This study aims to statistically assess the effects of training schedules and uncontrollable external factors on mental preparedness, as measured by trait emotional intelligence in 100-kilometer ultra-marathon runners. These factors include elevation gain, course surface type, average weekly mileage, and demographic characteristics. By examining these factors, the study seeks to offer a technical comprehension of the interplay between mental and physical readiness, which is essential for improving endurance sports performance. This work supports my own objective of utilizing data-driven techniques to maximize athletic training in sports like boxing and soccer, merging mental and physical fitness to create evidence-based performance gains. It also has practical importance for sports performance research.
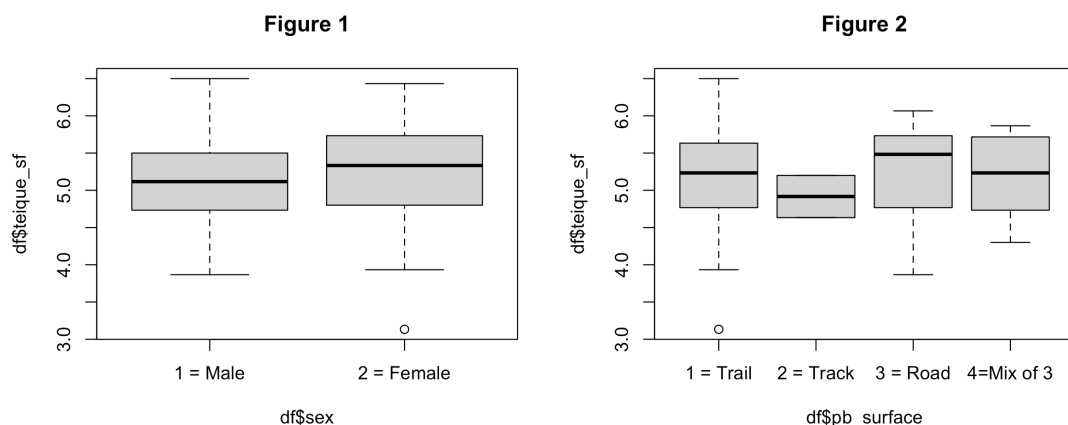
Eric Samtleben's (2023) study, which looks at the relationship between emotional intelligence and performance in 100-kilometer ultra-marathon runners, is the source of the dataset. It includes assessments of emotional intelligence together with important demographic, racial, and physical training factors based on self-reported data from 288 individuals. The study was approved by the Trent University Research Ethics Board, and participants were sourced online using sites including Facebook, Reddit, Strava, and Raceroster.com. This extensive dataset provides insightful information about the physical and psychological aspects of ultra-marathon performance.

This dataset is especially helpful for investigating how mental preparation, as determined by emotional intelligence, is impacted by practice and outside, uncontrolled events. In my exploratory data analysis (EDA), I considered predictors like age, sex, pb_surface, pb_elev, pb100k_dec, and avg_km, but I also concentrated on important factors like teique_sf, the main measure of trait emotional intelligence. These factors allow us to evaluate the potential interactions between an athlete's mental readiness and the physical demands of ultra-marathon
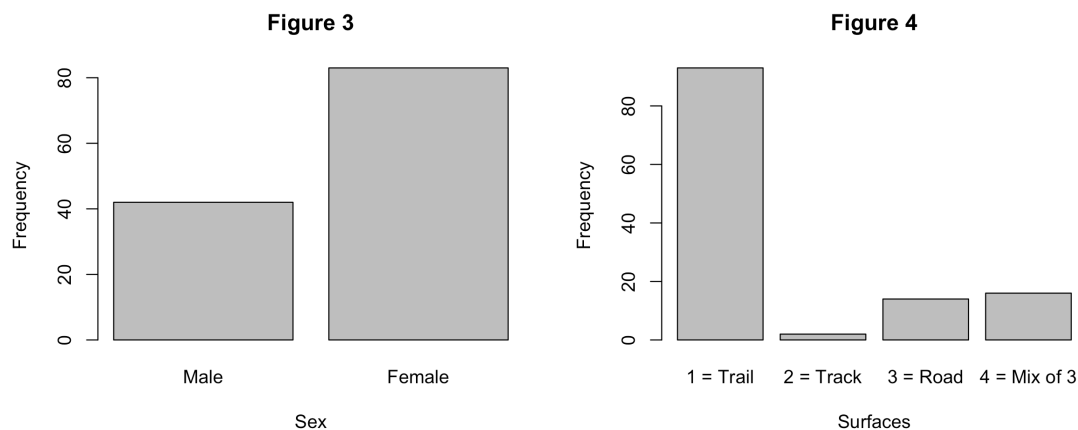
running. Clear definitions and ranges for every variable are provided by the comprehensive data dictionary and related documentation, guaranteeing that my analyses are based on a thorough comprehension of the underlying components. The dataset is extremely valuable for research attempting to create more comprehensive training and performance plans in endurance sports since it integrates psychological and physical metrics.

For my analysis, I started by importing the ultra-running dataset and eliminating any rows with missing values. Ten factors make up the dataset, which also includes assessments of emotional intelligence, physical performance metrics, and demographic data. In addition to predictors like age, sex, pb_surface (race surface type), pb_elev (elevation gain), pb100k_dec (personal best time in decimal hours), and avg_km (average kilometers ran each week), my main variables of interest are the trait emotional intelligence score (teique_sf).

I began by analyzing the distributions of the categorical variables, pb_surface and sex, to spot any possible outliers and have an idea of the distribution of these variables. A possible outlier among females is shown in the boxplot of teique_sf by sex (see Figure 1), which may be a result of data input errors or a real variance in the population. Likewise, Figure 2's boxplot of teique_sf by pb_surface raises the possibility that there is an outlier in the "trail" group with a lower emotional intelligence score.
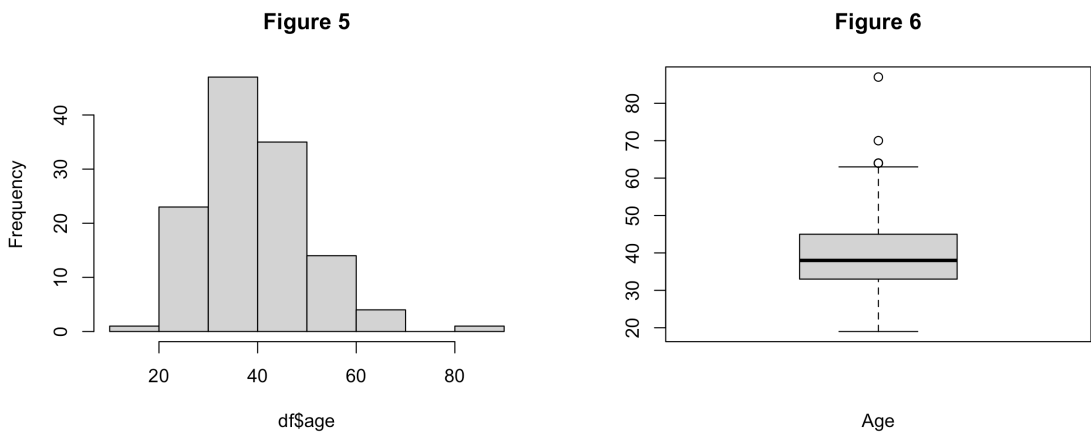


Figure 1

Figure 2

The frequency distributions of sex and pb_surface were also displayed using bar plots (Figures 3 and 4), demonstrating that the dataset is evenly represented in both categories. These illustrations assist us in determining if the final model's robust handling or additional transformation of the categorical components is necessary.

**Figure 3**



**Figure 4**

Then I looked at each numerical predictor separately. A possible upper outlier is indicated by age histograms and boxplots (Figures 5 and 6), which imply that some participants are significantly older than the bulk of the sample.
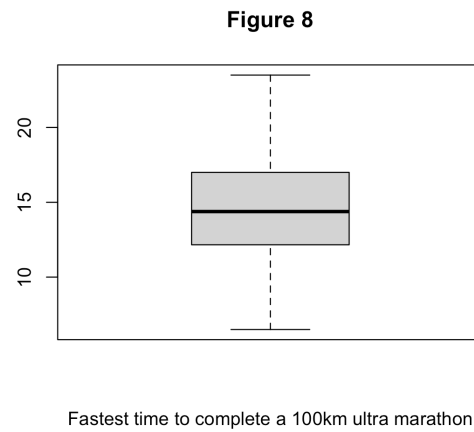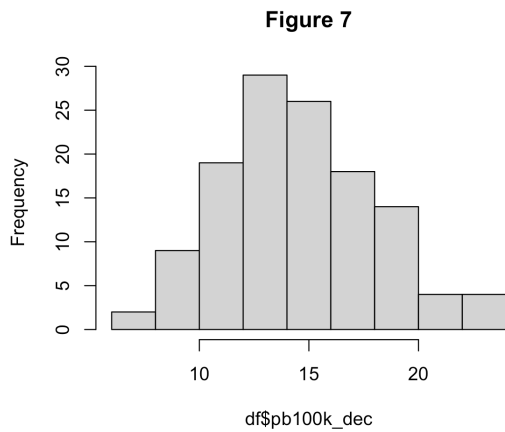
```
summary(df$age)
 Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 19.0    33.0    38.0    40.1    45.0    87.0
```

**Figure 5**



**Figure 6**

On the other hand, pb100k_dec seems to be regularly distributed, which is advantageous for regression modeling (Figures 7 and 8).

```
summary(df$pb100k_dec)
 Min. 1st Qu.  Median   Mean 3rd Qu.   Max.
 6.50   12.16   14.38  14.69   17.00  23.50
```

**Figure 7**



df$pb100k_dec

**Figure 8**



Fastest time to complete a 100km ultra marathon

However, pb_elev is noticeably right-skewed (Figures 9 and 10), with several high values on the top end. This skewness implies that any model that uses pb_elev must take this non-normality into account.

```
summary(df$pb_elev)
 Min. 1st Qu.  Median   Mean 3rd Qu.   Max.
    0     800    2200   2523    3658   9000
```
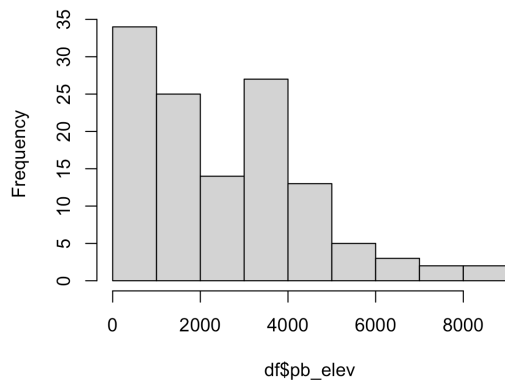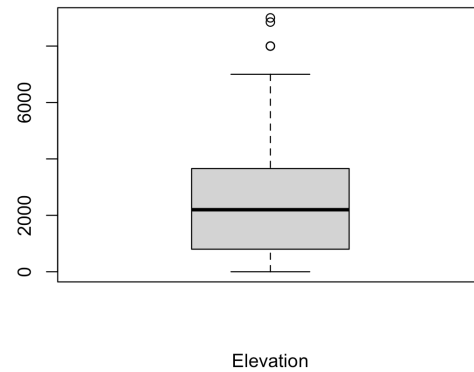
**Figure 9**

**Figure 10**

Although there are a few outliers at both ends, the average kilometers ran each week (avg_km) shows an approximately normal distribution (Figures 11 and 12).

```
> summary(df$avg_km)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
    5.0    56.0    70.0    72.1    80.0   160.0
```
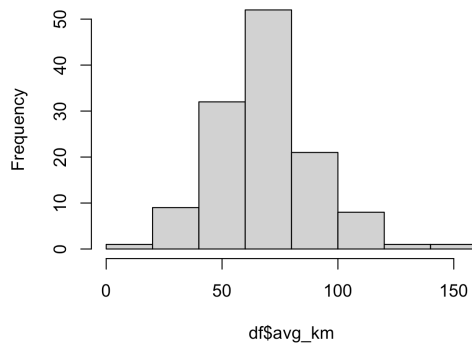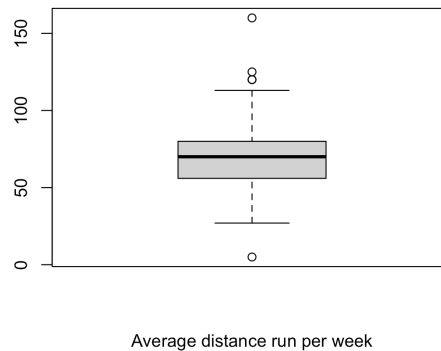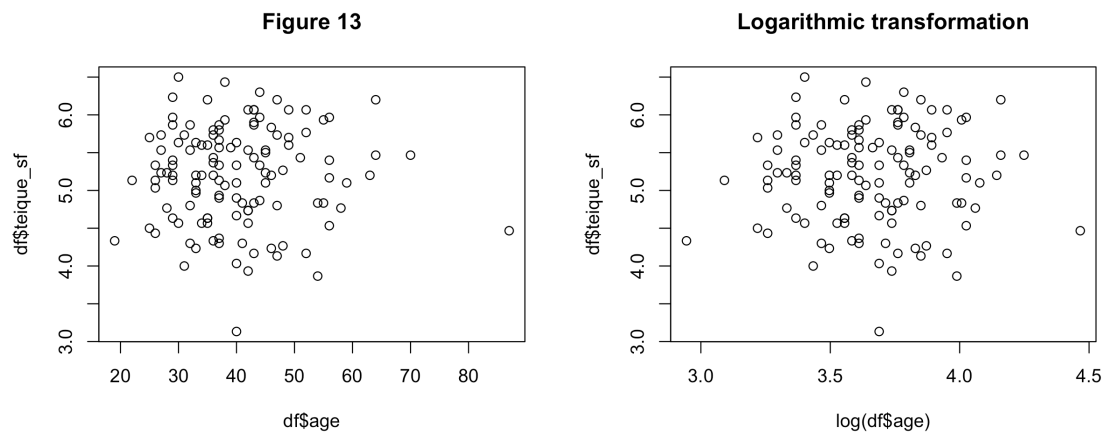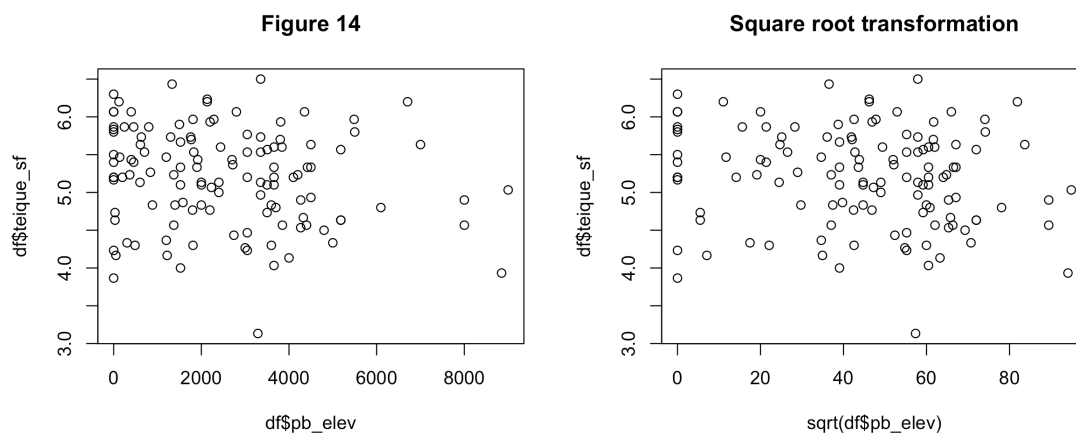
**Figure 11**

**Figure 12**

When comparing teique_sf to age in Figure 13, a single individual in their 80s jumps out as a possible outlier that may drag any fitted line, but there isn't a clear linear trend. The data span a

wide range of elevation values with several extreme points on the high end, and no

transformation significantly improved the linearity of the connection.

**Figure 13**     **Logarithmic transformation**



The data in Figure 14, which graphs teique_sf versus pb_elev spans a broad range of elevation

values, with several extreme spots on the high end. Applying a square root transformation seems

to provide a more balanced distribution and a little more linear scatter pattern, even if a square

root transformation did not significantly increase the linearity of the connection. This implies

that the square root of pb_elev could be better than the untransformed version, which might

improve residual diagnostics and model fit in subsequent studies.

**Figure 14**     **Square root transformation**

**Square root transformation on pb_elev**  **Square root transformation on pb_elev**


sqrt(df$pb_elev)


Elevation

Figure 15 shows a more pronounced linear trend between teique_sf and pb100k_dec, suggesting that faster timings may be linked to higher emotional intelligence. This variable also seems to be relatively normally distributed based on previous histograms. I tried doing a square root and logarithmic transformation in pb100k_dec to have more linearity, but the results produced a nearly similar scatter plot, indicating that it could be simpler to leave the variable untransformed.

**Figure 15**


df$pb100k_dec

**Square root transformation on pb_elev**   **Logarithmic transformation on pb_elev**



Figure 16, which displays Teique_sf vs. avg_km, shows a mostly linear trend with a few high-kilometrage outliers that might alter the slope. To obtain greater linearity as previously, I attempted a square root transformation, which just moved the data points to the right, and a logarithmic transformation, which tightened them even more while still shifting them to the right, indicating that neither transformation is helpful.

**Figure 16**

**Square root transformation on avg_km**  **Logarithmic transformation on avg_km**



The partial regression graphs for the model utilizing the untransformed pb_elev are displayed in

Figure 1. The age partial plot shows a largely horizontal band, but it is nevertheless impacted by

one older person who slightly pulls the fitted line. Elevation gain has a generally linear

relationship with teique_sf, though the residuals have a moderate dispersion around the fitted

line. Despite a few high-kilometrage data in avg_km, pb100k_dec and avg_km preserve

reasonably linear associations with teique_sf. Except for modest dispersion in the trail category,

the boxplots for the categorical variables pb_surface and sex demonstrate how each category

moves the response around the baseline level.



Component + Residual Plots

Following a square root transformation of pb_elev, the partial regression plots are shown in

Figure 2. A little tighter clustering of residuals and a wider spread is shown in the partial plot for

sqrt(pb_elev), indicating that this adjustment could enhance the model's capacity to represent the

impact of elevation gain on teique_sf.



No matter if I used the original pb_elev or its square root transformation, it seems from the

correlation matrices that none of the pairwise correlations between the numerical predictors have

values greater than 0.4. This implies that the regression model's collinearity is not likely to be a

significant issue. The overall pattern of low-to-moderate correlations suggests that the independent

variables are largely separate, even though there is a moderate negative correlation (around –0.33)

between pb100k_dec and avg_km. As a result, collinearity is most likely not an issue between the

independent variables.

```
> #Correlation between independent variables
> #No correlation is high
> cor(df[, -c(2, 3, 5, 8,9, 10)])
                 age     pb_elev  pb100k_dec      avg_km
age        1.0000000 -0.18024694 -0.1590007  0.20980204
pb_elev   -0.1802469  1.00000000  0.2593787  0.05380689
pb100k_dec -0.1590007  0.25937869  1.0000000 -0.32707374
avg_km     0.2098020  0.05380689 -0.3270737  1.00000000
> #Transformed pb_elev correlation between independent variables
> print(cor_transformed)
                  age sqrt_pb_elev pb100k_dec      avg_km
age         1.0000000  -0.18864963 -0.1590007  0.20980204
sqrt_pb_elev -0.1886496  1.00000000  0.2994052  0.04213896
pb100k_dec  -0.1590007  0.29940517  1.0000000 -0.32707374
avg_km       0.2098020  0.04213896 -0.3270737  1.00000000
```

Based on the patterns found in the exploratory study as well as domain knowledge, a few interactions now appear likely. For example, avg_km × sex may capture a possible variation in how weekly training volume effects trait across sexes, while age × sex may indicate if the effect of age on emotional intelligence varies between males and girls. Given that runners who contend with greater elevations could gain more, or less from increased training kilometers, sqrt_pb_elev × avg_km is another potential value. According to my exploratory investigation, the distributions seem acceptable following the necessary adjustments, and the predictors show primarily linear connections with the answer. A small number of outliers exist, but they don't appear to interfere with the general homoscedasticity or normalcy. As a result, at this point, no assumption is obviously broken, however diagnostics will be required for verification.

Based on theoretical considerations as well as findings from exploratory data analysis, the interaction term between sex and average weekly mileage (avg_km) was included. According to preliminary plots, male and female runners may have different relationships between training volume and trait emotional intelligence (teique_sf). Visualizations showed that the model would better represent these gender-specific differences if the interaction term was included. This

modification avoids assuming that the effect of avg_km is the same for both sexes. Therefore, by accounting for the differences in training effects between genders, the interaction term improves the model.

```
> final_model <- lm(teique_sf ~ age + I(sqrt(pb_elev)) + pb100k_dec + avg_km + pb_surface + sex + avg_km:sex, data = df)
> summary(final_model)

Call:
lm(formula = teique_sf ~ age + I(sqrt(pb_elev)) + pb100k_dec +
    avg_km + pb_surface + sex + avg_km:sex, data = df)

Residuals:
     Min      1Q  Median      3Q     Max
-1.99433 -0.39611  0.06311  0.41854  1.35548

Coefficients:
                   Estimate Std. Error t value Pr(>|t|)
(Intercept)       5.2726759  0.5114826  10.309   <2e-16 ***
age              -0.0003669  0.0056029  -0.065   0.9479
I(sqrt(pb_elev)) -0.0073620  0.0028354  -2.596   0.0106 *
pb100k_dec        0.0364524  0.0194031   1.879   0.0628 .
avg_km           -0.0056221  0.0048824  -1.151   0.2519
pb_surfacetrack  -0.3648192  0.4786327  -0.762   0.4475
pb_surfaceroad    0.0299353  0.2071706   0.144   0.8854
pb_surfacemix    -0.0736980  0.1788071  -0.412   0.6810
sexFemale        -0.5001546  0.4302743  -1.162   0.2475
avg_km:sexFemale  0.0100082  0.0058235   1.719   0.0884 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.6347 on 115 degrees of freedom
Multiple R-squared:  0.09603,   Adjusted R-squared:  0.02528
F-statistic: 1.357 on 9 and 115 DF,  p-value: 0.2157
```

Age, the square root transformation of elevation gain (sqrt(pb_elev)), sex, running surface (pb_surface), personal best time (pb100k_dec), and the avg_km:sex interaction term is included in the final model. It shows modest predictive accuracy with an in-sample Mean Squared Error (MSE) of 0.371 and an in-sample Mean Absolute Error (MAE) of 0.489. The model's R-squared of 9.6% and Adjusted R-squared of 2.5% indicate that the predictors only account for a minor amount of the variability in teique_sf. The error metrics show that the model's predictions are rather accurate on the sample data, even with this low explanatory power.

```
> cat("In-sample MAE:", mae_value, "\n")
In-sample MAE: 0.4894213
> cat("In-sample MSE:", mse_value, "\n")
In-sample MSE: 0.3705723
```

Using age, the square-root transformed total elevation gain (which indicates course difficulty), the fastest 100-kilometer time (pb100k_dec), the average weekly mileage (avg_km), the running surface (pb_surface), sex, and the interaction between avg_km and sex, our final model predicts trait emotional intelligence (teique_sf). In this model, a male runner on a track surface serves as the baseline, and the intercept shows the average degree of mental readiness for this population. A slower finishing time is somewhat related with greater mental readiness, but a larger square-root elevation increase is associated with a minor drop in teique_sf. Teique_sf is not much impacted by age, and variations in running surfaces result in negligible departures from the baseline. Increased training volume may help female runners' mental preparedness more than it does for male runners, according to the interplay between average weekly distance and sex.

Overall, the model explains around 9.6% of the variation in mental preparation, suggesting that this complex attribute is likely influenced by a wide range of other factors. In practical terms, factors like race performance and course difficulty are important, but they don't adequately account for an athlete's mental preparedness. The discovery that increased training mileage seems to support better mental readiness in female runners emphasizes how crucial it is to take gender variations in training methods into account. These revelations help us achieve our overarching objective of comprehending how psychological aspects of ultra-marathon performance interact with training routines and physical hurdles. To improve the model's prediction ability and find more variables that affect mental preparedness, more investigation is required.

According to the research, an ultra-runner's mental readiness is influenced by several training and performance aspects, although our dataset did not capture all of them. The need to

comprehend how physical demands and mental preparedness interact in endurance sports drove our research, and our findings show that factors like training volume, race time, and course difficulty (as determined by elevation gain) have quantifiable but slight effects, particularly when considering the differences between male and female runners. The intricacy of mental readiness is shown by the final model, which reveals that these variables only account for 9.6% of the variance in trait emotional intelligence. In other words, although our model explains certain elements of an ultra-runner's performance, it does not account for most factors that affect mental preparation. These results lend credence to the notion that future research should include variables not currently evaluated, such as psychological training, contextual circumstances, or other demographic traits.

This study has several significant drawbacks. First, the model's poor explanatory power suggests that many important factors influencing mental readiness are not taken into consideration. This might be because of the cross-sectional research design and the narrow range of data that is currently accessible. Furthermore, bias may be introduced using self-reported measures for psychological and performance qualities, and the sample size may restrict how broadly our findings may be applied. Longitudinal data collection, the inclusion of variables that capture other facets of mental preparation, and the investigation of more intricate modeling approaches might all lead to improvements. I gained important knowledge about data cleaning, model diagnostics, and the value of teamwork during this process, all of which will guide future studies targeted at enhancing endurance sports performance tactics.

| Project Team Member | Percentage of total effort |
|---|---|
|  |  |

| Victor Urdaneta | 100% |
|---|---|
|  |  |

Works Cited:

Samtleben, E. "Ultra-running the Upcoming Sport of the Endurance World: Is Emotional Intelligence Associated with Performance?" *Journal of Multidisciplinary Research at Trent*, vol. 3, no. 1, 2021, pp. 138–154. http://creativecommons.org/licenses/by-nc-sa/4.0/.https://causeweb.org/tshs/ultra-running/

"Ultra-Running and Emotional Intelligence Data Dictionary." 2021, http://creativecommons.org/licenses/by-nc-sa/4.0/. https://causeweb.org/tshs/datasets/ultrarunning_data_dictionary.pdf

R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, https://www.R-project.org/.