

# DAG C2M1 scripts

Video title	Isabel	Sean	Slides
<a href="#">L0V1 – Welcome to this course</a>	✓	✓	✓
<a href="#">L0V2 – Generative AI in this course</a>	✓	✓	✓
<a href="#">L0V3 – Module 1 introduction</a>	✓	✓	✓
<a href="#">L1V1 – Populations and sampling</a>	✓	✓	✓
<a href="#">L1V2 – Identifying the population</a>	✓	✓	✓
<a href="#">L1V3 – Probabilistic samples</a>	✓	✓	✓
<a href="#">L1V4 – Non-probabilistic samples</a>	✓	✓	✓
<a href="#">L1V5 – Types of bias</a>	✓	✓	✓
<a href="#">L2V1 – Histograms</a>	✓	✓	✓
<a href="#">L2V2 – Demo: plotting distributions</a>	✓	✓	✓
<a href="#">L2V3 – Central tendency, variability, and skewness</a>	✓	✓	✓
<a href="#">L2V4 – Central tendency: mean and mode</a>	✓	✓	✓
<a href="#">L2V5 – Central tendency: median</a>	✓	✓	✓
<a href="#">L2V6 – Demo: central tendency</a>	✓	✓	✓
<a href="#">L3V1 – Variability: range and interquartile range</a>	✓	✓	✓
<a href="#">L3V2 – Variability: variance and standard deviation</a>	✓	✓	✓
<a href="#">L3V3 – Skewness</a>	✓	✓	✓
<a href="#">L3V4 – Why use these measures?</a>	✓	✓	✓
<a href="#">L3V5 – Demo: variability and skewness</a>	✓	✓	✓
<a href="#">L3V6 – Box plots</a>	✓	✓	✓
<a href="#">L3V7 – Demo: LLMs for spreadsheet formulas &amp; errors</a>	✓	✓	✓
<a href="#">L4V1 – Correlation</a>	✓	✓	✓
<a href="#">L4V2 – Correlation and causation</a>	✓	✓	✓
<a href="#">L4V3 – Demo: correlations &amp; scatterplots in spreadsheets</a>	✓	✓	✓

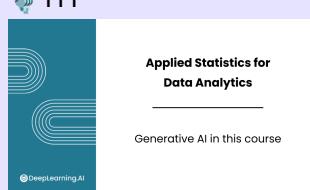
<a href="#">L5V1 – What is segmentation?</a>	✓	✓	✓
<a href="#">L5V2 – Demo: xlookup</a>	✓	✓	✓
<a href="#">L5V3 – Demo: pivot tables</a>	✓	✓	✓

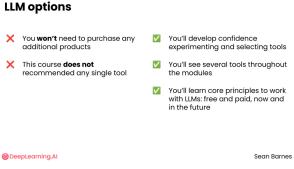
## Introduction

### L0V1 – Welcome to this course

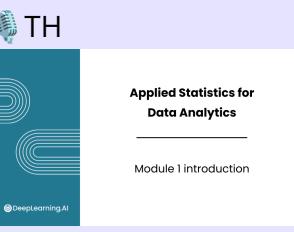
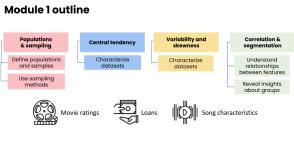
 [See DAG C1&C2 scripts for conversations with Andrew](#)

### L0V2 – Generative AI in this course

Visual	Script
 <p>TH Applied Statistics for Data Analytics Generative AI in this course DeepLearning.AI</p>	<p>One of the key elements of this course is learning to use generative AI, in particular large language models or LLMs like ChatGPT, Claude, Gemini, and so on. Effective use of LLMs will help you streamline your work and stand out.</p>
<p>In this course, you'll...</p> <ul style="list-style-type: none"> <li>Learn how to use LLMs to:           <ul style="list-style-type: none"> <li>Troubleshoot spreadsheet errors</li> <li>Create custom conditional formatting</li> <li>Design and run simulations</li> <li>Interpret inferential statistics</li> <li>Run statistical analysis for you</li> </ul> </li> </ul> <p>Explore LLMs' limitations, learning when to:     <ul style="list-style-type: none"> <li>Select the LLM for a particular task</li> <li>Use something else, like a spreadsheet</li> </ul> </p> <p>Sean Barnes</p>	<p>In this course, you'll <b>[CLICK]</b> learn how to use large language models, or LLMs, to</p> <ul style="list-style-type: none"> <li><b>[CLICK]</b> Troubleshoot spreadsheet errors</li> <li><b>[CLICK]</b> Create custom conditional formatting</li> <li><b>[CLICK]</b> Design and run simulations</li> <li><b>[CLICK]</b> Interpret inferential statistics</li> <li><b>[CLICK]</b> And run statistical analysis for you</li> </ul> <p>You'll also <b>[CLICK]</b> explore LLMs' key limitations, learning when to <b>[CLICK]</b> select the LLM for a particular task, and when to <b>[CLICK]</b> use something else, like a spreadsheet.</p> <p>LLMs are ever-evolving, and present challenges for both teaching and learning. I wanted to take a moment to share our team's philosophy about generative AI in this course.</p>
<p>LLMs in this course</p> <ul style="list-style-type: none"> <li>Demonstrates the most up-to-date capabilities as of 2024</li> <li>Evergreen principles</li> <li>How to think about and use generative AI regardless of product</li> <li>Develop a mindset of iteration and skepticism</li> </ul> <p>LLM progress since launch     <ul style="list-style-type: none"> <li>Abilities have advanced rapidly</li> <li>New features constantly being released</li> </ul> </p> <p>Changes you should expect     <ul style="list-style-type: none"> <li>More advanced and specialized features</li> <li>Cheaper tools</li> <li>Better tools</li> <li>Higher quality outputs overall</li> </ul> </p> <p>Sean Barnes</p>	<p>First, this course <b>[CLICK]</b> demonstrates the most up-to-date capabilities as of late 2024, and we expect changes in the coming months and years. This course is designed to teach <b>[CLICK]</b> evergreen principles: <b>[CLICK]</b> how to think about and use generative AI in your work regardless of which specific product you use. You will <b>[CLICK]</b> develop a mindset of iteration and skepticism.</p> <p><b>[CLICK]</b> New models and features are constantly being released, and here are some of the <b>[CLICK]</b> changes you should expect in the near future:</p> <ul style="list-style-type: none"> <li>First, genAI tools with <b>[CLICK]</b> more advanced and specialized features,</li> </ul>

	<p>like the ability to use apps for you</p> <ul style="list-style-type: none"> <li>• Expect [CLICK] cheaper tools</li> <li>• And [CLICK] faster tools</li> <li>• And [CLICK] higher quality outputs overall</li> </ul> <p>It's challenging to keep up with rapid progress in this field, but don't worry! In this course, you'll develop the metacognitive skills you need to harness those advancements in your work.</p>
 	<p>This course also demonstrates some paid features of LLMs, but you [CLICK] won't need to purchase any additional products to complete the assignments. It's important for you to see the available options, including paid options, so that you [CLICK] develop confidence experimenting and selecting the best tools in your work as a data analyst.</p> <p>This course [CLICK] does not recommend any single tool. You'll [CLICK] see several throughout the modules. And remember that the [CLICK] core principles you'll learn will prepare you to work with LLMs both free and paid, now and in the future.</p>
	<p>You'll encounter your first LLM demo and hands-on lab in Lesson 3 of this module. For now, join me in the next video to see all the exciting topics in this module. I'll see you there!</p>

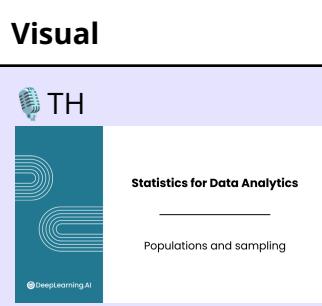
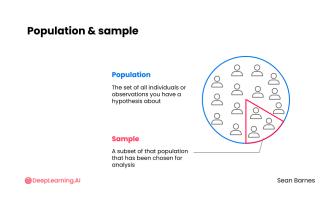
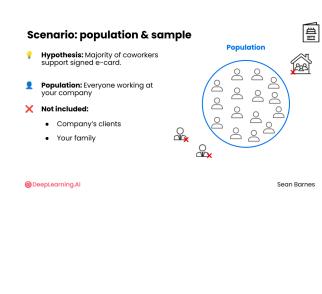
## LOV3 – Module 1 introduction

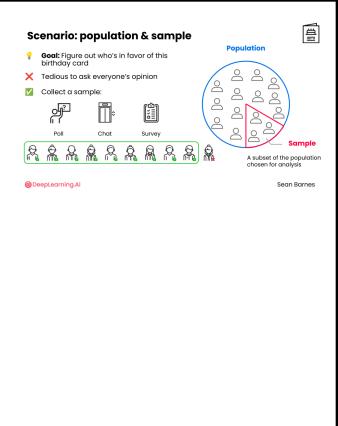
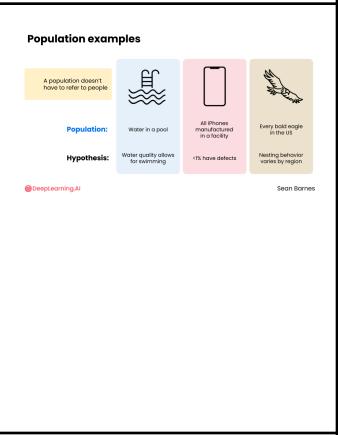
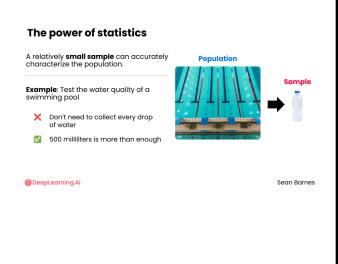
Visual	Script
	<p>Welcome to the first module of Applied Statistics for Data Analytics! Throughout this course, you will learn the fundamental statistical concepts, analyses, and visualizations that serve as the foundation for a career as a data analyst.</p>
	<p>In this module, you'll explore the essential building blocks of statistics that enable rigorous data analysis. You'll learn how to [CLICK] define populations, samples, and [CLICK] sampling methods; characterize datasets using measures of [CLICK] central tendency, [CLICK] variability, and skewness; use [CLICK] correlation to understand relationships between features; and employ segmentation to [CLICK] reveal insights about different groups within your data.</p> <p>You'll apply these concepts to real-world scenarios: analyzing [CLICK] movie ratings over time, identifying the most profitable [CLICK] loans, and analyzing [CLICK] song characteristics to create playlists. You'll also get hands-on</p>

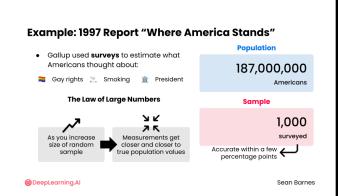
	practice with spreadsheet tools, building on the skills you learned in Data Analytics Foundations to make your analysis more efficient and effective.
	Whether you're new to statistics or looking to refresh your skills, this module will equip you with powerful techniques to extract meaningful insights from your data. By the end of this module, you will feel more confident and capable implementing rigorous statistical analyses in your career as a data analyst! Join me in the next video to get started with populations and sampling. I'll see you in the next video.

## Lesson 1

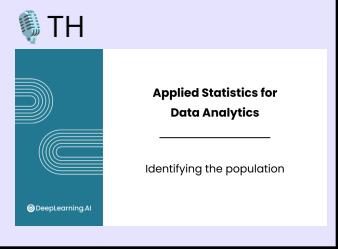
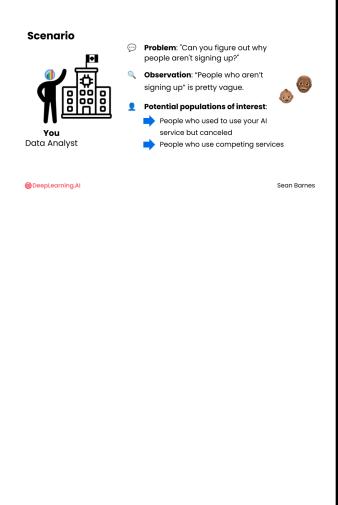
### L1V1 – Populations and sampling

Visual	Script
	<p>Say you're working at a small company of 50 people and you want to start a fun, new birthday tradition: have everyone at the company sign an electronic birthday card! You have a hypothesis that a majority of your coworkers would be on board, but how do you find out if that's true?</p> <p>You could go with the magic 8 ball [<b>Sean takes out magic 8 ball, then puts it away</b>] – just kidding. You'd probably ask a couple of coworkers, right?</p> <p>You're trying to figure out if your hypothesis about <b>everyone</b> is right by asking just a <b>few</b> people. In other words, you're taking a <b>sample</b> of a <b>population</b>. These two definitions – population and sample – underlie almost every analysis you'll carry out as a data analyst.</p>
	<p>Let's define these concepts more precisely. A [<b>CLICK</b>] population is the set of all individuals or observations that you have a hypothesis about, and a [<b>CLICK</b>] sample is a subset of that population that has been chosen for analysis.</p>
	<p>In the birthday card example, you hypothesize that the [<b>CLICK</b>] majority of your coworkers will support sending a signed e-card. So the [<b>CLICK</b>] population is everyone working at your company. That's who this new tradition will affect. Can you think of some people who are [<b>CLICK</b>] not included in this population? [<b>2 second pause</b>] The [<b>CLICK</b>] company's clients or [<b>CLICK</b>] your family would be good examples of people who aren't in this population. Their opinion ultimately won't affect the decision.</p>

 <p><b>Scenario: population &amp; sample</b></p> <p>Goal: Figure out who's in favor of this birthday card</p> <p>Tedious to ask everyone's opinion</p> <p>Collect a sample</p> <p>Population: A group of people</p> <p>Sample: A subset of the population chosen for analysis</p> <p>Sean Barnes</p>	<p>So you want to [CLICK] figure out who's in favor of this birthday card idea and who's against it. It would be [CLICK] tedious to ask everyone their opinion one-by-one in person. Instead, you could [CLICK] collect a sample!</p> <p>A [CLICK] sample is a subset of the population you've chosen for analysis. There are many ways you could sample this population: [CLICK] a meeting poll, [CLICK] elevator chat, [CLICK] a survey. Each of these samples is just a slice of the population, just a few coworkers. Say you end up asking [CLICK] 10 people, and [CLICK] 9 of them love the idea. That's a pretty good sign it could work!</p>		
 <p><b>Population examples</b></p> <p>Population: A population doesn't have to refer to people</p> <p>Hypothesis: Water in a pool allows for swimming</p> <p>Population: Water in a pool</p> <p>Hypothesis: Water quality allows for swimming</p> <p>Population: All iPhones manufactured in a factory</p> <p>Hypothesis: iPhone's have defects</p> <p>Population: Every bald eagle in the US</p> <p>Hypothesis: Nesting behavior varies by region</p> <p>Sean Barnes</p>	<p>Note that a population doesn't have to refer to people.</p> <p>Here are some [CLICK] hypotheses that relate to other types of [CLICK] populations</p> <ul style="list-style-type: none"> <li>The population might be the [CLICK] water in a swimming pool, if your hypothesis is that the [CLICK] water quality allows for swimming</li> <li>Or [CLICK] all the iPhones manufactured in a particular facility, if your hypothesis is that [CLICK] less than 1% have production defects</li> <li>Or, [CLICK] every bald eagle in the US, if your hypothesis is that their [CLICK] nesting behavior varies by region</li> </ul>		
 <p><b>Why this matters</b></p> <p>You won't typically have access to entire populations</p> <table border="1"> <tr> <td>Accessible <input checked="" type="checkbox"/> 50 coworkers</td> <td>Not accessible <input type="checkbox"/> 200 million Netflix users</td> </tr> </table> <ol style="list-style-type: none"> <li>Knowledge: Population might be unknown</li> <li>Practicality: Too much time or money required Might not be able to access its entirety</li> <li>Ethics: Testing could cause harm</li> </ol> <p>Sean Barnes</p>	Accessible <input checked="" type="checkbox"/> 50 coworkers	Not accessible <input type="checkbox"/> 200 million Netflix users	<p>This distinction between population and sample matters because [CLICK] you won't typically have access to your entire population. [CLICK] 50 coworkers, okay, you [CLICK] might be able to get a response from all of them. But surveying every Netflix user? There are over [CLICK] 200 million of them. Most populations can't or shouldn't be fully captured, and for three key reasons:</p> <ol style="list-style-type: none"> <li>First, [CLICK] knowledge. The population might just be [CLICK] unknown. For example, in a political poll, you can survey so-called "likely" voters, but you ultimately don't know who will actually go out and vote.</li> <li>Second, [CLICK] practicality. Reaching the entire population could take [CLICK] too much time or [CLICK] money, and you [CLICK] might not be able to access its entirety. For example, you might want to test every iPhone produced in a factory, but that may add an hour to the production of each phone.</li> <li>Finally, [CLICK] ethics. If you want to know if a certain fruit is toxic to dogs, well, [CLICK] testing that idea could cause harm.</li> </ol>
Accessible <input checked="" type="checkbox"/> 50 coworkers	Not accessible <input type="checkbox"/> 200 million Netflix users		
 <p><b>The power of statistics</b></p> <p>A relatively small sample can accurately characterize the population</p> <p>Example: Test the water quality of a swimming pool</p> <p>Don't need to collect every drop of water 500 milliliters is more than enough</p> <p>Population: A large swimming pool</p> <p>Sample: A small water bottle</p> <p>Sean Barnes</p>	<p>The good news is, through the power of statistics, even [CLICK] a relatively small sample, if collected properly, can accurately characterize the population. For example, say you want to [CLICK] test the water quality of a swimming pool. You [CLICK] don't need to collect every drop of water. Even for an Olympic size swimming pool, [CLICK] 500 milliliters is more than enough – about the [CLICK] size of a water bottle. [brief pause to absorb]</p>		

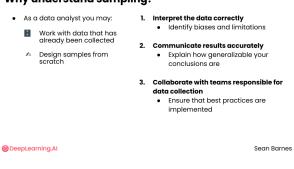
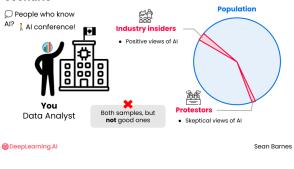
	<p>In 1997, the polling company Gallup released a report called Where America Stands. [CLICK] Gallup used surveys to estimate what all 187 million Americans thought about hundreds of issues like [CLICK] gay rights, [CLICK] smoking, and the [CLICK] president. How many people do you think they surveyed to get an accurate answer for all [CLICK] 187 million Americans? [pause for effect] [CLICK] 1,000. That's all!</p> <p>This surprising effectiveness is due to the power of random sampling and statistical theory, particularly a concept known as the [CLICK] Law of Large Numbers. This law states that as you [CLICK] increase the size of your random sample, your [CLICK] measurements from that sample tend to get closer and closer to the true values you would get if you could measure the entire population. When done correctly, a sample of 1,000 people can provide results [CLICK] accurate to within a few percentage points of the true values.</p>
 TH	<p>Think of the population as the truth: the way the world really is. Of your 50 coworkers, a certain number of them are in favor of sending birthday cards. Your sample is a window into that truth. You ask 10 coworkers and 9 are in favor. It's a glimpse at what everyone might be thinking, but it's not the whole truth.</p> <p>Join me in the next video to see how to identify the population you want to study.</p>

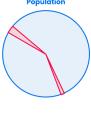
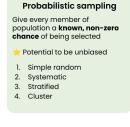
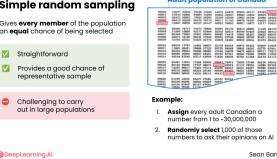
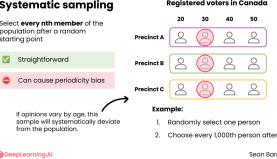
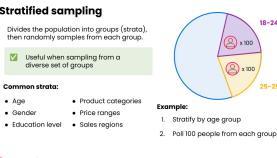
## L1V2 – Identifying the population

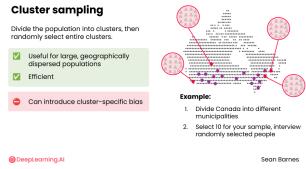
Visual	Script
	<p>Before you collect a sample, you'll need to know what you're trying to take a sample of. That may sound abstract, but the idea is that analyzing different populations can give you different insights.</p>
	<p>Let's explore an example. Suppose you're a [CLICK] data analyst at a company developing AI services in Canada. Your boss comes to you with a problem. She says, "We need to increase our user base. [CLICK] Can you figure out why people aren't signing up?"</p> <p>Of course you'll say, "Yeah, I can do that!"</p> <p>On looking deeper into the problem, though, you observe that [CLICK] "people who aren't signing up" is pretty vague. Everyone who doesn't use your service is a massive group – that's hundreds of millions of people from all over the [CLICK] globe, from [CLICK] babies to [CLICK] grandparents. [CLICK] Or maybe you're more interested in [CLICK] people who used to use your AI</p>

	<p>service but canceled. Or perhaps you want to focus on [CLICK] people who use competing services.</p> <p>Studying each of these populations would lead to very different insights. If you focus on [CLICK] former users, you might learn about retention. If you look at [CLICK] competitors' users instead, you might better understand what features you're missing.</p>
	<p>Here's a process you can use to nail down your relevant population:</p> <ol style="list-style-type: none"> <li>1. [CLICK] Start with your hypothesis or research question. In this case, it might be [CLICK] "Why aren't people signing up for our AI service?"</li> <li>2. Then, [CLICK] identify the key characteristics that define your population of interest. For an AI service, this might include factors like [CLICK] age (are we interested in all ages or just certain demographics?), [CLICK] geographic location (are we looking globally or just in certain markets?), and [CLICK] technology use (are we only interested in people who already use other AI services?).</li> <li>3. Next, [CLICK] consider practical constraints. While you might ideally want to collect data across the full population, that's probably not feasible. You might need to limit your population based on [CLICK] accessibility, [CLICK] budget, [CLICK] time, or [CLICK] ethical constraints.</li> <li>4. Then, [CLICK] think about generalizability. That means how broad do you want your conclusions to be? The [CLICK] [CLICK] more specific your population, the [CLICK] more precise your insights might be, but the [CLICK] less generalizable they'll be to other groups. A [CLICK] [CLICK] less specific population leads to [CLICK] less precise but [CLICK] more generalizable results.</li> <li>5. Finally, [CLICK] consult with stakeholders. [CLICK] Make sure your defined population aligns with the business goals. This step comes last so you can present a concrete idea for review.</li> </ol> <p>By following this process, you and your stakeholders might define the population as [CLICK] "internet-using adults aged 18-49 in Canada who currently use at least one AI service, but not ours." This clearly defined population can guide everything from your sampling method to how you interpret your results.</p>
TH	<p>Now you've seen how to identify your population. But as you've learned, you're unlikely to have access to the full extent of data. There are many ways to go about sampling your population. Join me in the next video to take a look at the gold standard: probabilistic methods.</p>

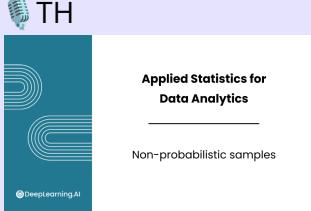
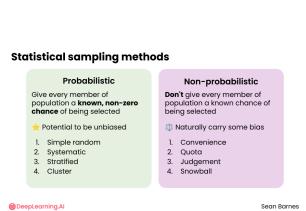
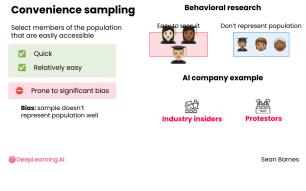
## LIV3 – Probabilistic samples

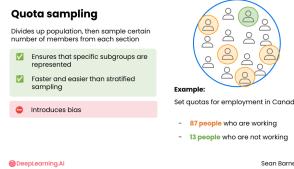
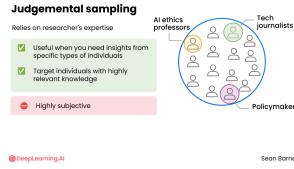
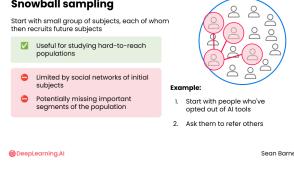
Visual	Script
 <p>TH Statistics for Data Analytics Probabilistic samples</p>	<p>Sampling is a critical skill for data analysts. Over the next few videos, you'll explore the tradeoffs involved in sampling as well as the biases it can introduce.</p>
 <p>Why understand sampling?</p> <ul style="list-style-type: none"> <li>As a data analyst you may:       <ul style="list-style-type: none"> <li>Work with data that has already been collected</li> <li>Design samples from scratch</li> </ul> </li> </ul> <ol style="list-style-type: none"> <li>Interpret the data correctly       <ul style="list-style-type: none"> <li>Identify biases and limitations</li> </ul> </li> <li>Communicate results accurately       <ul style="list-style-type: none"> <li>Explain how generalizable your conclusions are</li> </ul> </li> <li>Collaborate with teams responsible for data collection       <ul style="list-style-type: none"> <li>Ensure that best practices are implemented</li> </ul> </li> </ol> <p>Sean Barnes</p>	<p>[CLICK] As a data analyst, depending on the size of your team, you may [CLICK] work with data that has already been collected, or you may have to [CLICK] design a sample from scratch. Regardless, it's crucial for you to be able to understand the principles of sampling for a few key reasons:</p> <ul style="list-style-type: none"> <li>First, understanding sampling allows you to [CLICK] interpret the data correctly, [CLICK] identifying potential biases and limitations in the data</li> <li>It also allows you to [CLICK] communicate your results accurately. In particular, you will be able to [CLICK] explain how generalizable your conclusions are, meaning how accurately they reflect the population.</li> <li>You will also be better able to [CLICK] collaborate with teams responsible for data collection, [CLICK] ensuring that best practices are implemented.</li> </ul> <p>Let's explore sampling through an example.</p>
 <p>Scenario</p> <p>Problem: Understand Canadian public opinion on AI</p> <p>Population</p> <p>Industry insiders</p> <p>Protestors</p> <p>You Data Analyst</p> <p>Both samples, but not good ones</p> <p>Sean Barnes</p>	<p>Suppose you've been working at the Canadian AI company for a while and your CEO comes to you with a new problem. She wants to [CLICK] understand Canadian public opinion on AI to guide the company's roadmap. How would you go about collecting that data?</p> <p>Your first thought might be finding [CLICK] people who know what AI is. So you head to an [CLICK] AI conference, where you find a lot of enthusiastic [CLICK] industry insiders. You interview 100 attendees, and 95 of them express overwhelmingly [CLICK] positive views.</p> <p>Things seem promising, but as you're leaving the conference, you notice a group of [CLICK] protesters outside. Their signs display slogans about the potential negative impacts of AI. You interview 50 of them, and [CLICK] nearly all express skeptical views of AI.</p> <p>You've recorded two very different perspectives, but do either of these groups accurately represent overall public opinion on AI? Probably not. These are [CLICK] both samples, but they're not very good ones.</p>

 <p><b>Scenario</b> Goal: To analyze data that represents the entire population you're interested in. ↳ Adult population of Canada</p> <p><b>Sampling methods:</b> <input checked="" type="checkbox"/> Probabilistic    <input type="checkbox"/> Non-probabilistic</p> <p>©DeepLearning.AI    Sean Barnes</p>	<p>As a data analyst, your goal is to [CLICK] analyze data that represents the entire population you're interested in - in this case, [CLICK] the adult population of Canada – [CLICK] not just the vocal supporters and opponents. Statistics has a lot to say on how to do that, so let's talk about some [CLICK] specific sampling methods, which fall into two categories: [CLICK] probabilistic and [CLICK] nonprobabilistic. Let's look at the probabilistic ones first.</p>
 <p><b>Statistical sampling methods</b> Probabilistic sampling Give every member of population a known, non-zero chance of being selected Potential to be unbiased 1. Simple random 2. Systematic 3. Stratified 4. Cluster</p> <p>©DeepLearning.AI    Sean Barnes</p>	<p>Probabilistic sampling methods [CLICK] give every member of the population a known, non-zero chance of being selected. These methods have the [CLICK] potential to be unbiased, which means that your sample is truly representative of the population. That's what you're aiming for. The four most important probabilistic samples are:</p> <ul style="list-style-type: none"> <li>• [CLICK] Simple random,</li> <li>• [CLICK] Systematic,</li> <li>• [CLICK] Stratified, and</li> <li>• [CLICK] Cluster.</li> </ul>
 <p><b>Simple random sampling</b> Gives every member of the population an equal chance of being selected  <input checked="" type="checkbox"/> Straightforward  <input checked="" type="checkbox"/> Provides a good chance of representative sample  <span style="color: red;">☐ Challenging to carry out in large populations</span></p> <p>Example: 1. Assign every adult Canadian a number from 1 to ~30,000,000 2. Randomly select 1,000 of those numbers to call them up one on AI</p> <p>©DeepLearning.AI    Sean Barnes</p>	<p>Simple random sampling gives every member of the population an equal chance of being selected. It's [CLICK] straightforward and [CLICK] provides a good chance of having a representative sample, but it can be [CLICK] challenging to carry out in large populations.</p> <ul style="list-style-type: none"> <li>• [CLICK] In the AI opinion example from earlier, a simple random sample could look like [CLICK] assigning every adult Canadian a number from 1 to around 30 million, the total population of Canada, then [CLICK] randomly selecting 1,000 of those numbers. You can then call up each person whose number was chosen to ask them their opinion about AI.</li> </ul> <p>It sounds tough, right? Not every Canadian has a phone number and it might be expensive to place many phone calls this way.</p>
 <p><b>Systematic sampling</b> Select every nth member of the population after a random starting point  <input checked="" type="checkbox"/> Straightforward  <span style="color: red;">☐ Opinions vary by age; this sample will systematically deviate from the population</span></p> <p>Example: 1. Randomly select one person 2. Choose every 1000th person after</p> <p>©DeepLearning.AI    Sean Barnes</p>	<p>In systematic sampling, you select every nth member of the population after a random starting point. For example, you might have a [CLICK] list of all registered voters in Canada, [CLICK] randomly select one person, and [CLICK] then choose every 1,000th person after them on the list.</p> <p>This method is [CLICK] straightforward, but [CLICK] can cause periodicity bias if there's a pattern in your list. For example, if voters were sorted by [CLICK] precinct, [CLICK] then by age within that precinct, it would be possible to sample only people with similar ages. [CLICK] If opinions on AI vary by age, this sample will systematically deviate from the population.</p>
 <p><b>Stratified sampling</b> Divides the population into groups (strata), then randomly samples from each group  <input checked="" type="checkbox"/> Useful when sampling from a diverse set of groups</p> <p>Common strata:  <input type="checkbox"/> Age    <input type="checkbox"/> Product categories  <input type="checkbox"/> Gender    <input type="checkbox"/> Price ranges  <input type="checkbox"/> Education level    <input type="checkbox"/> Sales regions</p> <p>Example: 1. Stratify by age group 2. Pick 100 people from each group</p> <p>©DeepLearning.AI    Sean Barnes</p>	<p>Next up: stratified sampling. This method divides the population into groups, or strata, based on shared characteristics, then randomly samples from each group. [CLICK] It's useful when you want to make sure to sample from a diverse set of groups. [CLICK] Common strata include [CLICK] age, [CLICK] gender, and [CLICK] education level for people, or [CLICK] product categories,</p>

	<p><b>[CLICK]</b> price ranges, and <b>[CLICK]</b> sales regions for products.</p> <ul style="list-style-type: none"> <li>• <b>[CLICK]</b> In the AI example, you could <b>[CLICK]</b> stratify by age group – <b>[CLICK]</b> 18-24, <b>[CLICK]</b> 25-29, and so on – <b>[CLICK]</b> and poll 100 people from each age group.</li> </ul>
 Cluster sampling Divide the population into clusters, then randomly select entire clusters. ✓ Useful for large, geographically dispersed populations ✓ Efficient ❌ Can introduce cluster-specific bias  Example: 1. Divide Canada into different municipalities. 2. Select 10 of your sample. Interview randomly selected people.  ©DeepLearning.AI Sean Barnes	<p>Finally, you have cluster sampling, <b>[CLICK]</b> which is useful for large, geographically dispersed populations. You <b>[CLICK]</b> divide the population into clusters (often geographically), then randomly select entire clusters. <b>[CLICK]</b> It's efficient, but <b>[CLICK]</b> can introduce cluster-specific biases.</p> <ul style="list-style-type: none"> <li>• In the <b>[CLICK]</b> AI example, you could <b>[CLICK]</b> divide Canada up into different municipalities, <b>[CLICK]</b> then select 10 of them for your sample. You'll send a team to interview randomly selected people within each municipality. You can see how it might be easier to send teams out to just 10 areas, versus randomly sampling every Canadian adult, each of whom could live in the <b>[CLICK]</b> furthest corner of the country.</li> </ul>
 TH	<p>These four methods form the core of probabilistic sampling, helping to select a sample that's representative of the population you're trying to understand.</p> <p>Now, follow me to the next video to learn about <b>non</b>-probabilistic sampling methods.</p>

## L1V4 – Non-probabilistic samples

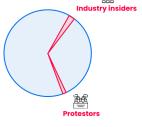
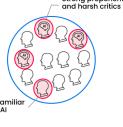
Visual	Script
 TH	<p>When budget, time, or workforce constraints make probabilistic samples infeasible, you can use a non-probabilistic sample. Non-probabilistic samples are generally less rigorous, but often more practical to carry out.</p>
 Statistical sampling methods Probabilistic Give every member of population a known, non-zero chance of being selected ★ Potential to be unbiased 1. Simple random 2. Systematic 3. Stratified 4. Cluster  Non-probabilistic Don't give every member of population a known chance of being selected ☐ Naturally carry some bias 1. Convenience 2. Quota 3. Judgement 4. Snowball  ©DeepLearning.AI Sean Barnes	<p>Nonprobabilistic sampling methods <b>[CLICK]</b> <b>don't</b> give every member of the population a known chance of being selected. These methods often arise from practical constraints and <b>[CLICK]</b> naturally carry some bias.</p> <p>The core non probabilistic methods you should know about as a data analyst are <b>[CLICK]</b> convenience, <b>[CLICK]</b> quota, <b>[CLICK]</b> judgment, and <b>[CLICK]</b> snowball samples. Let's get to it.</p>
 Convenience sampling Select members of the population that are easily accessible ✓ Quick ✓ Relatively easy  ❌ Prone to significant bias ❌ Sample doesn't represent population well  ©DeepLearning.AI Sean Barnes	<p>In convenience sampling, you <b>[CLICK]</b> select members of the population that are easily accessible. On the upside, it's <b>[CLICK]</b> quick and <b>[CLICK]</b> relatively easy to do. However, it's <b>[CLICK]</b> prone to significant bias. Bias just means your <b>[CLICK]</b> sample doesn't represent your population well. You'll learn more about it in the next video.</p>

	<p>A commonly cited example of convenience sampling is [CLICK] behavioral research. Because many researchers work at universities, they can [CLICK] easily recruit students to complete an experiment. But, these students [CLICK] don't represent the overall population very well. They're often younger than the average person, for example.</p> <ul style="list-style-type: none"> <li>In the [CLICK] AI company example from the previous video, surveying just the [CLICK] attendees and protesters at one conference is a convenience sample. Just one location, no randomness.</li> </ul>
 <p><b>Quota sampling</b> Divides up population, then samples certain number of members from each section  <input checked="" type="checkbox"/> Useful for studying hard-to-reach populations  <input checked="" type="checkbox"/> Faster and easier than stratified sampling  <input type="checkbox"/> Introduces bias</p> <p><b>Example:</b> Set quotas for employment in Canada      - 87 people who are working      - 13 people who are not working</p> <p>© DeepLearning.AI Sean Barnes</p>	<p>Quota sampling divides up the population based on certain characteristics, then samples a certain number of members from each section. You can think of it as the nonprobabilistic version of stratified sampling. It is often used to [CLICK] ensure that specific subgroups are represented. Like convenience sampling, it's [CLICK] faster and easier than stratified sampling, but [CLICK] introduces bias.</p> <ul style="list-style-type: none"> <li>[CLICK] In the AI example, you might set quotas to match the employment demographics of Canada, where the employment rate is 87%. So, you might interview [CLICK] 87 people who are employed, and [CLICK] 13 who are not. But without random sampling, interviewers may inadvertently skew the results based on who they choose to question.</li> </ul>
 <p><b>Judgemental sampling</b> Relies on researcher's expertise  <input checked="" type="checkbox"/> Useful when you need insights from specific types of individuals  <input checked="" type="checkbox"/> Target individuals with highly relevant knowledge  <input type="checkbox"/> Highly subjective</p> <p>© DeepLearning.AI Sean Barnes</p>	<p>Then you have judgmental sampling, which relies on the researcher's expertise to select the sample. [CLICK] It can be useful when you need insights from specific types of individuals. The benefit is that [CLICK] you can target individuals with highly relevant knowledge. The downside? [CLICK] It's highly subjective.</p> <ul style="list-style-type: none"> <li>For the [CLICK] AI company, using judgmental sampling might mean specifically choosing to interview [CLICK] AI ethics professors, [CLICK] tech journalists, and [CLICK] policymakers, based on the belief that their opinions are particularly valuable.</li> </ul>
 <p><b>Snowball sampling</b> Start with small group of subjects, each of whom then recruits future subjects  <input checked="" type="checkbox"/> Useful for studying hard-to-reach populations  <input type="checkbox"/> Limited by social networks of initial subjects  <input type="checkbox"/> Potentially missing important segments of the population</p> <p><b>Example:</b>      1. Start with people who've opted out of AI tools      2. Ask them to refer others</p> <p>© DeepLearning.AI Sean Barnes</p>	<p>Finally, take a look at snowball sampling. You [CLICK] start with a small group of subjects, each of whom then recruits future subjects from among their friends, family, and coworkers. Your sample size grows like a rolling snowball.</p> <p>This method is particularly [CLICK] useful for studying hard-to-reach populations. For instance, undocumented immigrants may be difficult to source for an interview, but you can start with a few contacts who then refer others from their community.</p> <p>The big drawback is that your sample is [CLICK] limited to the social networks of your initial subjects, [CLICK] potentially missing important segments of the population.</p> <ul style="list-style-type: none"> <li>In the [CLICK] AI example, snowball sampling might be used to study</li> </ul>

	<p>the opinions of people who've opted out of using AI tools. You could [CLICK] start with a few such individuals and [CLICK] ask them to refer others who share their habits.</p>
TH	<p>You've heard the term "bias" a few times now. It's one of the biggest issues with sampling. In the next video, you'll explore several common types of bias and how to mitigate them. I'll see you there.</p>

## L1V5 – Types of bias

Visual	Script
<p>TH</p> <p>Statistics for Data Analytics</p> <p>Types of bias</p> <p>What is bias? A systematic difference between sampling method and population • Happens in a predictable, likely preventable way</p> <p>Random occurrence</p> <p>Sean Barnes</p>	<p><b>Bias</b> in sampling happens when your sample doesn't represent your population of interest very well.</p> <p>Similar to how bias can negatively affect interactions between people, bias in sampling can lead to poor decision-making. Let's see how it happens.</p>
<p>What is bias? A systematic difference between sampling method and population • Happens in a predictable, likely preventable way</p> <p>Random occurrence</p> <p>Sean Barnes</p>	<p>Here's a formal definition for bias in sampling: a systematic difference between the sample and the population that creates an inaccurate depiction of reality. "Systematic" is important here – it means that [CLICK] the problem occurs in a predictable, likely preventable way.</p> <ul style="list-style-type: none"> <li>For example, if you're [CLICK] interviewing randomly selected Canadians about their opinion on AI, and one person happens to respond [CLICK] more positively than the average person, that's just a normal [CLICK] random occurrence. You expect some variation in opinion.</li> </ul>

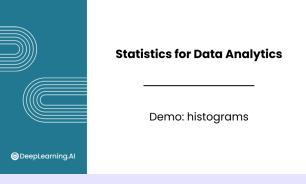
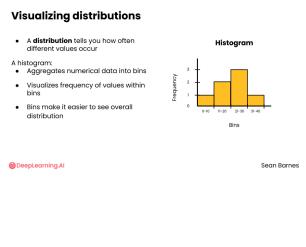
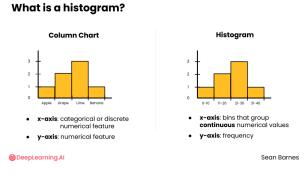
<p><b>Sampling bias</b></p> <p>Sample doesn't accurately represent the population of interest</p> <ul style="list-style-type: none"> <li>Unlikely those groups represent views of most Canadians</li> </ul>  <p>Sean Barnes</p> <p>DeepLearning.AI</p>	<p>Let's start with sampling bias. It's quite common, and it happens when a [CLICK] sample doesn't accurately represent the population of interest. This is the kind of bias that a researcher would run into if they interviewed [CLICK] people at an AI conference plus the [CLICK] protestors outside. [CLICK] It's unlikely those two extreme groups adequately represent the views of most Canadians.</p>
<p><b>Selection bias</b></p> <p>Sample is selected in a nonrandom way</p> <ul style="list-style-type: none"> <li>Mismatch between sample and population of interest</li> <li>Some groups are under or overrepresented</li> </ul> <ul style="list-style-type: none"> <li>Use probabilistic sampling methods</li> <li>Avoid over- or under-sampling any particular group</li> <li>Be transparent about limitations</li> </ul>  <p>Sean Barnes</p> <p>DeepLearning.AI</p>	<p>A common form of sampling bias is called selection bias. It happens when a [CLICK] sample is selected in a nonrandom way, causing a [CLICK] [CLICK] mismatch between the sample and the population of interest. As a result, [CLICK] some groups in the population are under or overrepresented.</p> <p>[CLICK] In the AI conference example, [CLICK] strong proponents and harsh critics are overrepresented, while people who [CLICK] aren't very familiar with AI are underrepresented.</p> <p>To avoid selection bias, [CLICK] use probabilistic sampling methods. [CLICK] Avoid over- or under-sampling any particular group. You should also [CLICK] be transparent about the limitations your sample may have. If you know your sample introduced selection bias, explain how that bias affects your conclusions.</p>
<p><b>Nonresponse bias</b></p> <ul style="list-style-type: none"> <li>Common in samples of people</li> <li>Not a random sample</li> <li>Typically really energized or really upset</li> </ul> <ul style="list-style-type: none"> <li>Consider sending follow-ups</li> <li>Include small incentives for participation</li> </ul>  <p>Sean Barnes</p> <p>DeepLearning.AI</p>	<p>Another common type of sampling bias is nonresponse bias. It's [CLICK] common in samples of people. For example, [CLICK] smartphone apps may ask a user to rate their app after using it for a while. But [CLICK] users can dismiss the notification rather than actually rating the app.</p> <p>By consequence, the [CLICK] people that do leave a review [CLICK] aren't a random sample. They [CLICK] often have a positive view of the app, since they're willing to accept the prompt and invest time in a review.</p> <p>To counteract nonresponse bias, [CLICK] consider sending follow-ups or [CLICK] even including small incentives for participation.</p>
<p><b>Instrument bias</b></p> <p>Comes from having faulty equipment or a poorly designed survey</p> <ul style="list-style-type: none"> <li>Use good quality tools</li> <li>Use well-worded surveys</li> <li>Take multiple measurements</li> </ul>  <p>Example: In 2020, Fitbit replaced smartwatches free due to issue with heart sensor</p> <p>Example: As simple as asking... "You agree that pineapple belongs on pizza, right?"</p> <p>Sean Barnes</p> <p>DeepLearning.AI</p>	<p>Next up: two key types of measurement bias. The first one, instrument bias, [CLICK] comes from having faulty equipment or a poorly designed survey. [CLICK] For example, in 2020, Fitbit [CLICK] replaced some of its smartwatches free of charge due to an issue with their heart rate sensor. But instrument bias can also be [CLICK] as simple as asking, [CLICK] "You agree that pineapple belongs on pizza, right?", which is a leading question. Make sure you're [CLICK] [CLICK] using good quality tools – [CLICK] and well-worded surveys – [CLICK] and take multiple measurements if you can.</p>

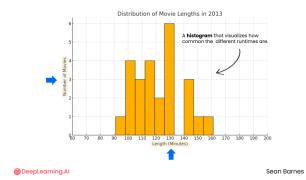
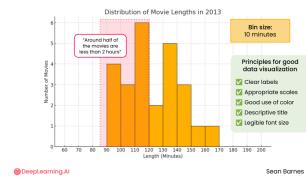
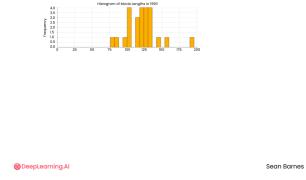
 <p><b>Observer bias</b> Occurs when the person measuring lets expectations affect what they see</p> <ul style="list-style-type: none"> <li>Might unconsciously focus on positive feedback</li> <li>Have multiple people take the measurement</li> </ul> <p><b>For example:</b></p> <ul style="list-style-type: none"> <li>Action: Rolling out dark mode</li> <li>Expectation: Feature will be well-reviewed</li> </ul> <p>DeepLearning.AI</p>	<p>Another type of measurement bias is observer bias. It <b>[CLICK]</b> occurs when the person doing the measuring lets their own expectations affect what they see.</p> <p><b>[CLICK]</b> For example, say your company is <b>[CLICK]</b> rolling out dark mode for one of its apps. If you as the data analyst <b>[CLICK]</b> expect this feature to be well-reviewed, then when you're interviewing a sample of users, you <b>[CLICK]</b> might unconsciously focus on their positive feedback. Try <b>[CLICK]</b> having multiple people take the measurement if possible.</p>
 <p><b>Response bias</b> Observed in samples of people</p> <ul style="list-style-type: none"> <li>People may not be completely forthcoming when answering questions</li> </ul> <p><b>Questions like:</b></p> <ul style="list-style-type: none"> <li>How much money do you make?</li> <li>When was the last time you dipped your french fries in mayonnaise?</li> </ul> <p><b>You can try:</b></p> <ul style="list-style-type: none"> <li>Emphasizing how important honest answers are</li> <li>Designing your questions to be as objective as possible</li> </ul> <p>DeepLearning.AI</p>	<p>Another type of bias is response bias, which is often <b>[CLICK]</b> observed in samples of people. Essentially, <b>[CLICK]</b> people may not be completely forthcoming when answering questions.</p> <p>So <b>[CLICK]</b> questions like <b>[CLICK]</b> "How much money do you make?" or <b>[CLICK]</b> "When was the last time you dipped your french fries in mayonnaise?" might elicit untruthful responses, even on an anonymous survey. Mitigating this type of bias isn't easy. <b>[CLICK]</b> You can try <b>[CLICK]</b> emphasizing how important honest answers are, or <b>[CLICK]</b> designing your questions to be as objective as possible.</p>
 <p><b>Analysis bias</b> The way you perform your analysis can lead to results that don't reflect the truth</p> <ul style="list-style-type: none"> <li>Confirmation bias</li> <li>Stay open-minded</li> <li>Let the data tell its own story</li> </ul> <p><b>Good:</b> A new feature is improving in user retention</p> <p><b>Problem:</b> Under some pressure to show that the feature was worth developing</p> <p><b>Analysis bias:</b> Focus on the two that show the most positive results</p> <p>DeepLearning.AI</p>	<p>Finally, let's consider analysis bias. You've sampled your population and collected your data, but the potential for bias isn't over! If you're not careful, <b>[CLICK]</b> the way you perform your analysis can lead to results that don't reflect the truth. The biggest pitfall is <b>[CLICK]</b> confirmation bias – looking for evidence for something you already believe, rather than looking at all the evidence objectively.</p> <p>For example, a <b>[CLICK]</b> product manager might <b>[CLICK]</b> hope that a new feature is improving user retention. They may be under some pressure to show that the feature was worth developing. They might <b>[CLICK]</b> analyze 10 metrics, but choose to <b>[CLICK]</b> focus on the two that show the most positive results.</p> <p>It's okay to have a hypothesis or a goal, of course, but <b>[CLICK]</b> stay open-minded and <b>[CLICK]</b> let the data tell its own story, even if it's not the one you were hoping to hear.</p>
	<p>Bias is often unavoidable, and most samples will contain some degree of it. However, by following best practices, in most cases you can mitigate its effects.</p> <p>That takes you to the end of this lesson! Great work so far. You've learned the core concepts of populations and sampling, which underlie all of statistics. After the practice assessment for this lesson, join me in the next module to learn how to characterize samples through measures of central tendency,</p>

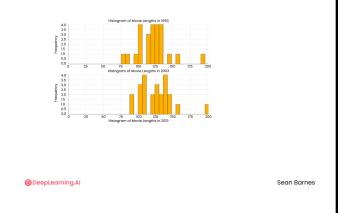
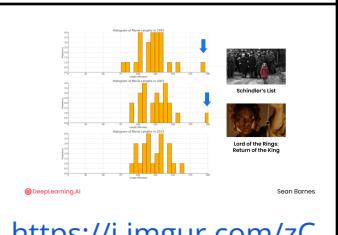
variability, and skewness. I'll see you in the next lesson.

## Lesson 2 – Central tendency

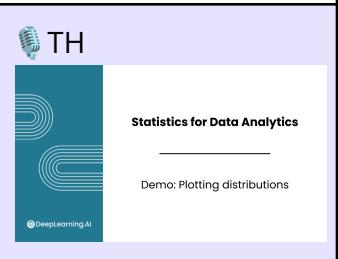
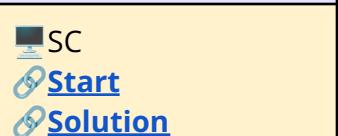
### L2V1 – Histograms

Visual	Script
 TH  	Once you've collected a sample, it's often useful to start your analysis by visualizing its distribution.
	<p><b>[CLICK]</b> A distribution tells you how often different values occur in a population or in your sample. Descriptive statistics is all about characterizing the distribution of your data.</p> <p>One useful visualization for distributions is <b>[CLICK]</b> the histogram.</p> <p>A histogram <b>[CLICK]</b> aggregates numerical data into bins and <b>[CLICK]</b> visualizes the frequency of the values within those bins. The purpose of the bins is to make it <b>[CLICK]</b> easier to see the overall distribution by showing how common ranges of values are.</p>
	Histograms are a special version of a <b>[CLICK]</b> column chart. In a column chart, you are trying to compare a numerical feature across different categories, so you display the <b>[CLICK]</b> categorical (or discrete numerical feature) on the x axis and <b>[CLICK]</b> the numerical feature on the y axis.
 TH	<p>In a histogram, you're interested in visualizing the frequency of different values in a continuous numerical feature. On the x-axis, it plots <b>[CLICK]</b> bins that group the <b>continuous</b> numerical values into categories, with the <b>[CLICK]</b> frequency on the y axis. Frequency means the number of observations in the data that fall into a particular bin. You may also see the proportion of observations plotted on the y axis, rather than the raw frequency.</p> <p>Let's work with some real data together.</p> <p>Do you feel like movies have been getting longer lately? <b>[pause]</b> Maybe it's just my attention span, but I've heard this hypothesis several times lately. Say you wanted to figure out how long movies typically are, and if they've been getting</p>

	<p>longer over time. You sample the 25 most popular movies from 2013.</p> <p>What can you do to describe this distribution of movie durations?</p>
 <a href="#">SC</a> <a href="#">Spreadsheet</a> <a href="#">Source data</a>	<p>So here's the data of movie durations of the top 25 most popular movies in the year 2013. If you'd like to explore the data yourself, you can find this spreadsheet in the downloads tab accompanying this video.</p> <p>The columns in the data set are the name of the movie, the year it was released – all 2013 – the rating out of 10 – this is from the International Movie Database or IMDB, the number of IMDB ratings, and the duration of the movie. My personal fav here was Iron Man 3, but maybe you liked The Great Gatsby. You learned a moment ago that a distribution represents how often values in the sample data occur. You have a sample of movies here. So what <b>values</b> should you examine to answer the question? <b>[pause for learner to think]</b> Those would be the durations of the movies!</p>
 <a href="https://i.imgur.com/Y15Hrj5.png">https://i.imgur.com/Y15Hrj5.png</a>	<p>Here's a <b>[CLICK]</b> histogram that visualizes how common the different movie durations are. On the x axis is the <b>[CLICK]</b> duration of the movie in minutes, and <b>[CLICK]</b> the y axis shows the number of movies with that duration.</p> <p>Note that duration is a continuous numerical feature, which has been grouped into <b>[CLICK]</b> 7-minute bins, giving this histogram <b>[CLICK]</b> 10 bins. Any movie with a length between 91 and 98 minutes is represented by <b>[CLICK]</b> this column.</p>
 <a href="https://i.imgur.com/0mqXBsD.png">https://i.imgur.com/0mqXBsD.png</a>	<p>Here's the same data but this time in histogram with <b>[CLICK]</b> bins of 10 minutes, which are a bit easier to interpret, though the overall picture of the data is similar. These bin sizes make it easier to say something like <b>[CLICK]</b> “around half of the movies are less than 2 hours.” Note that too few bins can overly simplify the data, and too many bins can make it difficult to identify any overall patterns.</p> <p>Aside from bin size, keep in mind your <b>[CLICK]</b> principles for good data visualization: <b>[CLICK]</b> clear labels, <b>[CLICK]</b> appropriate scales, <b>[CLICK]</b> good use of color, <b>[CLICK]</b> descriptive title. And make sure your <b>[CLICK]</b> font size is legible for your audience.</p>
 <a href="#">DeepLearning.AI</a> Sean Barnes	<p>Multiple histograms can also be plotted next to each other to compare distributions. Here are the distributions of movies in 1993 at the top,</p>

 @DeepLearning.AI Sean Barnes	<p>then 2003,</p>
 @DeepLearning.AI Sean Barnes	<p>then 2013 on the bottom. These all use the same bin size, to make comparison easier. What can you say about the change in durations over time? <b>[pause for learner to reflect]</b> It's a bit hard to say. Perhaps there was a small shift to the right? They look pretty consistent overall.</p>
<a href="https://i.imgur.com/zCxnuz.png">https://i.imgur.com/zCxnuz.png</a>	<p><b>[CLICK]</b> Check out those 200-minute movies though. <b>[CLICK]</b> Can you guess what they are? <b>[brief pause]</b> in 1993, <b>[CLICK]</b> that was Schindler's List; and in 2003, it was <b>[CLICK]</b> The Lord of the Rings: Return of the King.</p>
	<p>Histograms help visualize the distribution of a continuous numerical feature. Follow me to the next video to see how to create them in Google Sheets.</p>

## L2V2 – Demo: Plotting distributions

Visual	Script
 Statistics for Data Analytics	<p>Now that you've learned how to interpret histograms, let's take a look at how to plot distributions, using both histograms and column charts.</p>
 <u><a href="#">Source</a></u>	<p>The demos in this module use a subset of the Lending Tree Loan Dataset. The data set includes thousands of loans made through the lending tree platform, which allows individuals to lend to other individuals.</p> <ul style="list-style-type: none"> <li>Each row represents a particular loan, and each column includes features about the person who took out the loan, like their job title <b>[A]</b> and annual income <b>[F]</b>, as well as information about the loan like the amount <b>[K]</b>.</li> </ul> <p>Remember, if you'd like to follow along with the demo, you can find this spreadsheet and the solution in the downloads tab.</p> <p>Imagine you are thinking of becoming a lender on this platform. Before committing to anything, you want to better understand the risks involved. You can perform a statistical analysis on the existing loans to try to identify the level of risk of each loan, and which factors seem to affect risk.</p>
	<p>As a potential lender, you might be interested in the paid interest feature <b>[P]</b>. This feature essentially represents the amount of money the lender has made</p>

so far on the loan.

### Histogram

- Select column P
- Insert Chart → should be a histogram by default, or change it to one
- What if you want to customize the bar width?
- Customize → Histogram → Bucket size → 50
- Move to own sheet “paid\_interest\_chart”
- Title → Distribution of Paid Interest
- Comment on chart → looks like most loans have paid out around \$50 to \$600 dollars, with some loans paying a lot more than that

For discrete features you need to use a column chart rather than a histogram. You already learned this visualization in Data Analytics Foundations, but let's see how a column chart can show you the overall distribution of your data.

One categorical feature present in this dataset is the loan **grade**. The grade is the quality score given to each loan based on the borrower's credit history and some other factors. It's useful to you as a potential lender to better understand the distribution of good and bad loans on the platform

### Column chart

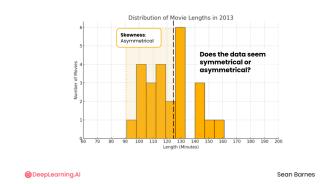
- Highlight the column N2:N
- Insert Chart - column chart
- Notice that the grades on the x axis are not in order, making it hard to visualize the distribution.
- You can sort the data to avoid this issue
  - Apply filter → sort A to Z
- Move to own sheet “grade\_chart”
- Add title “Distribution of Loan Grades”
- Notice the Aggregate box. This is on by default, and what it does is count the frequency of each grade in the dataset, rather than trying to display each value individually
- Great! Now you can see many loans are in the higher grades, A, B, and C. That may motivate you more as an investor, but it's worth investigating these terms more. Is a C significantly worse than an A? It's worth looking into



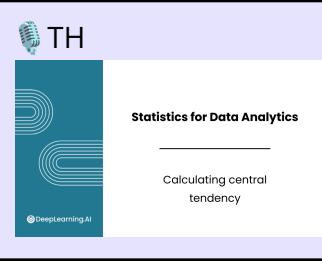
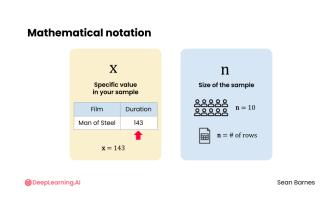
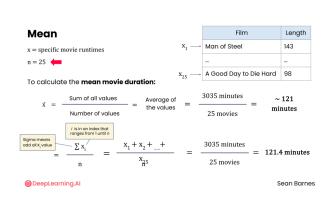
Great work on histograms! Join me in the next video to learn how to complement this visualization approach by calculating descriptive statistics.

## L2V3 – Central tendency, variability, and skewness

Visual	Script
<p>Once you've visualized the data in your sample, you'll want to describe it. But how? What are the most useful statistics for explaining what your sample is like?</p>	
<p>You can characterize a distribution through measures of [CLICK] central tendency, [CLICK] variability, and [CLICK] skewness. Here are the questions each type of descriptive statistic answers:</p> <ul style="list-style-type: none"> <li>For central tendency: [CLICK] where is the middle of the data? [CLICK] What values are most frequent? Common statistics include [CLICK] mean (or average), [CLICK] median, and [CLICK] mode.</li> <li>For variability: [CLICK] how different are the values in the data? [CLICK] Are they clustered together, or spread out? There are many [CLICK] measures of variability, including [CLICK] range, [CLICK] interquartile range, [CLICK] variance, and [CLICK] standard deviation.</li> <li>For skewness: [CLICK] is the distribution symmetrical, or does it lean to one side? Skewness characterizes the direction and extent of any asymmetry in the data. You can [CLICK] compare the [CLICK] mean and [CLICK] median to approximate skewness, or [CLICK] calculate it directly in a spreadsheet.</li> </ul> <p>You've learned many other useful descriptive statistics in Data Analytics Foundations, so don't forget about them! These include [CLICK] min, [CLICK] max, and [CLICK] frequency.</p>	<ul style="list-style-type: none"> <li>For central tendency: [CLICK] where is the middle of the data? [CLICK] What values are most frequent? Common statistics include [CLICK] mean (or average), [CLICK] median, and [CLICK] mode.</li> <li>For variability: [CLICK] how different are the values in the data? [CLICK] Are they clustered together, or spread out? There are many [CLICK] measures of variability, including [CLICK] range, [CLICK] interquartile range, [CLICK] variance, and [CLICK] standard deviation.</li> <li>For skewness: [CLICK] is the distribution symmetrical, or does it lean to one side? Skewness characterizes the direction and extent of any asymmetry in the data. You can [CLICK] compare the [CLICK] mean and [CLICK] median to approximate skewness, or [CLICK] calculate it directly in a spreadsheet.</li> </ul> <p>You've learned many other useful descriptive statistics in Data Analytics Foundations, so don't forget about them! These include [CLICK] min, [CLICK] max, and [CLICK] frequency.</p>
<p>Let's return to the movie durations data. In an earlier video, you saw this histogram of the durations. For now, let's explore the central tendency, variability, and skewness of this data by eye, then you'll see how to calculate these measures.</p> <ul style="list-style-type: none"> <li>For the mean, imagine you are trying to balance the distribution on the tip of your finger. You're finding the center of mass of the data. [CLICK] Where would you place the tip of your finger along the x-axis to balance the chart? [pause for learner to think] The center of mass looks to be [CLICK] somewhere around 120 or 130 minutes.</li> </ul>	<ul style="list-style-type: none"> <li>For the mean, imagine you are trying to balance the distribution on the tip of your finger. You're finding the center of mass of the data. [CLICK] Where would you place the tip of your finger along the x-axis to balance the chart? [pause for learner to think] The center of mass looks to be [CLICK] somewhere around 120 or 130 minutes.</li> </ul>
<ul style="list-style-type: none"> <li>What about variability? Look at how [CLICK] spread out the movie durations are on the x-axis. [CLICK] Are most of the movies clustered closely around the mean, or are they spread out? [pause for learner to think] If the mean is around 120, most of the durations are [CLICK] within 20 minutes to the left and right of the mean. That's % of the average duration, so they're relatively close together. There aren't any movies in this sample that are [CLICK] 60 minutes or 200 minutes long.</li> </ul>	<ul style="list-style-type: none"> <li>What about variability? Look at how [CLICK] spread out the movie durations are on the x-axis. [CLICK] Are most of the movies clustered closely around the mean, or are they spread out? [pause for learner to think] If the mean is around 120, most of the durations are [CLICK] within 20 minutes to the left and right of the mean. That's % of the average duration, so they're relatively close together. There aren't any movies in this sample that are [CLICK] 60 minutes or 200 minutes long.</li> </ul>

 <p><b>Sean Barnes</b></p>	<ul style="list-style-type: none"> <li>How about skewness? <b>[CLICK]</b> Does this data seem symmetrical or asymmetrical? <b>[pause for learner to think]</b> If you <b>[CLICK]</b> draw a line down the center here at the mean, you can see some asymmetry. There seem to be <b>[CLICK]</b> more movies on the left side, with shorter durations. However, the skewness isn't extreme.</li> </ul>
 TH <p><b>Sean Barnes</b></p>	<p>Now that you have an intuition about what these different statistics calculate, follow me to the next video to see how to calculate central tendency.</p>

## L2V4 – Central tendency: mean and mode

Visual	Script		
 <p><b>Sean Barnes</b></p>	<p>You have a few different tools to measure central tendency: mean, median, and mode. When should you use each one? Let's look at how to calculate and interpret these measures, starting with mean and mode.</p>		
 <p><b>Sean Barnes</b></p> <p><b>Mathematical notation</b></p> <table border="1"> <tr> <td style="text-align: center;"> <math>x</math>            Specific value in your sample            Film: Man of Steel   Duration: 143         </td> <td style="text-align: center;"> <math>n</math>            Size of the sample            10 rows   <math>n = 10</math>  <math>n = \# \text{ of rows}</math> </td> </tr> </table> <p><b>Sean Barnes</b></p>	$x$ Specific value in your sample Film: Man of Steel   Duration: 143	$n$ Size of the sample 10 rows   $n = 10$ $n = \# \text{ of rows}$	<p>First, you'll need a little bit of notation. Statisticians use mathematical notation to represent how the same calculation applies to <b>all</b> circumstances.</p> <p>It's common to use <b>[CLICK]</b> lowercase x to represent a <b>[CLICK]</b> specific value you found in your sample, like a <b>[CLICK]</b> specific movie duration. You'll see x used in many upcoming formulas.</p> <p>You'll also see <b>[CLICK]</b> lowercase n used to represent the <b>[CLICK]</b> size of the sample. So for example if you interview <b>[CLICK]</b> 10 people, <b>[CLICK]</b> n equals 10. In a sample dataset that you might have in a <b>[CLICK]</b> spreadsheet, <b>[CLICK]</b> n is the number of rows or observations.</p> <p>With these two terms now defined, we can now build up the formulas used to calculate these measures.</p>
$x$ Specific value in your sample Film: Man of Steel   Duration: 143	$n$ Size of the sample 10 rows   $n = 10$ $n = \# \text{ of rows}$		
 <p><b>Sean Barnes</b></p>	<p>Let's return to the movie duration example. You may already be familiar with how to calculate the mean, but you'll see it formalized using mathematical notation here.</p> <p>For this example, <b>[CLICK]</b> x represents specific movie durations. Note that the number of possible values for x is quite large, since movies can have many different durations, even if you round to the nearest minute. The sample size of the data, <b>[CLICK]</b> n, is 25 since you're working with the top 25 most popular movies from 2013.</p> <p>Let's walk through how to <b>[CLICK]</b> calculate the mean movie duration of your</p>		

sample. The sample mean is written as [CLICK]  $\bar{x}$ .

- The mean is calculated by [CLICK] summing up all of the values in the sample and [CLICK] dividing by the [CLICK] number of values, or the sample size. This calculation is [CLICK] an average of the values, which is why you may hear mean and average sometimes used interchangeably.
- Intuitively what's happening is that you're distributing the total sum of all values equally among the number of values in the sample. So in total all of these movies are [CLICK] 3,035 minutes, and [CLICK] there are 25 of them. If these movies were all the same duration, each one would be about [CLICK] 121 minutes.
- Using the notation you saw earlier, the mean would be the [CLICK] sum of all of the values  $x_i$ , [CLICK] where  $i$  is an index that ranges from 1 (the first data point in the sample) all the way until  $n$  (the last data point in the sample). So [CLICK]  $x_1$  would correspond to the duration of Man of Steel, which was 143 minutes, and [CLICK]  $x_{25}$  would correspond to the duration of A Good Day to Die Hard, which was 98 minutes. Remember, [CLICK]  $n$  is the sample size of the data, in this case 25. Don't be scared of this big Greek letter Sigma, it just [CLICK] means to add all of these  $x_i$  values up. And lastly, [CLICK] divide the sum by the sample size of [CLICK]  $n$ .
- Here's what that looks like all written out. You have [CLICK]  $x_1 + x_2$ , all the way to  $x_{25}$ , since you had 25 movies in your sample. The sum of all of these individual movie durations gives you the [CLICK] 3,035 total minutes that you've already seen. Then, you'll [CLICK] divide by the sample size of [CLICK] 25 to get the final result of [CLICK] 121.4 minutes. [pause for absorption]

#### Mode

The most common value in the sample data

- In the case of the movie lengths in 2013:

91	98	98	98	100	107	109	110	112	113	115	116	119	124	126
130	130	130	131	132	143	143	146	153	161					

DeepLearning.AI

Sean Barnes

Another common measure of centrality is the mode: the most common value in the sample data. [CLICK] In the case of the movie durations in 2013, you actually have two modes: [CLICK] 3 movies were 98 minutes long, and 3 movies were 130 minutes long.

#### Mode

Used most commonly for discrete numerical data or categorical data where the mean and median aren't really applicable

Candidate	% of Votes
A	55%
B	30%
C	15%

- There isn't a way to calculate the mean vote
- The mode is Candidate A who had the **most votes** and won the election

DeepLearning.AI

Sean Barnes

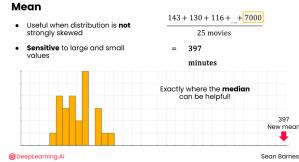
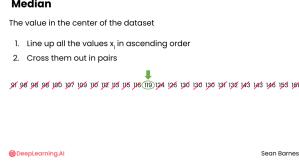
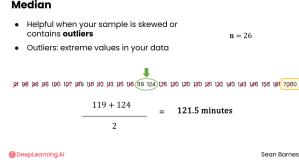
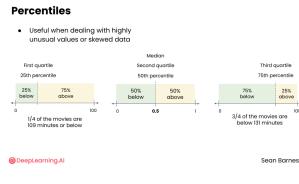
Mode is much more commonly [CLICK] used for discrete numerical data or categorical data where the [CLICK] mean and median aren't really applicable. Voting results are a great example. If there are [CLICK] three candidates, A, B, and C, with 55%, 30%, and 15% of the votes, respectively, [CLICK] there isn't a way to calculate the mean vote. What would you say is the "center of mass" of this data? [pause for thought] In this case, it would be [CLICK] candidate A, because they earned the most votes in the election. So, the [CLICK] mode being Candidate A is quite useful for describing this data.



TH

Now you've seen how to use mean and mode appropriately, but there's one more useful measure of central tendency. Follow me to the next video to calculate the median.

## L2V5 – Central tendency: median

Visual	Script									
 TH	<p>The median has a more complicated formal definition compared with mean and mode. For now, let's get an intuition of how it works.</p>									
 <p><b>Mean</b></p> <ul style="list-style-type: none"> <li>• Useful when distribution is not strongly skewed</li> <li>• Sensitive to large and small values</li> </ul> $\frac{143 + 130 + 116 + \dots + 7000}{25 \text{ movies}} = 397 \text{ minutes}$ <p>Exactly where the median can be helpful</p> <p>New mean ↓</p> <p>Sean Barnes</p>	<p>While the mean is especially [CLICK] useful when your distribution is not strongly skewed, it's quite [CLICK] sensitive to very large and very small values. Say the last movie in your [CLICK] 2013 dataset was actually the longest movie ever made, and it was [CLICK] 7000 minutes long! This very rare case would change the mean to [CLICK] 397 minutes. That's not the most useful way to think about the "middle" of the data. Looking at the [CLICK] histogram again, none of the observations are anywhere [CLICK] near that mark. This scenario is [CLICK] exactly when the <b>median</b> can be helpful!</p>									
 <p><b>Median</b></p> <p>The value in the center of the dataset</p> <ol style="list-style-type: none"> <li>1. Line up all the values <math>x_i</math> in ascending order</li> <li>2. Cross them out in pairs</li> </ol> <p>Sean Barnes</p>	<p>Median is calculated by selecting [CLICK] the value in the center of the dataset. One way to calculate it is to [CLICK] line up all the values <math>x_i</math> from your sample in ascending order. Then, [CLICK] cross them out in pairs. [CLICK] One on the left, [CLICK] one on the right. [CLICK] One on the left, [CLICK] one on the right, and [CLICK] eventually [CLICK] you'll [CLICK] come to the [CLICK] middle. Since there are 25 values in the movie duration dataset, there's one value left and that's [CLICK] 119.</p> <p>Of course computers are able to efficiently find the middle value, so you won't be crossing these out by hand.</p>									
 <p><b>Median</b></p> <ul style="list-style-type: none"> <li>• Helpful when your sample is skewed or contains outliers</li> <li>• Outliers: extreme values in your data</li> </ul> $\frac{119 + 124}{2} = 121.5 \text{ minutes}$ <p>Sean Barnes</p>	<p>The median is [CLICK] helpful when your sample is skewed or contains outliers. [CLICK] Outliers are extreme values in your data.</p> <p>Imagine adding the longest movie ever – the one that's [CLICK] 7000 minutes back into the dataset. Now you have [CLICK] 26 movies, so <math>n</math> is even. [CLICK] Repeat the process to calculate the median, and you can see that 7000 minutes gets thrown out right away. [CLICK] In the end, there are two remaining values: [CLICK] 119 and 124. To find the median, take the average of these two numbers by [CLICK] adding 119 and 124 together, [CLICK] then dividing the result by 2. This gives you a median of [CLICK] 121.5. That's pretty close to the original median of 119. Compare that 2.5 minute difference to the 276-minute difference between the original and new <b>means</b>. The median represents the center of this modified sample data much better.</p>									
 <p><b>Percentiles</b></p> <ul style="list-style-type: none"> <li>• Useful when dealing with highly unusual values or skewed data</li> </ul> <table border="1"> <thead> <tr> <th>First quartile</th> <th>Median</th> <th>Third quartile</th> </tr> </thead> <tbody> <tr> <td>25% below 100th percentile</td> <td>50% below 50th percentile</td> <td>75% below 100th percentile</td> </tr> <tr> <td>1/4 of the movies are 100 minutes or below</td> <td>0</td> <td>3/4 of the movies are below 100 minutes</td> </tr> </tbody> </table> <p>Sean Barnes</p>	First quartile	Median	Third quartile	25% below 100th percentile	50% below 50th percentile	75% below 100th percentile	1/4 of the movies are 100 minutes or below	0	3/4 of the movies are below 100 minutes	<p>The median is closely tied to the notion of percentiles. Percentiles are a powerful way to understand how values in your dataset are distributed. [CLICK] The median you just calculated is actually the 50th percentile - it's the value that's right in the middle, [CLICK] with 50% of the data below it and [CLICK] 50% above it. Some applications, like Google Sheets, have you select a number between [CLICK] 0 and 1 to calculate a percentile. So in that case,</p>
First quartile	Median	Third quartile								
25% below 100th percentile	50% below 50th percentile	75% below 100th percentile								
1/4 of the movies are 100 minutes or below	0	3/4 of the movies are below 100 minutes								

	<p>[CLICK] 0.5 would be the 50th percentile.</p> <p>But you're not limited to just the 50th percentile, the median. You can calculate any percentile from 0 (the minimum) to 100 (the maximum). For example, the [CLICK] 25th percentile, also known as the first quartile, is the value below which [CLICK] 25% of your data falls, in this case 109 minutes. So a [CLICK] quarter of the movies in this sample were 109 minutes or below.</p> <p>[CLICK] The second quartile is the median of 119, as you just saw. 50% above, 50% below.</p> <p>[CLICK] The 75th percentile, or third quartile, is the value below which 75% of your data falls, in this case 131. [CLICK] So three quarters of the movies are below 131 minutes.</p> <p>Percentiles are particularly [CLICK] useful when dealing with highly unusual values or skewed data. Remember the 7000-minute movie? It might dramatically affect the mean, but it has minimal impact on most percentiles. The 90th percentile, for instance, tells us about the longer movies in the sample without being swayed by this extreme value.</p>
TH	<p>Now you've seen how to select and calculate an appropriate measure of central tendency, one of the most crucial descriptive statistics you can use to describe a sample.</p> <p>Follow me to the next video to calculate these measures in a spreadsheet.</p>

## L2V6 – Demo: central tendency

Visual	Script
 TH  <b>Statistics for Data Analytics</b> <hr/> Demo: central tendency	<p>Let's see how you can use a spreadsheet to calculate central tendency in order to answer a business question. If you need a refresher on spreadsheet functions, I encourage you to revisit the Data Analytics Foundations course, where you can learn all the functions you'll see in this demo.</p>
 SC <a href="#"><u>Start here</u></a> <a href="#"><u>Solution</u></a>	<p>Recall that you previously created a histogram for the <b>paid interest</b> feature, essentially how much profit had accrued on the loan. Let's complement the histogram with some descriptive statistics, to help you get an idea of how much money you could be making from the interest on loans!</p>

And don't forget, if you'd like to follow along with the demo, you can find this spreadsheet and the solution in the downloads tab.

- Mean:
  - B2 → =AVERAGE(Data!Q:Q)
- Median:
  - B3 → =MEDIAN(Data!Q:Q)
- Percentile:
  - B4 → =PERCENTILE(Data!Q:Q,0.1).
    - Here you put the percentile you want, a number between 0 and 1.
    - If you want a different percentile you change the last value,
  - B5 → =PERCENTILE(Data!Q:Q,0.5)
    - Same values as the MEDIAN
  - B6 → =PERCENTILE(Data!Q:Q,0.99)
- Mode:
  - B7 → =MODE(Data!Q:Q)
  - This isn't really helpful, as your most common value may not be in the center

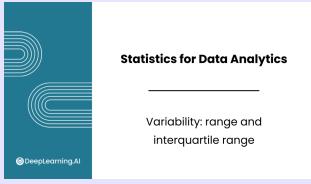
However, say you're interested in finding the center of mass of homeownership, to better understand the kind of collateral that is available to those who are taking out loans.

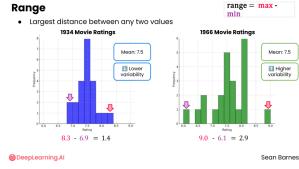
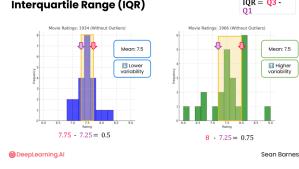
- Look at column D → It turns out you can't really find the median of this feature, what is the average of these?
- The mode function only uses numbers, so you can use homeownership numeric which has 1 for rent, 2 for mortgage and 3 for own.
  - Note that if you take the mean of this, you will get a number, but how would you interpret that? If the average is 1.5 for example, that doesn't correspond to a meaningful category. This operation also implies an order to the data that may not be accurate
- So, now take the MODE of this categorical column which will give you the center of mass
- Mode:
  - C7 → =MODE(Data!E:E)

	<ul style="list-style-type: none"> <li>○ So mortgage is the most common! Helps you get a better picture of where the people taking out loans are in life and the kinds of collateral they might have.</li> </ul>
TH	<p>By calculating central tendency for the paid interest, you're starting to understand the center of mass of the sample data, which gives you a sense of roughly what payout you can expect on average.</p> <p>Once you've completed the practice assessment and the practice lab for this lesson, join me in the next one to learn more about variability and skewness in your data.</p>

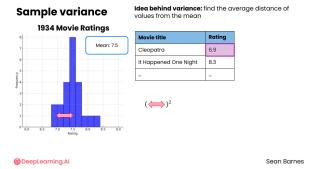
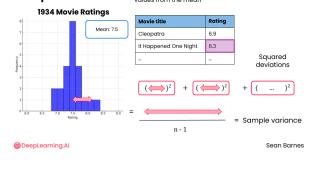
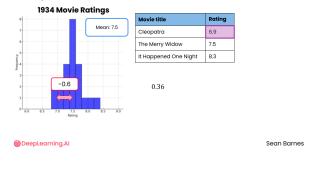
## Lesson 3 – Variability and skewness

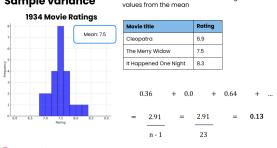
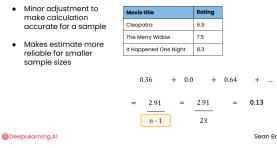
### L3V1 – Variability: range and interquartile range

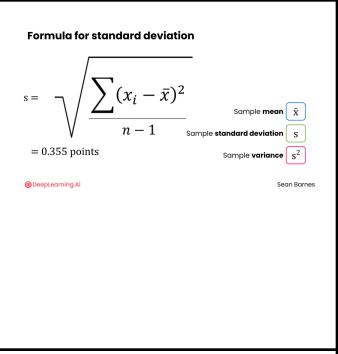
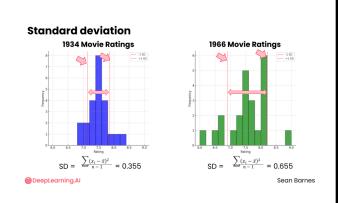
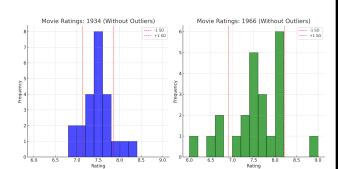
Visual	Script
  <u><a href="#">Source</a></u>	<p>In the 2020 Tokyo Olympics, elite swimmer Margaret Mac Neil won the women's 100 meter butterfly swim event by 0.05 seconds. [brief pause for effect] In fact, all 8 swimmers that competed in the finals achieved times within 1 and a half seconds of each other. That's a very, very small difference.</p> <p>Contrast that with your local high school swim tryouts. You might see the fastest swimmer finish in just 70 seconds, while others take 200 seconds or more. The range is much broader because the skill levels are more varied.</p> <p>This measure of spread across the data is known as <i>variability</i>: a measurement of how tightly or loosely data points cluster around the mean.</p>
 <a href="https://i.imgur.com/dLzfL5v.png">https://i.imgur.com/dLzfL5v.png</a>	<p>Let's start with some intuition. You saw earlier that the movie dataset also includes [CLICK] ratings, so you can compare how movies have been rated on the International Movie Database or IMDb, which is on a 1 to 10 scale.</p> <p>Let's zoom in on two years, [CLICK] 1934 and [CLICK] 1966. Take a look at [CLICK] these two histograms. What can you tell about the central tendency of these two years? [pause for thought] It turns out that on average, movies in 1934 and 1966 both have the [CLICK] same mean rating of 7.5. A quick note, by the way, that one outlier has been removed from each year, so there are only 24 movies here.</p> <p>Even though these two years have the same mean score, the distributions look quite different. What major difference can you spot? [pause for thought] Well, the distribution in 1934 looks [CLICK] tightly clustered around the mean</p>

	<p>of 7.5, while the distribution for 1966 is [CLICK] much more spread out. The values are more different from each other. In other words, 1966 has [CLICK] higher variability, while 1934 has [CLICK] lower variability.</p> <p>Let's look at the different tools you have available for calculating variability.</p>
	<p>First up: range. Range is simply the [CLICK] max minus the min. It answers the question, what is the [CLICK] largest distance between any two values? It's a simple yet useful measure of variability.</p> <p>In 1934, the highest rated movie was an [CLICK] 8.3, with the lowest a [CLICK] 6.9, making the range of 8.3 minus 6.9 a difference of 1.4.</p> <p>Contrast that year with 1966, which had a max rating of [CLICK] 9.0 and a minimum rating of [CLICK] 6.1, making the range [CLICK] 2.9, more than twice as large.</p> <p>Calculating the range is a quick way to see how spread out the ratings are, and whether a year had consistent or inconsistently rated movies. In this case, it looks like 1934 had more consistent ratings compared with 1966.</p>
	<p>A similar measure of variability you can use in conjunction with the median is the interquartile range, or IQR. It's calculated by [CLICK] subtracting the first quartile from the third quartile.</p> <p>Remember the first quartile defines the lowest 25% of the data and the third quartile defines the upper 75% of the data. That difference means that the IQR contains the middle 50% of the data.</p> <p>For 1934, the [CLICK] first quartile is 7.25, the median is 7.5, and the [CLICK] third quartile is 7.75, making the [CLICK] IQR 0.5.</p> <p>For 1966, the [CLICK] first quartile is also 7.25, the median is 7.6, so slightly higher than 1934, and the [CLICK] third quartile is 8, making the [CLICK] IQR 0.75, or 50% wider.</p> <p>So you can see IQR follows a similar pattern to the range: higher for 1966 compared with 1934, reflecting the greater variability in the data.</p>
<span style="color: #00AEEF;">TH</span>	<p>You've seen two common measures of variability: range and interquartile range. In the next video, you will learn about variance and standard deviation, two more useful measures of variability.</p>

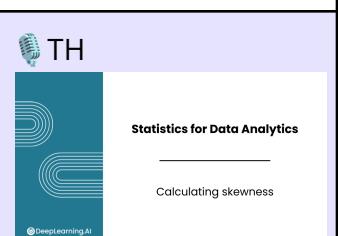
## L3V2 – Variability: variance and standard deviation

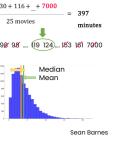
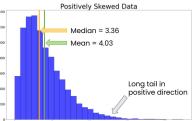
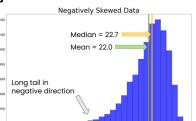
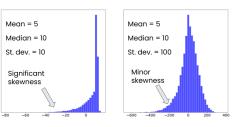
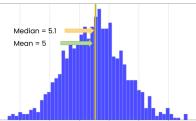
Visual	Script										
  <p>Statistics for Data Analytics</p> <p>Variability: variance and standard deviation</p>	<p>Now let's explore variance. Variance is a more complex calculation than range and IQR, and is the foundation of many other calculations in statistics, including the standard deviation.</p>										
 <p>Sample variance 1934 Movie Ratings</p> <p>Idea behind variance: find the average distance of values from the mean</p> <table border="1"> <thead> <tr> <th>Movie title</th> <th>Rating</th> </tr> </thead> <tbody> <tr> <td>Cleopatra</td> <td>6.9</td> </tr> <tr> <td>If Happened One Night</td> <td>8.3</td> </tr> </tbody> </table> <p><math>(\text{diff})^2</math></p> <p>Sean Barnes</p> <p>See cell G73 for values - <a href="https://docs.google.com/spreadsheets/d/1CVrznZlt1ieQB2Obr4f9RZFymMjp3plf33jPKaFcEuk/edit?gid=1892502851#gid=1892502851">https://docs.google.com/spreadsheets/d/1CVrznZlt1ieQB2Obr4f9RZFymMjp3plf33jPKaFcEuk/edit?gid=1892502851#gid=1892502851</a></p>	Movie title	Rating	Cleopatra	6.9	If Happened One Night	8.3	<p>Remember that the whole point of calculating variability is to quantify how spread out the values are.</p> <p>You can step through the calculation for variance using the [CLICK] movie ratings from 1934. The [CLICK] idea behind variance is roughly to find the average distance of the individual values from the mean. The process for calculating it involves looking at [CLICK] each value in your sample, in this case each movie rating, finding the [CLICK] difference between that value and the mean, [CLICK] squaring that difference, and then</p>				
Movie title	Rating										
Cleopatra	6.9										
If Happened One Night	8.3										
 <p>Sample variance 1934 Movie Ratings</p> <p>Idea behind variance: find the average distance of values from the mean</p> <table border="1"> <thead> <tr> <th>Movie title</th> <th>Rating</th> </tr> </thead> <tbody> <tr> <td>Cleopatra</td> <td>6.9</td> </tr> <tr> <td>If Happened One Night</td> <td>8.3</td> </tr> <tr> <td>The Merry Widow</td> <td>7.5</td> </tr> <tr> <td>If Happened One Night</td> <td>8.3</td> </tr> </tbody> </table> <p><math>(\text{diff})^2 + (\text{diff})^2 + (\text{diff})^2 = \text{Sample variance}</math></p> <p><math>n - 1</math></p> <p>Sean Barnes</p>	Movie title	Rating	Cleopatra	6.9	If Happened One Night	8.3	The Merry Widow	7.5	If Happened One Night	8.3	<p>repeat [CLICK] that [CLICK] for all [CLICK] the values. Then, you [CLICK] sum all of these squared differences and [CLICK] divide by one less than your sample size. We'll discuss the "one less" point in more detail shortly. The squared differences are also sometimes called [CLICK] "squared deviations", with deviation just meaning how different each value is from the mean.</p>
Movie title	Rating										
Cleopatra	6.9										
If Happened One Night	8.3										
The Merry Widow	7.5										
If Happened One Night	8.3										
 <p>Sample variance 1934 Movie Ratings</p> <p>Idea behind variance: find the average distance of values from the mean</p> <table border="1"> <thead> <tr> <th>Movie title</th> <th>Rating</th> </tr> </thead> <tbody> <tr> <td>Cleopatra</td> <td>6.9</td> </tr> <tr> <td>The Merry Widow</td> <td>7.5</td> </tr> <tr> <td>If Happened One Night</td> <td>8.3</td> </tr> </tbody> </table> <p>0.36</p> <p>Sean Barnes</p>	Movie title	Rating	Cleopatra	6.9	The Merry Widow	7.5	If Happened One Night	8.3	<p>In the case of the 1934 movies, you can start with Cleopatra, the lowest rated movie. Start with the movie's rating, [CLICK] 6.9. Subtract the mean of 7.5 from it giving you [CLICK] negative 0.6. Then square that, [CLICK] 0.36. Great, pause there.</p>		
Movie title	Rating										
Cleopatra	6.9										
The Merry Widow	7.5										
If Happened One Night	8.3										
 <p>Sample variance 1934 Movie Ratings</p> <p>Idea behind variance: find the average distance of values from the mean</p> <table border="1"> <thead> <tr> <th>Movie title</th> <th>Rating</th> </tr> </thead> <tbody> <tr> <td>Cleopatra</td> <td>6.9</td> </tr> <tr> <td>The Merry Widow</td> <td>7.5</td> </tr> <tr> <td>If Happened One Night</td> <td>8.3</td> </tr> </tbody> </table> <p>0.36 + 0.0</p> <p>Sean Barnes</p>	Movie title	Rating	Cleopatra	6.9	The Merry Widow	7.5	If Happened One Night	8.3	<p>Now you need to do the same for each other movie. So for example, The Merry Widow with a rating of [CLICK] 7.5, subtract the mean from that and you get [CLICK] 0, squared is still [CLICK] 0.</p>		
Movie title	Rating										
Cleopatra	6.9										
The Merry Widow	7.5										
If Happened One Night	8.3										

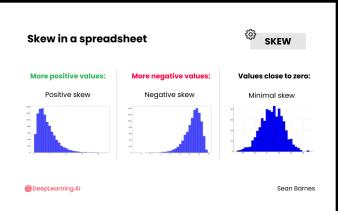
 <p>Sample variance 1934 Movie Ratings</p> <p>Idea behind variance: find the average distance of values from the mean</p> <table border="1"> <thead> <tr> <th>Movie Title</th> <th>Rating</th> </tr> </thead> <tbody> <tr> <td>Cleopatra</td> <td>6.9</td> </tr> <tr> <td>The Merry Widow</td> <td>7.5</td> </tr> <tr> <td>It Happened One Night</td> <td>8.3</td> </tr> </tbody> </table> <ul style="list-style-type: none"> <li>Purposes of squaring             <ul style="list-style-type: none"> <li>Positives and negatives don't cancel out</li> <li>Emphasizes larger deviations</li> </ul> </li> </ul> <p>Sean Barnes</p>	Movie Title	Rating	Cleopatra	6.9	The Merry Widow	7.5	It Happened One Night	8.3	<p>How about [mysterious voice] It Happened One Night, which received a rating of [CLICK] 8.3? Subtract the mean from 8.3 and you get [CLICK] 0.8, squared is [CLICK] 0.64. And you'll do this for [CLICK] each of the values. Here's how it looks.</p> <p>Notice how each of these values is either 0 or positive. That's one [CLICK] purpose of the squaring operation: to make sure that when you add these values all up, you [CLICK] don't have positive and negative values canceling each other out. Another effect of squaring is that it [CLICK] emphasizes larger deviations. So observations that are far away from the mean contribute much more significantly to the overall variance!</p>
Movie Title	Rating								
Cleopatra	6.9								
The Merry Widow	7.5								
It Happened One Night	8.3								
 <p>Sample variance 1934 Movie Ratings</p> <p>Idea behind variance: find the average distance of values from the mean</p> <table border="1"> <thead> <tr> <th>Movie Title</th> <th>Rating</th> </tr> </thead> <tbody> <tr> <td>Cleopatra</td> <td>6.9</td> </tr> <tr> <td>The Merry Widow</td> <td>7.5</td> </tr> <tr> <td>It Happened One Night</td> <td>8.3</td> </tr> </tbody> </table> <p>0.36 + 0.0 + 0.64 + ...      = 2.91      = 2.91 / 23      = 0.13</p> <p>Sean Barnes</p>	Movie Title	Rating	Cleopatra	6.9	The Merry Widow	7.5	It Happened One Night	8.3	<p>Your final step is to divide the [CLICK] total of these squared differences by the [CLICK] sample size <math>n - 1</math>, which in this case is [CLICK] 23. The result you get is [CLICK] 0.13.</p>
Movie Title	Rating								
Cleopatra	6.9								
The Merry Widow	7.5								
It Happened One Night	8.3								
 <p>Why <math>n - 1</math>?</p> <ul style="list-style-type: none"> <li>Minor adjustment to make calculation accurate for a sample</li> <li>Makes estimate more reliable for smaller sample sizes</li> </ul> <p>Idea behind variance: find the average distance of values from the mean</p> <table border="1"> <thead> <tr> <th>Movie Title</th> <th>Rating</th> </tr> </thead> <tbody> <tr> <td>Cleopatra</td> <td>6.9</td> </tr> <tr> <td>The Merry Widow</td> <td>7.5</td> </tr> <tr> <td>It Happened One Night</td> <td>8.3</td> </tr> </tbody> </table> <p>0.36 + 0.0 + 0.64 + ...      = 2.91      = 2.91 / 23      = 0.13</p> <p>Sean Barnes</p>	Movie Title	Rating	Cleopatra	6.9	The Merry Widow	7.5	It Happened One Night	8.3	<p>A quick note here, you might be wondering why you don't divide by just <math>n</math>, and calculate the average sum of squared deviations. That's great intuition! However, <math>n - 1</math> is [CLICK] a minor adjustment that helps make your calculation more accurate when working with a sample rather than the entire population. By using <math>n - 1</math>, you're applying a small correction that [CLICK] makes your estimate more reliable, especially for smaller sample sizes. Don't worry too much about the exact reasons for this adjustment right now; it's a bit of an advanced concept.</p>
Movie Title	Rating								
Cleopatra	6.9								
The Merry Widow	7.5								
It Happened One Night	8.3								
 <p>Sample variance 1934 Movie Ratings</p> <p>1966 Movie Ratings</p> <p>Sean Barnes</p>	<p>Compare the variance you just calculated with the [CLICK] variance for 1966, which is 0.42, almost three times higher. So this calculation is confirming your intuition that the movie ratings in 1966 were much more spread out.</p>								
<p>Formula for sample variance</p> $\sum_{n-1} (x_i - \bar{x})^2$ <p>Not in the same units as your data</p> <p>Sean Barnes</p>	<p>Now that you've calculated it through, let's formalize the definition for sample variance. Each value in your sample is [CLICK] <math>x</math> sub <math>i</math>, and you're subtracting the [CLICK] sample mean, <math>\bar{x}</math> from that. [CLICK] Square this up. Do you remember the notation for summing up all these individual values? [pause for learner to think] That would be [CLICK] sigma, the big E-looking symbol. And finally, [CLICK] divide by <math>n - 1</math>.</p> <p>Now one issue with interpreting variance is that, because you [CLICK] squared the differences between the mean and each value, the variance is [CLICK] not in the same units as your data. Right now, your variance is in the units points-squared. Awkward.</p>								

	<p>However, if you just [CLICK] take the square root, you'll get back to having a [CLICK] calculation in points, in this case 0.355. This square root of the variance is called the [CLICK] standard deviation. [brief pause]</p> <p>The sample standard deviation is written as [CLICK] s. So you already saw [CLICK] x-bar for the sample mean, now you have [CLICK] s for the standard deviation. Sample variance is written [CLICK] s-squared, since it's the standard deviation squared.</p>
  <a href="https://i.imgur.com/hvciIKO.png">https://i.imgur.com/hvciIKO.png</a>	<p>Let's compare with the [CLICK] standard deviation for movie ratings from 1966, which is 0.655. We have the same interpretation as before, that 1966 has a higher standard deviation.</p> <p>Now you can [CLICK] visualize these statistics on the histograms from before. The dashed red lines each indicate the mean minus one standard deviation on the left and the mean plus one standard deviation on the right. The difference on the left is [CLICK] much more narrow, while the difference on the right is [CLICK] wider. It looks like the values in the 1966 histogram are almost pulling those dashed lines outward, which mathematically, is exactly what happens!</p> <p>In practice, you'll see both standard deviation and variance, since variance is used to derive multiple other calculations and standard deviation is commonly used to characterize your distributions.</p>
	<p>Great work on calculating these measures of variability! Join me in the next video to better understand the third key measure of your distribution: skewness. I'll see you there.</p>

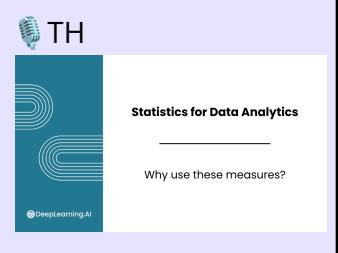
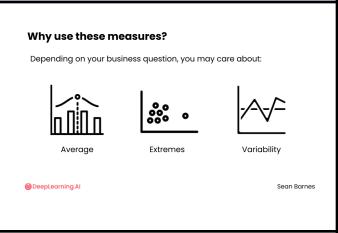
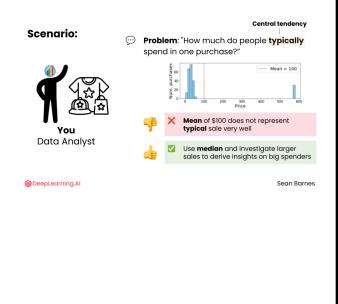
### L3V3 – Skewness

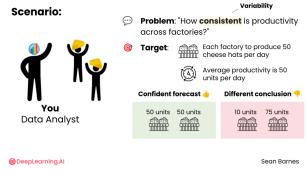
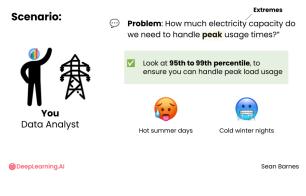
Visual	Script
 <p>Peek formula onto screen in whichever direction Sean sideyes:</p> $b_1 = \frac{m_3}{s^3} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left[ \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{3/2}}$	<p>Unlike mean, median, variance, and standard deviation ... skewness isn't regularly calculated manually. It has a rather complicated formula ... [formula peeks dramatically onto the screen, Sean sideyes it] No! Go away! It's quite a useful concept, though. Let me show you a quick shorthand for estimating it.</p>

<p><b>Skewness</b></p> <ul style="list-style-type: none"> <li>Mean: very sensitive to extreme values</li> <li>Median: relatively insensitive to unusual values</li> <li>One way to estimate skew: compare the mean and median</li> <li>Mean will move further towards the long tail of unusual values compared with the median</li> </ul>  <p>Sean Barnes</p> <p>Show histogram of positive skew (mean &gt; median)</p>	<p>Earlier you saw how the mean is [CLICK] very sensitive to extreme values. That's because the mean distributes the total of all the values over the sample size. You also saw that the median is [CLICK] relatively insensitive to these types of unusual values. So, [CLICK] one way to estimate skew is to compare the mean and median.</p> <p>When the values are asymmetrical, the [CLICK] mean will move much further towards the long tail of unusual values compared with the median.</p>
<p><b>Skewness</b></p>  <p>Positively Skewed Data</p> <p>Median = 3.38 Mean = 4.03</p> <p>Long tail in positive direction</p> <p>Sean Barnes</p>	<p>Here's what that looks like in a histogram:</p> <ul style="list-style-type: none"> <li>If the [CLICK] mean is greater than the [CLICK] median, the data is skewed towards the right. This is also called positive skew. You can remember this because this [CLICK] long tail of values goes up in the positive direction.</li> </ul>
<p><b>Skewness</b></p>  <p>Negatively Skewed Data</p> <p>Median = 22.7 Mean = 22.0</p> <p>Long tail in negative direction</p> <p>Sean Barnes</p>	<p>Meanwhile here's a different distribution. Here, the [CLICK] mean is less than the [CLICK] median, indicating a negative skew. You can see how the mean was pulled more strongly towards the left, towards the [CLICK] tail that goes out negatively from the center of the data.</p>
<p><b>Skewness</b></p>  <p>Significant skewness</p> <p>Minor skewness</p> <p>Mean = 5 Median = 10 St. dev. = 10</p> <p>Mean = 5 Median = 10 St. dev. = 100</p> <p>Sean Barnes</p>	<p>The larger this difference is, the more skewed your data is, because the mean will continue to be pulled further and further away from the median. The magnitude of this difference should be put in context with your standard deviation.</p> <ul style="list-style-type: none"> <li>For example, say you have a [CLICK] mean of 5, a [CLICK] median of 10, and a [CLICK] standard deviation of 10. The large difference between the mean and median compared with the standard deviation suggests [CLICK] significant skewness. [brief pause for learner to read graph]</li> <li>Now take the [CLICK] same mean and median, but the standard deviation is [CLICK] 100. The [CLICK] difference is not as impactful, but skewness is still present. [brief pause for reading]</li> </ul>
<p><b>Skewness</b></p>  <p>Median = 5.1 Mean = 5</p> <p>Sean Barnes</p>	<p>What if your data isn't skewed? Here's a [CLICK] histogram of data that is not skewed. In this case, the mean and median will be [CLICK] roughly equal, since the mean hasn't been influenced significantly by asymmetry in the data.</p>
<p><b>Quick test!</b></p> <ul style="list-style-type: none"> <li>s = 18</li> <li>Median = 10</li> <li>s = 15</li> </ul> <p><b>Hint:</b> Consider the direction the mean is being pulled</p> <p>Skewed positive</p> <p>Skewed negative</p> <p>Not skewed</p> <p>Sean Barnes</p>	<p>Quick test, say I have a [CLICK] mean of 18, a median of 10, and a standard deviation of 15. Would you say this data is [CLICK] skewed positive, [CLICK] skewed negative, or [CLICK] not skewed? [CLICK] Hint: consider the direction that the mean is being pulled. [second pause for learner to think]</p>

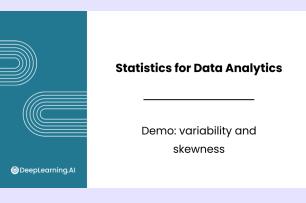
 <p>Skew in a spreadsheet</p> <p>More positive values: Positive skew More negative values: Negative skew Values close to zero: Minimal skew</p> <p>©DeepLearning.AI Sean Barnes</p>	<p>When calculating skew in a spreadsheet, you'll use the [CLICK] <code>SKEW()</code> function, which returns a number, so you'll need to know how to interpret it. [CLICK] More positive values indicate more [CLICK] positive skew, more [CLICK] negative values indicate [CLICK] negative skew, and [CLICK] values close to zero indicate [CLICK] minimal skewness.</p>
 <p>Skew in a spreadsheet</p> <p>Minimal skewness: <math>0 &lt; \text{Skewness} &lt; 0.5</math> Skewness = 0.07</p> <p>Moderate Skewness: <math>0.5 &lt; \text{Skewness} &lt; 1</math> Skewness = 0.78</p> <p>High Skewness: <math>\text{Skewness} &gt; 1</math> Skewness = 1.22</p> <p>©DeepLearning.AI Sean Barnes</p>	<p>Absolute values of skewness [CLICK] less than 0.5 indicate minimal skewness. Skewness between [CLICK] 0.5 and 1 is moderately skewed and skewness [CLICK] greater than 1 is highly skewed.</p>
 <p>TH</p>	<p>Now you're familiar with how to calculate and understand the key measures of central tendency, variability, and skewness. But how can these be used in data analytics? Follow me to the next video to better understand how to select each measure based on the business question you're trying to answer.</p>

## L3V4 – Why use these measures?

Visual	Script
 <p>Statistics for Data Analytics</p> <p>Why use these measures?</p> <p>©DeepLearning.AI</p>	<p>Why are these three measures – of central tendency, variability, and skewness – so important?</p>
 <p>Why use these measures?</p> <p>Depending on your business question, you may care about:</p> <p>Average      Extremes      Variability</p> <p>©DeepLearning.AI Sean Barnes</p>	<p>Depending on your business question, you may care <b>more</b> about <b>different aspects</b> of your distribution. Sometimes you care most about the [CLICK] average person, sometimes about the [CLICK] extremes, sometimes the [CLICK] variability. Let's look at a few examples of business questions and how descriptive statistics can help you answer them.</p>
 <p>Scenario:</p> <p>Problem: "How much do people typically spend in one purchase?"</p> <p>Central tendency</p> <p>You Data Analyst</p> <p>Mean of \$100 does not represent typical sale very well Use median and investigate larger sales to derive insights on big spenders</p> <p>©DeepLearning.AI Sean Barnes</p>	<p>A personalized merch company might wonder, [CLICK] how much do people <b>typically</b> spend in one purchase? In this case, [CLICK] you're interested in central tendency. Basically, where's the middle? You have to be mindful of how you characterize the typical purchase amount. The average could be misleading if you have a skewed distribution. [CLICK] For example, if your mean purchase is \$100, but you have lots of small purchases and a few very large ones, that [CLICK] mean of \$100 doesn't represent your typical sale very</p>

	<p>well. In this case, you might want to [CLICK] use the median, and maybe investigate those large sales further to see if you can derive insights into your big spenders.</p>
 <p>Scenario: You Data Analyst Problem: How consistent is productivity across factories? Target: Each factory to produce 50 cheese hats per day Confident forecast: Average productivity is 50 units per day Different conclusion: 50 units vs 75 units</p> <p>Sean Barnes @DeepLearning.AI</p>	<p>Now say you're working for a [CLICK] cheese hat manufacturer and they want you to figure out, [CLICK] how <b>consistent</b> is productivity across factories? In this case, [CLICK] variability is quite valuable. Say your [CLICK] target is for [CLICK] each factory to produce 50 cheese hats per day, and the [CLICK] average productivity is 50 units per day. It looks good on the surface, but a mean value with no variance means something very different than a mean value with a lot of variance.</p> <p>If [CLICK] each factory produces almost exactly 50 cheese hats on average per day, you can be pretty [CLICK] confident in your production forecasts. But if a factory produces [CLICK] 10 units in one day and [CLICK] 75 on another, you'll come to a [CLICK] different conclusion. It's possible to reach up to 75 cheese hats per day, and you may want to implement the practices used on these high productivity days company-wide.</p>
 <p>Scenario: You Data Analyst Problem: How much electricity capacity do we need to handle peak usage times? Extremes: Look at 95th to 99th percentile to ensure you can handle peak load usage Hot summer days Cold winter nights</p> <p>Sean Barnes @DeepLearning.AI</p>	<p>Lastly, consider an electricity company that may investigate: [CLICK] how much electricity capacity do we need to handle <b>peak</b> usage times? For this question, you're not just interested in the average electricity usage, [CLICK] but also the extremes. You'd want to look at the highest usage times, perhaps the [CLICK] 95th or 99th percentile, to ensure you can handle peak loads without blackouts. This analysis is especially important for utilities planning for [CLICK] hot summer days or [CLICK] cold winter nights when usage spikes.</p>
	<p>Cool! You've seen how to match up these different measures with the different insights they best support. Follow me to the next video to see conduct these analyses in a spreadsheet environment.</p>

## L3V5 – Demo: variability and skewness

Visual	Script
 <p>TH Statistics for Data Analytics Demo: variability and skewness</p> <p>DeepLearning.AI</p>	<p>Now let's see how you can use a spreadsheet to calculate measures of variability and skewness in the Lending Tree data. These measures complement your previous exploration of the central tendency of the sample data.</p>
 <a href="#">Start</a> <a href="#">Solution</a>	<p>Let's take another look at the paid interest feature. You know that the mean value is around \$617 and the median is around \$456. But that doesn't tell you much about the variability in the data. So, while you may be interested in</p>

seeing the average payout, you're probably also interested in understanding the variability too – how common is it to make less than the average or more than average? Is the amount of interest paid spread out over hundreds of dollars or is it tightly clustered around the mean?

Remember that if you'd like to follow along with the demo, you can find this spreadsheet and the solution in the downloads tab.

- You can start with the range, which tells you the difference between the maximum and minimum payments. This will help you understand the difference between the largest and smallest possible payouts
- Range
  - B9 → =MAX(Data!Q:Q)
  - B10 → =MIN(Data!Q:Q)
  - B11 → =B9-B10
- That tells you the maximum possible distance between values but you have more tools available to you. Variance accounts for all the data points, not just these two extremes
- Variance
  - B13 → =VAR(Data!Q:Q)
  - Reduce decimal points to 2
  - This calculation is useful, but it's not in dollars, it's in dollars squared which is weird. Let's convert to standard deviation to have it be more interpretable
- SD
  - B14 → =STDEV(Data!Q:Q)
  - Show that  $STDEV = \sqrt{VAR}$ 
    - B14 → =SQRT(B13)
  - Make standard deviation into dollars
- Skewness
  - Based on the histogram you saw earlier, what would you say is the skew of this feature?
  - You can also look at the mean, median and mode. Do you think that the distribution is symmetrical or skewed? [pause for thought]
    - The mean > median, which indicates Positive Skewness
  - B16 → =SKEW(Data!Q:Q)
  - Reduce to 2 decimal points
  - Skewness is positive sign as expected

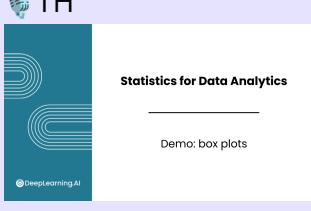
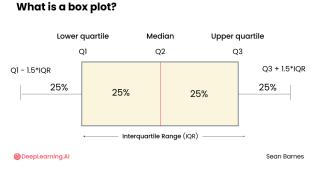
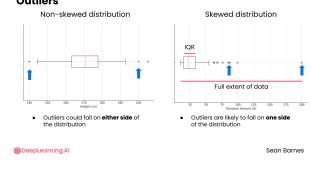
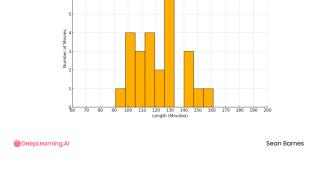


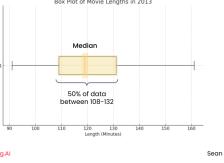
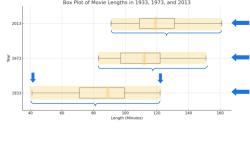
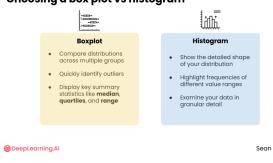
TH

Even just by calculating central tendency, variability, and skewness, you learn a

lot about the potential interest payout. In the next video, you'll wrap up this discussion of distributions by interpreting another common visualization: the box plot.

## L3V6 – Box plots

Visual	Script
	<p>Another great visualization tool for distributions is the box plot, also sometimes called a box and whiskers plot by cat enthusiasts [Sean puts on cat ears]. They're a bit less common than histograms, but quite useful. You should know how to interpret them.</p>
	<p>A box plot visualizes the quartiles of the data, [CLICK] including the median, which you saw previously is the second quartile. [CLICK] The first and [CLICK] third quartiles form the box, and the whiskers typically extend to the [CLICK] smallest and [CLICK] largest values that are within <math>1.5 * \text{IQR}</math> from the first and third quartiles. [CLICK] The box covers the length of the interquartile range.</p>
	<p>Box plots also help visualize outliers. Outliers are any values that lie beyond the whiskers, and are often represented as individual markers.</p> <p>In [CLICK] non-skewed distributions, [CLICK] outliers could fall on either side of the center of the distribution. For example, there are outliers on either side of the distribution for human height – [CLICK] people who are unusually short, and [CLICK] unusually tall.</p> <p>For [CLICK] skewed distributions, however, most of the [CLICK] outliers are likely to fall on one side of the distribution. This is because most of the values in skewed distributions are highly concentrated in one region of the distribution, which causes the [CLICK] interquartile range to be small relative to the [CLICK] full extent of the data. For example, in this distribution of charity donations, most donations are 10 to 20 dollars, but there are a [CLICK] few larger donations. So, most outliers are likely to be large donations, not really small ones.</p>
	<p>As a quick reminder, you saw this histogram of movie durations in the year 2013 in an earlier lesson. You guessed that the mean movie duration was somewhere between 120 and 130.</p> <p>Now that you're familiar with IQR, you can construct a box plot for this same data to compare the two visualizations.</p>

 <p>Box Plot of Movie Lengths in 2013</p> <p>Median</p> <p>50% of data between 108-132</p> <p>Length (Minutes)</p> <p>Sean Barnes</p> <p><a href="https://i.imgur.com/IF2b8uw.png">https://i.imgur.com/IF2b8uw.png</a></p>	<p>Here's a box plot of the same movie durations data. It's easier to tell that the [CLICK] median is a bit below 120 minutes and that [CLICK] 50% of the data is roughly between 108 and 132 minutes.</p> <p>Box plots are especially good for comparing multiple distributions, since you're directly comparing the medians, quartiles, and potential outliers, as well as the variability, whereas a histogram does not enable those comparisons as directly, even when you control for bin size and axis scale.</p>		
 <p>Box Plot of Movie Lengths in 1933, 1973, and 2013</p> <p>Length (Minutes)</p> <p>Sean Barnes</p>	<p>Let's see how box plots enable better comparison for the movie duration data. [CLICK] Here are box plots for movie durations from [CLICK] 1933 on the bottom, [CLICK] 1973 in the middle, and [CLICK] 2013 on the top [CLICK]. What can you tell about the distributions? [pause for learner to think] The median value for 2013 as you have seen is around [CLICK] 120, while for 1973 it's close to [CLICK] 115. That's not a massive difference, and the [CLICK] variability in those two years is similar, considering both the range and IQR.</p> <p>1933 is a bit of a different story. Its median duration is around [CLICK] 90 minutes – what a time! And it seems movies regularly were as [CLICK] short as 40 minutes! The range only barely reaches [CLICK] 120, which was the median length in 2013. It seems between 1933 and 2013 there was a noticeable shift in movie durations, but that in the latter half of this period, the change had pretty much run its course.</p>		
 <p>Choosing a box plot vs histogram</p> <table border="1"> <tr> <td style="background-color: #f2e0aa; padding: 5px;"> <b>Boxplot</b> <ul style="list-style-type: none"> <li>• Compare distributions across multiple groups</li> <li>• Quickly identify outliers</li> <li>• Display key summary statistics like the median, quartiles, and range</li> </ul> </td> <td style="background-color: #d9eaf7; padding: 5px;"> <b>Histogram</b> <ul style="list-style-type: none"> <li>• Show the detailed shape of your distribution</li> <li>• Highlight frequencies of different value ranges</li> <li>• Examine your data in granular detail</li> </ul> </td> </tr> </table> <p>Sean Barnes</p>	<b>Boxplot</b> <ul style="list-style-type: none"> <li>• Compare distributions across multiple groups</li> <li>• Quickly identify outliers</li> <li>• Display key summary statistics like the median, quartiles, and range</li> </ul>	<b>Histogram</b> <ul style="list-style-type: none"> <li>• Show the detailed shape of your distribution</li> <li>• Highlight frequencies of different value ranges</li> <li>• Examine your data in granular detail</li> </ul>	<p>So when should you choose a [CLICK] box plot versus a [CLICK] histogram?</p> <p>Histograms are a great choice when you want to [CLICK] show the detailed shape of your distribution, [CLICK] highlight frequencies of different value ranges in your data, or [CLICK] examine your data in granular detail.</p> <p>Box plots on the other hand are ideal when you want to [CLICK] compare distributions across multiple groups, [CLICK] quickly identify outliers, or [CLICK] display key summary statistics like the median, quartiles, and range at a glance.</p>
<b>Boxplot</b> <ul style="list-style-type: none"> <li>• Compare distributions across multiple groups</li> <li>• Quickly identify outliers</li> <li>• Display key summary statistics like the median, quartiles, and range</li> </ul>	<b>Histogram</b> <ul style="list-style-type: none"> <li>• Show the detailed shape of your distribution</li> <li>• Highlight frequencies of different value ranges</li> <li>• Examine your data in granular detail</li> </ul>		
 <p>TH</p>	<p>Histograms and box plots are both powerful tools for data visualization. Before you complete the upcoming ungraded lab, I hope you'll join me in the next video to see how to use LLMs to help with spreadsheet errors and formulas. I know it will be so helpful to you. I'll see you there!</p>		

## L3V7 – Demo: LLMs for spreadsheet errors & formulas

Visual	Script
 <p>TH</p>	<p>You've encountered a lot of complex formulas and spreadsheet tasks so far. A</p>

 <p><b>Statistics for Data Analytics</b></p> <hr/> <p>Demo: LLMs for spreadsheet formulas &amp; errors</p>	<p>large language model can serve as a helpful thought partner for writing formulas as well as troubleshooting any errors you may encounter.</p>
<p><b><u>Start</u></b></p> <p><b><u>Solution</u></b></p>	<p>All right, so you've seen this data set before. This is the data set of the top 25 movies in each year, their rating on IMDB, the number of ratings they have, and their duration. Let's say you're working with this data in a spreadsheet and you've set up a few calculations that you want to do.</p> <p>Let's say you want to see how many movies are over two hours. You know you want to use COUNTIF, so... count the number of movies if it's greater than 120 minutes. So, you select your E column, and you say greater than 120.</p> <p>But you get this error. Frustrating, because it just says formula parse error. It really doesn't tell you that much. Let's see if we can figure out the problem with an LLM.</p> <p><b>PROMPT:</b></p> <p><i>Hey, I got a formula parse error, can you help me fix it? Here is my formula: =COUNTIF(E:E,&gt;120).</i></p> <p><i>Here is the first few rows of my data:</i></p> <p>I'm just going to demonstrate ChatGPT 4-oh here. I'm going to start a new conversation and I'm going to use this prompt right here: "Hey, I got a formula parse error, can you help me fix it? Here is my formula." I'm just pasting the formula here. And add a few rows of the data.</p> <p><b>PROMPT:</b></p> <p><i>Hey, I got a formula parse error, can you help me fix it? Here is my formula: =COUNTIF(E:E,&gt;120).</i></p> <p><i>Here is the first few rows of my data:</i></p> <p><i>Film Year Rating Number_Ratings Duration_Minutes</i></p> <p><i>M 1931 8.5 71033 10</i></p> <p><i>City Lights 1931 8.6 65709 87</i></p> <p><i>Frankenstein 1931 8 37746 70</i></p>

So I want to give the LLM as much context as possible. Maybe it's not relevant, but maybe it can help me figure out what the problem is.

Okay, it looks like the issue with your formula is related to the syntax. In Google Sheets, you need to enclose a criterion in quotation marks. Okay, so, I forgot about my quotation marks. Now, let's say I'm not 100 percent sure that this is correct. Remember, you should have a mindset of skepticism. So let's just ask,

**PROMPT:**

*Are you sure the greater than sign needs to be within the quotes?*

And just double check. It says yes, I'm sure. just by asking again, we're just giving it a chance to rethink. So, let's go ahead and copy this function right here, go back to our sheet, and paste it in. And there we have 638 movies that are greater than two hours. Okay. Now, say I want to see the percent of movies that are greater than two hours.

If I want to count the total number of movies here, I can scroll down to the bottom and see 2,051 (minus 1 for the header row). That's kind of tedious. Maybe the number of movies will change if I add more years.

So in this case, I'm just going to ask,

**PROMPT:**

*How can I count the number of rows in my dataset with a formula? Give me some options.*

Okay, so it gives me a couple of options here. [READ AND COMMENT ON SUGGESTIONS - FIND AN INCORRECT ONE AND DO IT]

countA for the A column, rows for the A column. Here's another one using countIf, but this is only counting rows where the value in column A is greater than 120, so not the most useful suggestion, that isn't what I asked. And counting rows with multiple criteria, okay, but I don't care about these criteria. So let's go ahead and borrow the first option that it gave us. And I'll put under total movies, COUNTA for this column. Well, 2052, okay, that's clearly not correct. It looks like this is counting the header row. Let's go back to the model and just say,

**PROMPT:**

*I tried the first option, but it isn't working for me. It counts the header row as well.*

Okay, so now it suggests [READ AND COMMENT ON SUGGESTIONS]

This is interesting because I happen to know that the count function is also a good option. So you can just do equals COUNT, this counts the number of numeric values in a data set, and you can just select any numeric column: year, rating, anything like that, and that will give you 2051, which we already checked manually is correct. But these other options work as well. So we could do COUNTA for column A minus one, that works as well. To get percent, we can just divide these two results. And that'll give us 31.11%. Okay, interesting. So about a third of movies are greater than two hours.

Let's say that you want to see – where do these 2 hour movies fall in the distribution of all movies? So first we'll calculate the median of duration. It's about 109 minutes.

So it seems like the 2 hour movies will be above the median. Possibly above the 75th percentile. Let's check. So we'll do percentile, and we'll select the duration column, and I want to know the 75th percentile. But I get an error. Okay. "Error. Function percentile parameter 2 value 75 is out of range"

Okay, it tells me that it's out of range, but it doesn't tell me what the range is. I can head to ChatGPT and ask

**PROMPT:**

*What is causing a parameter 2 value 75 is out of range error. =percentile(E:E, 75)?*

what is causing the range error. And just give it the function here. So the error occurs because the percentile function expects a percentile to be between 0 and 1, and it gives us the corrected formula. Okay, now I remember. Okay, that's helpful. Alright, so 120 is just below the 75th percentile. And we can go ahead and fill in the same for 25th percentile if we're curious.

Okay, one cool trick I want to show you is using LLMs for conditional formatting. Conditional formatting is actually really, really flexible, So, when you go to Format and Conditional Formatting, under Format Rules, one of the options is "Custom Formula Is...".

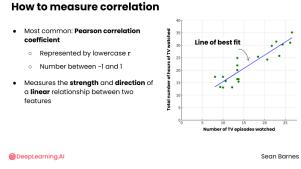
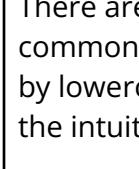
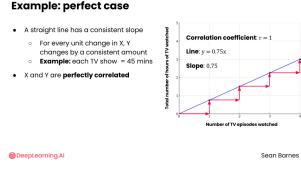
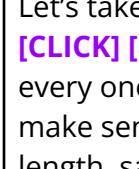
	<p>And so that gives you a lot of flexibility, but it is kind of a pain to remember how to format custom formulas.</p> <p>Let's say I want to highlight each row from column A to column E if and only if the rating of that movie is in the 75th percentile or higher. So let's just go ahead and calculate the 75th percentile.</p> <p>Okay, so 75th percentile rating is 7.9. So let's go over to chat GPT</p> <p>Here is my prompt:</p> <p><b>PROMPT:</b></p> <p><i>How can I use conditional formatting to highlight the entire row from column A to column E green, if the rating is in the 75th percentile or higher. My 75th percentile is in cell H6.</i></p> <p>So, this is important because now I'm giving the LLM the ability to reference that cell in the formula rather than having it write a formula and then fix that formula later myself. So let's see what it says.</p> <p>Okay, so calculate 75th percentile, check. Select the range of cells, for example, select A2, okay. Open conditional formatting. Set the rule using custom formula is. Okay,</p> <p>Select the range. Make sure not to select the header row, because that can cause some issues. And then change to "custom formula is...". And paste in that value that the LLM just gave you. Okay, so right away I see that it's applied. Let's double check. It should only be the rows with a rating greater than or equal to 7.9 that are highlighted. So this 7.9 is highlighted. 8.0 is highlighted here. It looks like it's working really well. So this conditional formatting helps me see at a glance which of these movies are in the 75th percentile or higher based on their duration.</p> <p>So, it's pretty cool. LLMs can help you discover stuff that you didn't even know existed. Great work!</p>
--	--

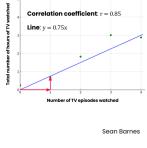
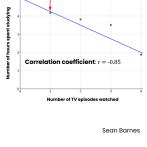
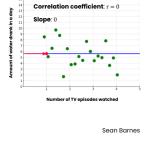
 TH	<p>LLMs can help you discover spreadsheet functionality and enjoy working with them more.</p> <p>You've learned a ton so far in this module. Coming up, you'll take the practice</p>
--	--

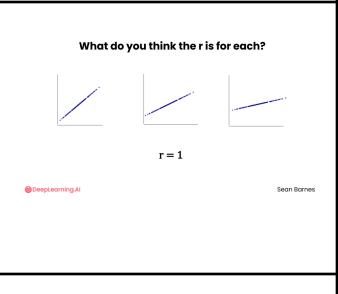
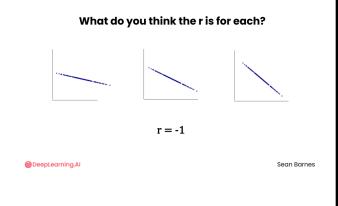
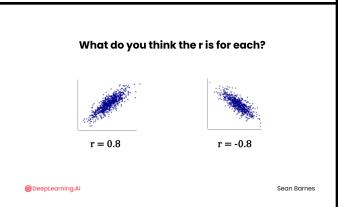
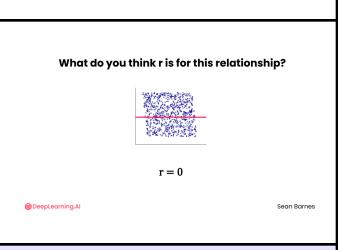
	<p>assessment for this lesson, as well as complete two practice labs. In the first, you'll practice calculating and interpreting variability and skew using the dataset of Spotify songs from the previous lesson. In the second, you'll work with a large language model to explore new spreadsheet functionality and troubleshoot errors.</p> <p>Once you're done, join me in the next lesson to explore correlation.</p>
--	---

## Lesson 4 – Correlation

### L4V1 – Correlation

Visual	Script
  <b>Correlation</b>	<p>Correlation is a way to quantify the relationship between two numerical features. You've seen scatter plots that visually represent how two features are related. Now you'll formalize that relationship.</p>
  <b>Sean Barnes</b>	<p>There are multiple ways of measuring correlation, but [CLICK] the most common is the Pearson correlation coefficient, which is [CLICK] represented by lowercase r. r is always a [CLICK] number between -1 and 1. Let's start with the intuition behind this measure.</p> <p>Say you have a [CLICK] scatterplot of two numerical features. On the x axis you have [CLICK] number of tv episodes watched, and on the y axis the [CLICK] total number of hours of tv watched.</p> <p>You can think of r as measuring how well you can fit a [CLICK] straight line through these points, called a [CLICK] line of best fit. More technically, r [CLICK] measures the strength and direction of a linear relationship between two features.</p>
  <b>Sean Barnes</b>	<p>Let's take the perfect case. Say you have these points for your scatter plot. [CLICK] [CLICK] A straight line has a consistent slope, meaning [CLICK] for every one-unit change in X, Y changes by a consistent amount. That would make sense for a tv show where [CLICK] every episode is the exact same length, say 45 minutes. If you know how many episodes of tv someone watched, you also know exactly how long they've been watching tv for.</p> <p>In this example, you have a line at [CLICK] <math>y = 0.75x</math>, meaning that each episode is three quarters of an hour. So this line has a [CLICK] slope of</p>

	<p>positive 0.75, and [CLICK] for every change of 1 in X, [CLICK] Y changes by 0.75. [next 6 clicks before the next line of script] [CLICK] [CLICK] [CLICK] [CLICK] [CLICK] [CLICK]</p> <p>In this case, the correlation coefficient [CLICK] <math>r</math> equals positive 1. You can say that [CLICK] X and Y are perfectly correlated because changes in X produce exactly predictable changes in Y.</p>
<p><b>Example: less perfect case</b></p> <ul style="list-style-type: none"> <li>Each point is either a little above or a little below the line</li> <li>X goes up, Y goes up</li> <li>Episodes aren't the same length</li> <li>Line of best fit is the same, but correlation isn't as strong</li> </ul>  <p>Sean Barnes</p> <p>@DeepLearning.AI</p>	<p>Now imagine your data is less perfect, with this set of points. It's not as easy to fit a line through the middle, but it is [CLICK] definitely possible to get a good approximation. [CLICK] Each point is either a little above or a little below the line, but as [CLICK] X goes up, [CLICK] Y also goes up, and it's relatively predictable how much it will go up. You can tell that the more tv episodes you watch, the more time you've watched tv, but the [CLICK] episodes aren't all the same length. The correlation coefficient might be close to positive [CLICK] 0.85. In this scatter plot, the line of best fit is the [CLICK] same as the previous one, but the correlation isn't as strong.</p>
<p><b>Example: negative correlation</b></p> <ul style="list-style-type: none"> <li>X goes up, Y goes down</li> <li>Line of best fit has a negative slope</li> <li>Correlation is equally as strong as the previous example</li> <li>Correlation is called "negative"</li> </ul>  <p>Sean Barnes</p> <p>@DeepLearning.AI</p>	<p>Here's another example where as [CLICK] X goes up, Y goes down. In this case, X is the [CLICK] number of tv episodes watched, and y is the [CLICK] number of hours spent studying. So the more tv a person watches, the less time they might have to study.</p> <p>In this case, [CLICK] the line of best fit has a negative slope. However, [CLICK] the correlation is equally as strong as the previous example. The correlation coefficient would be around negative 0.85. Remember, you should think of correlation as essentially about how well you can fit a straight line through the data, regardless of whether the relationship is positive or negative. So, when [CLICK] X goes up, [CLICK] Y goes down, and it goes down a roughly predictable amount. This correlation is as strong as the previous example, [CLICK] it's just called "negative", like the slope of the line of best fit.</p>
<p><b>No correlation</b></p> <ul style="list-style-type: none"> <li>Line of best fit has a slope of 0</li> <li>X goes up, Y doesn't go up or down in any predictable way</li> <li>Line of best fit just predicts the mean value of Y</li> </ul>  <p>Sean Barnes</p> <p>@DeepLearning.AI</p>	<p>Here's what a scatter plot would look like with no correlation, which means the [CLICK] correlation coefficient is equal to zero. This graph represents the [CLICK] number of tv episodes watched and the [CLICK] amount of water each person drank in a day. You can drink water whether you watch tv or not, so there doesn't appear to be any relationship.</p> <p>[CLICK] The line of best fit has a slope of zero, because [CLICK] as X goes up, [CLICK] Y doesn't go up or down in any predictable way. So, [CLICK] the line of best fit just predicts the mean value of Y no matter the X value – that's the best possible guess, and it's not very good.</p>

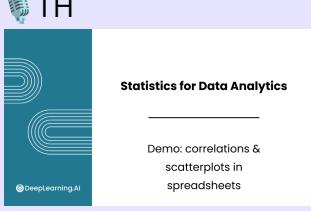
 <p><b>How to interpret r</b></p> <table border="1"> <thead> <tr> <th colspan="2">Strength of the correlation</th> </tr> <tr> <th>Absolute value of r</th> <th>Strength</th> </tr> </thead> <tbody> <tr> <td>0 - 0.3</td> <td>Weak</td> </tr> <tr> <td>0.3 - 0.7</td> <td>Moderate</td> </tr> <tr> <td>0.7 - 1</td> <td>Strong</td> </tr> </tbody> </table> <table border="1"> <thead> <tr> <th colspan="2">Direction of the correlation</th> </tr> <tr> <th colspan="2">Positive: both features increase or decrease together</th> </tr> <tr> <th colspan="2">Negative: one feature increases, the other decrease</th> </tr> </thead> </table> <p>©DeepLearning.AI      Sean Barnes</p>	Strength of the correlation		Absolute value of r	Strength	0 - 0.3	Weak	0.3 - 0.7	Moderate	0.7 - 1	Strong	Direction of the correlation		Positive: both features increase or decrease together		Negative: one feature increases, the other decrease		<p>Calculating r by hand can get complicated, but computers nowadays make it much easier. Let's examine how to interpret r, which tells you two important things about the relationship between two features:</p> <ul style="list-style-type: none"> <li>First, [CLICK] the strength of the correlation. Here's a general guide for interpreting r based on its [CLICK] absolute value, or distance from zero.             <ul style="list-style-type: none"> <li>[CLICK] Positive or negative values between 0 and 0.3 indicate a weak correlation.</li> <li>[CLICK] Between 0.3 and 0.7, positive or negative suggests a moderate correlation.</li> <li>[CLICK] Between 0.7 and 1 positive or negative indicates a strong correlation.</li> </ul> </li> <li>Second, r tells you [CLICK] the direction of the correlation. A [CLICK] positive r means both features tend to increase or decrease together. A [CLICK] negative r means that as one feature increases, the other tends to decrease.</li> </ul>
Strength of the correlation																	
Absolute value of r	Strength																
0 - 0.3	Weak																
0.3 - 0.7	Moderate																
0.7 - 1	Strong																
Direction of the correlation																	
Positive: both features increase or decrease together																	
Negative: one feature increases, the other decrease																	
 <p>What do you think the r is for each?</p> <p>r = 1</p> <p>©DeepLearning.AI      Sean Barnes</p>	<p>Here are three scatterplots. What do you think the r is for each? Here's a hint: see how they all look like a straight line? [pause for learner's answer] [CLICK] These each have an r of 1. You can perfectly fit a line through this data, and as one feature goes up, so does the other. The slope of the line is positive, no matter how big or small, and the relationship is highly predictable.</p>																
 <p>What do you think the r is for each?</p> <p>r = -1</p> <p>©DeepLearning.AI      Sean Barnes</p>	<p>What about these three scatter plots, what r do you think they have? [pause for learner's answer] [CLICK] These each have an r of -1. You can fit a line perfectly through them, however as X goes up, Y actually goes down, making the correlation negative.</p>																
 <p>What do you think the r is for each?</p> <p>r = 0.8      r = -0.8</p> <p>©DeepLearning.AI      Sean Barnes</p>	<p>Here are two scatter plots with correlations of 0.8 and -0.8. Can you tell which is which? [pause for learner's answer] [CLICK] The one on the left is positive 0.8, [CLICK] while the one on the right is -0.8.</p>																
 <p>What do you think r is for this relationship?</p> <p>r = 0</p> <p>©DeepLearning.AI      Sean Barnes</p>	<p>Okay last one, what's your guess for the correlation of this scatter plot? [pause for learner's answer] [CLICK] It's 0! There's no discernable relationship here, and the [CLICK] straight line here would just be flat.</p>																
 <p>TH</p>	<p>Great work studying correlations! It's exciting to calculate a number that maps onto your intuition about the data. Join me in the next video to learn more about the difference between correlation and causation.</p>																

## L4V2 – Correlation and causation

Visual	Script
TH	<p>Correlation is easy to misinterpret because it's often confused with a related concept – causation.</p>
	<p>Correlation and causation might sound similar, but they refer to two different aspects of the relationship between two features.</p> <p>Take a look at [click] this plot of [click] ice cream sales on the x axis and the [click] number of hours spent sunbathing on the y axis. You can see that as ice cream sales go [click] up, [click] so do hours spent sunbathing. The r value might be around [click] positive 0.8. Based on that information, can you conclude that buying ice cream causes people to sunbathe more? [pause for thought] [click] No. [click] Although these two occurrences are correlated, one does not cause the other.</p> <p>Have a look at [click] this very similar plot of number of [click] sunny hours in the day and number of [click] hours spent sunbathing. You can see that the plot is quite similar to the one you saw a moment ago. The r value might again be around [click] positive 0.8. But is there causation? [pause for thought] [click] Yes there is, because the availability of sunlight directly affects people's ability to sunbathe. You would be less likely to sunbathe on a cloudy day. Causation means that [click] one event is the result of another; [click] there's a cause-and-effect relationship.</p>
	<p>To sum up the difference, [click] correlation indicates [click] a relationship or association between two features, but [click] does not imply that one feature causes the other. It [click] is positive when the features move in the same direction or [click] negative when the features move in opposite directions.</p> <p>[click] Causation, on the other hand, implies that [click] one feature directly affects the other and it [click] establishes a cause-and-effect relationship. You can't establish causation using a scatterplot or the Pearson correlation coefficient r. [click] It requires more rigorous testing than correlation.</p> <p>To establish whether there is a causal relationship between the features, you should note that</p> <ul style="list-style-type: none"> <li>• [click] There might be a third feature influencing both features of interest, giving the illusion that they directly affect each other. In the</li> </ul>

	<p>case of the ice cream and sunbathing example, while these two correlate, they are both caused by sunnier days, a third feature.</p> <ul style="list-style-type: none"> <li>Another potential pitfall is [click] the directionality problem. Even if a causal relationship exists, it might be challenging to determine which feature is the cause and which is the effect without experimental evidence. It's clear that [click] sunnier days [click] cause more [click] sunbathing and [click] not vice versa. But how about if you're interested in the relationship between loneliness and social media usage? Do lonely people use more social media, or does social media make people feel more lonely?</li> </ul>
TH	<p>While correlation can be established using scatterplots and the Pearson correlation coefficient, causation requires a higher level of rigor. Be mindful of how you interpret correlation.</p> <p>Follow me to the next video to see how to run correlation analysis in a spreadsheet!</p>

## L4V3 – Demo: correlation in spreadsheets

Visual	Script
	<p>Let's put correlation into practice on the lending tree Loan Dataset. You can choose any pair of numerical features to see how they correlate!</p>
<a href="#"><u>Start</u></a> <a href="#"><u>Solution</u></a>  <a href="#"><u>Dataset</u></a>	<p>As a reminder, this dataset includes thousands of loans made through the lending tree platform, which allows individuals to lend to other individuals.</p> <ul style="list-style-type: none"> <li>Each row represents a particular loan, and each column includes features about the person who took out the loan, like their job title [A] and annual income [F], as well as information about the loan like the amount [K].</li> </ul> <p>And don't forget, if you'd like to follow along with the demo, you can find this spreadsheet and the solution in the downloads tab.</p> <p>So far, you've taken a look at the distribution of the paid_interest feature, and examined its central tendency, variability, and skewness. This is an interesting feature from a lender's point of view, but it's even more useful with more</p>

context. You're also interested in which factors contribute to higher interest and loan amounts.

Let's consider the correlation between **paid\_interest** and **installment**.

Remember the installment is the monthly payment the applicant should be making.

- How do you expect these features to be correlated? Positive or negative? Strong, moderate or weak? **[pause for thought]**
- Create the scatter plot.
  - Select the columns → Insert Chart
  - Chart Setup → Chart type Scatter
  - Move to new sheet
  - Add x label
    - Customize → Chart and axis title → Horizontal axis title → installment
  - Add y label
    - Customize → Chart and axis title → Vertical axis title → paid\_interest
  - Add the trendline: Customize → Series → Select trendline at the bottom. Make it thicker and change the color so it stands out more.
  - What kind of correlation do you see in the scatterplot? **[pause for thought]**
    - Trendline has positive slope
    - For small values of paid\_interest, you get small values of installment, and you have a lot of these points
    - For larger values of paid\_interest you get bigger values of installment, but with a lot more dispersion, but less density of points.
    - It suggests positive correlation, probably between moderate and strong.
- Now find the actual correlation between features
  - Pick a cell → =CORREL( → Select paid\_interest → CTRL → installment →ENTER
    - Mention you can swap the columns and it is the same, because correlation is "symmetric".
    - You get a correlation of 0.69, which is right on the edge between moderate and strong correlation. Remember moderate is between 0.3 and 0.7, and strong is above 0.7

Let's see some other examples!

Let's take a look at how the **paid\_interest** correlates with the **annual\_income**.

- How do you expect these features to be correlated? Positive or negative? Strong, moderate or weak? **[pause for thought]**
- Create scatter plot with trendline (same as before)
  - Select the columns → create plot → Scatter plot
  - Add axis labels
  - Add trendline (fix color, opacity and width)
- What kind of correlation do you see in the scatterplot? **[pause for thought]**
- Interpret:
  - There is a higher density of points with low paid\_interest and low annual\_income, but the points are pretty dispersed
  - You also get some points pretty far away, that will affect correlation
  - The points very loosely fit the trendline, meaning that there is likely not a strong linear relation between the features.
  - This suggests a weak positive correlation
- Get the actual correlation
- This gives you a correlation of 0.206 which is below 0.3, so it's weak but positive. The more people make in a year, the more interest they pay, but there are a lot of other factors that explain the variation in interest paid. It's difficult to predict the exact value of interest paid based on income.

Let's look now at two new features: **annual\_income** and **debt\_to\_income**.

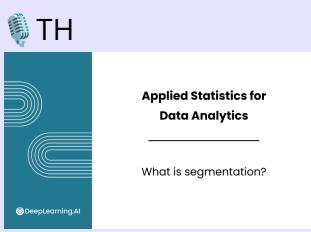
Debt to income is the ratio between how much money you owe on loans and how much money you make in a year. A lower debt to income ratio means a person has proportionally less debt and is more likely to be able to pay off a new loan. This correlation helps answers the question what's the relationship between how much money these borrowers make & the relative amount of debt they have?

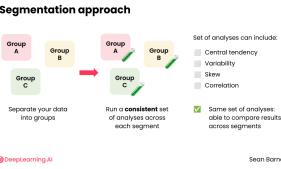
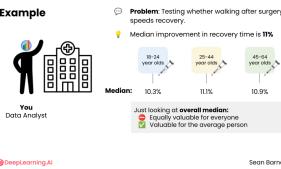
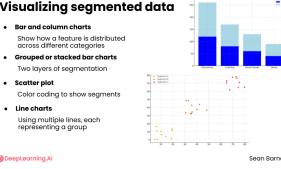
- How do you expect these features to be correlated? Positive or negative? Strong, moderate or weak? **[pause for thought]**
- Create scatter plot
  - Select the columns → create plot → Scatter plot
  - Add axis labels
  - Add trendline (fix color, opacity and width)
- What kind of correlation do you see in the scatterplot? **[pause for thought]**

	<ul style="list-style-type: none"> <li>Interpret: <ul style="list-style-type: none"> <li>Trendline has negative slope</li> <li>You have a lot of points where you get high values of annual_income paired with low debt_to_income, and vice versa, this suggests negative correlation</li> <li>The scatter plot shows a pretty non linear relationship between features. This will impact correlation value</li> </ul> </li> <li>Find the actual correlation <ul style="list-style-type: none"> <li>-0.177</li> <li>This is a weak negative correlation. As income goes up, people taking out loans tend to have a lower debt to income ratio.</li> </ul> </li> </ul>
TH	<p>Great work using scatter plots spot the sign and magnitude of a correlation, as well as interpreting the output of the CORREL function.</p> <p>That takes you to the end of this lesson. Next up, you'll complete the practice assessment as well as the practice lab for this lesson. In the practice lab, you'll explore the correlations between different features of music in order to build an even better playlist.</p>

## Lesson 5 – Segmentation

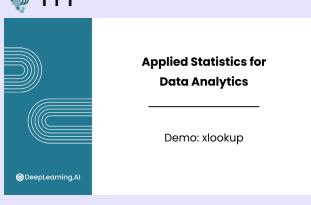
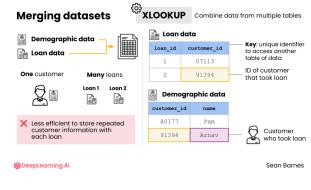
### L5V1 – What is segmentation?

Visual	Script									
 <p>What is segmentation?</p>	<p>Segmentation is a powerful technique that allows you to develop insights for subsets of your data. You saw several examples of segmentation in Data Analytics Foundations, but now let's formalize this concept.</p>									
<p><b>What is segmentation?</b></p> <p>Dividing data into meaningful groups and analyzing each separately</p> <ul style="list-style-type: none"> <li>Analyzing feature or outcome across groups</li> <li>Make different decisions based on characteristics</li> </ul> <table border="1"> <tr> <th>Company</th> <th>Segment by</th> <th>In order to</th> </tr> <tr> <td>Streaming service</td> <td>Viewing habits</td> <td>Tailor content recommendations</td> </tr> <tr> <td>Healthcare provider</td> <td>Risk factors</td> <td>Determine care methods</td> </tr> </table> <p>Sean Barnes</p>	Company	Segment by	In order to	Streaming service	Viewing habits	Tailor content recommendations	Healthcare provider	Risk factors	Determine care methods	<p>At its core, segmentation is about [CLICK] dividing your data into meaningful groups and analyzing each group separately. Generally, it's useful when you're interested in [CLICK] analyzing a particular feature or outcome across different groups. It's also valuable when you might want to [CLICK] make different decisions based on group characteristics. Common segmentations include [CLICK] age group, [CLICK] geographical region, or [CLICK] habits.</p> <p>For instance, [CLICK] a streaming service might segment its users by [CLICK] viewing habits to [CLICK] tailor content recommendations. Or [CLICK] a</p>
Company	Segment by	In order to								
Streaming service	Viewing habits	Tailor content recommendations								
Healthcare provider	Risk factors	Determine care methods								

	<p>healthcare provider might segment patients by [CLICK] risk factors to [CLICK] determine appropriate interventions.</p>
 <span style="font-size: small;">@DeepLearning.AI</span> <span style="float: right;">Sean Barnes</span>	<p>The approach is straightforward: [CLICK] you separate your data into these groups and [CLICK] then run a consistent set of analyses across each segment. [CLICK] These analyses often [CLICK] include descriptive statistics, like measures of [CLICK] central tendency, [CLICK] variability, or [CLICK] skew, and [CLICK] correlation analyses. Using the same set of analyses is important because you want to be [CLICK] able to compare results across your segments.</p>
 <span style="font-size: small;">@DeepLearning.AI</span> <span style="float: right;">Sean Barnes</span>	<p>Let's look more closely at the healthcare example. Say you're [CLICK] testing whether walking after surgery speeds recovery. You know the [CLICK] median improvement in recovery time is 11%. But you might want to compare the effect across age groups. Is walking after surgery more effective for the young, the old? So, [CLICK] you segment your data by age group and [CLICK] measure the median recovery time for each group. Let's say you find that for 18-24 year olds, the median improvement in recovery time was [CLICK] 10.3%, for 25-44 year olds, it was [CLICK] 11.1%, and for 45-64 year olds, it was [CLICK] 10.9%.</p> <p>So walking after surgery is beneficial no matter the age! If you just looked at the overall median, you [CLICK] wouldn't be able to say for certain whether it was equally valuable for everyone, just that [CLICK] it's valuable for the average person.</p> <p>Segmentation is crucial in the medical field. Drugs and other interventions often affect men and women differently and require different dosages for children.</p>
<p><b>Understanding segmentation</b></p> <ul style="list-style-type: none"> <li>• Doesn't have to have a formula</li> <li>• Describes an approach of: <ul style="list-style-type: none"> <li>◦ Breaking down the data</li> <li>◦ Trying to learn how subsets compare to each other</li> </ul> </li> </ul> <span style="font-size: small;">@DeepLearning.AI</span> <span style="float: right;">Sean Barnes</span>	<p>There's no magical "segmentation method." [CLICK] It doesn't have a formula like the Pearson Correlation Coefficient does. [CLICK] It describes an approach of [CLICK] breaking down the data and [CLICK] trying to learn about how different subsets compare to each other.</p>
 <span style="font-size: small;">@DeepLearning.AI</span> <span style="float: right;">Sean Barnes</span>	<p>For visualizing segmented data, you have several options:</p> <ul style="list-style-type: none"> <li>• [CLICK] Bar and column charts are [CLICK] great for showing how a single feature is distributed across different categories.</li> <li>• You can also further segment your data using [CLICK] grouped or stacked bar charts. So, [CLICK] two layers of segmentation!</li> <li>• [CLICK] Scatter plots can be [CLICK] enhanced with color coding to show different segments within the data. You can also show several smaller scatter plots in a grid, one for each segment.</li> <li>• [CLICK] Line charts can be segmented by [CLICK] using multiple lines,</li> </ul>

	<p>each representing a different group.</p> <p>As always, when creating these visualizations, don't forget the principles of good design that you've already learned. Clarity, efficiency, and context are still crucial when displaying segmented data.</p>
TH	Great work with segmentation. Join me in the next video to see a powerful spreadsheet technique for working with data from multiple files, which is common in segmentation tasks.

## L5V2 – Demo: xlookup

Visual	Script
	<p>Often, customer data and product data are stored separately. They may be collected and analyzed in different ways. However, if you want to segment product usage patterns based on customer features, you'll need to merge these datasets. It's a common task in data analytics.</p>
	<p>In the case of the Lending tree Loan dataset you've seen throughout this module, the [CLICK] customer demographic data, like their income, and the [CLICK] loan data, like the interest paid, were likely originally stored in separate files. [CLICK] One customer may have [CLICK] many loans, so it's [CLICK] less efficient to store this repeated customer information together with each loan.</p> <p>Here's what typically happens. The loans dataset has a feature called [CLICK] customer_id, and its value for each loan is the [CLICK] id of the customer that took it out. So, if you take that id and [CLICK] search the dataset of customers, you'll come back with [CLICK] one unique result: the [CLICK] customer who took out the loan. The main advantage of this approach is that you can store the customer and loan data separately, while still having access to both.</p> <p>The feature customer id is called a [CLICK] key, a unique identifier that allows you to access another table of data.</p> <p>Now, in order to perform segmentation analysis, you'll need to [CLICK] merge these two datasets. That way you can answer questions like how does income correlate with interest paid?</p> <p>Spreadsheets have a powerful function called [CLICK] XLOOKUP that allows you to [CLICK] combine data from multiple tables, in this case to create detailed profiles of customers with the most profitable loans. Let's see how it works.</p>

[Start](#)  
 [Solution](#)

- Check out the data – demographics in one sheet, loan info in the other
  - [Use this sheet](#). A **customer\_id** column was added
  - I randomly sorted the customer card info so it wasn't just like, copy over the rows
- Which one helps us match?
- Point out **customer\_id** is in both,
- Do it the manual way first
  - Copy first **customer\_id**, ctrl f, look up in second sheet
  - **Grade** is A
- We need to do this 3,000 times, let's do it programmatically
- Start with one column, grade
  - Add the column name in demographic sheet
  - Use =XLOOKUP(), look at help menu
    - Search key = loan id
    - lookup\_range= where you're looking (loan\_id in second sheet)
    - Result\_range = which column you want returned (grade in the second sheet)
- Can manually drag handle down, let's use arrayformula though
  - =ARRAYFORMULA(XLOOKUP(A2:A,applicants\_credit\_history!A:A, applicants\_credit\_history!M:M))
- Let's do one more column
  - =ARRAYFORMULA(XLOOKUP(\$A2:\$A,applicants\_credit\_history!\$A:\$A, applicants\_credit\_history!J:J))
  - You can see that the result\_range changed



TH Great work!

Follow me to the next and final video of this module, where you'll explore how to apply descriptive statistics to segmented data using pivot tables.

### L5V3 – Demo: pivot tables

Visual	Script
TH Statistics for Data Analytics Demo: pivot tables ©DeepLearning.AI	In addition to charts, you'll typically want to examine the descriptive statistics for your segmented data. Spreadsheets power users employ pivot tables for this task. Let's see how they work.
<a href="#"> Start here</a> <a href="#"> Solution</a>	Different kinds of loans have different risk levels for the lender, which reflects the likelihood that the borrower will not be able to pay the loan. and the interest rate reflects those risk levels. Higher risk loans have higher interest rates, so these loans can be more profitable as long as the borrower is able to

repay. You're interested in investigating the differences in interest rates to understand your potential loans better.

Start segmenting by **grade**. Remember that grade is the classification system that involves assigning a quality score to a loan, where grade A represents the lowest risk and G is the riskiest.

Remember to check the downloads tab for this spreadsheet and the solution.

- Select the entire dataset
- Go to **Insert → Pivot table**.
- A pop up menu appears, select insert to New sheet. Alternatively you could use an existing sheet for this, and select the cell where you want the pivot table to be in.
- On the right you will get the **Pivot table editor**, here is where you set up the pivot table.
  - Rows: Add the **grade** feature (this is the one you use for segmentation)
  - **Values:** Add the feature you want to aggregate
    - Click on **Add under Values**.
    - Select **interest\_rate** .
    - Now choose the summarization method from **SUM** to other options like **AVERAGE**. This should show you the mean interest rate for each card category
  - Observations: As expected, grade A loans pay the smallest interest rate, and it gets consistently higher for riskier grades
- If you get out of the pivot table and need to edit, don't worry, just click on the pencil button on the lower left corner, and you get all the options again.
- **Columns:** If you want to pivot over two features, this is the place to add it, but let's skip this for now.
- Add more statistics
  - Now add another statistic, like standard deviation. Values → Add → Select **interest\_rate** → change the summarization method to STDEV
  - Observations: what happened with the standard deviation of grade G? **[pause for thought]** It probably has just one observation, so the variability is 0!
- You can also add more features, it's the same step.

- Values → Add → Select **total\_credit\_utilized** → change the summarization method to AVERAGE
- You can keep adding as many features and statistics as you like!

Another cool thing is that you can create categories within a feature to do the segmentation. Let's see how

- Select the data table
- Insert → Pivot table
- Select Insert to Existing Sheet → Choose a cell below the previous pivot table (Say A9)
- For Rows select **emp\_length**
- For Values select **paid\_interest** and summarize by **AVERAGE**.
- Observe that there seems to be some difference between the **emp\_length = 1,2,3** segments and the rest. Let's segment with this criteria
  - Right Click anywhere in the **emp\_length** column in the pivot table
  - Select Create pivot group Rule
  - Minimum Value: 4(so it groups 1, 2 and 3)
  - Interval Value: 7 (so you cover all values above 3 up to 10, which is the maximum value)
- You can add MODE, VAR, STDEV or whichever statistics you like!



Excellent work with segmentation. You learned how to assemble two datasets into one, visualize the segments in your data, and calculate descriptive statistics for each one using pivot tables.

Coming up, you'll take the practice assessment for this lesson. You'll also complete the two graded items for this module, including the graded assessment and the graded lab. In the lab, you'll use all the skills you learned in this module to help the Portuguese National Park Service prevent forest fires.

After you've completed these items, you'll move on to the next module: Probability and Simulation. Once you've taken a sample of the population and described the distribution of that sample, you can apply rules of probability and statistics to estimate features of the entire population. I'll see you in the next module to learn more.