

## Module1 Syllabus

### Module1:

#### Chapter1:

**Introduction to Data Science:** What is Data Science? Big Data and Data Science hype – and getting past the hype, Why now? – Datafication, Current landscape of perspectives, Skill sets.

#### Chapter2:

**Needed Statistical Inference:** Populations and samples, Statistical modelling, probability distributions, fitting a model.

#### Question Bank:

- Questions from Previous Question Papers
- Additional Questions
- Multiple Choice Questions with Answers

## Chapter1:

### Introduction to Data Science: What is Data Science?

Authors, Cathy and Rachel were unsure and puzzled about all the excitement around data science and Big Data in the media. They discussed their confusion over breakfast, wondering if this new trend might be important and match their skills. Instead of ignoring it, they decided to learn more. Eventually, Rachel started teaching a data science course at Columbia University, Cathy wrote about it on her blog, and now you're reading a book based on their experiences.

*[ Book Details: Rachel Schutt and Cathy O'Neil Copyright © 2014 Rachel Schutt and Cathy O'Neil. All rights reserved. Printed in the United States of America. Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472. ]*

#### 1. Big Data and Data Science Hype:

Here are five reasons why Big Data and Data Science have caused a lot of excitement:

1. **Confusing Terminology:** People are unsure about what "Big Data" and "data science" really mean. These terms are used in different ways, making them hard to understand and leaving many questions unanswered.
2. **Lack of Recognition for Previous Work:** Researchers in fields like statistics and computer science have been studying data for a long time, but sometimes the media makes it seem like these ideas are brand new. This overlooks the decades of work that laid the groundwork for today's data science.
3. **Exaggerated Hype:** Some people use exaggerated phrases like "Masters of the Universe" to describe data scientists, which makes the field seem more impressive than it really is. This hype can make it harder to see the real value of data science.
4. **Confusion with Statistics:** Statisticians feel like data science is just a rebranding of what they already do. They worry that their work is being misrepresented as something new when it's actually part of their field.
5. **Debates about Science vs. Craft:** Some people question whether data science is really a science or just a skilled craft. This debate adds to the confusion about what data science truly represents in the broader context.

These issues contribute to the hype around Big Data and Data Science, making it challenging to understand their true impact and potential.

## 2. Getting Past the Hype

**Rachel** studied statistics in school and then started working at Google. Initially skeptical about data science hype, Rachel found her job at Google revealed its validity. School laid a foundation, but **coding, data visualization, and domain knowledge were crucial** for her role, highlighting a significant gap between academia and industry.

This made her curious about a new field called **data science**. Data science combines **statistics and computer science** to solve real-world problems using data.

Rachel investigated this field by talking to people at Google, start-ups, and universities.

From those meetings she started to form a clearer picture of the new thing that's emerging. She ultimately decided to continue the investigation by giving a course at Columbia called "Introduction to Data Science," which Cathy covered on her blog.

The author believes data science is a real, new field that brings together different subjects and has a new way of working with data

Now, Rachel and Cathy want to share their knowledge about data science with many more people through this book.

## 3. Why now the Data Science has Emerged?

The two key factors that explain the significance of the current time period for the emergence of data science are:

### 1. Availability of massive amounts of data:

- Our online activities like shopping, communicating, reading news, listening to music, searching for information, and expressing opinions are being tracked and generating a lot of data.
- Additionally, there is increasing "datafication" of our offline behaviors as well, mirroring the online data collection.

- Data is being collected across various sectors like finance, healthcare, pharmaceuticals, bioinformatics, social welfare, government, education, retail, and many others.

## **2. Abundance of inexpensive computing power:**

- Along with the availability of massive amounts of data, there is also an abundance of inexpensive computing power to process and analyze this data.

The combination of these two factors – massive data and affordable computing resources – has created an opportunity for data science to thrive. The data itself is becoming the building block for creating data products, such as recommendation systems, credit ratings, trading algorithms, personalized learning, and data-driven policies.

The content highlights that we are witnessing the beginning of a "massive, culturally saturated feedback loop" where our behavior changes the product, and the product changes our behavior, enabled by technology and infrastructure for large-scale data processing, increased memory, and bandwidth.

The author emphasizes the need to think seriously about the ethical and technical responsibilities of conducting this feedback loop, considering its potential impact, which is one of the goals of the book.

## 4. Datafication and its Implications:

Datafication is defined as the process of taking all aspects of life and turning them into data. It involves quantifying and recording various activities and behaviors, both online and offline, for later examination and analysis.

**Examples of datafication include:**

- **"Liking"** something on social media platforms
- Google's augmented reality glasses **recording what people look at**
- Twitter capturing **people's thoughts**
- LinkedIn **capturing professional** networks

**Spectrum of intentionality in datafication:**

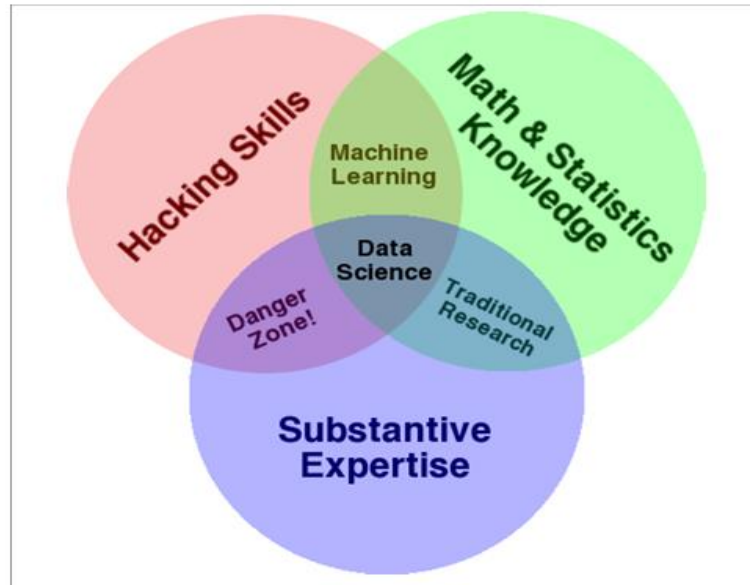
- People intentionally participate in some forms of datafication, like social media interactions.
- People unintentionally or passively get datafied through activities like browsing the web (via cookies) or walking in public spaces (via sensors and cameras).

"Once we datafy things, we can transform their purpose and turn the information into new forms of value."

Here "we" likely refers to modelers and entrepreneurs who monetize the data for purposes like targeted advertising and increased efficiency through automation. It also raises the concern that if "we" is meant to refer to people in general, this perspective of datafication creating value primarily for commercial interests goes against the broader interests of society.

## 5. What is Data Science? The Current Landscape (with a Little History)/Different Perspectives

- What data science is and whether it is distinct from statistics/machine learning.
- The perspective from the Quora discussion and blog posts that view data science as encompassing skills like statistics, data munging (parsing/formatting data), coding, visualization, etc. Suggesting it is a broad set of skills and techniques.
- Driscoll then refers to Drew Conway's Venn diagram of data science from 2010, shown in Venn Diagram below:



- Cosma Shalizi's perspective that data science is just a rebranding and "unwelcome takeover" of the field of statistics. He argues that any good statistics department already does everything described as data science.
- ASA President Nancy Geller's perspective defending the importance of statistics for making sense of data across fields, but her examples don't fully capture data science in the high-tech industry.
- The perspective that "data scientist" emerged as a new job title in 2008 at companies like LinkedIn and Facebook to describe a hybrid skillset combining statistics, computer science, curiosity and persistence for working on data problems.
- The media's role in fueling the hype around data science after the Harvard Business Review called it the "sexiest job of the 21st century".
- William Cleveland's 2001 position paper viewing "data science" as an action plan to expand the field of statistics, suggesting data science existed conceptually before the job title.
- The questions raised about whether data science should be defined by what data scientists do, who has the authority to define the field, and whether it's semantics or a distinct domain.
- The implication that much of data science's development is happening in industry rather than academia.

## 6. The Role of the Social Scientist in Data Science

1. Providing expertise in understanding and analyzing human/user behavior data, which was particularly relevant for early data science roles at social network companies like LinkedIn and Facebook.
2. Contributing the "substantive expertise" component required for data science problems that involve social phenomena, as represented in Drew Conway's data science Venn diagram.[Fig1]
3. Asking good investigative questions and having strong inquiry skills, which are valuable qualities for a data scientist.
4. Bringing a combination of quantitative, programming, and social science skills, which can make them well-suited for data science roles focused on social science-related problems.
5. Potentially being a part of the emerging field of "computational social sciences," which is described as a subset of data science focused on social phenomena.

In essence, the content highlights that social scientists can play a crucial role in data science when the problems involve understanding and analyzing human behavior, social interactions, and other social phenomena. Their expertise in social science concepts and methods, combined with quantitative and programming abilities, can be valuable for certain types of data science applications and domains.

## 7. Data Scientist Jobs [ During 2010 – 2014]

The growing demand and job opportunities for data scientists during 2010 - 2014:

1. Columbia University deciding to start an Institute for Data Sciences and Engineering with Bloomberg's help, indicating recognition of data science as an important field worth investing in.
2. The specific mention of 465 job openings for data scientists in New York City alone at the time the content was written. This number is described as "a lot", highlighting the significant demand for data science roles in just one major city.
3. The statement "So even if data science isn't a real field, it has real jobs" directly acknowledges the existence of real job opportunities with the

"data scientist" title, regardless of debates about the legitimacy of the field itself.

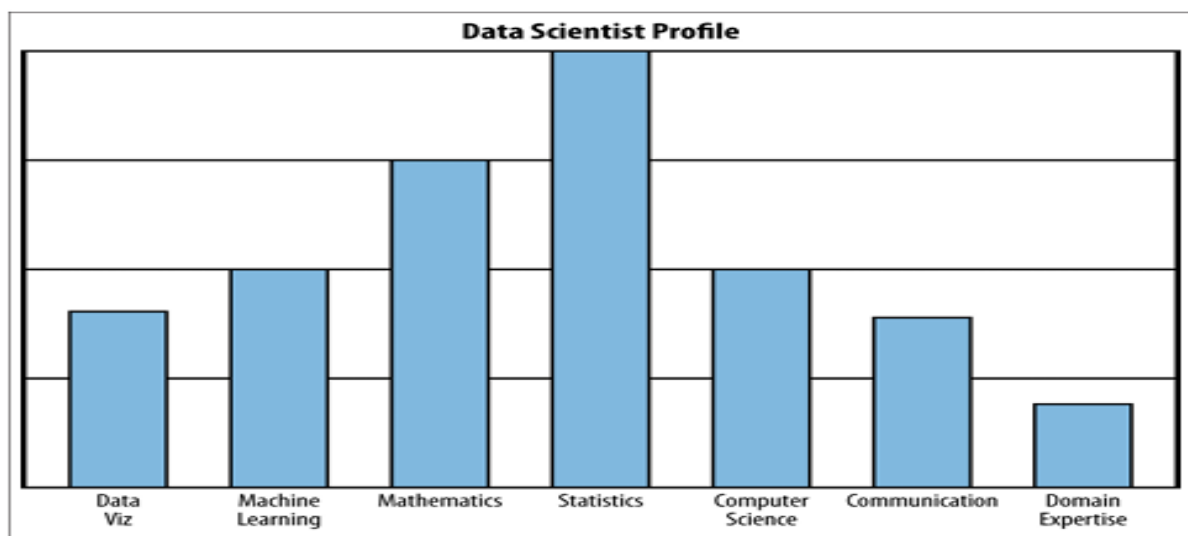
4. The observation that most job descriptions for data scientists require expertise in multiple areas such as computer science, statistics, communication, data visualization, and extensive domain expertise. This implies a high demand for professionals with a diverse and multidisciplinary skillset.

The author does not provide any other specific numbers or details about the job demand beyond the example of 465 openings in New York City. However, author clearly portrays data science as an area with growing recognition and a substantial number of job opportunities, even if its status as a distinct field is still being debated.

## 8. Data Scientist Profile

The key elements that define the Data Scientist Profile are:

1. Computer science skills
2. Mathematical skills
3. Statistical skills
4. Machine learning skills
5. Domain expertise (subject matter knowledge)
6. Communication and presentation skills
7. Data visualization skills



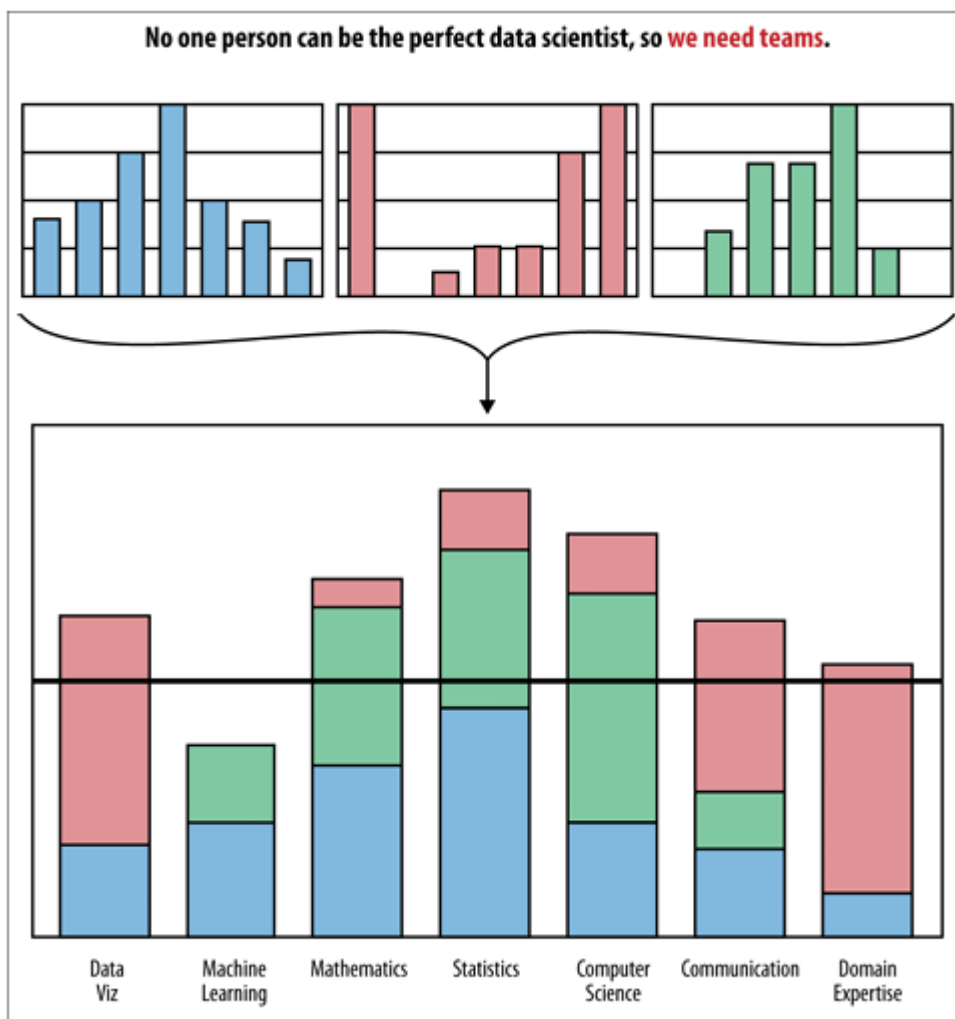
**Fig:** Rachel's data science profile, which she created to illustrate trying to visualize oneself as a data scientist; she wanted students and guest lecturers to "riff" on this—to add buckets or remove skills, use a different scale or visualization method, and think about the drawbacks of self-reporting



This highlights that a well-rounded data scientist is expected to have a combination of technical skills (computer science, math, statistics, machine learning), subject matter expertise, and soft skills (communication, presentation, visualization).

The data scientists can have diverse educational and professional backgrounds, but need to develop a multidisciplinary skillset.

Importantly, the authors raise the idea that rather than defining an idealized "data scientist" who excels at everything, it may be more practical to assemble "data science teams" where individuals with complementary strengths collectively cover all the necessary skills.



**Fig:** Data science team profiles can be constructed from data scientist profiles; there should be alignment between the data science team profile and the profile of the data problems they try to solve

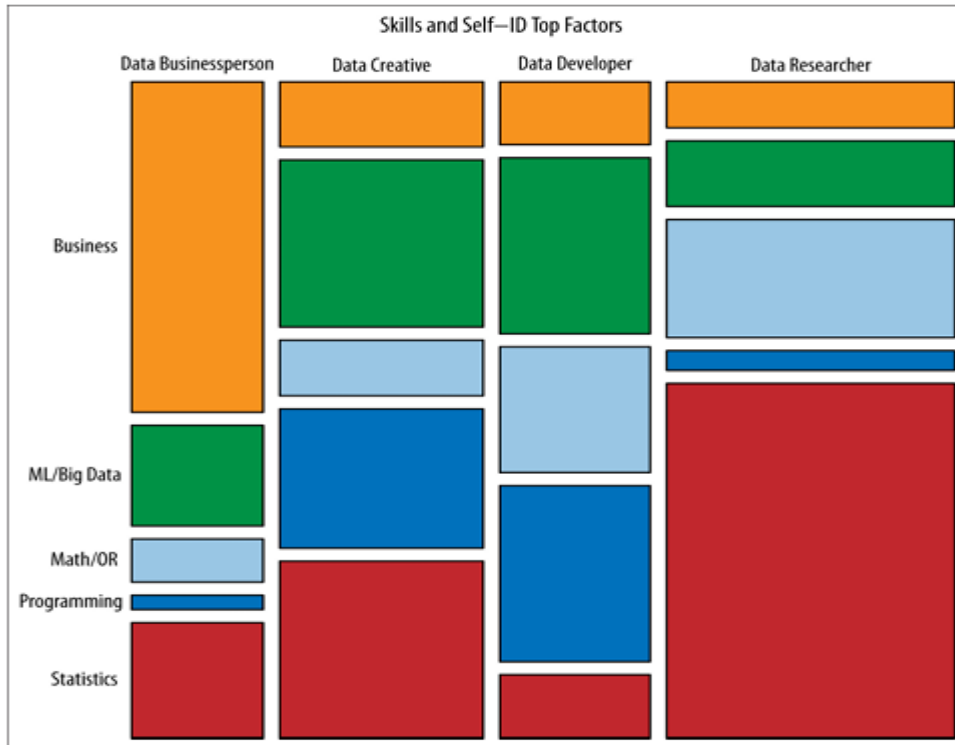
So in summary, the data scientist profile encompasses a broad range of technical, analytical, domain, and communication abilities, but teamwork and diversity of skills within a data science team is also emphasized as a viable approach.

## 9. Thought Experiment:

The "thought experiment" explores the idea of using data science itself to define what data science is - a meta-definition of sorts. Here are the key points regarding this thought experiment:

- The experiment poses the open-ended question: "Can we use data science to define data science?"
- One approach suggested is text mining - performing a Google search for "data science" and applying text mining models on the results. However, this would rely on the masses/public usage to define the term, rather than an authoritative source.
- An alternative is to look at how practitioners of data science describe what they do (e.g. via word clouds) and compare that to how other professionals like statisticians or economists describe their work.
- This data could then be fed into clustering algorithms or other models to see if the model can accurately predict which field a person belongs to based on the "stuff they do."

The content refers to a recent study by Harlan Harris, who used a survey and clustering to define sub-fields within data science (illustrated in Figure ).



In essence, the thought experiment explores using data-driven techniques like text mining, natural language processing, and clustering to analyze descriptions

of what data scientists and other professionals do, in order to derive a definition of data science in a bottom-up manner from the data itself. It proposes applying the methods of data science to tackle the meta-problem of defining the nascent field.

## 10. Who is Data Scientist in the context of academia and industry?

### In Academia:

- Currently, no one directly calls themselves a "data scientist" in academia, except perhaps as a secondary title when affiliated with a data science institute or grant.
- In Academic students interested in becoming data scientists come from diverse backgrounds like statistics, applied math, computer science, social sciences, journalism, biomedical informatics, etc. They are interested in using data to solve important real-world problems.
- It proposes defining an academic data scientist as: ***"A scientist, trained in anything from social science to biology, who works with large amounts of data, and must grapple with computational problems posed by the structure, size, messiness, and the complexity and nature of the data, while simultaneously solving a real-world problem."***
- The rationale is that across disciplines, the computational and data-related challenges have commonalities, so researchers from different departments could collaborate to solve problems in various domains.

### In Industry:

- For the internet/online industry where the term originated, it describes data scientists at different levels of seniority.
- A chief data scientist sets the company's data strategy, including data collection infrastructure, privacy concerns, user-facing data products, data-driven decision making, managing teams of engineers/scientists/analysts, communicating with leadership, patenting solutions, and setting research goals.
- In general, a data scientist extracts meaning from data using statistics, machine learning, and domain knowledge. They spend significant time collecting, cleaning, and wrangling messy data, which requires persistence, statistics, and software engineering skills.

- Key responsibilities include exploratory data analysis, visualization, finding patterns, building models/algorithms to understand product usage, designing experiments, and communicating insights clearly to teams and leadership through data visualizations.
- The data scientist plays a critical role in data-driven decision making and prototyping data products to be integrated into the company's offerings.

In summary, it describes data scientists in academia as interdisciplinary researchers working on computational challenges with large, complex datasets to solve real-world problems across domains. In industry, especially tech, data scientists are depicted as strategic decision-makers and hands-on practitioners who wrangle data, build models, and communicate insights to drive data-driven products and decisions.

tocx10

## Chapter 2:

**Needed Statistical Inference:** What is Big Data? Populations and samples, Statistical modelling, probability distributions, fitting a model.

### 2.1 What is Big Data?

**Big Data** refers to datasets whose size or complexity is beyond the ability of traditional database software tools to **capture, store, manage and analyze**. The different perspectives on what constitutes "**Big Data**":

1. "**Big**" is a relative term, not an absolute threshold like 1 petabyte. Data qualifies as "Big Data" when its size outgrows the current computational capabilities (memory, storage, processing speed) available to handle it effectively.
2. "**Big Data**" refers to datasets that can't fit or be processed on a single machine. Once data exceeds what a single computer can handle, new tools and methods are required to work with it, marking it as "Big Data."
3. **Big Data is also described as a cultural phenomenon** - it highlights how massive amounts of data have become pervasive in our lives due to rapid technology advances.
4. The **4 Vs are / Key characteristics of Big Data include:**
  - **Volume (huge quantities)** : Big Data involves massive amounts of data, ranging from terabytes to petabytes or even more. This massive volume is generated from various sources like social media, sensors, digital transactions, etc.
  - **Variety (diverse formats - structured, unstructured etc.):** Big Data comes in different formats like structured data (databases), semi-structured (XML, JSON) and unstructured (text, audio, video, etc). This heterogeneity of data types makes processing complex.
  - **Velocity (high speed of data generation/processing):** The speed at which Big Data is generated, accessed, processed and analyzed is extremely high and ever-increasing. Handling high velocity data for timely decision making is a challenge.

- **Value (analyzing for extracting valuable insights):** The core reason to deal with Big Data is to uncover hidden insights, patterns and correlations to extract substantial value from the ocean of data.

In essence, "**Big Data**" is a relative concept describing data that is too large, too varied, too rapidly changing or too complex for traditional data processing systems. It requires new age computational tools and represents the data-intensive modern era we live in.

### 2.1.1 Definition of Big Data According to Steve Lohr

Steve Lohr defines Big Data as follows:

1. **"A bundle of technologies"** - This aligns with the technological aspect of Big Data, referring to the various tools, frameworks and platforms like Hadoop, Spark, NoSQL databases etc. that enable storage, processing and analysis of massive, complex datasets.
2. **"A potential revolution in measurement"** - This points to the transformative impact Big Data can have on measurement and metrics across industries/domains by allowing organizations to capture, quantify and analyze more data points than ever before.
3. **"A philosophy about how decisions will/should be made"** - This highlights the paradigm shift Big Data brings in terms of data-driven decision making. With insights extracted from large datasets, decisions can potentially be made based more on hard data/evidence rather than just intuition or experience alone.

So in essence, **Lohr** covers the technological enablers of Big Data, its potential to enable more rigorous measurement/quantification, and the philosophical/cultural shift it could drive in favoring data/analytics-based decision making over traditional approaches. It's a holistic view that captures Big Data's technological, measurement and decision-making facets concisely.

## **2.1.2 Here are some examples of Big Data across different domains:**

### **1. Social Media:**

- Facebook processes over 500 terabytes of data every day from user activities like posts, comments, likes, shares, etc.
- Twitter generates over 12 terabytes of data every day from tweets, retweets, and other user interactions.

### **2. E-commerce:**

- Amazon processes millions of transactions every day and captures vast amounts of data related to customer preferences, browsing histories, and purchase patterns.
- Walmart handles over 1 million customer transactions per hour, generating huge volumes of data related to inventory, sales, and customer behavior.

### **3. Healthcare:**

- Electronic Health Records (EHRs) store massive amounts of patient data, including medical histories, test results, imaging data, and treatment plans.
- Genomic sequencing projects like the Human Genome Project generate petabytes of data from DNA sequencing and analysis.

### **4. Internet of Things (IoT):**

- Smart cities generate massive data streams from sensors monitoring traffic, weather, air quality, and infrastructure.
- Industrial IoT devices in manufacturing plants produce continuous data about machine performance, quality control, and supply chain operations.

### **5. Scientific Research:**

- The Large Hadron Collider at CERN generates around 30 petabytes of data annually from particle collisions and experiments.
- Astronomical projects like the Sloan Digital Sky Survey have captured terabytes of data from telescopic observations of galaxies and celestial objects.

### **6. Telecommunications:**

- Telecom companies handle billions of call detail records, user locations, and network data every day from mobile devices and network infrastructure.

### **7. Finance and Banking:**

- Financial institutions process millions of transactions, trades, and market data every minute, generating massive volumes of data for analysis and compliance.

### **8. Transportation and Logistics:** Fleet management systems track real-time location, performance, and maintenance data from thousands of vehicles, generating terabytes of data.

These are just a few examples that illustrate the diverse sources and massive scales of data that constitute Big Data in various industries and sectors. The ability to capture, store, and analyze these vast and complex datasets is driving valuable insights and innovations across domains.

## 2.2 Statistical Thinking in the Age of Big Data

When developing skills as a data scientist one should have the following foundational skills:

- Statistics,
- Linear algebra, and
- Programming

Data scientists need to develop several interdependent skill sets in parallel, such as

- Data preparation,
- Modeling,
- Coding,
- Visualization, and
- Communication.

Statistics is one of the "**foundational pieces**" that need to be in place when developing skills as a data scientist. Even for readers who are already "**awesome software engineers**" and can code well, should possess the solid grasp of statistical knowledge to become data scientists. In the age of Big Data, classical statistics methods need to be revisited and reimaged in new contexts.

### 2.2.1 Statistical Inference

The world we live in is **complex, random, and uncertain**, but it is also a "data-generating machine" through various processes and activities. Some of the **examples of potential data-generating processes are**: Counting people passing by, collecting email data, or Analyzing DNA samples, etc.

These real-world data-generating processes, should be understood and make sense of them, either for scientific curiosity or to solve problems.

The **randomness** is inherent in the data generating process itself and the **uncertainty is associated** with the data collection methods.

Once data is collected, you cannot simply look at it and understand the underlying process that generated it. You need to simplify and capture the data in a more comprehensible way, using mathematical models or statistical estimators.



*The statistical inference is a discipline/field that deals with understanding and making sense of the complex, random, and uncertain real-world processes by collecting data, developing statistical procedures, methods, models and theorems for extracting meaningful information from the data generated by stochastic (random) processes.*

### 2.2.2 Population and Samples

In statistics, population and sample are two important concepts that are used to gather and analyze data. Here's an explanation of these terms with examples:

#### Population:

A population refers to ***the entire group of individuals, objects, or items*** that a researcher is interested in studying. It is the *complete set of elements that share some common characteristics*, which the researcher wants to draw conclusions about.

**Example:** If you want to study the average height of all adults in a particular city, the population would be all the adults living in that city.

#### Sample:

A sample is a subset of the population that is selected for observation and analysis. It is a smaller, manageable group that represents the characteristics of the larger population.

**Example:** If the population is all adults in a city, a sample could be 500 randomly selected adults from different neighbourhoods within that city.

The main reasons for using a sample instead of studying the entire population are:

1. **Cost and time efficiency:** Studying an entire population can be expensive and time-consuming, especially when the population is large.
2. **Accessibility:** In some cases, it may not be possible or practical to study the entire population due to geographical constraints or other limitations.
3. **Destructive testing:** If the study involves destructive testing, it is not feasible to test the entire population.

There are different types of sampling methods, such as simple random sampling, stratified sampling, cluster sampling, and systematic sampling, each with its own advantages and disadvantages.

The goal of sampling is to select a representative sample from the population, so that the characteristics observed in the sample can be generalized to the larger population with a known degree of accuracy and precision.

**Example:** To estimate the average income of households in a city, a researcher might randomly select **1000 households** from different neighborhoods and use **their income data as a sample** to make inferences about the entire population of households in that city.

In summary, the population is the entire group of interest, while a sample is a smaller, manageable subset selected from the population for the purpose of studying and making inferences about the population as a whole. Here are 5 examples to explain population and samples in statistics:

**Example 1: Studying student performance in a school district**

**Population:** All students enrolled in the school district

**Sample:** A random selection of 500 students from different schools within the district

**Example 2: Estimating customer satisfaction for a retail chain**

**Population:** All customers who have shopped at the retail chain in the last year

**Sample:** A random sample of 2,000 customers from the company's customer database

**Example 3: Measuring public opinion on a political issue**

**Population:** All eligible voters in a particular country

**Sample:** A representative sample of 1,500 voters from different regions, age groups, and demographic backgrounds

**Example 4: Testing the effectiveness of a new medication**

**Population:** All patients diagnosed with a specific medical condition

**Sample:** A random sample of 300 patients who meet the inclusion criteria for the clinical trial

### Example 5: Analysing consumer preferences for a new product

**Population:** All potential customers in the target market

**Sample:** A focus group of 20 individuals representative of the target demographic

In each of these examples, the population represents the entire group of interest, while the sample is a smaller subset selected from the population. The sample is chosen to be representative of the population so that the observations and conclusions drawn from the sample can be generalized to the larger population with a certain degree of accuracy and confidence.

The selection of an appropriate sample size and sampling method is crucial to ensure the validity and reliability of the statistical analysis and findings. Additionally, proper randomization and representative sampling techniques are employed to minimize bias and increase the generalizability of the results to the target population.

### 2.2.3 Big Data Can Mean Big Assumptions

The major assumptions and pitfalls associated with Big Data analysis are as follows:

#### 1. The claim that with Big Data we have "N=ALL" i.e. all the data:

This is almost never true. We often miss out on important perspectives from groups who don't actively contribute data (e.g. non-voters, marginalized communities). Examples show even seemingly comprehensive data like **exit polls** or **internet surveillance** miss key segments.

Assuming N=ALL risks excluding vital voices and leads to incomplete/biased conclusions.

#### 2. Data is not Objective:

It's wrong to believe data is inherently objective or that "**data speaks for itself.**" Data can **reflect historical biases and discrimination**. Blindly trusting data without **context perpetuates problems**. The example of a hiring algorithm shows how data alone can discriminate against women if underlying causes like workplace bias aren't considered.

### 3. Ignoring causation:

Models focusing only on **correlations in data** while ignoring root causes and reasons behind patterns can entrench existing societal issues rather than solve them.

### 4. User-level modelling (n=1):

- Traditionally a sample of **1 is too** small to make inferences.
- But with Big Data, we can now gather extensive data on a single person's actions/events.
- This allows making observations about that individual based on their comprehensive personal data footprint.

In essence, one should be cautious while making overconfident claims with Big Data. One must account for **missing data sources, underlying reasons behind patterns**, and **provide full context** - not treat data simplistically as automatic truth.

## 2.2.3 Modeling

Models are **simplified representations or abstractions** of reality that humans create to understand the world around them. Whether *architectural blueprints, molecular visualizations, or statistical functions*, models attempt to capture the essence of the underlying phenomena or processes generating the observed data.

The term "**model**" can have different meanings in different contexts, like **data models for databases** versus **statistical/mathematical models**.

Models necessarily involve **removing or abstracting** away **extraneous details** from the **full complexity of reality**. This abstraction process is a key characteristic of models.

However, after **analysing a model**, one must pay attention to the details that were **abstracted** away, as they may have been important and overlooked.

However, one can disagree with **Wired article by Chris Anderson** that with enough data, models are obsolete and just finding correlations is sufficient. It argues models are still essential representations.

**Statistical models** specifically aim to capture the uncertainty and randomness inherent in data-generating processes through mathematical functions expressing the shape and structure of the data.

However, statistical modelers must be cautious about mistakenly excluding key variables, including irrelevant ones, or assuming an unrealistic mathematical structure divorced from reality.

Modelling refers to the process of creating simplified representations or abstractions of reality in an attempt to understand, explain and make predictions about real-world phenomena or processes. In data science specifically, modeling involves creating mathematical or statistical representations of these phenomena using data.

The key aspects of modeling in data science include:

1. **Data Representation:** Capturing data in mathematical/statistical forms like equations, functions or algorithms to analyse it systematically.
2. **Abstraction and Simplification:** Models necessarily abstract away extraneous details from the full complexity of reality to focus on the essential features and relationships relevant to the problem.
3. **Variable Selection:** Identifying the relevant variables/features believed to influence the phenomenon based on domain knowledge.
4. **Assumption Making:** Statistical/machine learning models make assumptions about data distributions, variable relationships and underlying processes.
5. **Model Fitting:** Estimating the model parameters that best explain the observed data.
6. **Evaluation:** Assessing the model's performance, accuracy and generalizability using techniques like cross-validation.
7. **Interpretation:** Well-constructed models provide insights into underlying processes, variable importance and relationships within the data.
8. **Prediction:** Once validated, models can make predictions about new, unseen data or future scenarios.

**Modeling** is fundamentally the human effort to represent and understand the nature of reality through particular lenses like architectural, biological or mathematical. Statistical models specifically aim to capture uncertainty and randomness in data-generating processes.

However, care must be taken as models are inherently abstracted representations - modelers must consider abstracted details that may have been overlooked. The

goal is creating an accurate yet interpretable model balancing simplicity with underlying assumptions and constraints.

### 2.2.4 What is Model?

A **model** is a simplified representation or abstraction of reality that humans create in an attempt to understand the nature of the world around them. It is an **artificial construction** that captures the essence of a phenomenon or process through a particular lens or perspective, such as architectural, biological, or mathematical.

#### Key points about models:

1. Models are human efforts to understand and represent reality by focusing on specific aspects through a chosen **viewpoint or framework**.
2. They involve **removing or abstracting** away **extraneous details** from the full complexity of reality. Only the **most relevant attributes** are included in the model.
3. **Models in different domains serve this purpose** - architects use **blueprints**, molecular biologists use **3D visualizations**, statisticians use **mathematical functions**.
4. The act of abstracting away details is **inherent to models**. However, care must be taken to **not overlook** important details that were abstracted.
5. Models make assumptions and simplifications **about the underlying reality**. For example, a protein model may not account for **quantum mechanics** governing **electron behaviour**.
6. **Statistical models** aim to capture the uncertainty and randomness in data-generating processes using **mathematical functions** representing the data's shape and structure.
7. However, **statistical modelers** can mistakenly exclude key variables, include irrelevant ones, or assume an **unrealistic mathematical structure**.

In essence, a model is a purposeful simplification that retains **only the most essential elements of a process or phenomenon from a particular perspective**, while intentionally disregarding other complexities, in order to facilitate understanding, analysis and reasoning.

### 2.2.5 Statistical Modeling:

Statistical modeling is the process of representing the underlying data-generating process or phenomenon using mathematical or statistical expressions. It involves the following key aspects:

1. **Conceptualizing the Process:** Before diving into data analysis, it is useful to first conceptualize and draw a picture or diagram of what the underlying process might be, including the relationships between different variables and how they influence or cause each other.
2. **Mathematical Representation:** Statistical modelers express these conceptualized relationships mathematically using equations or expressions. These expressions contain parameters (represented by Greek letters like  $\beta$ ) whose values are initially unknown and need to be estimated from the data.
3. **Variable Relationships:** The mathematical expressions aim to capture how the variables relate to each other. For example, if there is a hypothesized linear relationship between two variables  $x$  and  $y$ , it could be expressed as  $y = \beta_0 + \beta_1 x$ , where  $\beta_0$  and  $\beta_1$  are the parameters to be estimated.
4. **Visualizing Data Flow:** Some modelers prefer to first visualize the data flow and relationships using diagrams with arrows, showing how variables affect each other or how processes evolve over time. This provides an abstract picture before translating it into mathematical equations.
5. **Parameter Estimation:** Once the mathematical expressions are defined, the next step is to estimate the values of the unknown parameters (like  $\beta_0$  and  $\beta_1$ ) using the available data, typically through techniques like maximum likelihood estimation or least squares regression.

In summary, statistical modeling involves conceptualizing the underlying process, representing the hypothesized relationships between variables mathematically using expressions with unknown parameters, and then estimating those parameter values from data to obtain a quantitative model that captures the essence of the data-generating phenomenon.

### 2.2.6 Key points: How to build a Model?

- **Model building** is not straightforward, and there are no global standards or obvious starting points. It requires making assumptions about the underlying structure of reality.
- **Exploratory data analysis (EDA)**, which involves plotting and visualizing the data, can help build intuition about the dataset and guide the modeling process.
- **The recommended approach** is to start simple and then gradually increase the complexity of the model. Beginning with the simplest possible model, even if it's wrong, can help understand the data and refine the assumptions.

- **Writing down assumptions** in the form of equations and code forces one to think critically about whether the assumptions make sense and what could be improved.
- **There is a trade-off between simplicity and accuracy in modeling.** Simple models may be easier to interpret and build, and sometimes they can get you most of the way to the desired accuracy.
- **Building an arsenal of potential models and understanding probability distributions** are essential components of the modeling process.

### 2.2.7 Probability Distributions

Probability distributions are the foundation of statistical models. They describe the probability of different outcomes or values of a random variable. Many common probability distributions like **normal, Poisson, Weibull, etc.** were named after scientists who observed recurring patterns in real-world phenomena. For example, the **normal or Gaussian distribution** was named after **Gauss** who noticed human heights followed a **bell-shaped curve**.

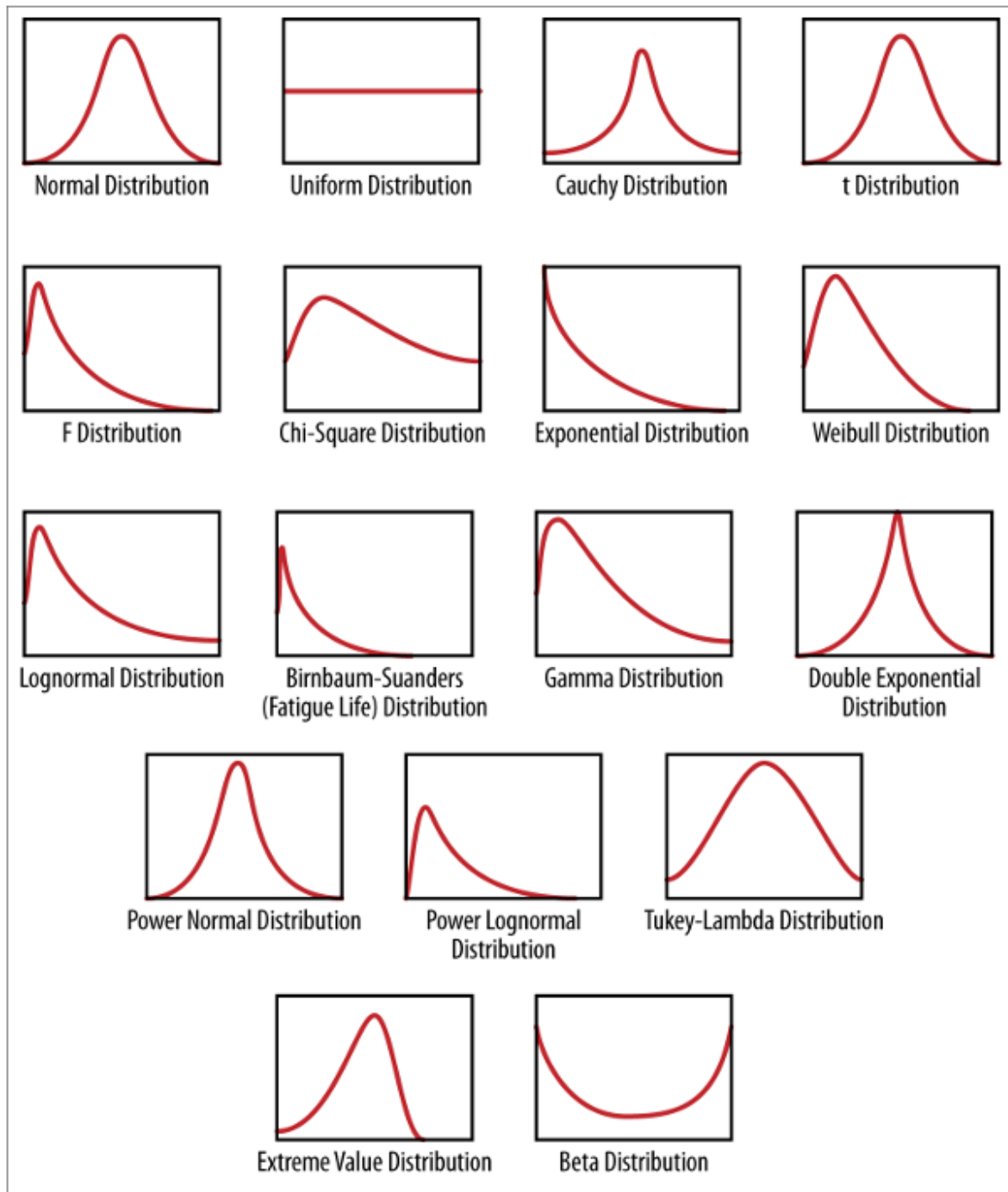
Natural processes tend to generate data that can be approximated by mathematical functions called **probability distributions**, with parameters that can be estimated from the data. For instance, if **X** represents human height, it can be modeled by a normal distribution:

$$N(x|\mu, \sigma) \sim \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Where  $\mu$  is the mean and  $\sigma$  is the standard deviation, estimated from data.

Figure below illustrates a bunch of continuous density functions (aka probability distributions)





Probability distributions assign probabilities to subsets of possible outcomes and have corresponding density functions. The density  $f(x)$  must integrate to 1 over the entire range:

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

For example, let  $X$  be the time until the next bus arrives. If  $X$  follows an exponential distribution with rate  $\lambda$ , the density is:

$$f(x|\lambda) = \lambda e^{-\lambda x}$$

To find the probability the bus arrives between 12 and 13 minutes:

$$P(12 \leq X \leq 13) = \int(12 \text{ to } 13) \lambda e^{-\lambda x} dx \quad [\text{Here } \lambda \text{ is rate/waiting time is } 2]$$

There are also **multivariate joint distributions** for multiple random variables. If X and Y are random variables, the joint density  $f(x,y)$  must integrate to 1 over the X-Y plane.

**Conditional distributions** describe the probability of one variable given another. For X and Y:

$$f(x|y) = f(x,y) / f(y)$$

If X is money spent and Y is items viewed before purchase,  $f(x|y>5)$  is the distribution of money spent given more than 5 items viewed.

When working with a dataset of n rows and k columns/variables, each row represents an observed realization from the k-variable joint distribution.

So probability distributions, univariate, multivariate, conditional, and their parameters are core components of statistical modeling real-world data.

## Definitions and Examples

**1. Probability Distribution:** A probability distribution refers to a mathematical function that describes the likelihood of different outcomes (or values) of a random variable. It specifies how the total probability of 1 is distributed among all possible values that the random variable can take. Probability distributions can be either discrete or continuous.

A probability distribution  $P(X)$  assigns probabilities to different outcomes or values of a random variable X. The general equation involves the probability mass function (for discrete random variables) or probability density function (for continuous random variables) that specifies how probabilities are distributed across the possible values of X.

For discrete X, the probability distribution is described by:

$$P(X = x) = f(x)$$

where  $f(x)$  is the probability mass function. For continuous X, the probability distribution is described by:

**Dr.Thyagaraju G S**

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

where  $f(x)$  is the probability density function.

**Example 1 (Discrete):** Consider a fair six-sided die roll  $X$ . The probability distribution  $P(X)$  can be described as:

$$P(X = x) = \frac{1}{6} \quad \text{for } x = 1, 2, 3, 4, 5, 6$$

Here,  $f(x) = \frac{1}{6}$  for each outcome  $x$ .

**Example 2 (Continuous):** Suppose  $X$  follows a normal distribution  $N(\mu, \sigma^2)$ , where  $\mu$  is the mean and  $\sigma^2$  is the variance. The probability distribution  $P(X)$  is given by its probability density function:

$$P(a \leq X \leq b) = \int_a^b \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

$$\text{Here, } f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

2. **Univariate Probability distribution:** A univariate probability distribution describes the probabilities associated with a single random variable. It deals with the distribution of outcomes or values of this single variable, either in a discrete or continuous context.

**a)** Let a single variable  $X$  represent the **heights** of adult males.  $X$  can be modeled by a normal distribution, with parameters  $\mu$  (mean height) and  $\sigma$  (standard deviation of heights) estimated from data.

**b)** Let a single variable  $X$  represent the **waiting times** between customer arrivals at a bank.  $X$  can be modeled by an exponential distribution with rate parameter  $\lambda$ , where  $1/\lambda$  is the average waiting time

3. **Multivariate Probability Distribution:** A multivariate probability distribution describes the probabilities associated with multiple random variables simultaneously. It specifies the joint probabilities of these variables taking certain combinations of values. Multivariate distributions can be used to model dependencies and relationships between different variables.

a) Let  $X$  represent the heights of individuals and  $Y$  represent their weights. The joint distribution of  $X$  and  $Y$  can be modeled by a bivariate normal distribution, with parameters  $\mu_x, \mu_y$  (means for height and weight),  $\sigma_x, \sigma_y$  (standard deviations for height and weight), and  $\rho$  (correlation coefficient between height and weight).

b) Let  $X$  represent the daily stock returns of company A and  $Y$  represent the daily stock returns of company B. The joint distribution of  $X$  and  $Y$  can be modeled by a multivariate normal distribution, with parameters  $\mu_x, \mu_y$  (mean returns for companies A and B),  $\sigma_x, \sigma_y$  (volatilities of returns for companies A and B), and  $\rho$  (correlation between returns of companies A and B).

4. **Conditional Probability Distribution:** conditional probability distribution refers to the probability distribution of one or more random variables given specific values or conditions of other random variables. It describes how the probabilities of certain outcomes are affected or conditioned by the values of other variables.

a) Let  $X$  represent the SAT scores of students. The condition is whether a student took a prep course or not, represented by a binary variable  $C$  (1 for took prep course, 0 for didn't take prep course). The distribution of  $X$  can be modeled by conditional normal distributions, with parameters  $\mu_1, \sigma_1$  (mean and standard deviation of scores for students who took the prep course, i.e.,  $C=1$ ) and  $\mu_0, \sigma_0$  (mean and standard deviation of scores for students who didn't take the prep course, i.e.,  $C=0$ ).

b) Let  $X$  represent the insurance claims of drivers. The condition is whether the driver is over 65 years old or not, represented by a binary variable  $A$  (1 for over 65, 0 for under 65). The distribution of  $X$  can be modeled by conditional distributions (e.g., gamma or log-normal), with parameters  $\theta_1$  (parameters of the distribution for drivers over 65, i.e.,  $A=1$ ) and  $\theta_0$  (parameters of the distribution for drivers under 65, i.e.,  $A=0$ ).

### 2.2.8 Fitting a Model:

Fitting a model involves estimating the parameters of a mathematical model using observed or given data. The goal is to use the data to approximate the underlying real-world mathematical process that generated the data. This step is crucial in statistical analysis and modeling to understand relationships and make predictions. Fitting the model often involves **optimization methods** and algorithms, such as maximum likelihood estimation, to help estimate the best parameter values.

**Estimating Parameters:** When fitting a model, you estimate parameters (like coefficients in a regression model) using optimization methods such as maximum likelihood estimation. These parameters are statistical estimators that represent functions of the observed data. The aim is to find the best-fitting model that explains the relationship **between variables based on the assumption of a certain mathematical form (e.g., linear, exponential)**.

**Expressing the Model:** Once the model is fitted, it can be expressed in a mathematical form like  $y = 7.2 + 4.5x$ . This equation represents the best estimated relationship between the variables based on the observed data and the assumed model structure.

**Coding the Model:** Implementing the fitted model involves coding the specified mathematical form (e.g.,  $y=7.2+4.5x$ ) into programming languages like **R or Python**. The code reads in the data and utilizes built-in optimization algorithms to compute the most likely parameter values given the data.

**Optimization Methods:** Optimization methods are used to find the best values of model parameters that maximize the likelihood of observing the given data. While these methods are typically built into statistical software like R or Python, understanding the concept of optimization is important for interpreting and validating model results.

**Sophistication and Expertise:** As you become more experienced or specialized in modeling, you may explore and customize optimization methods. Initially, however, it's essential to grasp the underlying optimization process even if you rely on existing software functions to perform the computations.

In summary, fitting a model involves the iterative process of **parameter estimation, model specification, coding, and optimization** to develop a statistical representation that explains the relationship within the observed

data. This process forms the foundation of statistical modeling and inference in data analysis.

### 2.2.9 Overfitting

Overfitting is the term used to mean that you used a dataset to estimate the parameters of your model, but your model isn't that good at capturing reality beyond your sampled data.

**Overfitting** in model development refers to the phenomenon where a model learns not only the underlying patterns in the training data but also captures noise, random fluctuations, or outliers that are specific to the training dataset. This results in a model that performs extremely well on the training data but fails to generalize effectively to new, unseen data.

#### Implications of Overfitting on Generalization:

- When a model is overfitted, it may exhibit poor performance when applied to new data that was not part of the training dataset.
- Overfitted models tend to be too complex, capturing spurious relationships and details that are not reflective of the true underlying patterns in the data.
- Overfitting can lead to misleading conclusions and inaccurate predictions when deployed in real-world scenarios.

#### Strategies to Mitigate Overfitting and Improve Model Generalizability:

##### 1. Simplify the Model:

- Use simpler models with fewer parameters to reduce complexity.
- For example, in linear regression, use fewer predictors or lower-degree polynomial features.

##### 2. Regularization:

- Apply regularization techniques (e.g., L1/L2 regularization) to penalize large coefficients and prevent the model from fitting noise.
- Regularization helps in reducing model complexity and improving generalization.

##### 3. Cross-Validation:

- Use cross-validation techniques to evaluate the model's performance on unseen data.

- Cross-validation helps assess how well the model generalizes to new data by splitting the dataset into multiple training and validation sets.

#### **4. Feature Selection:**

- Carefully select relevant features and avoid using unnecessary or noisy variables that do not contribute to the model's predictive power.
- Feature selection helps in reducing the risk of overfitting by focusing on the most informative attributes.

#### **5. Early Stopping:**

- Monitor the model's performance on a validation dataset during training.
- Stop training when the validation error starts to increase, indicating that the model is beginning to overfit.

### **Real-World Examples Demonstrating Overfitting:**

#### **Example 1: Predictive Modeling in Finance**

An algorithm trained on historical stock market data may overfit by capturing noise or specific market events that are not indicative of future market behavior. The overfitted model may fail to generalize when applied to new market conditions, leading to unreliable investment predictions.

#### **Example 2: Healthcare Diagnostics**

A medical diagnostic model trained on a small dataset with rare disease cases may overfit by memorizing specific patient profiles instead of learning general disease patterns.

The overfitted model may perform poorly when tested on a larger, more diverse dataset of patients, leading to inaccurate diagnosis outcomes.

In summary, overfitting poses a significant challenge in model development by compromising the model's ability to generalize to new data. Employing strategies such as model simplification, regularization, cross-validation, feature selection, and early stopping can help mitigate overfitting and improve the robustness and generalizability of machine learning and statistical models in real-world applications.

## Module 1: Question Bank

### Questions from Previous Year Question Papers

1. What is Data Science? Explain the Venn Diagram of data Science.
2. Draw the Venn diagram of data science and discuss in brief.
3. What is Data Science? Explain in detail about BigData and Data Science Hype. Why Hype?
4. What is Data Science? Explain in detail about BigData and Data Science Hype. How to get past this hype?
5. Explain the Data Science Profile.
6. Explain the work of the Data Scientist in academia and industry.
7. Describe the reasons for BigData and Data Science Hype.
8. "The Data Scientists play a major role in Academics and industry". Justify.
9. Explain the concept of Datafication with example.
10. What is Datafication? Discuss in detail about the current landscape of perspectives and skill sets.
11. What is Datafication? Why it is required in the current scenario?
12. "There are varieties of types of data that a Data Scientist has to deal", explain the statement with examples.
13. Explain the following concepts involved in building models:
  - a. Statistical modeling.
  - b. Probability distributions
  - c. Fitting and overfitting
14. What is Statistical Thinking? Discuss in details statistical inferences and also explain about population and Samples
15. What is Statistical Modelling? Explain in detail about probability distribution and fitting a model. Write about over fitting.
16. Explain the following concepts with examples
  - a. Statistical inference
  - b. Population
  - c. Samples
  - d. Types of Data
17. Explain the Probability Distribution
18. What is Statistical Modelling? How do you build a model?
19. Explain probability distributions with examples.



## Additional Question Bank

### Chapter1:

1. What are some reasons contributing to the excitement and confusion around Big Data and Data Science, as outlined in the book?
2. Describe Rachel's transition from skepticism to interest in data science during her time at Google. How did her experiences shape her perspective on the field?
3. Explain why the current era is significant for the emergence of data science, focusing on the availability of massive data and affordable computing power. How do these factors create opportunities for data science?
4. Define datafication and discuss its implications in modern society. Provide examples of datafication and its role in transforming various aspects of life into data.
5. What are the different perspectives on the definition of data science presented in the book? How do statisticians and practitioners view the relationship between data science and their respective fields?
6. What role do social scientists play in data science, according to the content? How does their expertise contribute to solving data science problems related to social phenomena?
7. During the period from 2010 to 2014, what observations were made regarding the demand for data scientists in New York City? How did this reflect the growing recognition of data science as a field?
8. Outline the key elements that define the profile of a data scientist according to the authors. What technical, analytical, and soft skills are essential for a data scientist?
9. Explain the concept of the "thought experiment" proposed in the book regarding defining data science using data science itself. What methods are suggested to derive a definition of data science from practitioners' descriptions?
10. Compare the role of data scientists in academia and industry based on the book's content. How does their work and expertise differ in these two contexts?
11. Discuss the challenges and benefits associated with the data scientist profile described in the book. How do multidisciplinary skills contribute to the effectiveness of data scientists in various domains?

12. Examine the role of domain expertise in the profile of a data scientist. Why is subject matter knowledge considered essential for data science roles?
13. What are some examples of datafication mentioned in the book, and how do these illustrate the transformation of real-world activities into data?
14. Describe the importance of communication and presentation skills for data scientists. How do these skills complement technical expertise in the role of a data scientist?
15. Explain why computer science skills are considered fundamental for data scientists. How do these skills enable data scientists to work effectively with large datasets and complex computational challenges?
16. How does the content describe the evolution of data science within the academic and industry settings? What are the key differences in how data science is perceived and practiced in these environments?
17. Discuss the concept of datafication and its implications for privacy and society. How does the transformation of everyday activities into data raise ethical concerns?
18. Examine the role of data visualization in data science, as discussed in the book. Why is data visualization considered a critical skill for data scientists?
19. Describe the interdisciplinary nature of data science based on the book's content. How do data scientists collaborate across different domains to solve complex problems?
20. What role does statistical expertise play in the skillset of a data scientist, according to the authors? How do statistical methods contribute to data analysis and modeling in data science?

## Chapter2:

- 1.Explain "Big Data" using the 4 Vs (Volume, Variety, Velocity, Value). Provide examples for each characteristic illustrating datasets that qualify as Big Data.
- 2.Discuss how "Big Data" is a cultural phenomenon impacting modern society, industries, and decision-making processes. Highlight technological advancements and challenges associated with handling massive datasets.
- 3.Why is statistics essential for data scientists dealing with Big Data, regardless of programming proficiency? How are classical statistical methods being adapted to address challenges posed by Big Data? Give examples showcasing statistical thinking in modern data science practices.
- 4.Define statistical inference's relevance in understanding complex data-generating processes. How does it aid in extracting meaningful insights from data? Use specific examples like counting people passing by or analyzing DNA samples.
- 5.Explain population and sample in statistics, emphasizing the distinction and importance of sampling. Provide examples illustrating populations and samples in different scenarios (e.g., estimating city income, studying student performance).
- 6.Discuss the reasons for using sampling methods in statistical analysis. How do sampling methods address challenges like cost, time efficiency, and accessibility? Highlight types of sampling (e.g., random, stratified) and their applications.
- 7.Discuss major assumptions and pitfalls in Big Data analysis, and why "N=ALL" is misleading. Provide examples showing incomplete or biased conclusions due to excluded perspectives. How can researchers mitigate risks with large datasets?
- 8.Explain data objectivity's critical role in Big Data analysis and its challenges (e.g., historical bias). Use the hiring algorithm example to illustrate pitfalls of data reliance. Discuss causation versus correlation importance.
- 9.Define modeling in data science and its aspects aiding in understanding real-world phenomena. Discuss model development steps (e.g., data representation, variable selection). Provide examples of model evaluation and interpretation.

- 10.** Describe the concept of a model's significance across domains (e.g., architecture, biology). Discuss model characteristics, assumptions, and pitfalls. Provide examples showing models' simplified representations and challenges.
- 11.** Explain statistical modeling's process and components, including conceptualizing data relationships. Discuss role of mathematical representation and parameter estimation. Use examples like linear regression for illustration.
- 12.** Discuss the importance of parameter estimation in statistical modeling and common techniques (e.g., maximum likelihood). Explain how estimation refines relationship representations in models. Provide insights into data-driven model development.
- 13.** Explain model building challenges, the role of assumptions, and exploratory data analysis (EDA). Discuss the approach of starting with a simple model and increasing complexity iteratively. Highlight trade-offs and assumptions' critical evaluation.
- 14.** Discuss writing assumptions in equations/code during model building and their impact on critical thinking. Explain simplicity versus accuracy trade-offs and building an arsenal of models. Use examples to illustrate effective model building.
- 15.** Describe probability distributions in statistical modeling and their role in complex models. Discuss how distributions are used to model real-world phenomena and parameter estimation techniques (e.g., maximum likelihood).
- 16.** Define univariate probability distribution, its components, and properties. Explain how it models single random variables using examples.
- 17.** What is a conditional probability distribution? Provide a scenario demonstrating its use in modeling relationships between dependent random variables. Discuss importance in statistical modeling and applications.
- 18.** Distinguish between discrete and continuous probability distributions, with examples. Explain how probabilities are assigned differently and implications for modeling real-world phenomena.
- 19.** Explain multivariate probability distribution using examples. Discuss its role in modeling joint behavior of multiple random variables and capturing dependencies or correlations.

- 20.** Describe the process of fitting a model, including parameter estimation and optimization methods. Explain how coding facilitates model implementation and understanding of optimization. Provide examples illustrating model fitting in statistical analysis.
- 21.** Define overfitting in model development and its implications on generalization to new data. Discuss strategies to mitigate overfitting and improve model generalizability. Provide real-world examples demonstrating overfitting's impact on data analysis.

tocx10

## Multiple Choice Questions

### Chapter 1: Introduction to Data Science

1. What is one of the reasons for the hype around Big Data and Data Science mentioned in the text?
  - a. Exaggerated claims about the importance of data scientists
  - b. Lack of recognition for previous work in related fields
  - c. Confusion with existing fields like statistics
  - d. **All of the above**
2. Which of the following is NOT one of the key factors that explains the significance of the current time period for the emergence of data science?
  - a. Availability of massive amounts of data
  - b. Abundance of inexpensive computing power
  - c. **Advancements in quantum computing**
  - d. Datafication of offline behaviors
3. What does "datafication" refer to?
  - a. The process of converting data into different formats
  - b. **The process of quantifying and recording various activities and behaviors**
  - c. The process of cleaning and preprocessing data
  - d. The process of visualizing data
4. According to the Venn diagram by Drew Conway, what are the three overlapping skill sets that define data science?
  - a. Statistics, Machine Learning, and Visualization
  - b. Hacking Skills, Math, and Business Knowledge
  - c. **Statistics, Hacking Skills, and Domain Expertise**
  - d. Coding, Math, and Visualization
5. What role can social scientists play in data science, according to the text?
  - a. **Contributing expertise in understanding and analyzing human/user behavior data**
  - b. Developing new statistical techniques for data analysis
  - c. Building data storage and processing infrastructure
  - d. Designing user interfaces for data products

6. During the period of 2010-2014, what example is given to illustrate the growing demand for data scientists?
- a. The creation of the Data Science Institute at Columbia University
  - b. The launch of the Data Science degree program at MIT
  - c. The mention of 465 job openings for data scientists in New York City**
  - d. The announcement of the Data Science Certification by Google
7. Which of the following is NOT part of the data scientist profile mentioned in the text?
- a. Computer science skills
  - b. Mathematical skills
  - c. Design skills**
  - d. Machine learning skills
8. What is the "thought experiment" proposed in the text to define data science?
- a. Using surveys and interviews with data scientists
  - b. Analysing job descriptions for data scientist roles
  - c. Using data science techniques like text mining and clustering on descriptions of what data scientists do**
  - d. Conducting a literature review of academic papers on data science
9. In the context of academia, how is an academic data scientist defined in the text?
- a. A scientist who works with large amounts of data and computational problems across various domains**
  - b. A researcher who develops new statistical methods and algorithms
  - c. A faculty member who teaches data science courses
  - d. A computer scientist who builds data storage and processing infrastructure
10. In the context of industry, what is the role of a chief data scientist described in the text?
- a. Setting the company's data strategy and managing teams of data professionals**
  - b. Developing algorithms and models for specific data products
  - c. Collecting and cleaning data from various sources
  - d. Presenting data-driven insights to leadership and stakeholders

## Chapter 2: Needed Statistical Inference

1. Which of the following is NOT considered one of the key characteristics of Big Data?
  - a. Volume
  - b. Variety
  - c. Velocity
  - d. Veracity**
2. According to Steve Lohr's definition, what does "Big Data" refer to as a "potential revolution"?
  - a. A revolution in measurement**
  - b. A revolution in computing power
  - c. A revolution in data storage
  - d. A revolution in data visualization
3. What is the purpose of statistical inference, according to the text?
  - a. To collect and store large amounts of data
  - b. To understand and make sense of complex, random, and uncertain real-world processes**
  - c. To develop new mathematical models
  - d. To visualize data in meaningful ways
4. What is the difference between a population and a sample in statistics?
  - a. A population is a subset of a sample
  - b. A sample is a subset of a population**
  - c. A population and a sample are the same thing
  - d. There is no difference between a population and a sample
5. Which assumption or pitfall related to Big Data analysis is mentioned in the text?
  - a. The assumption that Big Data is always accurate and error-free
  - b. The assumption that Big Data can only be analyzed using specialized software
  - c. The claim that with Big Data we have "N=ALL" or all the data**
  - d. The assumption that Big Data is only useful for large corporations
6. What is the recommended approach for building a model, according to the text?
  - a. Start with the most complex model and simplify it if necessary
  - b. Start with a simple model and gradually increase complexity**



- c. Always use the most complex model available
  - d. Use only linear models for simplicity
7. Which of the following is NOT a key component of the modeling process mentioned in the text?
- a. Understanding probability distributions
  - b. Writing down assumptions in the form of equations
  - c. **Developing new statistical methods**
  - d. Estimating model parameters from data
8. What is the purpose of fitting a model, according to the text?
- a. **To estimate the parameters of the model using observed data**
  - b. To visualize the data in different ways
  - c. To clean and preprocess the data
  - d. To develop new mathematical models
9. What is the difference between a univariate and a multivariate probability distribution?
- a. **A univariate distribution describes a single random variable, while a multivariate distribution describes multiple random variables**
  - b. A univariate distribution is always continuous, while a multivariate distribution is always discrete
  - c. A univariate distribution is used for small datasets, while a multivariate distribution is used for large datasets
  - d. There is no difference between univariate and multivariate distributions
10. What is the purpose of a conditional probability distribution?
- a. **To describe the probability of one variable given the value of another variable**
  - b. To describe the probability of multiple variables simultaneously
  - c. To describe the probability of a single variable
  - d. To describe the probability of all possible outcomes