

The Regression Cookbook

Now with Machine Learning and Stats Flavours!

G. Alexi Rodríguez-Arelis

2024-07-30

This book aims to set a common ground between machine learning and statistics regarding linear regression techniques, using `Python` and `R`, under two perspectives: **inference** and **prediction**.

Table of contents

Preface	5
Audience and Scope	6
1 Getting Ready for Regression Cooking!	8
1.1 The ML-Stats Dictionary	10
1.2 A Quick Review on Probability and Frequentist Statistical Inference	12
1.2.1 Basics of Probability	14
1.2.2 Basics of Frequentist Statistical Inference	15
1.2.3 What is Maximum Likelihood Estimation?	15
1.2.4 Supervised Learning and Regression Analysis	15
1.3 The Data Science Workflow	15
1.3.1 Study Design	15
1.3.2 Data Collection and Wrangling	17
1.3.3 Exploratory Data Analysis	18
1.3.4 Data Modelling	18
1.3.5 Estimation	18
1.3.6 Goodness of Fit	18
1.3.7 Results	18
1.3.8 Storytelling	18
1.4 Mindmap of Regression Analysis	18
2 Ordinary Least-squares	27
References	28
Appendices	29
A The ML-Stats Dictionary	29
D	29
Dependent variable	29
O	29
Outcome	29
Output	29

R	30
Response	30
T	30
Target	30

Preface

Let the regression cooking begin!

Throughout my journey as a postdoctoral fellow in the [Master of Data Science \(MDS\)](#) at the University of British Columbia, I became aware of the fascinating overlap between **machine learning** and **statistics**. Many data science students usually come across common machine learning/statistics concepts or ideas that might only differ in names. For instance, simple terms such as weights in **supervised learning** (and their equivalent statistical counterpart as regression coefficients) might be misleading for students starting their data science formation. On the other hand, from an instructor’s perspective in a data science program that subsets its courses in machine learning in **Python** and statistics in **R**, regression courses in **R** also demand the inclusion of **Python**-related packages as alternative tools. In my MDS teaching experience, this is especially critical for students whose career plans lean towards industry where **Python** is more heavily used.

As a data science educator, I view this field as a **substantial synergy** between machine learning and statistics. Nevertheless, many gaps between both disciplines still need to be addressed. Thus, closing these critical gaps is imperative in a domain with accelerated growth, such as data science. In this regard, the [MDS Stat-ML dictionary](#) inspired me to write this textbook. It basically consists of **common ground** between **foundational supervised learning models** from machine learning and **regression models** commonly used in statistics. I strive to explore **linear modelling approaches** as a primary step while highlighting different terminology found in both fields. Furthermore, this discussion is more comprehensive than a simple conceptual exploration. Hence, the second step is **hands-on practice** via the corresponding **Python** packages for machine learning and **R** for statistics.

Fun fact!

While thinking about possible names for this work, I was planning to name it “*Machine Learning and Statistics: A Common Ground*.” Nevertheless, it was quite plain and boring! That said, this whole textbook idea sounded analogous to a cookbook, given its heavily applied focus.

Hence, the cookbook name idea!

Audience and Scope

This book mainly focuses on regression analysis and its supervised learning counterpart. Thus, it is not introductory statistics and machine learning material. Instead, the following topics are suggested as prerequisites:

- **Multivariable differential calculus and linear algebra.** Certain sections of each chapter pertain to modelling estimation. Therefore, topics such as partial derivatives and matrix algebra are a great asset. You can find helpful learning resources on the [MDS webpage](#).
- **Basic Python programming.** When necessary, Python `{pandas}` library will be used to perform data wrangling. The MDS course [DSCI 511 \(Programming for Data Science\)](#) is an ideal example of this prerequisite.



Figure 1: Image by *Lubos Houska* via *Pixabay*.

- **Basic R programming.** Knowledge of data wrangling and plotting through R `{tidyverse}` is recommended for hands-on practice via the examples provided in each one of

the chapters of this book. The MDS courses [DSCI 523 \(Programming for Data Manipulation\)](#) and [DSCI 531 \(Data Visualization I\)](#) are ideal examples of this prerequisite.

- **Foundations of probability and basic distributional knowledge.** The reader should be familiar with elemental discrete and continuous distributions since they are a vital component of any given regression or supervised learning model. The MDS course [DSCI 551 \(Descriptive Statistics and Probability for Data Science\)](#) is an ideal example of this prerequisite.
- **Foundations of frequentist statistical inference.** One of the data science paradigms to be covered in this book is statistical inference, i.e., identifying relationships between different variables in a given population or system of interest via a sampled dataset. I only aim to cover a frequentist approach using inferential tools such as parameter estimation, hypothesis testing, and confidence intervals. The MDS course [DSCI 552 \(Statistical Inference and Computation I\)](#) is an ideal example of this prerequisite.
- **Foundations of supervised learning.** The second data science paradigm to be covered pertains to prediction, which is core in machine learning. The reader should be familiar with basic terminology, such as training and testing data, overfitting, underfitting, cross-validation, etc. The MDS course [DSCI 571 \(Machine Learning I\)](#) provides these foundations.
- **Foundations of feature and model selection.** This prerequisite also relates to machine learning and its corresponding prediction paradigm. Basic knowledge of prediction accuracy and variable selection tools is recommended. The MDS course [DSCI 573 \(Feature and Model Selection\)](#) is an ideal example of this prerequisite.

A further remark on probability and statistical inference

In case the reader is not 100% familiar with probabilistic and inferential topics, as discussed above, I will provide a quick refresher Chapter 1 with crucial points that are needed to follow along the **statistical way** each one of the chapters is delivered (more specifically for **modelling estimation/training** matters!).

Furthermore, this refresher will be integrated into **the three big pillars** that will be fully expanded in this book: a **data science workflow**, the right **workflow flavour** (**inferential** or **predictive**), and a **regression toolbox**.

1 Getting Ready for Regression Cooking!

It is time to prepare for the different regression techniques we will use in each of the subsequent chapters of this book. That said, there is a strong message I want to convey across all this work:

Different modelling estimation techniques in regression analysis can be smoothly grasped when we develop a fair probabilistic and inferential intuition on our populations or systems of interest.



Figure 1.1: Image by [Lucas Israel](#) via [Pixabay](#).

The above statement has a key statistical foundation on how data is generated and can be modelled via different regression approaches. More details on the concepts and ideas associated with this foundation will be delivered in Section 1.2.

Then, once we have reviewed these statistical concepts and ideas, we will move on to the three big pillars I previously pointed out:

1. The use of an ordered **data science workflow** in Section 1.3,

2. choosing the proper workflow flavour according to either an **inferential** or **predictive** paradigm as shown in Figure 1.3, and
3. the correct use of an **appropriate regression model** based on the response or outcome of interest as shown in the mind map from Section 1.4 (analogous to a **regression toolbox**).

💡 The Rationale Behind the Three Pillars

Each data science-related problem that uses regression analysis might have distinctive characteristics considering inferential (statistics!) or predictive (machine learning!) inquiries. Specific problems would implicate using outcomes (or responses) related to survival times (e.g., the time until one particular equipment of a given brand fails), categories (e.g., a preferred musical genre in the Canadian young population), counts (e.g., how many customers we would expect on a regular Monday morning in a national central bank), etc. Moreover, under this regression context, our analyses would be expanded to explore and assess how our outcome of interest is related to a further set of variables (the so-called features!). For instance, following up with the categorical outcome of the preferred musical genre in the Canadian young population, we might analyze how specific age groups prefer certain genres over others or even how preferred genres compare each other across different Canadian provinces in this young population. The sky is the limit here!

Therefore, we might be tempted to say that we should approach each regression problem should have its own workflow, given that the regression model to use would implicate particular analysis phases. **However, it turns out that is not the case to a certain extent**, and we have a regression workflow in Figure 1.3 to support this bold statement as a proof of concept for thirteen different regression models (i.e, thirteen subsequent chapters in this book). The workflow aims to homogenize our data analyses and make our modelling process more transparent and smoother. We can deliver exactly concluding insights as data storytelling while addressing our initial main inquiries. Of course, when depicting the workflow as a flowchart, there will be decision points that will turn it into **inferential** or **predictive** (the second pillar). Finally, where does the third pillar come into play in this workflow? This pillar is contained in the **data modelling stage**, where the mind map from Figure 1.12 will come in handy.

Now, before delving into probability and frequentist statistical inference, let us establish a convention on the use of admonitions beginning Section 1.2 and subsequent chapters in this textbook:

! Definition

A formal statistical and/or machine learning definition. This admonition aims to untangle the significant amount of jargon and concepts that both fields have. Furthermore, alternative terminology will be brought up when necessary to indicate the same definition across both fields.

i Heads-up!

An idea (or ideas!) of key relevance for any given modelling approach, specific workflow stage or data science-related terminology. This admonition also extends to crucial statistical or machine learning topics that the reader would be interested in exploring more in-depth.

💡 Tip

An idea (or ideas!) that might be slightly out of the scope of the topic any specific section is discussing. Still, I will provide significant insights on the matter along with further literature references to look for.

The core idea of the above admonition arrangement is to allow the reader to discern between ideas or concepts that are key to grasp from those whose understanding might not be highly essential (but still interesting to explore to check out in external references!).

1.1 The ML-Stats Dictionary

The above admonition for a **definition** will pave the way to a complimentary aspect of this textbook that I have had in mind since I started teaching statistics (and, more specifically, regression analysis) in a data science program. Machine learning and statistics usually overlap across many subjects, and regression modelling is no exception. Topics we teach in an utterly regression-based course, under a purely statistical framework, also appear in machine learning-based courses such as fundamental supervised learning, but often with different terminology. On this basis, the Master of Data Science (MDS) Program at the University of British Columbia (UBC) provides the [MDS Stat-ML dictionary](#) (Gelbart 2017) under the following premises:

This document is intended to help students navigate the large amount of jargon, terminology, and acronyms encountered in the MDS program and beyond.

This section covers terms that have different meanings in different contexts, specifically statistics vs. machine learning (ML).

Indeed, both disciplines have a tremendous amount of jargon and terminology. Furthermore, as I previously emphasized in the **Preface**, machine learning and statistics construct a **substantial synergy** that is reflected in data science. Even with this, people in both fields could encounter miscommunication issues when working together. This should not happen if we build solid bridges between both disciplines. Hence, a comprehensive **ML-Stats dictionary** (*ML* stands for *Machine Learning*) is imperative, and this textbook offers a perfect opportunity to build this resource. Primarily, this dictionary clarifies any potential confusion between statistics and machine learning regarding terminology within supervised learning and regression analysis contexts.

i Note 1: Heads-up on terminology highlights!

Following the spirit of the **ML-Stats dictionary**, throughout the book, all **statistical** terms will be highlighted in **magenta** whereas the **machine learning** terms will be highlighted in **orange**. This colour scheme strives to combine this terminology so we can switch from one field to another in an easier way. With practice and time, we should be able to jump back and forth when using these concepts.

Finally, this Appendix **A** will be one of the appendices of this book where the reader can find all those **statistical** and **machine learning-related** terms in alphabetical order. Specific terms (either statistical or machine learning-related) will include an admonition identifying which terms (again, either statistical or machine learning-related) are equivalent. Take as an example the statistical term **dependent variable**:

In supervised learning, it is the main variable of interest we are trying to **learn** or **predict**, or equivalently, the variable we are trying **explain** in a statistical inference framework.

Then, the above definition will be followed by this admonition:

⚠ Equivalent to:

Response, **outcome**, **output** or **target**.

Note we have identified four equivalent terms for the term **dependent variable**. Furthermore, according to our already defined colour scheme, these terms can be **statistical** or **machine learning-related**.

i Note 2: Heads-up on the use of terminology!

Throughout this book, I will interchangeably use specific terms when explaining the different regression approaches in each subsequent chapter. Whenever confusion arises about using these interchangeable terms, it is highly recommended to consult their definitions

and equivalences in Appendix A.

Now, let us proceed to a quick review of probability and statistics in a frequentist framework. This review will be especially essential to understanding the philosophy of modelling parameter estimation, specifically in relation to statistical inference.

1.2 A Quick Review on Probability and Frequentist Statistical Inference

Back in the old times, when I was an undergraduate student and took my very first course in probability and statistics (inference included!) in an industrial engineering context, I used to feel quite overwhelmed by the large amount of jargon and formulas one had to grasp and use regularly for primary engineering fields such as quality control in a manufacturing facility. *Population parameters, hypothesis testing, tests statistics, significance level, p-values, and confidence intervals (do not worry, our *statistical/machine learning* scheme will come in later in this section)* were appearing here and there! And to my frustration, I could never find a statistical connection between all these inferential tools! Instead, I relied on mechanistic procedures when solving assignments or exam problems.

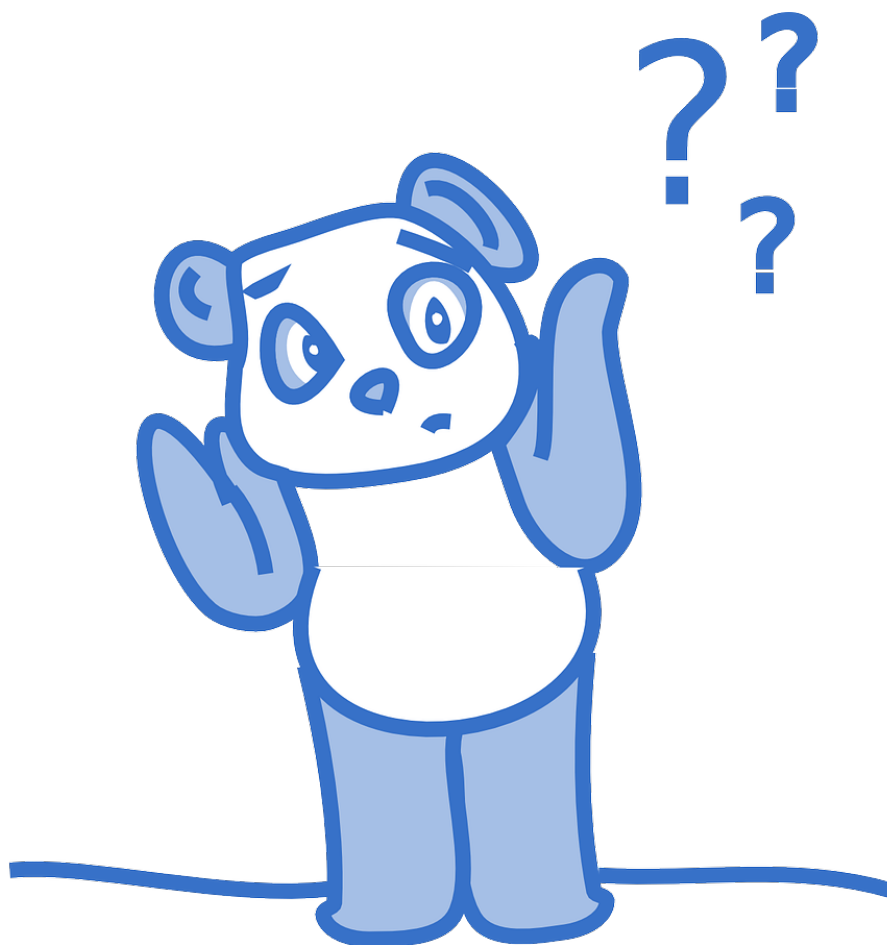


Figure 1.2: Image by [OpenClipart-Vectors](#) via [Pixabay](#).

For instance, when performing hypothesis testing for a two-sample t -test, I never reflected on what the hypotheses were trying to indicate for the corresponding population parameters(s) nor how the test statistic was related to these hypotheses. Moreover, my interpretation of the resulting p -value and/or confidence interval was purely mechanical with the inherent:

With a significance level $\alpha = 0.05$ we reject (or fail to reject, if that is the case!) the null hypothesis given that...

Honestly, I am not proud of this whole mechanical way of doing statistics now that I reflect on it after many years of practice, teaching, and research. Then, of course, the above situation should not happen when we learn key statistical topics for the very first time as undergraduate students. That is why we will dig into a more intuitive way of viewing probability and its crucial role in statistical inference. This matter will help us deliver more coherent storytelling when presenting our results in practice during any regression analysis. Note that the role of

probability also extends to model training when it comes to supervised learning and not just in regard to statistical inference.

Having said all this, it is time to introduce a statement that I ended up reasoning the very first time I taught hypothesis testing in the introductory statistical inference course in the MDS program at UBC (after reflecting on the basics of hypothesis testing):

In statistical inference, everything always boils down to randomness and how we can control it!

I know! That is quite a bold statement. Nonetheless, once one starts teaching statistical topics to audiences not entirely familiar with the usual field jargon, the idea of randomness always persists across many different tools. And, of course, regression analysis is not an exception at all since it also involves inference on population parameters of interest! This is why I have allocated this section in the textbook to explain core probabilistic and inferential concepts to pave the way to its role in regression analysis.

On the other hand, even though this book has prerequisites related to the basics of probability via different distributions and the fundamentals of frequentist statistical inference as stated in **Audience and Scope**, we will retake essential concepts as follows:

- The role of *random variables* and *probability distributions* and the governance of *population (or system) parameters* (i.e., the so-called Greek letters we usually see in statistical inference and regression analysis). Section 1.2.1 will explore these topics more in detail while connecting them to the subsequent inferential terrain under a *frequentist context*.
- Section 1.2.2 will explore how a *random sample* is connected to the basics of *hypothesis testing* and its intrinsic components such as *null* and *alternative hypotheses*, *type I* and *type II* errors, *test statistic*, *standard error*, *p-value*, and *confidence interval*.
- When delving into supervised learning and regression analysis, we might wonder how randomness is incorporated into *model fitting*. That is quite a fascinating aspect, implemented via a crucial statistical tool known as *maximum likelihood estimation*. This tool is heavily related to the concept of *loss function* in supervised learning. Section 1.2.3 will explore these matters in more detail.
- Finally, Section 1.2.4 will briefly discuss the connections between supervised learning and regression analysis regarding terminology.

1.2.1 Basics of Probability

Under this foundation, our data is coming from this given population or system of interest. Moreover, the population or system is assumed to be governed by parameters. In practice, we will never know the parameters

1.2.2 Basics of Frequentist Statistical Inference

1.2.3 What is Maximum Likelihood Estimation?

1.2.4 Supervised Learning and Regression Analysis

1.3 The Data Science Workflow

It is time to review the so-called data science workflow. Each one of these three pillars is heavily connected since a general Data Science workflow is applied in each one of these regression models, which aims to help in our learning (i.e., we would be able to know what exact stage to expect in our data analysis regardless of the regression model we are being exposed to). Therefore, a crucial aspect of the practice of Regression Analysis is the need for this systematic Data Science workflow that will allow us to solve our respective inquiries in a reproducible way. Figure 1.3 shows this workflow which has the following general stages (I briefly define each one of them; note a broader delivery will be done in subsequent subsections):

1. **Study design:**
2. **Data collection and wrangling:**
3. **Exploratory data analysis:**
4. **Data modelling:**
5. **Estimation:**
6. **Goodness of fit:**
7. **Results:**
8. **Storytelling**

i What if there is no formal structure in our regression analysis?

Since very early learning stages in data analysis, it is crucial to

Now, suppose we do not follow a predefined workflow in practice. In that case, we might be at stake in incorrectly addressing our inquiries, translating into meaningless results outside the context of the problem we aim to solve. This is why the formation of a Data Scientist must stress this workflow from the very introductory learning stages.

1.3.1 Study Design

The first stage of this workflow is heavily related to the *main statistical inquiries* we aim to address throughout the whole data analysis process. As a data scientist, it is your task to primarily translate these inquiries from the stakeholders of the problem as *inferential* or *predictive*. Roughly speaking, this primary classification can be explained as follows:

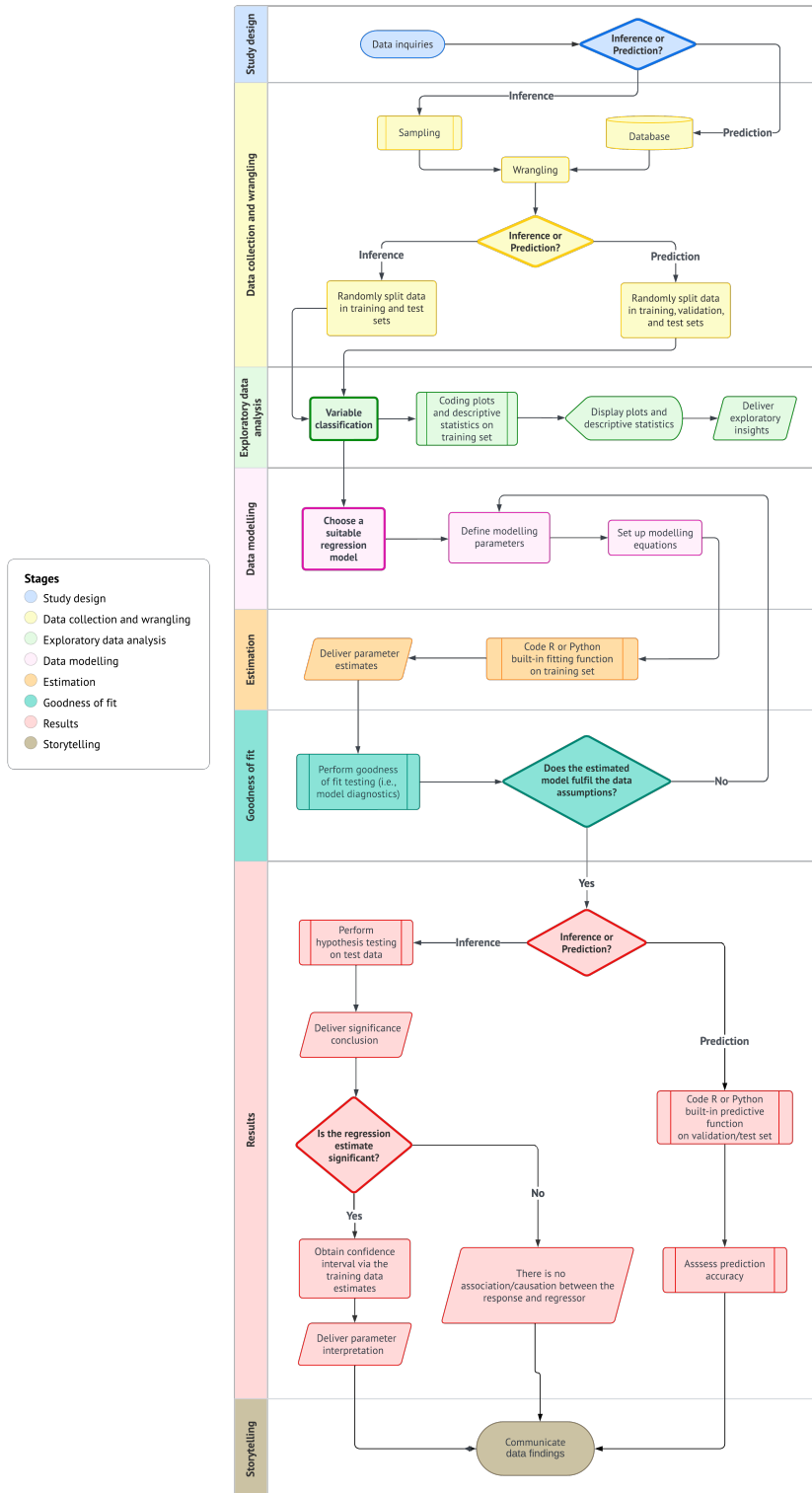


Figure 1.3: Data science workflow for *inferential* and *predictive* inquiries in regression analysis and supervised learning, respectively.

- **Inferential.** The main objective is to untangle relationships of *association* or *causation* between the regressors (i.e., explanatory variables) and the corresponding response in the context of the problem of interest. Firstly, we would assess whether there is a statistical relationship between them. Then, if significant, we would quantify by how much.
- **Predictive.** The main objective is to deliver response predictions on further observations of regressors, having estimated a given model via a current training dataset. Unlike inferential inquiries, assessing a statistically significant association or causation between our variables of interest is not a primary objective but *accurate predictions*. **This is one of the fundamental paradigms of machine learning.**

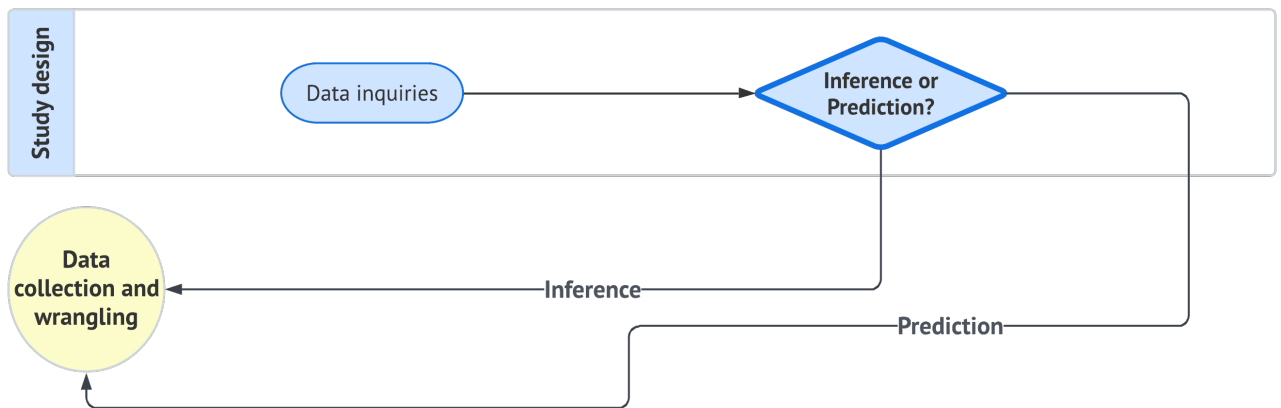


Figure 1.4: *Study design* stage from the data science workflow in Figure 1.3. This stage is directly followed by *data collection and wrangling*.

1.3.2 Data Collection and Wrangling

Once we have defined our main statistical inquiries, it is time to collect our data. Note we have to be careful about the way we collect this data since it might have a particular impact on the quality of our statistical practice:

- Regarding inferential inquiries, recall we are approaching populations or systems of interest governed by *unknown and fixed distributional parameters*. Thus, via sampled data, we aim to *estimate* these distributional parameters. This is why **a proper sampling method** on this population or system of interest is critical to obtaining representative data for *appropriate hypothesis testing*.

💡 Tip 1: A Quick Debrief on Sampling!

This stage is coloured in gray in {numref}data-science-workflow, unlike the other ones coloured in yellow. This is because sampling topics are out of the scope of this course and MDS in general. Nevertheless, we still need to stress that a proper sampling method is also key in inferential inquiries to assess association and/or causation between the regressors and your response of interest. That said, depending on the context of the problem, we could apply either one of the following methods of sampling:

- **Simple random sampling.**
- **Systematic sampling.**
- **Stratified sampling.**
- **Clustered sampling.**
- Etc.

As in the case of Regression Analysis, statistical sampling is a vast field, and we could spend a whole course on it. If you are more interested in these topics, *Sampling: design and analysis* by Lohr offers great foundations.

- In practice, regarding predictive inquiries, we would likely have to deal with databases given that our trained models will not be used to make inference and parameter interpretations.

1.3.3 Exploratory Data Analysis

1.3.4 Data Modelling

1.3.5 Estimation

1.3.6 Goodness of Fit

1.3.7 Results

1.3.8 Storytelling

1.4 Mindmap of Regression Analysis

Having defined the necessary statistical aspects to execute a proper supervised learning analysis, either *inferential* or *predictive* across its seven sequential phases, we must dig into the different approaches we might encounter in practice as regression models. The nature of our outcome of interest will dictate any given modelling approach to apply, depicted as clouds in

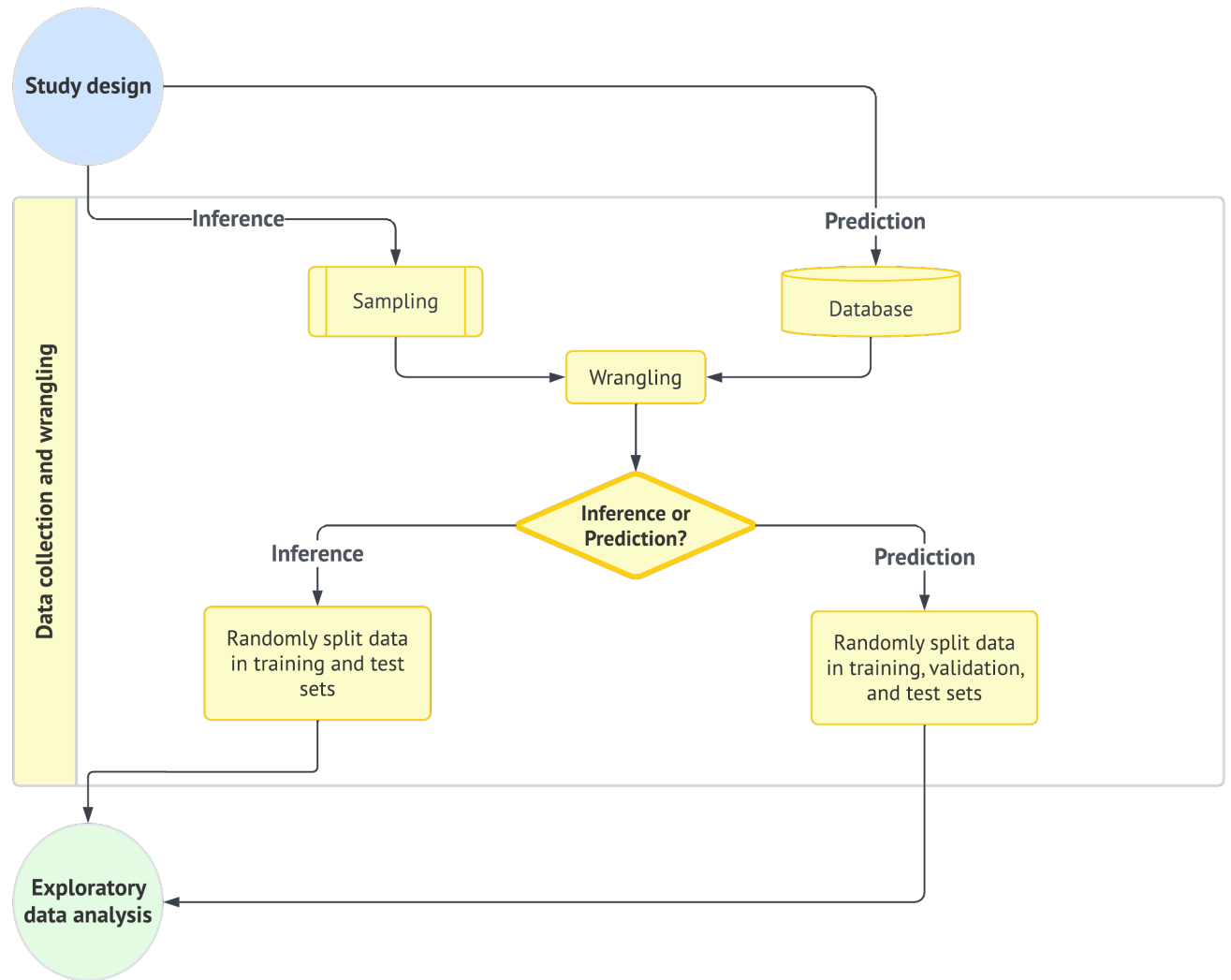


Figure 1.5: *Data collection and wrangling* stage from the data science workflow in Figure 1.3. This stage is directly followed by *exploratory data analysis* and preceded by *study design*.

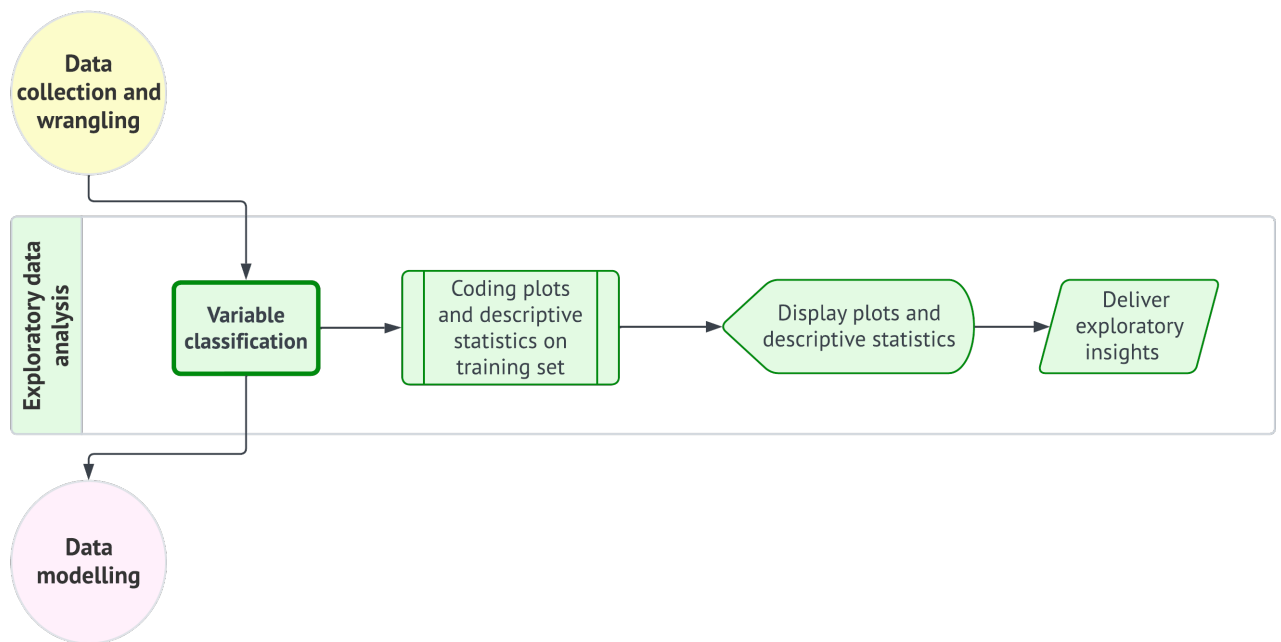


Figure 1.6: *Exploratory data analysis* stage from the data science workflow in Figure 1.3. This stage is directly followed by *data modelling* and preceded by *data collection and wrangling*.

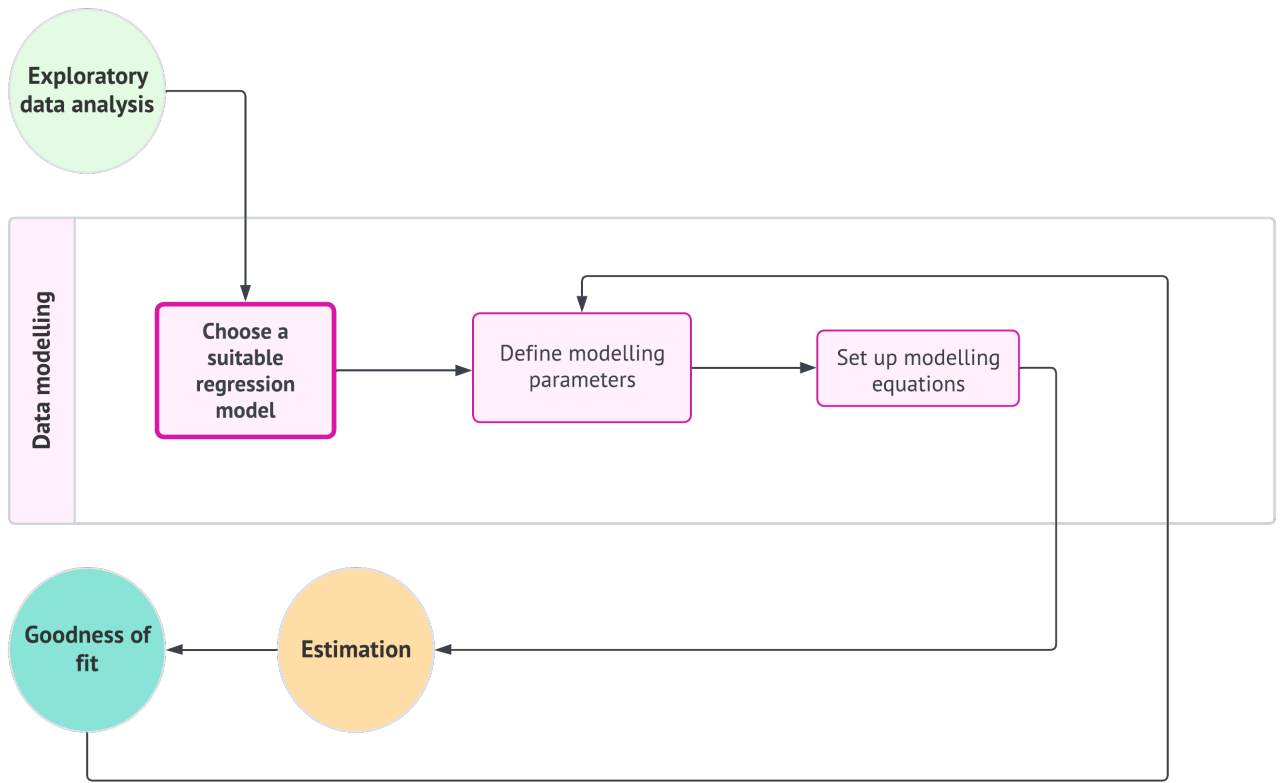


Figure 1.7: *Data modelling* stage from the data science workflow in Figure 1.3. This stage is directly preceded by *exploratory data analysis*. On the other hand, it is directly followed by *estimation* but indirectly with *goodness of fit*. If necessary, the *goodness of fit* stage could retake the process to *data modelling*.

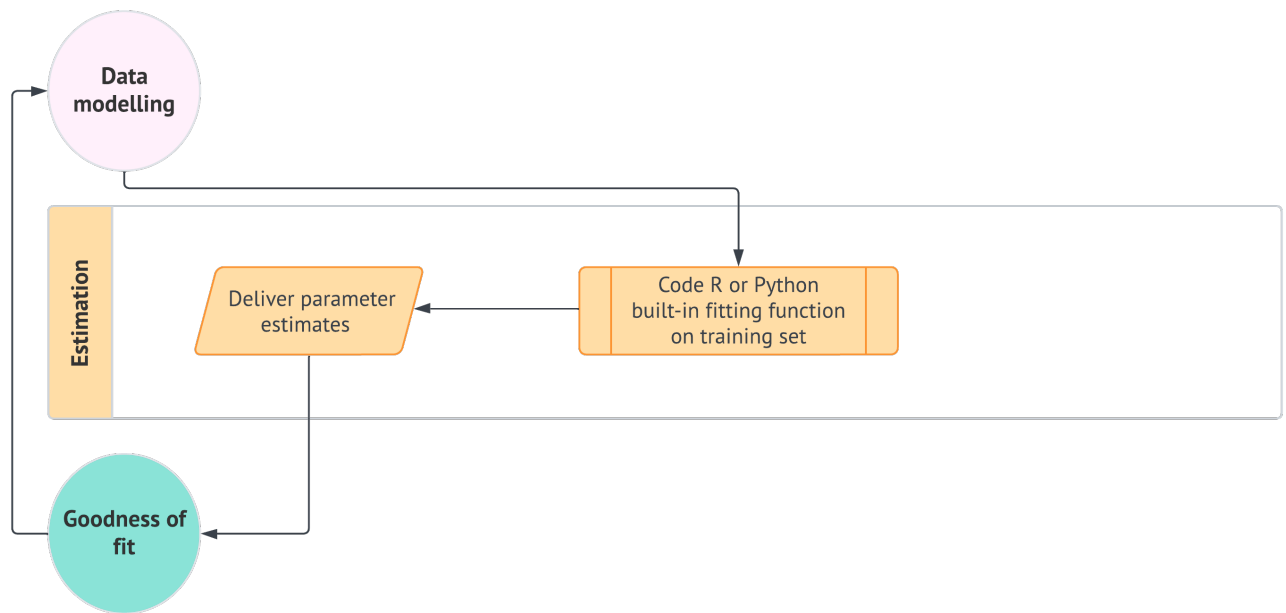


Figure 1.8: *Estimation* stage from the data science workflow in Figure 1.3. This stage is directly preceded by *data modelling* and followed by *goodness of fit*. If necessary, the *goodness of fit* stage could retake the process to *data modelling* and then to *estimation*.

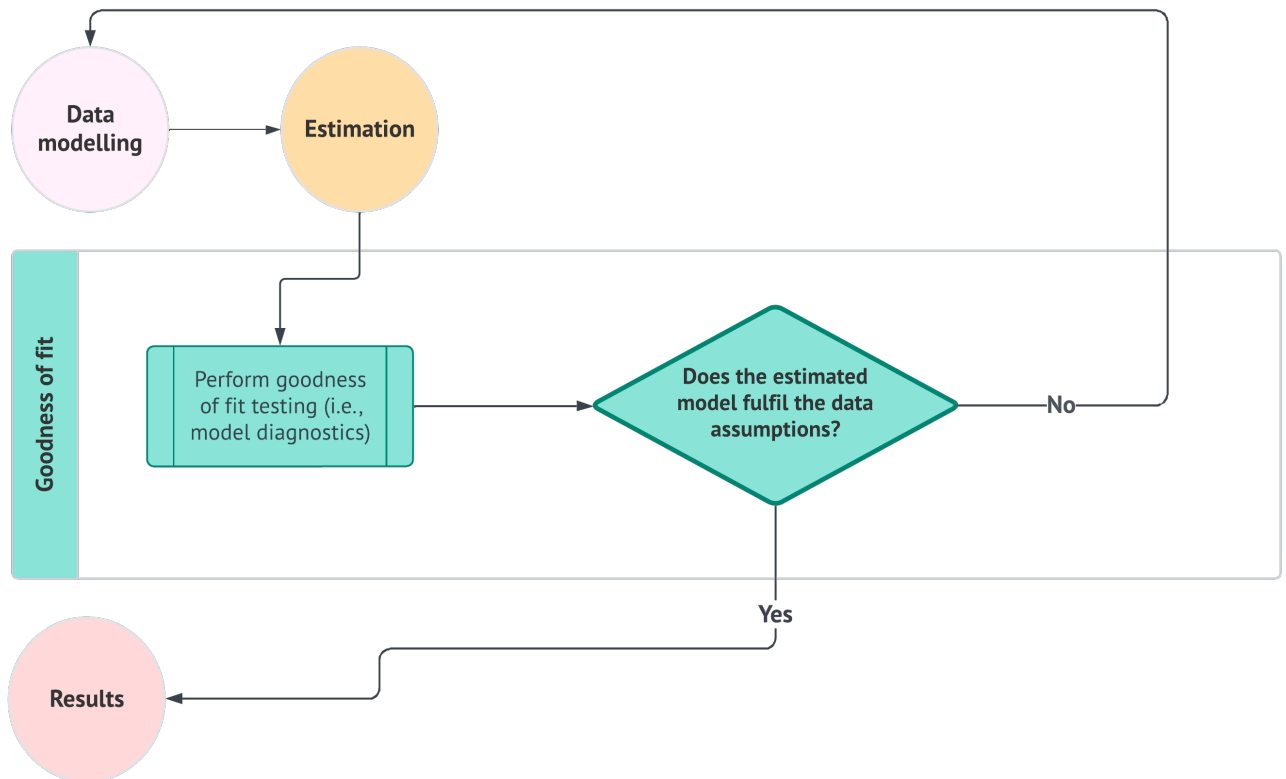


Figure 1.9: *Goodness of fit* stage from the data science workflow in Figure 1.3. This stage is directly preceded by *estimation* and followed by *results*. If necessary, the *goodness of fit* stage could retake the process to *data modelling* and then to *estimation*.

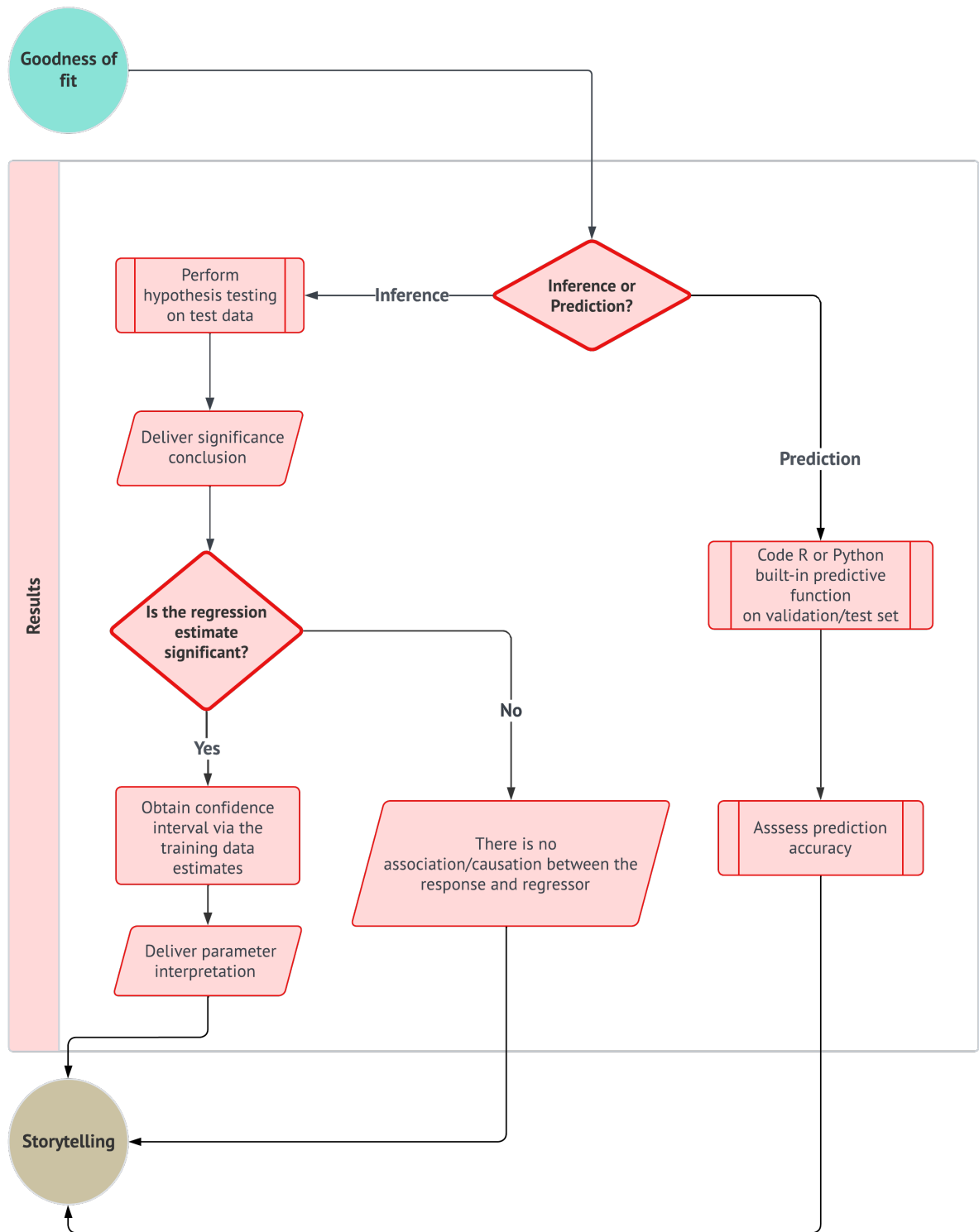


Figure 1.10: *Results* stage from the data science workflow in Figure 1.3. This stage is directly followed by *storytelling* and preceded by *goodness of fit*.

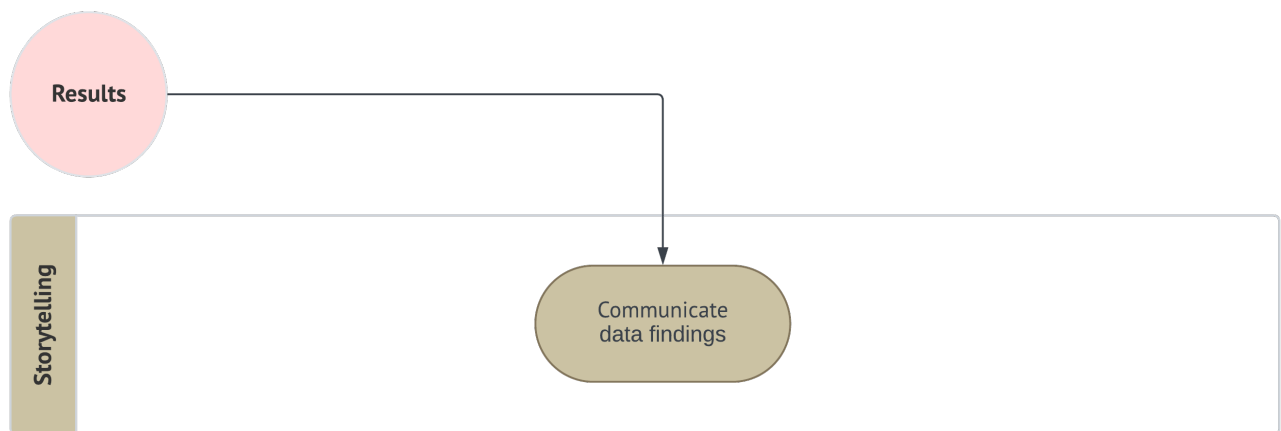


Figure 1.11: *Storytelling* stage from the data science workflow in Figure 1.3. This stage preceded by *results*.

Figure 1.12. Note these regression models can be split into two sets depending on whether the outcome of interest is *continuous* or *discrete*. Therefore, under a probabilistic view, identifying the nature of a given random variable is crucial in regression analysis.

That said, we will go beyond OLS regression and explore further regression techniques. In practice, these techniques have been developed in the statistical literature to address practical cases where the OLS modelling framework and assumptions are not suitable anymore. Thus, throughout this block, we will cover (at least) one new regression model per lecture.

As we can see in the clouds of Figure 1.12, there are 13 regression models: 8 belonging to discrete outcomes and 5 to continuous outcomes. Each of these models is contained in a chapter of this book, beginning with the most basic regression tool known as ordinary least-squares in Chapter 2. We must clarify that the current statistical literature is not restricted to these 13 regression models. The field of regression analysis is vast, and one might encounter more complex models to target certain specific inquiries. Nonetheless, I consider these models the fundamental regression approaches that any data scientist must be familiar with in everyday practice.

Even though this book comprises 13 chapters, each depicting a different regression model, we have split these chapters into two major subsets: those with *continuous* outcomes and those with *discrete* outcomes.

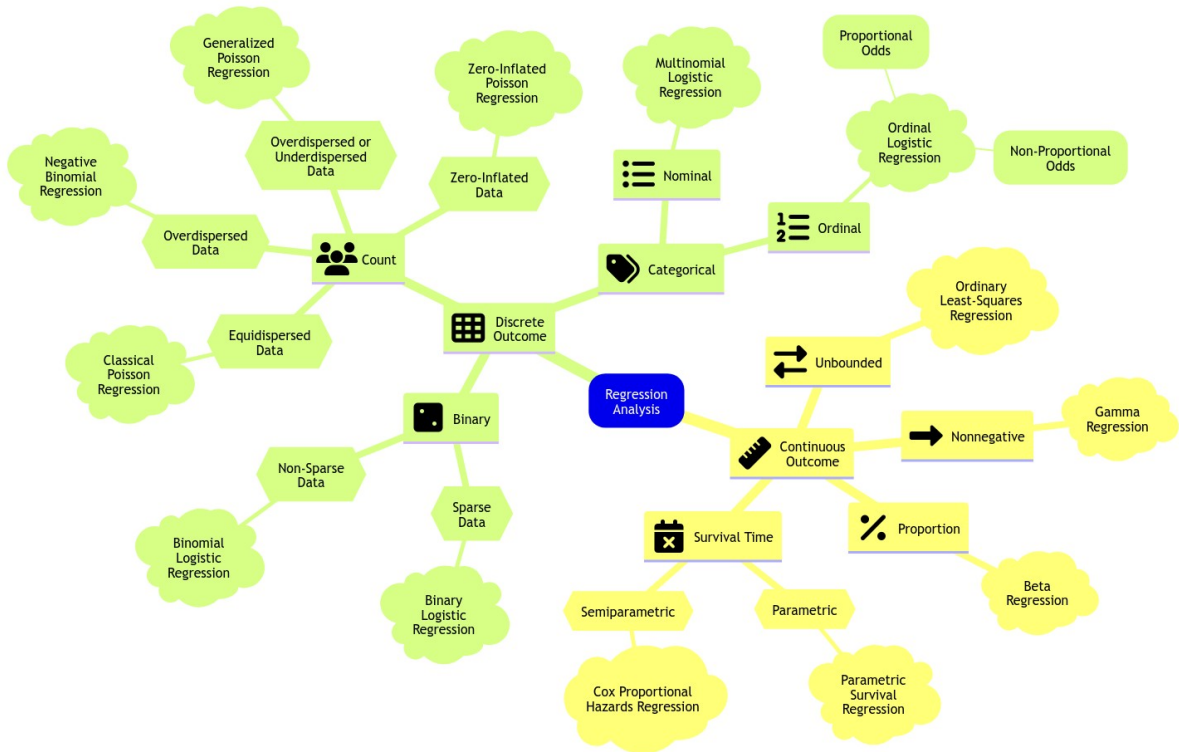


Figure 1.12: Regression analysis mindmap depicting all modelling techniques to be explored in this book. These techniques are split into two big sets: *continuous* and *discrete* outcomes.

2 Ordinary Least-squares

References

Gelbart, Michael. 2017. “Data Science Terminology.” *UBC MDS*. Master of Data Science at the University of British Columbia. https://ubc-mds.github.io/resources_pages/terminology/.

A The ML-Stats Dictionary

D

Dependent variable

In supervised learning, it is the main variable of interest we are trying to **learn** or **predict**, or equivalently, the variable we are trying **explain** in a statistical inference framework.

⚠ Equivalent to:

Response, outcome, output or target.

O

Outcome

In supervised learning, it is the main variable of interest we are trying to **learn** or **predict**, or equivalently, the variable we are trying **explain** in a statistical inference framework.

⚠ Equivalent to:

Dependent variable, response, output or target.

Output

In supervised learning, it is the main variable of interest we are trying to **learn** or **predict**, or equivalently, the variable we are trying **explain** in a statistical inference framework.

⚠ Equivalent to:

Dependent variable, response, outcome or target.

R

Response

In supervised learning, it is the main variable of interest we are trying to **learn** or **predict**, or equivalently, the variable we are trying **explain** in a statistical inference framework.

⚠ Equivalent to:

Dependent variable, outcome, output or target.

T

Target

In supervised learning, it is the main variable of interest we are trying to **learn** or **predict**, or equivalently, the variable we are trying **explain** in a statistical inference framework.

⚠ Equivalent to:

Dependent variable, response, outcome or output.