompdf-vs-Binomcdf

USING CALCULATOR TI83: binompdf(n,p,x)

\boxed{DIST} key in yellow ($\boxed{2nd} > \boxed{VARS}$) > Scroll to 10 “binompdf” or scroll to A “binomcdf”

- Binompdf is when we want exactly x trials to be successful so this is binomial distribution pdf. Thus this is 1-valued random variables.
- Binomcdf is when we want multiple values of x to be true. It is defined as

$$\text{binomcdf}(n, p, x) = P(X \leq x) \quad (2)$$

Notice the sneaky “ \leq ” less than or equal to sign in the binomcdf. This means:

$$\begin{aligned}\text{binomcdf}(n, p, x) &= P(X \leq x) = P(X = 0, 1, 2, \dots, x) \\ &= P(X = 0) + P(X = 1) + P(X = 2) + \dots + P(X = x)\end{aligned}$$

This can help us with “at most” and “at least” type of problems

Visualizing a Binomial Distribution

Activity 1: Visualizing-Binomial-Distribution

Let X be the number of heads that turn up after flipping a coin five times. Then $n = 5$ and x can be 0, 1, 2, 3, 4, 5. We can calculate the probability of zero heads turning up with $P(X = 0)$, one head turning up with $P(X = 1)$, etc. Using our calculator check that:

$$P(X = 0) = \frac{1}{32}, P(X = 1) = \frac{5}{32}, P(X = 2) = \frac{10}{32}, P(X = 3) = \frac{10}{32}, P(X = 4) = \frac{5}{32}, P(X = 5) = \frac{1}{32}$$

- (a) Plot a histogram for the random variable X probability distribution.

To do this, on the horizontal axis scale from $x = 0, 1, 2, 3, 4, 5$ and the vertical axis scale from 0 to $10/32$ with $1/32$ intervals.

- (b) Describe any interesting features from your histogram.

- (c) Sketch what you think the histogram would look like for the same random variable X but the number of trials is $n = 100$.

Activity 2: Visualizing-Binomial-Distribution

Let our experiment be shooting free-throws. Assume that the probability of making a freethrow is 70% and that these are independent events. Let X be the number of made in taking six shots.

- (a) Use your calculator to find $P(X = x)$ for $x = 0, 1, 2, 3, 4, 5, 6$.

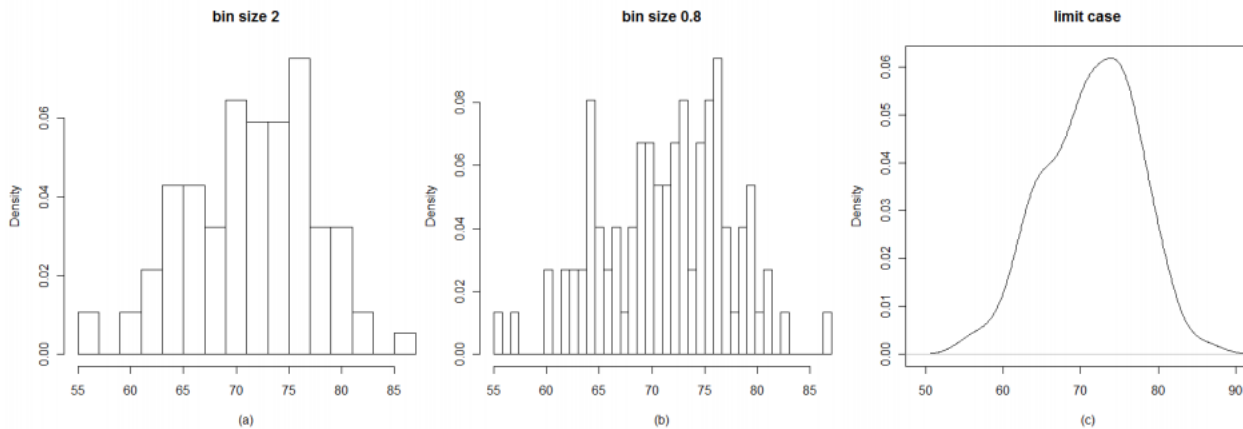
- (b) Plot a histogram for the random variable X probability distribution.

- (c) Describe any interesting features from your histogram.

Probability Density Functions

Let's consider the random variable X to be the temperature in a room. We want to know the probability that the temperature is in any given interval. For example, what's the probability for the temperature between 70° and 80° ?

Ultimately, we want to know the probability distribution for X . One way to do that is to record the temperature from time to time and then plot the histogram. However, when you plot the histogram, it's up to you to choose the bin size. But if we make the bin size finer and finer (meanwhile we need more and more data), the histogram will become a smooth curve which will represent the probability distribution for X .



Definition 4: Probability-Density-Functions

Recall a random variable can be either discrete or continuous. If we want to know the probability of a random variable in certain range, then we

- **Probability Density Functions:** Let X be a continuous random variable. Then a **probability density (or probability distribution) function (pdf) of X** is a function $f(x)$ if:

1. f is a continuous function on \mathbb{R}
2. f is nonnegative, that is, $f(x) \geq 0$
3. The probability that X takes on a value in the interval $[a, b]$ is the area above this interval and under the graph of the density function:

$$P(a \leq X \leq b) = \text{Area under the curve between } x = a \text{ and } x = b \quad (3)$$

The graph of $f(x)$ is often referred to as the density curve.

4. The total area under the graph of $f(x)$ is 1. This corresponds to $P(S) = 1$.

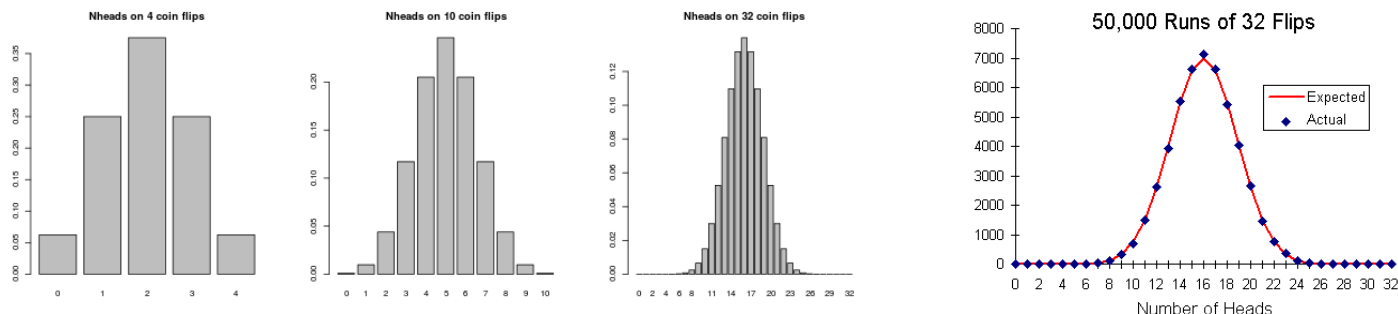
NOTE: for continuous pdf, $P(X = x) = 0$ why?

Example: X = temperature in the room. What is $P(X = 70)$, that is, the temperature in the room is *exactly* 70° ? Well, logically, $P(X = 70) = 0$ right? Notice it is NOT the y -coordinate of the point $(70, 0.056)$. The key to using a PDF is that you must know a RANGE of values for X , not insist on exact values. This is because we are dealing with infinite possibilities unlike the discrete case.

Normal Distribution

The normal distribution is the most widely known and used of all distributions. Because the normal distribution approximates many natural phenomena so well, it has developed into a standard of reference for many probability problems.

If you recall Activity 1, I asked you to visualize the histogram of the probability of flipping a coin 100 times. What if we did it for 1,000,000 times? What if we imagine $n \rightarrow \infty$. We go from a discrete random variable X to a continuous random variable X in the limit. That means we can talk about its PDF $f(x)$. What does the graph of $f(x)$ look like for flipping a coin infinitely many times?



Definition 5: Normal-Distribution

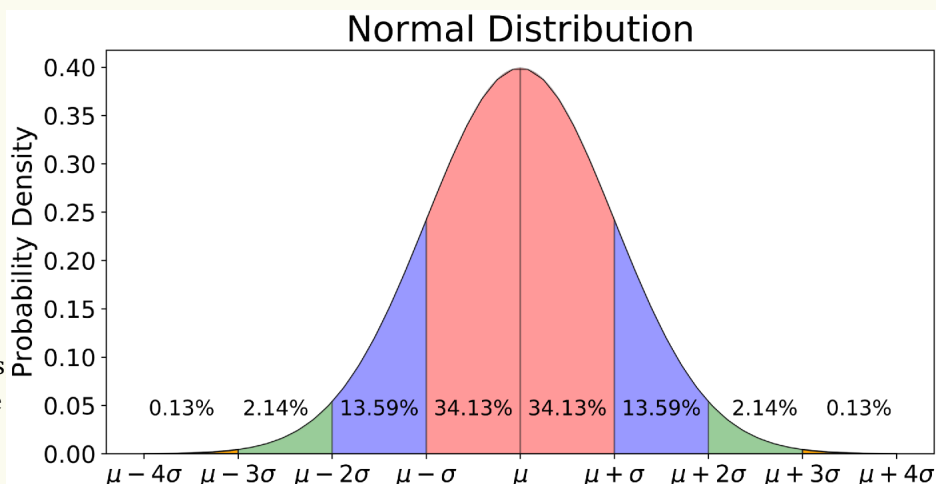
- **Normal Distribution:** Is the PDF for a very special function, and its graph is given by:

Properties:

- Symmetric, bell shaped
- Continuous on all of \mathbb{R}
- Mean: $\bar{x} = \mu$
- Standard deviation: σ
- The associated PDF

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)}$$
- Two parameters, μ and σ .

Note that the normal distribution is actually a family of distributions since μ and σ determine the shape of the distribution.



- **Empirical Rule (68-95-99.7 Rule):** is a shorthand used to remember the percentage of values that lie within a band around the mean in a normal distribution with a width of two, four and six standard deviations, respectively. More accurately, 68.27%, 95.45% and 99.73% of the values lie within one, two and three standard deviations of the mean, respectively.

In mathematical notation, these facts can be expressed as follows, where X is an observation from a normally distributed random variable, $\mu (= \bar{x})$ is the mean of the distribution, and σ is its standard deviation:

Within 1 standard deviation of mean: $P(\mu - \sigma \leq X \leq \mu + \sigma) \approx 68\%$

Within 2 standard deviations of mean: $P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%$

Within 3 standard deviations of mean: $P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99.7\%$

- Calculating probabilities outside of empirical rule: use table or calculator!

USING CALCULATOR TI83: `normalcdf(min, max, μ , σ)`

`DIST` key in yellow (`2nd` > `VARS`) > Scroll to 2. “normalcdf”

Why is the standard normal distribution useful?

- Many things actually are normally distributed, or very close to it. For example, height and intelligence are approximately normally distributed; measurement errors also often have a normal distribution
- The normal distribution is easy to work with mathematically. In many practical cases, the methods developed using normal theory work quite well even when the distribution is not normal.
- There is a very strong connection between the size of a sample N and the extent to which a sampling distribution approaches the normal form. Many sampling distributions based on large N can be approximated by the normal distribution even though the population distribution itself is definitely not normal.

Activity 3: Normal-Distribution

Weight (in grams) of bags of sugar from a factory are normally distributed, with a mean of 1000g, and standard deviation of 13g. Find the following.

- The probability that a randomly selected bag of sugar weighs in between 974g and 1000g.
- The percentage of bags whose weight is above 1026g.

Activity 4: Normal-Distribution

The time it takes employees to get to work from home (in minutes) is normally distributed with a mean of 30 minutes, and a standard deviation of 5 minutes. Find the following.

- The percentage of employees that take between 20 and 40 minutes to get to work. Do this without a calculator.
- The percentage of employees that take between 28 and 37 minutes to get to work. Do this with a calculator.

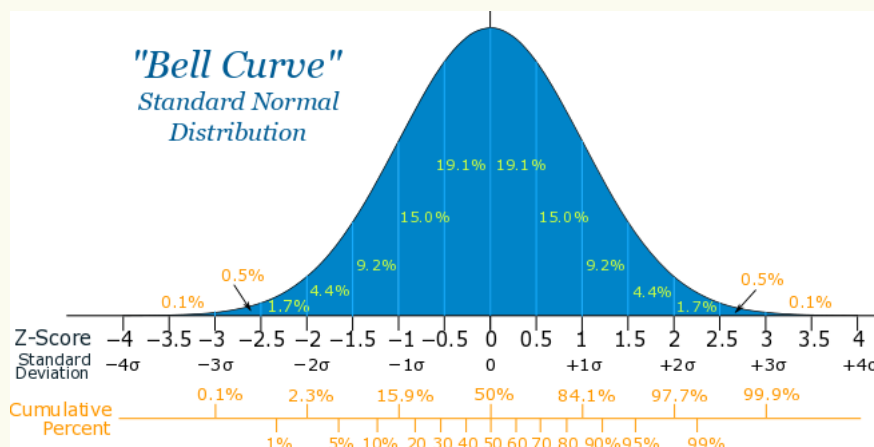
Standard Normal Distribution

Definition 6: Standard-Normal-Distribution

- **Standard Normal Distribution:** Is the PDF for a very special function, and it's graph is given by:

Properties:

- Symmetric, bell shaped
- Continuous on all of \mathbb{R}
- Mean: $\bar{x} = \mu = 0$
- Standard deviation: $\sigma = 1$
- The associated PDF $f(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$



- Calculating probabilities outside of empirical rule: use table or calculator!

USING CALCULATOR TI83: `normalcdf(min, max)` – NOTICE: when using Standard Normal we do not need to put μ and σ since the calculator already has is programed for $\mu = 0$ and $\sigma = 1$

`DIST` key in yellow (`2nd` > `VARS`) > Scroll to 2. "normalcdf"

Definition 7: Using-Normal-Distribution-Z-Scores

- General Procedure: As you might suspect from the formula for the normal density function, it would be difficult and tedious to do the calculus every time we had a new set of parameters for μ and σ . So instead, we usually work with the standardized normal distribution, where $\mu = 0$ and $\sigma = 1$. That is, rather than directly solve a problem involving a normally distributed variable X with mean μ and standard deviation σ , an indirect approach is used.
- We first convert the problem into an equivalent one dealing with a normal variable measured in standardized deviation units, called a **standardized normal variable**. To do this,

$$Z = \frac{X - \mu}{\sigma}$$

- A table of standardized normal values (Appendix C) can then be used to obtain an answer in terms of the converted problem.
- If necessary, we can then convert back to the original units of measurement. To do this, simply note that, if we take the formula for Z , multiply both sides by σ , and then add μ to both sides, we get

$$X = Z\sigma + \mu$$

- The interpretation of Z values is straightforward. Since $\sigma = 1$, if $Z = 2$, the corresponding X value is exactly 2 standard deviations above the mean. If $Z = -1$, the corresponding X value is one standard deviation below the mean. If $Z = 0$, $X =$ the mean, i.e. μ .

KEY: suffices to know the standard normal since we can go back and forth between normal and standard normal using the formulas

Activity 5: Conver-Z-values

Convert each of the following between x and z values.

- (a) $x = 35$ where $\mu = 40, \sigma = 2$
- (b) $x = 130$ where $\mu = 100, \sigma = 12$
- (c) $z = -0.57$ where $\mu = 14, \sigma = 1.5$

Activity 6: Standard-Normal-Distribution

Find the area under the standard normal curve between

- (a) $Z = 0$ and $Z = 1.2$
- (b) $Z = -0.68$ and $Z = 0$
- (c) $Z = -0.46$ and $Z = 2.21$
- (d) to the right of $Z = -1.28$

Activity 7: Standard-Normal-Distribution

The mean weight of 500 male students at a certain college is 151 lb and the standard deviation is 15 lb. Assuming that the weights are normally distributed, find without using a calculator how many students weigh

- (a) between 120 and 155 lb
- (b) more than 185 lb.

Definition 8: Skewness

Skewness is asymmetry in a statistical distribution, in which the histogram (or curve) appears distorted or skewed either to the left or to the right. Skewness can be quantified to define the extent to which a distribution differs from a normal distribution.

- **Positively Skewed:** the “tail” of the distribution is to the right of the mean
- **Negatively Skewed:** the “tail” of the distribution is to the left of the mean

Activity 8: Frequency-Skewness

The following is a list of prices (in dollars) of birthday cards found in various drug stores:

1.45	2.20	0.75	1.23	1.25
1.25	3.09	1.99	2.00	0.78
1.32	2.25	3.15	3.85	0.52
0.99	1.38	1.75	1.22	1.75

- Organize this data with intervals of 50 cents (i.e. .50-0.99, 1.00-0.49, and so on) using create a frequency distribution table.
- Draw a Histogram of the data. State the skewness of the data.