# On-Field Insights from Tokyo 2020 Olympic Swimming Events

**Ranice TAN Hui Qi**

# Contents

# 1. Overview

Sports analytics involves using data and statistical analytics to drive decision making. On-field analytics involves studying the performance, strategy or health and fitness of an athlete or team. The examination of these data can potentially enable athletes and coaches to make better decisions, allowing them to stand a competitive edge.

In this study, we plan to review data from the Tokyo 2020 Olympics data for Swimming events, to discover interesting on-field insights and use statistical analysis to proof these theories.

# 2. Objective

The objective of this study is to understand the factors affecting the reaction time of Olympic swimmers, as well as determine if reaction time affects the performance of the swimmers. The reaction time across all swimming events will be examined, to identify any relationship between the events, physical attributes or performance and reaction time.

JMP Pro 16 statistical analysis was used for data preparation, analysis and insights discovery using Interactive Data Exploration and Analysis.

# 3. Data

## 3.1. Data Used

The data used is the official result report, accessible from the official Tokyo 2020 Olympics webpage. The dataset fetched and combined the reports from all the Swimming events. The list of datasets can be found in Appendix B: Datasets.

## 3.2. Data Preparation

The dataset was imported into JMP Pro to ensure all fields are correctly filled, and all columns are appropriately formatted. A preliminary inspection of the data summary statistics and distribution was also conducted to remove entries with missing data, extreme outliers, and redundant data. The data preparation log is accessible in  Finally, a separate study can be done for relay team performance.

Appendix .

## 3.3. Data Quality

### 3.3.1. Inconsistent Formatting for 'Finals_Time' field

Finals_Time field contains inconsistent formatting for data <60 secs, and ≥ 60 secs. Hence, unable to recognise data and model type as numeric and continuous, and unable to sort the data by ascending order.



*Figure 1: Inconsistency in Formatting in Finals_Time field*

First, formula was written into new column Finals_Time_2 to standardise the formatting of all data to mm:ss, 2 decimal places.

```
If( Length( :Finals_Time ) < 6,
    :Finals_Time_2 = "00:" || :Finals_Time,
    :Finals_Time_2 = :Finals_Time
)
```

*Figure 2 Code to standardise formatting in Finals_Time_2 field*

Next, Finals_Time_2 was recoded to remove missing ('NA') data into new unlocked column Finals_Time_Formatted.



*Figure 3 Recoding of 'NA' data from Finals_Time_2 field*

### 3.3.2. Inaccurate Data and Modelling Types

The variables from the dataset which has incorrect data or modelling types are identified and corrected in the data table below. It is important to input the correct model type so that JMP can infer the appropriate statistic for analysis, without which, may result in poor and errors in analysis.

*Table 1 Changes in Data and Model Types for Variables*

| Variable Names | Default | | Change To | |
|---|---|---|---|---|
| | Data | Model | Data | Model |
| Place | Numeric | Continuous | Numeric | Ordinal |
| Lane | Numeric | Continuous | Numeric | Nominal |
| Finals_Time | Character | Nominal | Numeric | Continuous |
| DQ | Numeric | Continuous | Numeric | Nominal |
| Exhibition | Numeric | Continuous | Numeric | Nominal |
| Split_100, 150, 200, 300, 350, 400, 500, 550, 600, 700, 750, 800 | Character | Nominal | Numeric | Continuous |

3

### 3.3.3. Large numbers of categories for Categorical Data

Too many categorical variables may hinder analysis, hence, for variables important for analysis, some changes were made as follows:

*Table 2 Changes made to columns with large categorical variables*

| Columns | N Categories | Change |
|---|---|---|
| Name | 746 | Information not relevant |
| Team | 214 | Removed duplicate entries for teams e.g. BRA and BRA-Brazil. |
| Event | 35 | Recoded into Gender, Stroke and Distance columns, of 3, 7 and 8 levels respectively. |
| Relay_Swimmer_1 | 125 | Information not relevant |
| Relay_Swimmer_2 | 116 | Information not relevant |
| Relay_Swimmer_3 | 119 | Information not relevant |
| Relay_Swimmer_4 | 110 | Information not relevant |



*Figure 4 Summary Statistic for Categorical Columns with large levels*



*Figure 5 Recoding from Event field to Gender, Stroke and Distance fields*

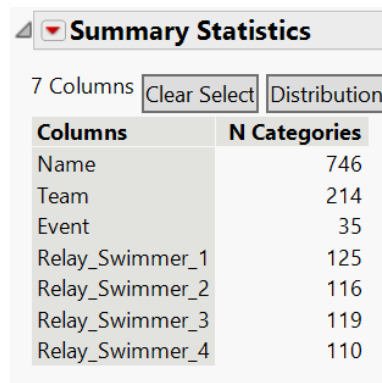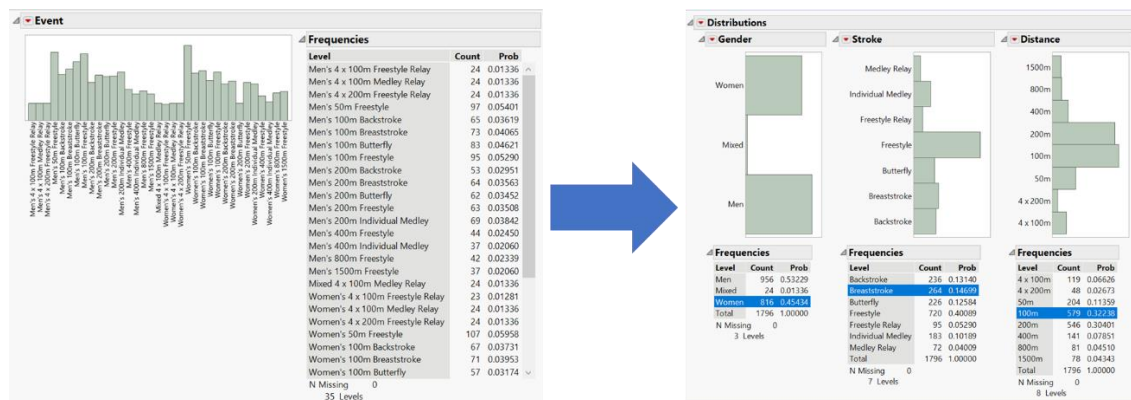*Figure 6 Team categorical data contains duplicate entries*

### 3.3.4. Incorrect Information Pre-Filled

2 rows by swimmer Rodriguez Villanueva were pre-filled with incorrect information and manually corrected as follows:

| Place | Heat | Lane | Name | Team | Reaction_Time | Finals_Time |
|---|---|---|---|---|---|---|
| 4 | Heat_1 | 3 | RODRIGUEZ VILLANUEVA Byanca Melissa MEX | 0.66 | 33.19 | 0:01.84 |
| 6 | Heat_3 | 7 | RODRIGUEZ VILLANUEVA Byanca Melissa MEX | 0.65 | • | 0:01.94 |

| Place | Heat | Lane | Name | Team | Reaction_Time | Finals_Time |
|---|---|---|---|---|---|---|
| 6 | Heat_3 | 7 | RODRIGUEZ VILLANUEVA Byanca Melissa | MEX | 0.65 | 1:08.76 |
| 4 | Heat_1 | 3 | RODRIGUEZ VILLANUEVA Byanca Melissa | MEX | 0.66 | 2:26.82 |

*Figure 7 Manual correction of incorrect data entry*

### 3.3.5. Missing Data

192 out of 1796 of reaction time data is missing from disqualified swimmers and relay events. Hence, limiting study to individual performances only.



*Figure 8 Summary Statistics for Reaction_Time Column*

### 3.3.6. Inconsistent definition for split timing data between events

Definition for some split time is inconsistent throughout the events, hence may lead to confusion if the data is used for calculations or analysis.

*Table 3 Inconsistencies in Split time data*

| Event | Inconsistencies |
|---|---|
| Individual Events | Time taken for each 50m interval i.e. Split_50: 0-50m, Split _100: 50m-100m |
| 4 X 100m Relay Events | Time taken for each swimmer's 1st 50m interval, and 1st 100m interval, with some entries missing data. i.e. Split_50: 0-50m, Split _100: 0m-100m |
| 4 X 200m Relay Events | Similar to "4 X 100m Relay Events", but missing data for 100-200m intervals. |

5

| Event | Split_50 | Split_100 | Split_150 | Split_200 | Split_250 | Split_300 | Split_350 | Split_400 | Split_450 |
|---|---|---|---|---|---|---|---|---|---|
| Women's 4 x 100m Medley Relay | 29.17 | • | 31.85 | • | 26.43 | 57.67 | 25.78 | 54.89 | • |
| Women's 4 x 100m Medley Relay | 30.29 | • | 31.81 | • | 27.13 | 58.99 | 26.09 | 54.43 | • |
| Women's 4 x 200m Freestyle Relay | 26.99 | 56.02 | | • | 26.88 | 56.81 | | • | 26.81 |
| Women's 4 x 200m Freestyle Relay | 27.7 | 56.7 | | • | 27 | 56.17 | | • | 26.81 |

*Figure 9 Examples of Relay Events with Definition Discrepancies in Split Timing*

### 3.3.7. Outliers in Reaction Time data

Reaction Time data contains 7 outlier points, which accounts for <0.5% of data, hence may not drastically impact the analysis when removed.



*Figure 10 Statistics of Reaction Time Data showing few outliers*

### 3.3.8. Missing parameters

Some parameters that are useful to the study and are missing include time of event, and age of swimmer. Time – morning, afternoon and night may be used to gauge the alertness of the swimmer. Age data may be used to analyse any correlation between reaction time or performance to identify peak performance for an athlete.

# 4. Data Analysis

## 4.1. General Insights

### 4.1.1. Analysis 1: Olympic swimmers have similar mean and median reaction time.

The mean and median reaction times for Olympic swimmers is 0.66s and 0.66 respectively.



*Figure 11 Statistics of Reaction Time*

### 4.1.2. Analysis 2: Lanes 3-4 swimmers have a lower reaction time than swimmers in the other lanes

The reaction time of swimmers in Lanes 3, 4 have a median of 0.65s and their means are < 0.66s, whereas swimmers in the other lanes have medians of 0.66s and means > 0.66s. Hence, swimmers in Lanes 3 and 4 seem to be able to react faster than the other lanes.



*Figure 12 Boxplot of Reaction Time across Lanes*

### 4.1.2.1. Hypothesis 1: Swimmers in Lanes 3 and 4 react faster than the swimmers in the other lanes

We are interested to test if the swimmers in lanes 3 and 4 are truly faster in reacting then the swimmers in the other lanes. A normality test was conducted to determine the appropriate hypothesis testing method. In this hypothesis, we will use Shapiro-Wilk and Anderson-Darling Goodness-of-Fit Test with confidence level of 99% to test:

- $H_0$: Reaction Time is normally distributed
- $H_1$: Reaction Time is not normally distributed

*Figure 13 Normal Distribution Test for Reaction Times across Swim lanes*

Based on the results above, we accept that reaction time for all lanes is normally distributed as smallest p-value 0.0480 (Lane 6) > critical value 0.01. Hence, the hypothesis testing will be conducted using parametric ANOVA and Welch, with 99% confidence level. Our hypothesis is defined as:

- $H_0$: Reaction Time is same for all lanes
- $H_1$: Reaction Time is not the same for all lanes



*Figure 14 Welch's Test for Reaction Time across Lanes*



*Figure 15 Pair comparison tests using student's t and Tukey-Kramer method*

The results have a p-values (smallest 0.0383) greater than critical value of 0.01. This statistical evidence can be used to accept the null hypothesis and conclude that the reaction time is consistent across all lanes. Hence, the hypothesis Lanes 3 and 4 swimmers react faster is not true.

## 4.2. Insights from Gender

### 4.2.1. Analysis 3: Male swimmers react faster than female swimmers

Across all events, men had lower mean and median reaction times than women. Women swimmers also demonstrated higher variability in their reaction time, as shown by the higher standard deviation.

*Table 4 Reaction Times Statistics across Gender*

| Gender | Reaction Times | | | | |
|--------|------|------|--------|---------|------|
| | N | Mean | Median | Std Dev | IQR |
| Men | 872 | 0.645 | 0.64 | 0.053 | 0.07 |
| Women | 732 | 0.680 | 0.68 | 0.057 | 0.07 |



*Figure 16 Boxplot of Reaction Time across Genders*

#### 4.2.1.1. Hypothesis 2: Males have faster reaction time than Female swimmers

We are interested to test if there is evidence to suggest that the reaction times are different among male and female swimmers. A normality test was conducted to determine the appropriate hypothesis testing method. In this hypothesis, we will use Shapiro-Wilk and Anderson-Darling Goodness-of-Fit Test with confidence level of 95% to test:

- $H_0$: Reaction Time is normally distributed
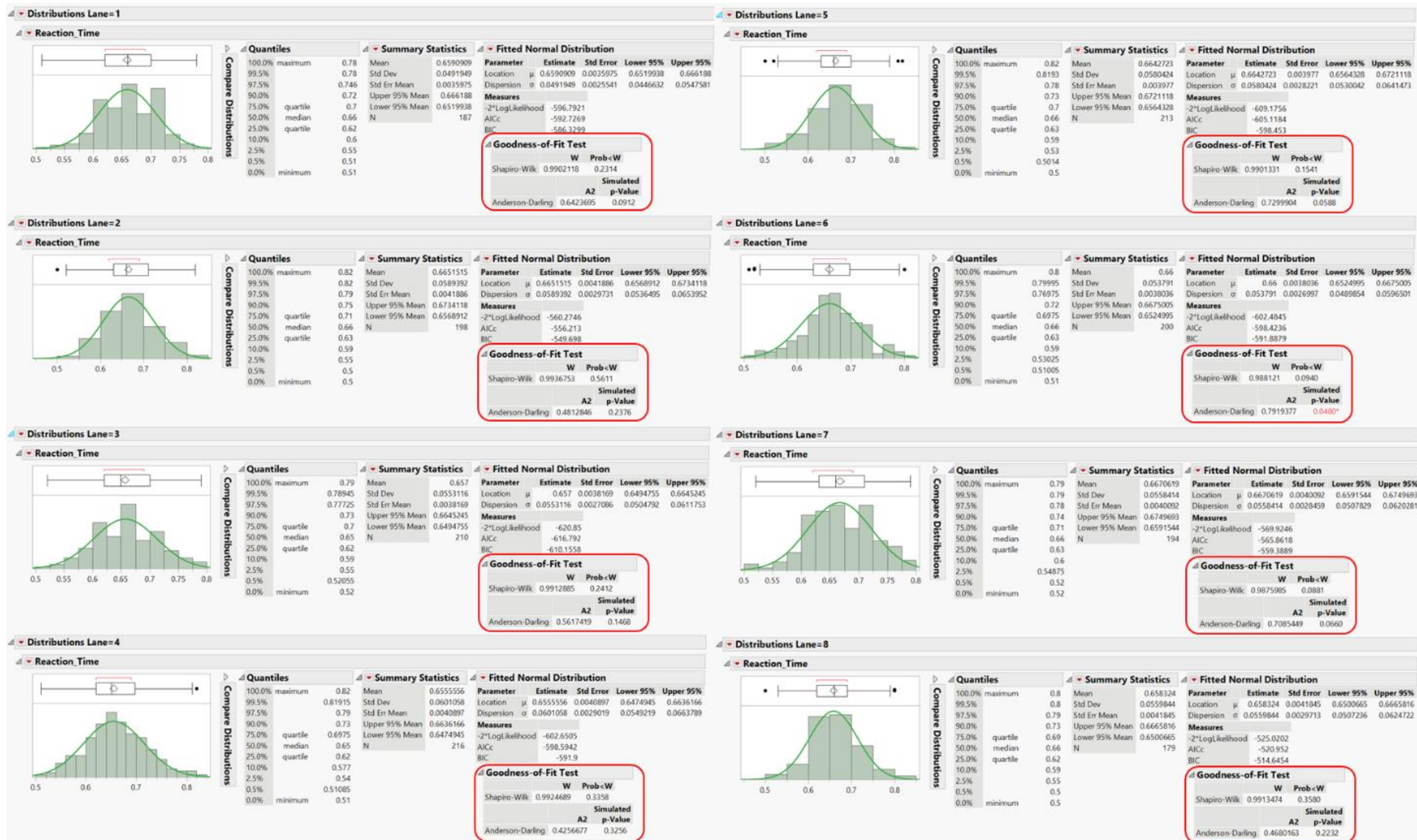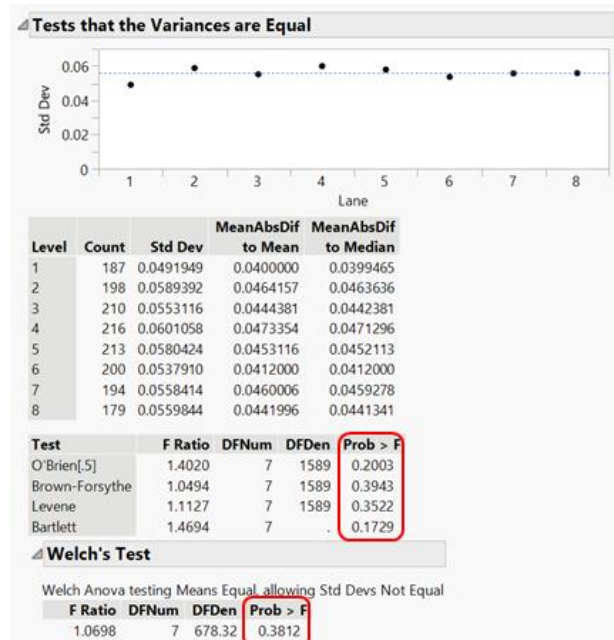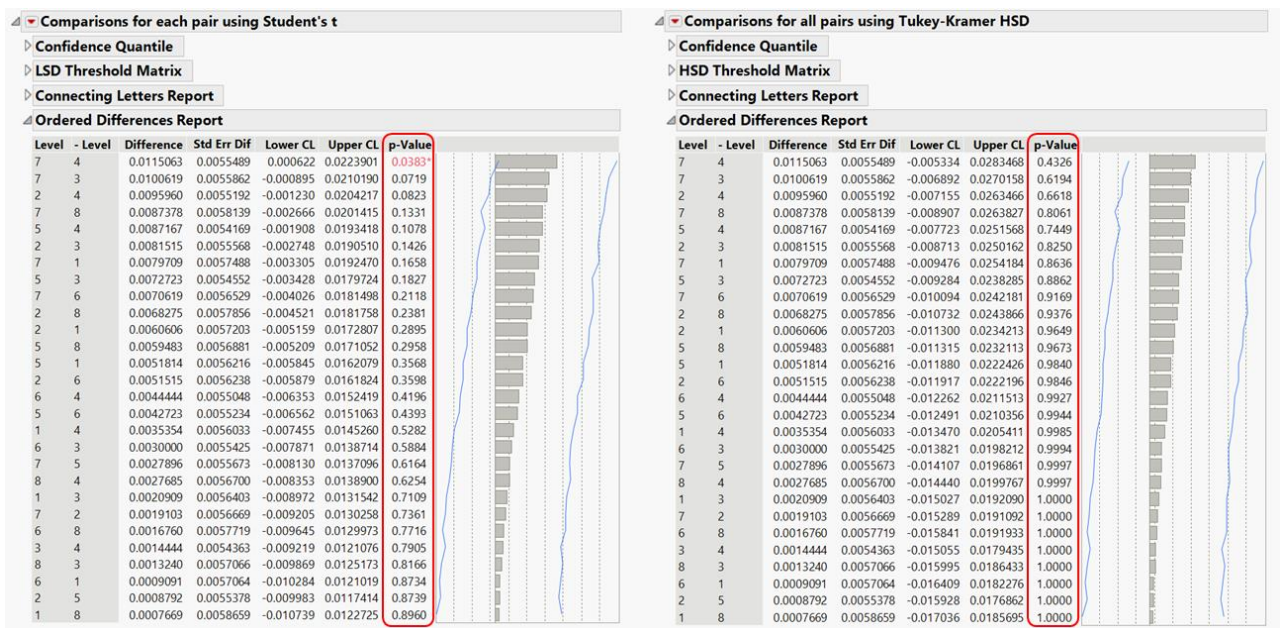- $H_1$: Reaction Time is not normally distributed

*Figure 17 Normal distribution test for Reaction Time across Genders*

Based on the results above, we reject that reaction time for all lanes is normally distributed as p-values < 0.0001 which are less than critical value of 0.05. Hence, the hypothesis testing will be conducted using non-parametric Kruskal-Wallis test, with 95% confidence level. Our hypothesis is defined as:

- $H_0$: Reaction Time is same for both genders
- $H_1$: Reaction Time is not the same for both genders



*Figure 18 Kruskal-Wallis test for Reaction Time across Genders*

The results have a p-values less than critical value of 0.01. We have evidence to reject the the null hypothesis and conclude that the reaction time is consistent not the same for both genders. Hence, the hypothesis that male swimmers have faster reaction time than female swimmers it True.

11

### 4.2.2. Analysis 4: Most of the male swimmers are faster than female swimmers for Freestyle 200m event.

From the boxplot below, male swimmers are about 11s faster than their female counterparts for the Freestyle 200m race. Other than 1 outlier, all other male swimmers are faster than the fastest female swimmer. Women freestyle swimmers have higher variability in their final timing. The distribution of final time for both genders are skewed right, with the means higher than the median.
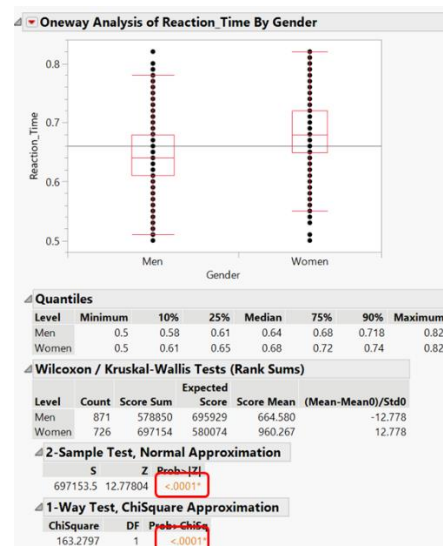


*Figure 19 Boxplot of Final Time across Gender for 200m Freestyle Event*

## 4.3. Insights from Events

### 4.3.1. Analysis 5: Longer distance events tend to have longer reaction times

Events with shorter distances (<400m) tend to have lower mean and median reaction times than the longer distance events (≥ 400m).

*Table 5  Statistics of reaction time across Event's distance*

| Distance | Reaction Times | | | | |
|---|---|---|---|---|---|
| | N | Mean | Median | Std Dev | IQR |
| 50m | 202 | 0.667 | 0.66 | 0.055 | 0.07 |
| 100m | 562 | 0.642 | 0.64 | 0.052 | 0.07 |
| 200m | 543 | 0.653 | 0.66 | 0.052 | 0.07 |
| 400m | 140 | 0.699 | 0.70 | 0.050 | 0.07 |
| 800m | 79 | 0.716 | 0.72 | 0.047 | 0.07 |
| 1500m | 77 | 0.717 | 0.72 | 0.059 | 0.07 |

*Figure 20 Boxplot of Reaction Time across Distance*

### 4.3.2. Analysis 6: Backstroke swimmers react faster than the other swimmers

Backstroke swimmers have the lowest mean and median reaction times, which indicates they have the fastest reaction time. Breaststroke reaction time distribution is the narrowest inter-quartile range and variability, which may represent more consistency in their reaction.

*Table 6 Statistics of Reaction Time Across Event's Stokes*

| Stroke | Reaction Times | | | | |
|---|---|---|---|---|---|
| | N | Mean | Median | Std Dev | IQR |
| Backstroke | 232 | 0.598 | 0.60 | 0.053 | 0.07 |
| Breaststroke | 255 | 0.667 | 0.67 | 0.038 | 0.05 |
| Butterfly | 222 | 0.653 | 0.65 | 0.047 | 0.07 |
| Freestyle | 713 | 0.680 | 0.68 | 0.055 | 0.08 |
| Individual Medley | 182 | 0.669 | 0.66 | 0.050 | 0.06 |


*Figure 21 Boxplot of Reaction Time across Strokes*

13

### 4.3.3. Analysis 7: Freestyle swimmers are the fastest swimmers

Freestyle swimmers have the lowest mean and median final timings, indicating they are faster than the other strokes. Despite backstroke swimmers having the fastest reaction time, they end up with mid-range final times.



*Figure 22 Boxplot of Final Times across Strokes (Limit Distance = 200m)*

### 4.3.4. Analysis 8: Swimmers in Semis event tend to have lower reaction times

Semis event swimmers have the lowest mean and median reaction time of 0.65s, showing that swimmers tend to react slightly faster for the semi events, compared to heats and finals. An unexpected observation is that the reaction time for the finals event has the widest IQR and does not have outliers.



*Figure 23 Boxplot of Reaction Time across Heat types*

14

## 4.4. Insights from Performance

### 4.4.1. Analysis 9: Reaction time does not have much impact on a swimmer's placings

The boxplot below shows reaction time fluctuates across swimmers for all placings, with 1st-3rd place swimmers having similar reaction times to 7th-8th place. This may mean that the reaction time may not have a large impact on a swimmer' performance. Though, in the finals event (2nd boxplot), Gold medalists have the lowest mean (0.649s) and median (0.65s) reaction times, they also have a moderate IQR.



*Figure 24 Boxplot of Reaction Time across Placing for Heats, Semis and Finals*
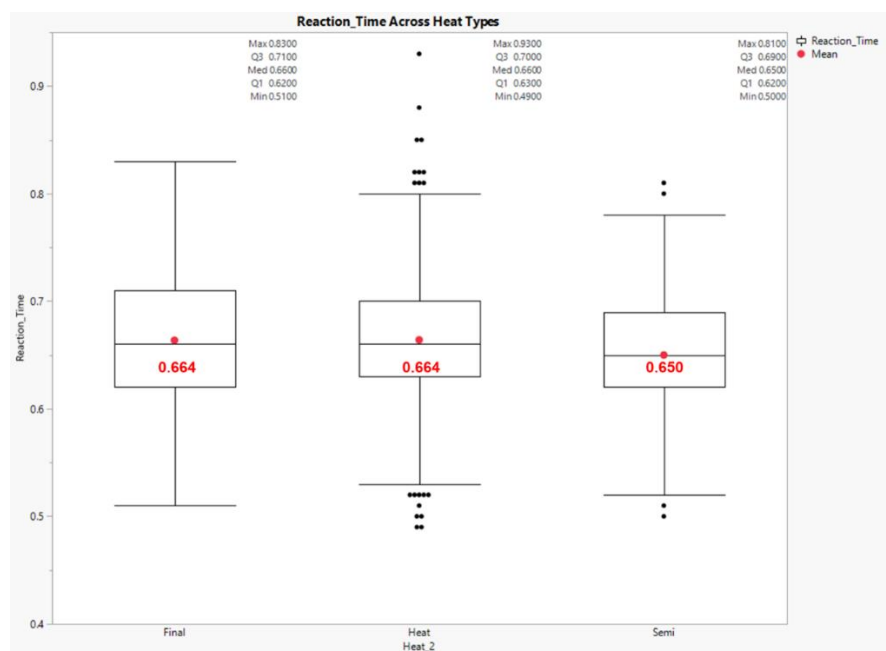


*Figure 25 Boxplot of Reaction Time across Placings for Finals events only*

#### 4.4.1.1. Hypothesis 3: Swimmers with lower reaction times have better performance

We are interested to test if reaction time impacts the performance of the swimmers. A normality test was conducted to determine the appropriate hypothesis testing method. We will use Goodness-of-Fit Test with confidence level of 99% to test:

- $H_0$: Reaction Time is normally distributed
- $H_1$: Reaction Time is not normally distributed

*Figure 26 Normality Test for Reaction Time across Placings*

Based on the results above, we accept that reaction time for all placings is normally distributed as smallest p-value 0.0356 (5$^{th}$ place) > critical value 0.01. Hence, the hypothesis testing will be conducted using parametric Welch test, then compare using student's t and Tukey-Kramer with 99% confidence level. Our hypothesis is defined as:

- $H_0$: Reaction Time is same for all performance
- $H_1$: Reaction Time is not the same for all performance



Figure 27 Welch's test for Reaction Time across performances



Figure 28 Student's t and Tukey-Kramer pair test for Reaction Time across placings

The results have a p-values (smallest 0.2142) greater than critical value of 0.01. This statistical evidence can be used to accept the null hypothesis and conclude that the reaction time is consistent across placings. Hence, the hypothesis that swimmers with better performance have lower reaction time is false. This evidence suggests that reaction time does not affect the performance of the swimmers to a large extent.

### 4.4.2. Analysis 10: Swimmers in the centre lanes tend to have higher place finishes

The dot plot below shows that swimmers with better performance seem to be placed in the centre Lanes 3-4, with lane 4 having the highest concentration of Gold medalists. Most podium finishers (1st, 2nd, 3rd) were also competing in Lanes 3-5.



*Figure 29 Dotplot of Swimmer's Placing vs Lane*

### 4.4.3. Analysis 11: 96% of Olympic finalists in Lane 4 end up with podium finishes

The heatmap below shows that of the swimmers who compete in Lane 4 finals, 96% end up in the Top 3, and receive medals. There is also a high proportion of Olympic medalists in Lanes 3-5, where >50% of the finalists end up with podium finishes.



*Figure 30 Heatmap of Medalists across Lanes for Olympic Swimming Finalists*

### 4.4.4. Analysis 12: There is little correlation between the Final Time and Reaction Time

The regression between the Final Time and Reaction Time for different events, with distance of 100m, has poor line of fit. The highest $R^2$ is 0.283, from the Butterfly 100m event.

*Figure 31 Regression of Final Time and Reaction Time for Event Distance = 100m*

### 4.4.4.1. Hypothesis 4: Final Time increases as Reaction Time increases

In this hypothesis testing, we shortlist the Freestyle 100m event data, to test if there is any relationship between the Final Time (Performance) and Reaction Time of the swimmers. Bivariate fit and analysis, with a 95% confidence level is used to test the hypothesis:

$H_0$: True slope coefficient is 0.
$H_1$: True slope coefficient is not 0.



*Figure 32 Bivariate Fit Test for Final Time and Reaction Time for Freestyle 100m Event*

Based on the test result, Final Time and Reaction Time has a p-value (<0.001) less than critical value. Hence, the null hypothesis is rejected. Since the variables are positively correlated ($R^2$ = 0.15), we can conclude that the Final Time increases as Reaction Time increases, and the hypothesis is true.

# 5. Interpretation of Analysis Results

## 5.1. Male swimmers have faster reaction time than female swimmers

Male swimmers have a mean and median reaction time of 0.645s and 0.64s respectively, whereas female swimmers have a mean and median reaction time of 0.680s and 0.68s respectively (Analysis 3: Male swimmers react faster than female swimmers. This is further reiterated in hypothesis testing, whereby the null hypothesis of male and female swimmers having the same reaction time was rejected. (Hypothesis 2: Males have faster reaction time than Female swimmers). This may be due to physical attributes of the human body.

## 5.2. Reaction time increases across events with longer distances

Swimmers who participate in events with shorter distances (<400m) tend to have lower mean and median reaction 0.64s – 0.67s, as compared to the longer distance events (≥ 400m) of 0.70s – 0.72s Analysis 5: Longer distance events tend to have longer reaction times. This may be because there are other factors that may affect performance for longer distance events.

## 5.3. Backstroke swimmers have the fastest reaction time but does not final time

Backstroke swimmers have the lowest mean and median reaction times of 0.598s and 0.60s respectively Analysis 6: Backstroke swimmers react faster than the other swimmers . However, amongst all strokes in 200m events, they only have mid-range mean and median final times Analysis 7: Freestyle swimmers are the fastest swimmers. This shows that there are other stroke-related factors impact a swimmers' performance.

## 5.4. Although the centre lanes tend to churn out more medallists and higher placed finishers, the reaction time across all lanes are consistent

A higher promotion of well performing athletes tend to stem from the center Lanes 3-5 (Analysis 10: Swimmers in the centre lanes tend to have higher place finishes), even with > 50% of finalists in Lanes 3-5 ending up with podium finishes (Analysis 11: 96% of Olympic finalists in Lane 4 end up with podium finishes). However, the hypothesis testing shows that the athletes in the centre do not have faster reaction time, it is consistent across all lanes 1-8 Hypothesis 1: Swimmers in Lanes 3 and 4 react faster than the swimmers in the other lanes. Hence, the reason centre lane swimmers perform better is unlikely due to better reaction time.

## 5.5. Reaction time does not significantly affect the performance of the swimmers

From the 100m Freestyle event analysis, the Final Time and Reaction Time of a swimmer is positively correlated but has a poor cor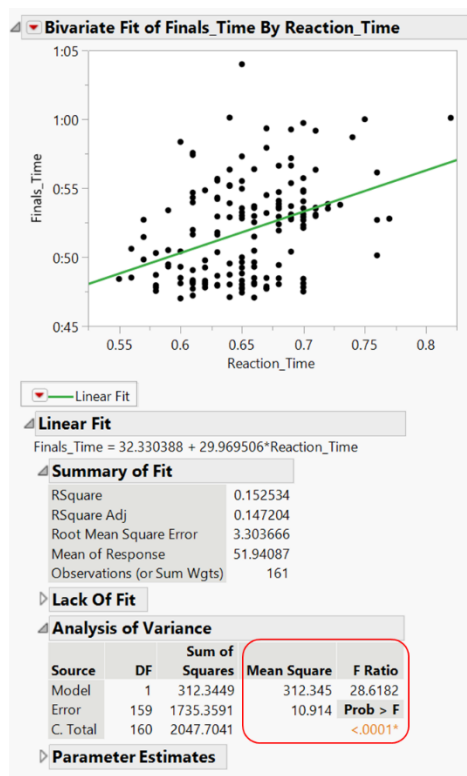relation value and line of fit (Analysis 12: There is little correlation between the Final Time and Reaction Time, Hypothesis 4: Final Time increases as Reaction Time increases).

Reaction time also fluctuates across swimmers for all placings, with 1st-3rd place swimmers having similar reaction times to 7th-8th place Analysis 9: Reaction time does not have much impact on a swimmer's placings. The hypothesis testing also provides statistical evidence that swimmers with better placing have lower reaction time is false Hypothesis 3: Swimmers with lower reaction times have better performance. This evidence suggests that reaction time does not impact the performance of the swimmers to a large extent.

# 6. Recommendation

The results from this study show that generally, the impact of reaction time on the performance of an Olympic swimmer is minor. However, that does not mean that the impact of the reaction time can be entirely neglected, it is dependent on individual performance and other factors. The recommendation for an athlete/coaches would be to compare their individual performance across the data in this study to determine if there are room of improvement.

To improve the study, a suggestion is to widen the scope to collect and analyse other factors like age, start time, stroke distance, stroke speed, turning time and its impact on performance. Finally, a separate study can be done for relay team performance.

# 7. Appendix

## 7.1. Appendix A: Data Preparation Change Log

Dataset: Swimming

| Item | Variable Name | Issue | Action |
|------|---------------|-------|--------|
| 1 | Place | Should not be continuous modelling type. | Change to nominal modelling type. |
| 2 | Lane | Should not be continuous modelling type. | Change to nominal modelling type. |
| 3 | DQ | Should not be continuous modelling type. | Change to nominal modelling type. |
| 4 | Split_100, 150, 200, 300, 350, 400, 500, 550, 600, 700, 750, 800 | Should not be in character data and nominal modelling type. | Change to numeric data type and continuous modelling type. |
| 5 | Finals_Time | Inconsistent format for data < 1 min (ss) and data > 1min (min:ss). | Add in new column with formula, Finals_Time_2 to standardise data format (min:ss, 2 decimal place). |
| 6 | Finals_Time_2 | Column is locked and contains 'NA' information for disqualified swimmers. | Recode into new unlocked column, Finals_Time_recoded and remove 'NA' data. |
| 7 | Finals_Time_recoded | Should not be in character data, and nominal modelling type. | Change to numeric data, and continuous modelling type. |
| 8 | Finals_Time | Intermediate variable and cannot be used for analysis. | Change naming to Finals_Time_1 |
| 9 | Finals_Time_recoded | Final variable to be used for analysis. | Change naming to Finals_Time |
| 10 | Finals_Time Finals_Time_1 Finals_Time_2 | Duplicate columns | Hide and exclude columns for Finals_Time_1 and Finals_Time_2. |
| 11 | Event | Contains too many categories, will be difficult for analysis. | Recode and categorised into 'Gender', 'Distance' and 'Stroke' columns. |
| 12 | Heat | Contains too many categories, will be difficult for analysis. | Recode into 'Heat_Type', categorised into 'Heat', 'Semi' and 'Final' |
| 13 | Team | Field data contains duplicate naming for some teams i.e. 'AUS' and 'AUS-Australia' | Recode and standardise naming. |

| 14 | Name, Team, Reaction_Time, Finals_Time, Split | 2 rows by swimmer 'VILLANUEVA Byanca Melissa' contains data incorrectly | Refer to raw pdf to update the fields. |
|---|---|---|---|
| 15 | Reaction_Time | Contains missing data for disqualified swimmers and relay events. | Exclude rows with missing Reaction_Time data. |
| 16 | Exhibition | Does not contain much information, all rows showing same value. Data is redundant. | Hide and exclude exhibition field. |
| 17 | Split | Too many columns for split timing. | Hide columns |
| 18 | Relay_Swimmer Relay_Swimmer_Gender | Since relay events are missing reaction time data and not included in this analysis, fields are redundant. | Hide and exclude Relay_Swimmer and Relay_Swimmer_Gender fields |

## 7.2. Appendix B: Datasets

| Item | Dataset Name | Brief Description | Usage |
|---|---|---|---|
| 1 | swimming_cleaned | Main clean dataset | Used for majority of reaction time analysis, that does not need |
| 2 | swimming_cleaned_reaction and final time statistics | Statistical mean and median data extracted from main swimming_cleaned dataset. | Used for brief comparison and extraction of numbers. |
| 3 | swimming_cleaned_100m Freestyle | Extracted only 100m Freestyle data from main swimming_cleaned data set | Used for bivariate analysis between Final Time and Reaction Time |