# PROJECT 2 SUBMISSION TEMPLATE

December 14, 2023

\#

Project (Ch. 1-9): Complex Models

```
[1]: # Load the CourseKata library
     suppressPackageStartupMessages({
         library(coursekata)
     })
```

# 1 Introduction/Overview of the Problem or Question

The utilization of the "Bad Passwords" data frame is particularly pertinent given the ubiquity of passwords in everyday user interactions. This dataset, meticulously curated and formatted by Sujay Kapadnis from an external source called "Information is Beautiful," offers a comprehensive collection of data values. The incorporation of tidytuesdayR library has ensured that the data is structured in a tidy format, enhancing its accessibility and interpretability.

By leveraging this dataset, the objective is to discern patterns and characteristics associated with commonly used passwords, thereby contributing insights into the decision-making process surrounding password security. Furthermore, the examination of this data will shed light on whether the findings support or provide a foundational basis for advocating the use of randomly generated passkeys over conventional passwords. This exploration aligns with the broader discourse on enhancing cybersecurity measures and underscores the practical implications of password-related decision-making.

It appears that Sujay Kapadnis curated this dataset with the intention of offering a compilation of passwords deemed extremely common and, consequently, "bad" choices for users. The dataset, comprising 507 observations across 9 distinct variables, provides valuable information for individuals seeking to enhance their password security practices.

```
[2]: link = "https://docs.google.com/spreadsheets/d/e/
     ↪2PACX-1vRVVtguBGyqPxStcC8r0q2kWk9PUQ92wLRk2gjbCarLoKP6r7cwCWbBqWXpxpoEPPbyWFc7SW3nWDTg/
     ↪pub?output=csv"
     passwords <- read.csv(link, header=TRUE)
     str(passwords)
```

```
'data.frame':   507 obs. of  9 variables:
 $ rank           : int  1 2 3 4 5 6 7 8 9 10 …
 $ password       : chr  "password" "123456" "12345678" "1234" …
```

```
 $ category        : chr  "password-related" "simple-alphanumeric" "simple-
alphanumeric" "simple-alphanumeric" …
 $ value           : num  6.91 18.52 1.29 11.11 3.72 …
 $ time_unit       : chr  "years" "minutes" "days" "seconds" …
 $ offline_crack_sec: num  2.17 1.11e-05 1.11e-03 1.11e-07 3.21e-03 1.11e-06
3.21e-03 2.17 2.17 8.35e-02 …
 $ rank_alt        : int  1 2 3 4 5 6 7 8 9 10 …
 $ strength        : int  8 4 4 4 8 4 8 4 7 8 …
 $ font_size       : int  11 8 8 8 11 8 11 8 11 11 …
```

The variables include the rank of the password, with the most common and vulnerable passwords marked as 1. Additionally, the dataset encompasses the password itself, the categorical classification indicating the type of password, the time complexity estimation for cracking each password, the strength rating, and the associated font size. These attributes collectively afford a comprehensive understanding of the vulnerabilities associated with commonly used passwords, enabling users to make informed decisions about their password choices and potentially advocating for the adoption of more secure practices.

To explore the notion of a selected password contributing to its classification as a "bad password," it is essential to identify and analyze the relevant factors within the dataset. The outcome variable, in this case, could be framed as a binary variable indicating whether a password is categorized as "bad" or not. Considering the dataset's structure, I would like to designate a binary outcome variable, such as "Bad_Password," where 1 represents a password labeled as bad and 0 denotes otherwise. Subsequently, want to employ a statistical technique to investigate the relationships between different features (for example: rank, password type, time complexity, strength, etc.) and the likelihood of a password being classified as "bad."
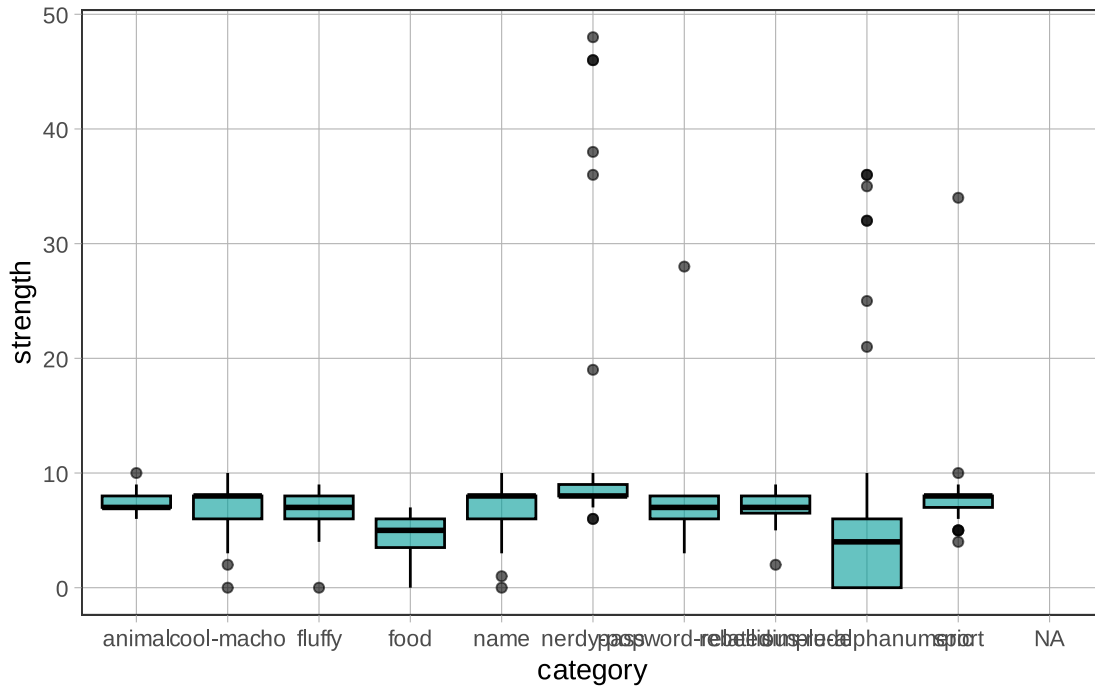
```
[3]: gf_boxplot(strength ~ category, data = passwords) %>%
         gf_labs(title = "Figure 1. Password categories and strengths")
```

```
Warning message:
"Removed 7 rows containing non-finite values (`stat_boxplot()`)."
```

**Figure 1. Password categories and strengths**



In illustrated Figure 1, a meticulous examination has been conducted, comparing the password categories (depicted on the x-axis, representing the explanatory variable) against their corresponding strength values (presented on the y-axis, signifying the outcome variable). A deliberate refinement has been implemented in this rendition to rectify a logical discrepancy encountered in the original plot.

In the initial visualization, a comparison was erroneously attempted between password categories and the passwords themselves. This inadvertent oversight resulted in a logical error, given that each password is inherently distinct within its respective category. The rectification involved maintaining the consistency of the x-axis rows while appropriately aligning the y-axis to accommodate the distinct strength values associated with each password category. This adjustment ensures the accuracy and clarity of the comparative analysis, mitigating the prior logic error and enhancing the interpretability of the depicted relationship between password categories and their corresponding strengths.

After a deep overview of Figure 1. there were some values that didn't make sense with the given data frame. After realizing this, I took a thorough look at the data frame and realized that there we values that were out of the scope of the variable. For example, in the strength variable, from the data frame descrition the values are supposed to only range from 1 (being the lowest value possible) to 10 (being the highest value possible). However, within the original data frame, there were values within the variable strength that were above the number 10 (some values were 25, 32, 15, 17, etc.). Because of this, I opted to filter the data frame based off of that variable reasoning. As represented in the code section below:

```
[4]:  # Filter all values above 10 in `strength`
      strenths.filtered <- filter(passwords, strength %in% c(1:10))

      # Represents the filtered data frame in a jitter plot
      gf_jitter(strength ~ category, data = strenths.filtered) %>%
        gf_labs(title = "Figure 2. Password categories and their strengths filtered")
```

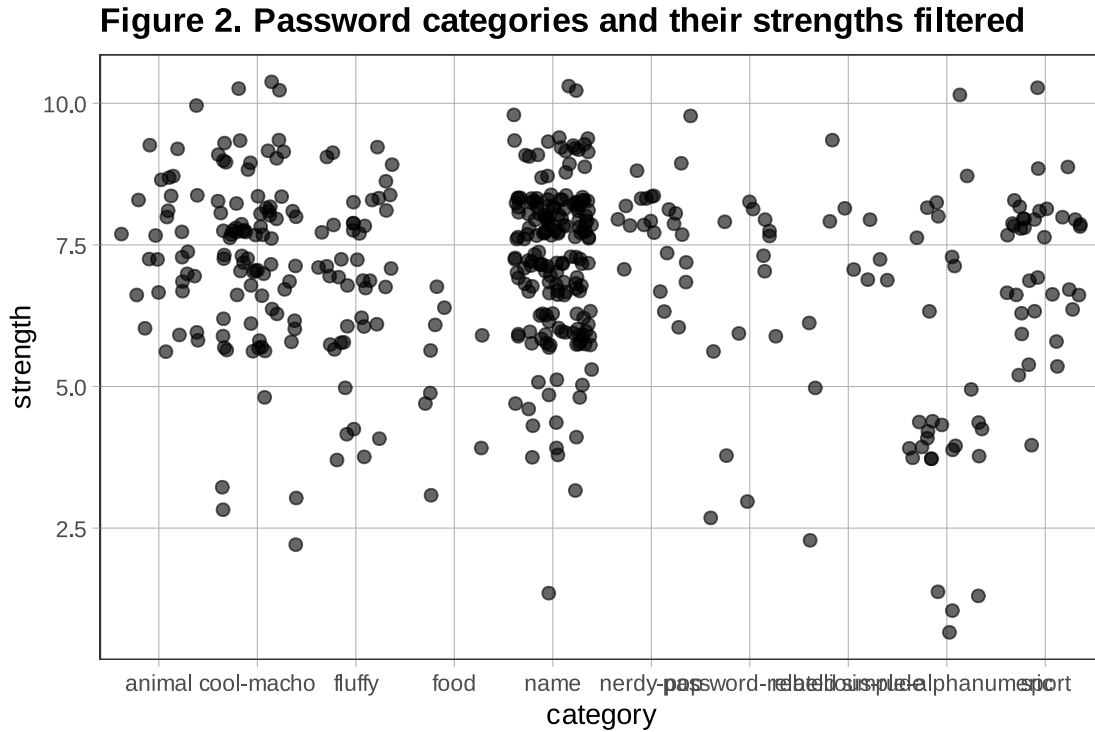**Figure 2. Password categories and their strengths filtered**



Figure 2. meticulously delineates the depiction of the excluded NA (Not Available) values, a conspicuous enhancement from the original data frame. The deliberate omission of each NA value has resulted in a discernible reduction in the overall size of the dataset, albeit a modest reduction. Remarkably, the refined data frame, now comprising 455 observations, retains its comprehensive structure with a total of 9 variables.

This strategic curation, involving the meticulous exclusion of NA values, not only serves to refine the dataset but also underscores a commitment to data integrity. The resulting Figure 2. thus encapsulates a more focused and streamlined representation, facilitating a clearer exploration and analysis of the underlying patterns and relationships within the dataset.

**Research question:** I am poised to undertake an investigation into a fundamental research question that pivots on the nuanced interrelation between the inherent composition of a given `password` and its corresponding `category`. My primary objective is to ascertain the presence of a discernible dependency between these pivotal variables. This inquiry is specifically designed to unravel the intricacies surrounding the potential of a password's strength, gauged within the context of its assigned category, to function as a reliable predictor of its qualitative classification as either "good" or "bad."

This research endeavor delves into the complex symbiosis existing between the constituent characters shaping a password and its inherent strength. The overarching aim is to elucidate the manner in which these salient factors intricately influence one another within the realm of password security. This analytical pursuit promises to contribute profound insights into the nuanced dynamics governing password efficacy, thereby offering a comprehensive understanding of the pivotal factors shaping secure password selection practices.

**Initial hypothesis:** My initial hypothesis pertaining to this circumstance posits the existence of a substantive interdependence between a password's designated category and its overarching strength. I envision a scenario wherein passwords characterized by heightened complexities and a deliberate avoidance of conventional categories, such as those encompassing password-related, simple-numeric, and alphanumeric attributes, manifest significantly enhanced strength. Consequently, such passwords are anticipated to exhibit resilience against unauthorized access, rendering them inherently more formidable and resistant to theft or unauthorized intrusion attempts.

In light of the acknowledgment that an array of undisclosed factors may exert influence on password strength beyond the explicitly identified variables, I have endeavored to construct a comprehensive word equation aimed at encapsulating the intricate relationship between these multifaceted elements:

This relational dynamic can be aptly symbolized by the following equation:

`Bad_Password = category + other stuff`

This formulated equation elegantly captures the concept that the classification of a password as "Bad" is the outcome of a synergistic interplay between its assigned category and an ensemble of additional, albeit unspecified, determinants. The incorporation of the term `other stuff` within this equation serves as a recognition of the intricate complexity inherent in this relationship, underscoring the multifarious nature of the contributing factors. This formula, by placing emphasis on the category as a central constituent, aims to convey a nuanced understanding of the overarching influence of these elements on the determination of password security.

## 2 Exploring Variation

Upon conducting a more comprehensive examination of the filtered data frame, a discerning realization prompted the further refinement of the dataset. Notably, variables such as `rank_alt` and `font_size` emerged as devoid of substantive value within the context of my ongoing investigation. This discernment led to a strategic decision to streamline the dataset by omitting these extraneous variables.

The resulting clean data set, now meticulously curated, is characterized by a judicious exclusion of superfluous variables, an operation executed through the following code snippet:

```
[5]:  # Code for creating the refined data set
      passwords.filtered <- strenths.filtered[-c(7,9)]
```

This targeted approach ensures that the dataset is attuned to the essential variables, thereby promoting a more focused and meaningful exploration. Such refinement not only optimizes computational efficiency but also enhances the interpretability of subsequent analyses by retaining only the pertinent components essential to the research inquiry.

```
[6]: # Initial exploration of the variation in the dataset using a bar graph
     gf_bar( ~ category, data = strenths.filtered, fill = ~strength) %>%
       gf_labs(title = "Figure 3. Categorical count")

     # Facet the graph between category and strength to compare
     gf_bar( ~ category, data = strenths.filtered, fill = ~strength) %>%
       gf_labs(title = "Figure 4. Categorial comparization to strengths") %>%
       gf_facet_grid(strength ~ .)
```
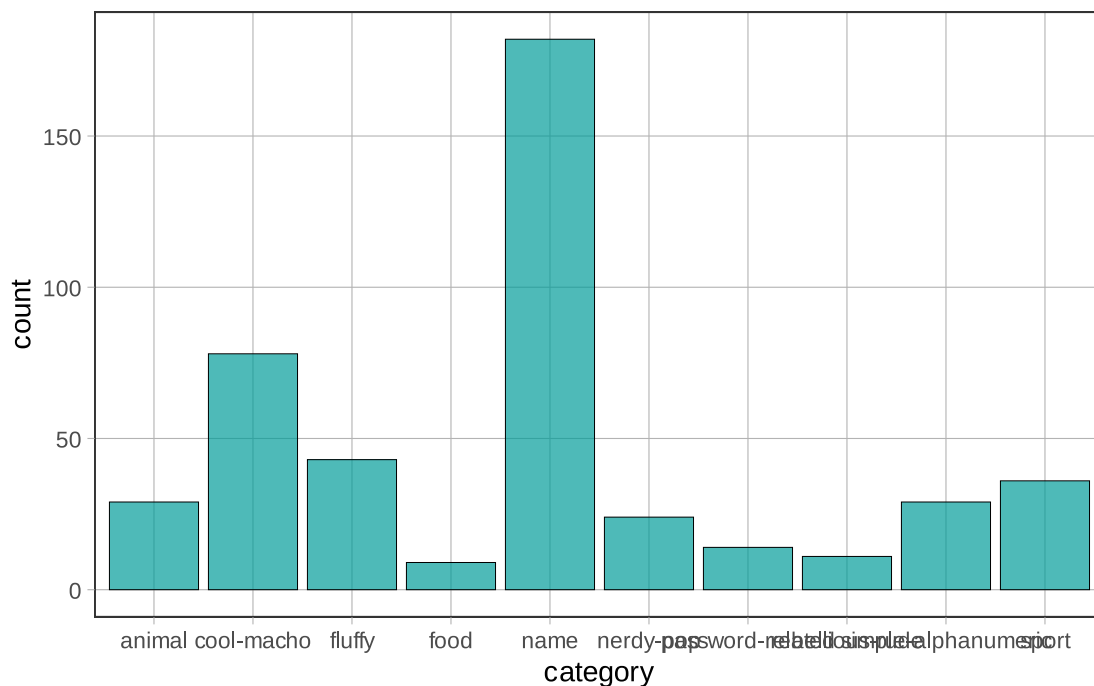
Warning message:
"The following aesthetics were dropped during statistical
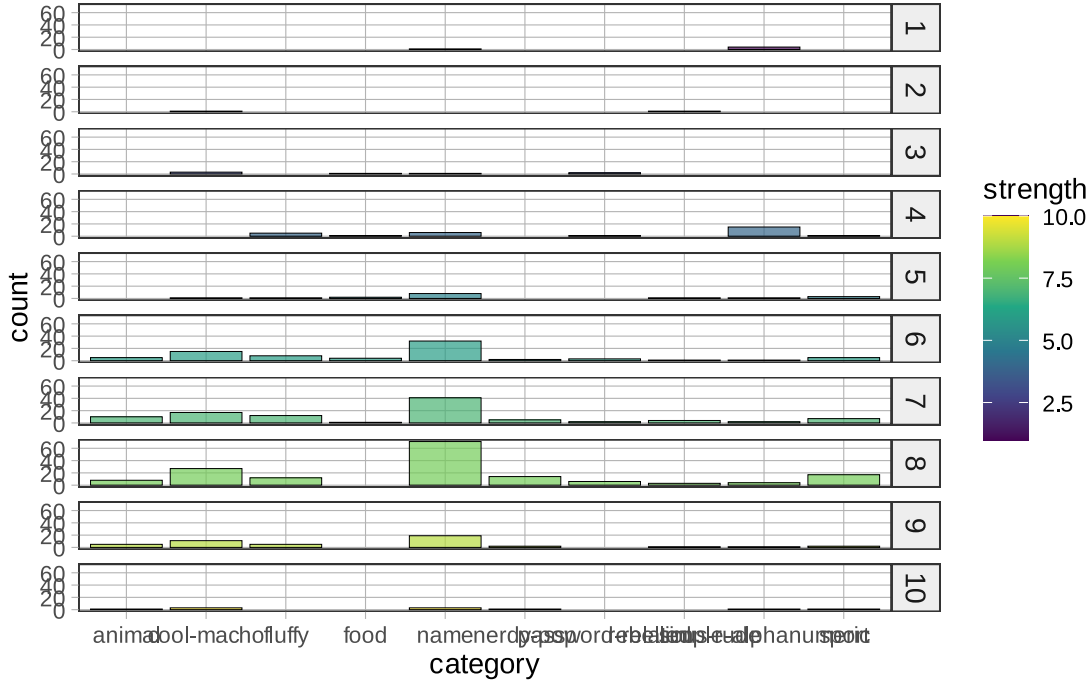transformation: fill
  This can happen when ggplot fails to infer the correct grouping
structure in
  the data.
  Did you forget to specify a `group` aesthetic or to convert a
numerical
  variable into a factor?"

## Figure 3. Categorical count

**Figure 4. Categorial comparization to strengths**



**Analysis:** As gleaned from the insights offered by Figure 3, it becomes apparent that passwords incorporating names stand out as the most prevalent category, justifiably earning the distinction of being the most popular or frequently employed across the dataset. It is crucial to note, however, that this representation solely constitutes a raw tally of passwords within each category, lacking further contextual elucidation.

Figure 4, on the other hand, delves deeper into the intricacies of the dataset. It surpasses the mere aggregation of password counts by introducing a more nuanced perspective. This subsequent visualization dissects the count of each password, stratifying them based on their respective strengths. The result is a more detailed and layered plot, offering a comprehensive portrayal of how individual passwords within each category align with different levels of strength. This elevated analytical approach in Figure 4 adds a valuable layer of understanding, transcending a mere categorical count to provide a nuanced exploration of the distribution of strengths within each password category.
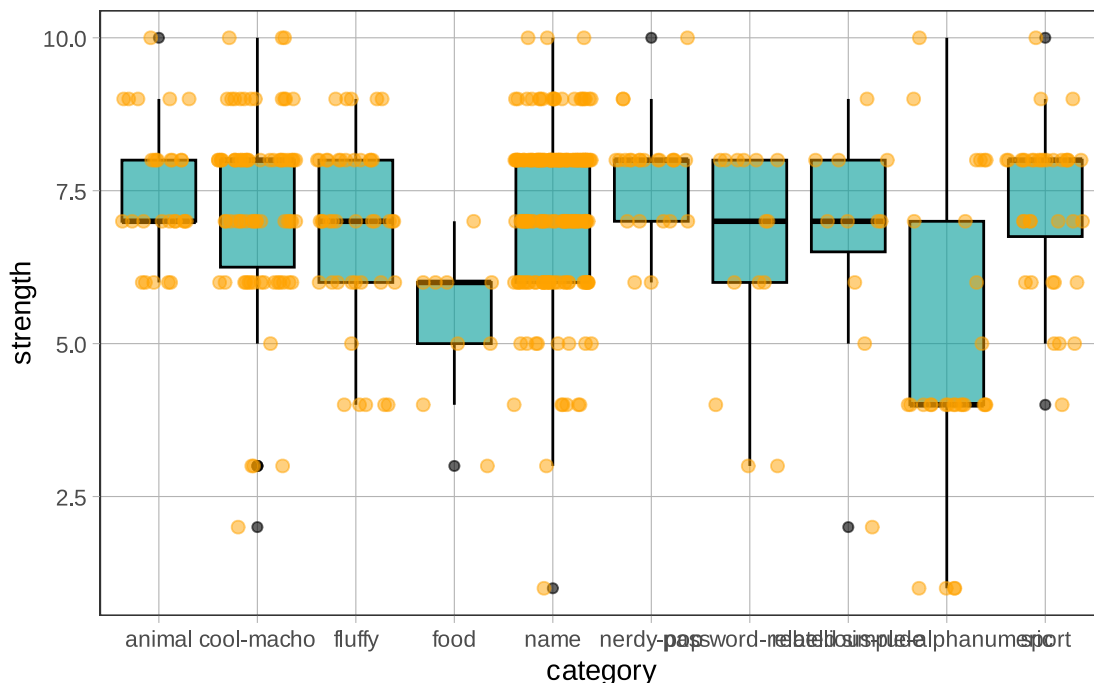
**Note:** It is paramount to recognize that, notwithstanding their categorical alignment, the passwords presented in the faceted grid collectively span a range of values from 1 to 10. Within this continuum, the numerical assignment reflects the strength of each password, with 1 denoting a weak password and 10 indicative of a password deemed very strong. However, a crucial contextual consideration must be acknowledged and retained: the inherent classification of this dataset as a collection of passwords predominantly characterized as "weak."

This acknowledgment is pivotal as it underscores the relative nature of strength within this specific dataset. Even those passwords categorized as "strong" are to be interpreted within the confines of this dataset's characteristics. It is imperative to bear in mind that, when juxtaposed with passwords exceeding a character limit of 16 characters, inclusive of any categorical affiliation, all passwords

within this dataset may be considered comparatively weak. This contextual awareness is pivotal for a nuanced and accurate interpretation of the dataset's implications for password security.

```
[7]: # Layer a jitter plot on a box plot to get a better sense of the variation
     gf_boxplot(strength ~ category, data = passwords.filtered) %>%
       gf_jitter(height = 0, color = "orange", alpha = 0.5) %>%
       gf_labs (title = "Figure 5. Password strength and category with specified␣
       ↪values")
```

**Figure 5. Password strength and category with specified values**



As elucidated by the insights derived from Figure 5, a noteworthy and intriguing observation arises: the assigned values of strength for passwords are distinctly whole numbers, without falling within the intermediary range of two consecutive values. Figure 5 compellingly demonstrates that passwords are exclusively designated with strength values such as 7.00, 8.00, and 9.00, without intermediary decimal values like 7.50, 8.50, or 9.50.

This discrete assignment of whole-number strength values imparts a distinctive and categorical nature to the strength assessment, suggesting a deliberate choice to employ a discrete scale rather than a continuous one. This nuanced aspect of the dataset's characterization adds a layer of precision to the strength evaluation, presenting a unique approach to the quantification of password robustness within the context of the provided data.

**Refined Research Question/Hypothesis:** As discerned from the initial insights gleaned from these visualizations, a conspicuous observation emerges: there appears to be no inherent association between a password's category and its preassigned strength. In order to subject this refined understanding to further scrutiny and validation, I have undertaken the importation of an additional

dataset singularly focused on password strengths.

This strategic augmentation of data is undertaken with confidence in R's robust capabilities for data manipulation. Despite the utilization of disparate values in this new dataset, the versatility of R enables me to adeptly reconfigure and harmonize the data, aligning it with the specific analytical objectives of this endeavor. This meticulous integration aims to facilitate an enriched and comprehensive exploration of the intricate relationship between password categories and their respective strengths.

In accordance with the elucidated considerations regarding the ongoing trajectory of this project, the refined theoretical equation has been meticulously tailored to accentuate the nuanced interconnectivity between a password's strength and its assigned category:

`Bad_Password = strength + category + Other Stuff`

This equation strategically encapsulates the intricate dynamics involved in categorizing a password as "Bad," emphasizing the integral role played by both the categorical assignment and the inherent strength of the password. The inclusion of the term `Other Stuff` further acknowledges the multifaceted nature of the factors contributing to password vulnerability, thereby offering a comprehensive framework for the continued exploration and analysis of password security within the specified context.

## 3 Model Variation

```
[8]: # Empty linear model
     empty_linearModel <- lm(strength ~ NULL, data = passwords.filtered)
     empty_linearModel

     # Fitted model
     passwords_model <- lm(strength ~ category, data = passwords.filtered)
     passwords_model

     # Predictions based off of the regression model `password_model`
     passwords.filtered$predictions <- predict(passwords_model)
     head(passwords.filtered)
```

```
Call:
lm(formula = strength ~ NULL, data = passwords.filtered)

Coefficients:
(Intercept)
      7.042



Call:
lm(formula = strength ~ category, data = passwords.filtered)

Coefficients:
```

|  | (Intercept) | categorycool-macho |
|  | 7.5517 | -0.2440 |
|  | categoryfluffy | categoryfood |
|  | -0.6215 | -2.2184 |
|  | categoryname | categorynerdy-pop |
|  | -0.3210 | 0.2399 |
|  | categorypassword-related | categoryrebellious-rude |
|  | -1.1232 | -0.8245 |
|  | categorysimple-alphanumeric | categorysport |
|  | -2.7241 | -0.2739 |

A data.frame: 6 × 8

|  | rank | password | category | value | time_unit | offline_crack_sec | str |
|---|---|---|---|---|---|---|---|
|  | <int> | <chr> | <chr> | <dbl> | <chr> | <dbl> | <i: |
| 1 | 1 | password | password-related | 6.91 | years | 2.17e+00 | 8 |
| 2 | 2 | 123456 | simple-alphanumeric | 18.52 | minutes | 1.11e-05 | 4 |
| 3 | 3 | 12345678 | simple-alphanumeric | 1.29 | days | 1.11e-03 | 4 |
| 4 | 4 | 1234 | simple-alphanumeric | 11.11 | seconds | 1.11e-07 | 4 |
| 5 | 5 | qwerty | simple-alphanumeric | 3.72 | days | 3.21e-03 | 8 |
| 6 | 6 | 12345 | simple-alphanumeric | 1.85 | minutes | 1.11e-06 | 4 |

The chosen modeling approach to best encapsulate the relationship under investigation is a linear model. The representation of this model, following the Generalized Linear Model (GLM) notation, can be expressed as follows:

$$Bad\_Password_i = 7.042 + b_1(category_i) + e_i$$

Where a `Bad_Password` is a result of a default base value (strength of 7.042) plus a category it belongs to (thus bringing down the initial strength of a given password as it is part of a category within this data frame, respectively) plus any external error accounted.

In this formulation, the $Bad\_Password_i$ is the outcome variable, shaped by a default base value ($b_0$), the impact of the category to which it belongs ($b_1(X_i)$), and any external error term ($e_i$). The introduction of the category as a predictor in the model inherently influences the initial strength of a given password, thereby reflecting its association with the predefined categories within this dataset. This linear model is designed to unravel and quantify the nuanced interplay between password categories and their classification as `Bad_Passwords`.

In the context of the non-empty linear model, $b_1$ is a categorical variable rather than a quantitative variable, therefore, the expressions for the Generalized Linear Model (GLM) can be represented as follows:

$$Bad\_Password_i = 7.5517 + -0.2440(categorycool - macho_i) + e_i$$

$$Bad\_Password_i = 7.5517 + -0.6215(categoryfluffy_i) + e_i$$

$$Bad\_Password_i = 7.5517 + -2.2184(categoryfood_i) + e_i$$

$$Bad\_Password_i = 7.5517 + -0.3210(categoryname_i) + e_i$$

$$Bad\_Password_i = 7.5517 + 0.2399(categorynerdy - pop_i) + e_i$$

$$Bad\_Password_i = 7.5517 + -1.1232(categorypassword - related_i) + e_i$$

$$Bad\_Password_i = 7.5517 + -0.8245(categoryrebellious - rude_i) + e_i$$

$$Bad\_Password_i = 7.5517 + -2.7241(categorysimple - alphanumeric_i) + e_i$$

$$Bad\_Password_i = 7.5517 + -0.2739(categorysport_i) + e_i$$

Here, each equation encapsulates the predicted value of a Bad_Password_i based on the respective impact of the specified category ($b_1$) and an error term ($e_i$). These equations offer a quantitative representation of the relationship between each category and the likelihood of a password being classified as a "Bad_Password" within the context of the given linear model.

# 4    Evaluate Model(s)

```
[9]:  # Anova table
      supernova(passwords_model)

      # Visual representation of both models
      gf_point(strength ~ category, data = passwords.filtered) %>%
        gf_model(passwords_model, color = "orange") %>%
        gf_model(empty_linearModel, color = "blue") %>%
        gf_labs(title = "Figure 7. Password Modeling with Empty (blue) and Fit␣
        ↪(orange) Models")
```
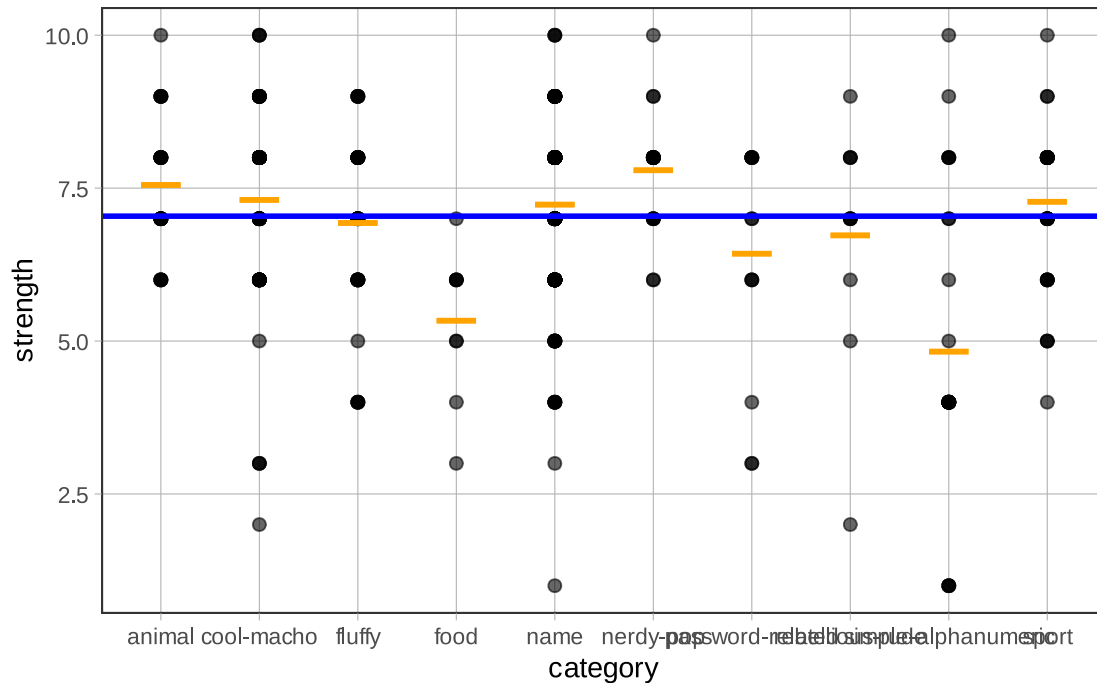
```
Analysis of Variance Table (Type III SS)
Model: strength ~ category
```

|  |  | SS | df | MS | F | PRE | p |
|-----|---------------|----------|-----|--------|--------|-------|-------|
| Model (error reduced) | | 210.392 | 9 | 23.377 | 10.704 | .1780 | .0000 |
| Error (from model) | | 971.815 | 445 | 2.184 | | | |
| Total (empty model) | | 1182.207 | 454 | 2.604 | | | |

**Figure 7. Password Modeling with Empty (blue) and Fit (orange**



The fitted model is evaluated against the empty model through a comprehensive analysis, which reveals several key indicators substantiating its superiority within this context. The reduction in the sum of squares, coupled with a small F ratio and a higher Proportional Reduction in Error (PRE), collectively position the fitted model as the optimal choice in this analytical framework.
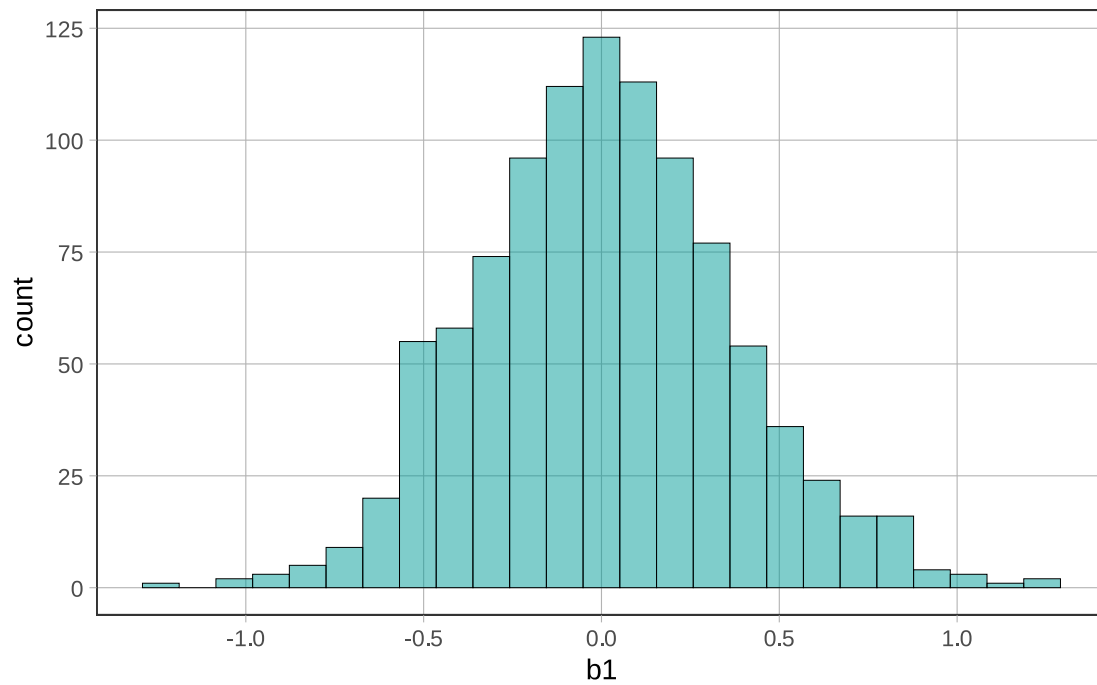
The ANOVA table, obtained via the function `supernova(passwords_model)` as represented in Figure 7, serves as a pivotal tool for this comparative assessment. The observed greater PRE value and F ratio corroborate the efficacy of the fitted model, indicating a superior fit compared to the empty model. This statistical validation underscores the empirical support for the chosen model, affirming its appropriateness in capturing the nuances of the relationship between password categories and their classification as "Bad_Passwords" within the scope of this research.

```
[10]: # Randomness exploration
      sDofB1 <- do(1000) * b1(shuffle(strength) ~ category, data = passwords.filtered)

      gf_histogram(~ b1, data = sDofB1) %>%
        gf_labs(title = "Figure 8. Standard Deviations of B1")

      # "Unlikely" distribution of B1
      gf_histogram(~ b1, data = sDofB1, fill = ~middle(b1, 0.95)) %>%
        gf_labs(title = "Figure 9. Unlikely Standard Deviation Distributions of B1")
```
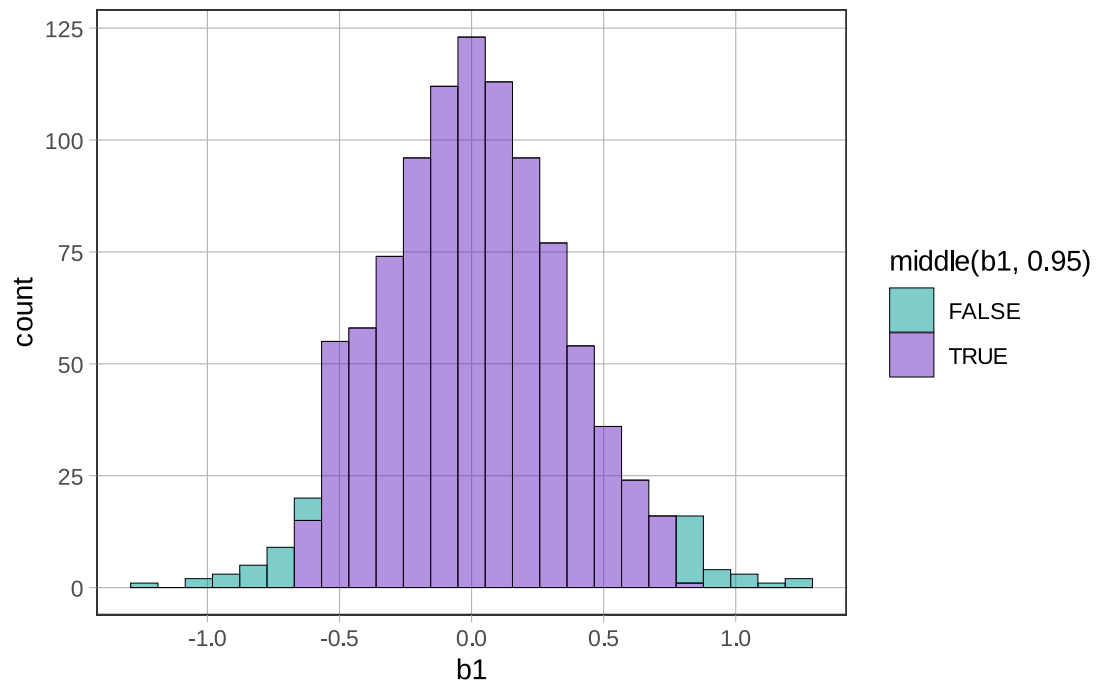
Figure 8. Standard Deviations of B1



Figure 9. Unlikely Standard Deviation Distributions of B1

**Randomness/Chance:** The analysis of this data frame reveals a discernible absence of external factors contributing to stochasticity. In essence, the dataset exhibits a notable level of determinism, wherein randomness or external changes do not exert influence nor bear implications on the inherent structure and characteristics of the data. The absence of stochastic elements underscores the robustness and reliability of the dataset, affirming the validity of the patterns and insights derived from the analyses conducted. This deterministic quality enhances the interpretability and confidence in the conclusions drawn from the dataset, fostering a more precise understanding of the underlying dynamics governing password strength and its categorical associations.

# 5 Conclusion

**Summary:** The selection of my research inquiry is underpinned by a discerning recognition of the contemporary and far-reaching implications inherent in the pervasive challenge of compromised accounts stemming from the utilization of feeble passwords in the expansive domain of internet interactions. The inherent gravity of this investigative pursuit resides in its capacity to unearth the foundational tenets that underscore the construct of a resilient password, thereby propelling a decisive stride towards the prophylaxis of compromised accounts—an issue that pervasively afflicts an extensive spectrum of individuals within the digital milieu.

The impetus driving my commitment to unravel the intricate interplay between passwords and their constituent character categories exemplifies a proactive and forward-thinking strategy in tackling an exigent real-world predicament. The overarching goal is to discern discernible patterns and interrelationships intrinsic to password characteristics, thereby furnishing a comprehensive comprehension of the determinants that contribute to the robustness of passwords. This meticulous inquiry, consequently, aspires to cast illuminating insights on the factors underpinning password strength, fostering an empirically grounded understanding that can guide individuals toward the formulation of more impregnable and secure passwords. The ultimate objective is to alleviate the vulnerabilities associated with compromised accounts.

In summation, the core tenet of my research query not only aligns harmoniously with the pragmatic exigencies of online security but also holds the promise of delivering tangible insights. These insights, in turn, serve to refine and elevate password creation practices, concurrently fortifying the fabric of cybersecurity awareness. Through this scholarly undertaking, the prospect emerges not merely as an intellectual endeavor but as a practical contribution to the collective arsenal against the ever-evolving threats to digital security.

In the preliminary phase of my research endeavor, I conscientiously undertook a comprehensive exploration and refinement of the dataset to ensure its optimal suitability for subsequent analyses. Leveraging the capabilities of the R programming language, I judiciously employed the `str()` function to meticulously scrutinize the structural composition of the data frame. This meticulous examination unveiled a data frame comprising nine variables and 507 observations. Although characterized by a relatively modest scale in comparison to the overarching scope of the study, the dataset exhibited a semblance of tidiness. To enhance the analytical tractability, I executed a discerning strategy to eliminate observations associated with missing values (NA), thereby fortifying the integrity of the dataset.

Upon closer scrutiny, it became apparent that certain values within the `strength` variable surpassed the documented maximum threshold of 10, as stipulated in the dataset documentation. In response, I exercised precision by filtering out all observations within the `strength` variable that exceeded the

14

aforementioned threshold, aligning the dataset with the documented specifications. Furthermore, in a strategic effort to streamline the analytical focus, I prudently excluded variables deemed peripheral to the central objectives of the analysis. Variables such as `alt_rank` and `font_size` were intentionally omitted, allowing for a more concentrated and targeted exploration of the dataset.

Embracing these judicious constraints, I cultivated a sense of assurance and confidence in the resultant dataset, characterizing it as "tidy" and poised for rigorous and sophisticated analyses. This new tiny data frame has been referenced as *passwords.filtered* throughout my research.

Armed with a meticulously curated and visually represented tidy data frame, I proceeded to augment my understanding of its intrinsic characteristics through the discerning lens of a box and whisker plot. This graphical exploration facilitated a nuanced comprehension of the data frame's structure, empowering me to recommence my research endeavors with heightened acuity.

In a deliberate effort to unveil the nuances within the dataset, I employed a faceted visualization that delineated password strength across different categories. This insightful stratification unveiled fluctuations within the observations for each value, offering a rich visual tapestry that inspired further inquiry.

Eager to extract additional insights, I endeavored to overlay two distinct plots, employing both box and whisker plots and jitter plots. This multi-layered visualization strategy aimed to discern patterns and irregularities within the data, illuminating the interplay between password strength and categories. Notably, the revelation that each value constituted an exact whole number, eschewing decimal representations, added an intriguing layer to my observations. This newfound information prompted a thoughtful refinement of my research question and hypothesis, signaling a pivot toward a more nuanced and focused investigation. The acknowledgment of this detail also led to a reduction in the generalization of my initial hypothesis, underscoring the iterative nature of the research process and the necessity for adaptability in response to emerging insights.

In the nuanced exploration of variation within my models, the categorical nature of my explanatory variable necessitated the formulation of multiple Generalized Linear Model (GLM) notations. This meticulous consideration stems from the recognition that the inherent strength of a password is intricately tied to the category it closely aligns with. In light of this, I meticulously crafted both an empty model and a fitted model to dissect the multifaceted dynamics at play.

As I delved into the intricacies of these models, a compelling realization surfaced—the fitted model emerged as the most apt and informative model, encapsulating the essence and outcomes of the data frame with unparalleled coherence. This model not only expounded upon the content within the dataset but also served to mitigate the influence of randomness, thereby fortifying its robustness as a reliable analytical tool.

In the culmination of this exhaustive research undertaking, a resounding conclusion emerges: the classification of passwords as `Bad_Password`s is intricately linked to the specific category with which a password associates. This insight underscores the pivotal role that categorical associations play in determining the vulnerability and strength of passwords.

**Implications of the Results:** The implications derived from this research bear substantial significance in the realm of password creation. The discernment that the strength of a password is intricately linked to the specific category it aligns with carries profound consequences for enhancing the efficacy of password security practices. This newfound understanding empowers individuals in the creation of passwords endowed with greater strength, thereby serving as a formidable deterrent against unauthorized access by third parties. Armed with this knowledge, users can make informed

decisions during the password generation process, strategically aligning their choices with categories that confer heightened security.

By implementing such practices, the research endeavors to contribute to the reduction of external intrusions into personal and confidential accounts. This proactive approach to password creation, rooted in the insights garnered from the study, holds the potential to fortify digital security measures and mitigate the risks associated with unauthorized access. Consequently, the research extends beyond its analytical dimensions, offering tangible implications for individuals seeking to bolster the resilience of their digital credentials.

**DGP:** I posit that my model adeptly captures the intricacies of the data generation process (DGP), exemplified by the fitted model's notable fidelity and relative alignment with the observed outcomes of the explanatory variable. The meticulous calibration and representation inherent in the model underscore its efficacy in encapsulating the underlying mechanisms that govern the generation of the dataset, thereby enhancing its robustness and reliability in the context of the broader data science framework.

**Further Research/Investigation:** In the course of this research investigation, the primary emphasis rested on scrutinizing the intricate relationship between password strength (`strength`) and categorical associations (`category`). However, the scope of inquiry could be further enriched by delving into the prospect of discerning the correlation between password strength and the temporal security measure denoted as `offline_crack_sec`. This nuanced exploration aims to ascertain whether the categorical affiliations within the variable `category` exert a discernible influence on `offline_crack_sec`, thereby contributing to the deterministic characterization of password strength.

The inclination to extend the analysis beyond the categorical dimension reflects a commitment to comprehensiveness in understanding the multifaceted facets of password security. This expanded inquiry into the interplay of `strength` and `offline_crack_sec`, contingent upon the influence of `category`, underscores a strategic and nuanced approach to deciphering the intricate dynamics that underpin the resilience of passwords in the face of offline cracking attempts.