

```
In [ ]: library(readr)
```

```
In [ ]: df <- read_csv("/content/combined_ecg_data.csv")
```

```
Rows: 140000 Columns: 13
— Column specification —
Delimiter: ","
dbl (13): V1, V2, V3, V4, V5, V6, V7, V8, V9, V10, V11, V12, Athlete

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
In [ ]: head(df)
```

A tibble: 6 × 13

V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	Athlete
<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
10251	-1096	-10267	-3724	9391	-5395	13580	11410	14721	16103	6662	-3806	1
8643	-2558	-10829	-1862	8973	-6448	13331	11096	14093	15416	5897	-4548	1
5427	-3776	-9985	743	6469	-6711	13829	10991	13644	14815	5460	-5105	1
5427	-4507	-10829	1116	7304	-7501	13331	10467	13016	14128	4696	-5848	1
6231	-4751	-11673	1116	8138	-8027	13331	10049	12477	13355	3822	-6962	1
5427	-5481	-11954	1860	8138	-8554	12832	9630	11938	12926	3385	-7704	1

Twenty-eight healthy athletes were recruited for this study. 19 (68%) of the participants were men and 9 (32%) were women. Participant's ages ranged from 20 to 43 years (Mean = 25 years, standard deviation = 4.7 years). The distribution among sports was 24 rowers (86%), 2 kayakers (7%) and 2 cyclists (7%). The average amount of training hours for 2017 was 822 hours with a standard deviation of 117 hours, in 2018 the average amount of training was 820 hours with a standard deviation of 113 hours and in 2019 the average amount of training was 798 hours with a standard deviation of 171 hours.

```
In [ ]: # Number of observations per athlete
obs_per_athlete <- nrow(df) / 28

# Create a vector for Gender, repeating "Male" and "Female" the necessary number of times
gender_vector <- rep(c(rep("Male", obs_per_athlete * 19), rep("Female", obs_per_athlete * 9)), times = 1)

# Assign the gender vector to the dataframe
df$Gender <- gender_vector

# Similarly for Sport, adjust the numbers as per your distribution
sport_vector <- rep(c(rep("Rowing", obs_per_athlete * 24), rep("Kayaking", obs_per_athlete * 2), rep("Cycling",
df$Sport <- sport_vector
```

```
In [ ]: # Assuming a normal distribution and that df has a row for each athlete for each year
df$Training2017 <- rnorm(nrow(df), mean = 822, sd = 117)
df$Training2018 <- rnorm(nrow(df), mean = 820, sd = 113)
df$Training2019 <- rnorm(nrow(df), mean = 798, sd = 171)

# Then, calculate an average of these simulated values for a more nuanced estimate
df$Simulated_Avg_Training <- rowMeans(df[, c('Training2017', 'Training2018', 'Training2019')])

# View the first few rows to verify
head(df,10)
```

A tibble: 10 × 19

V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	Athlete	Gender	Sport	Training2017	Training2018
<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<chr>	<dbl>	<dbl>
10251	-1096	-10267	-3724	9391	-5395	13580	11410	14721	16103	6662	-3806	1	Male	Rowing	935.3347	903.6270
8643	-2558	-10829	-1862	8973	-6448	13331	11096	14093	15416	5897	-4548	1	Male	Rowing	686.9353	939.8613
5427	-3776	-9985	743	6469	-6711	13829	10991	13644	14815	5460	-5105	1	Male	Rowing	826.7363	843.8690
5427	-4507	-10829	1116	7304	-7501	13331	10467	13016	14128	4696	-5848	1	Male	Rowing	781.6344	921.4977
6231	-4751	-11673	1116	8138	-8027	13331	10049	12477	13355	3822	-6962	1	Male	Rowing	752.8014	845.6427
5427	-5481	-11954	1860	8138	-8554	12832	9630	11938	12926	3385	-7704	1	Male	Rowing	811.6301	947.5145
4623	-5969	-11954	2605	7721	-8817	12832	9316	11669	12410	2730	-8076	1	Male	Rowing	721.2272	724.3908
4221	-6943	-12517	3722	7721	-9606	12583	9107	11400	11981	2074	-9004	1	Male	Rowing	824.9267	749.8709
3416	-7187	-12235	4467	6886	-9606	12334	8897	11131	11723	1637	-9190	1	Male	Rowing	820.1125	763.6142
4221	-7430	-13079	4094	8138	-10133	12085	8792	11041	11637	1419	-9375	1	Male	Rowing	819.8759	929.0274

```
In [ ]: set.seed(42) # For reproducibility
```

```
# Generate ages based on the provided distribution
ages <- rnorm(n = 28, mean = 25, sd = 4.7)

# Round ages, ensure they are within the specified range
ages <- round(ages)
ages <- ifelse(ages < 20, 20, ages)
ages <- ifelse(ages > 43, 43, ages)

# Create a dataframe for athletes and their ages
athletes_ages <- data.frame(Athlete = 1:28, Age = ages)

# Merge the age information with the main dataset
df<- merge(df, athletes_ages, by = "Athlete", all.x = TRUE)

# Display the first few rows to check
head(df)
```

A data.frame: 6 × 20

	Athlete	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	Gender	Sport	Training2017	Training201
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<chr>	<dbl>	<dbl>
1	1	10251	-1096	-10267	-3724	9391	-5395	13580	11410	14721	16103	6662	-3806	Male	Rowing	935.3347	903.627
2	1	8643	-2558	-10829	-1862	8973	-6448	13331	11096	14093	15416	5897	-4548	Male	Rowing	686.9353	939.861
3	1	5427	-3776	-9985	743	6469	-6711	13829	10991	13644	14815	5460	-5105	Male	Rowing	826.7363	843.869
4	1	5427	-4507	-10829	1116	7304	-7501	13331	10467	13016	14128	4696	-5848	Male	Rowing	781.6344	921.497
5	1	6231	-4751	-11673	1116	8138	-8027	13331	10049	12477	13355	3822	-6962	Male	Rowing	752.8014	845.642
6	1	5427	-5481	-11954	1860	8138	-8554	12832	9630	11938	12926	3385	-7704	Male	Rowing	811.6301	947.514

```
In [ ]: library(dplyr)
library(ggplot2)
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

filter, lag

The following objects are masked from 'package:base':

intersect, setdiff, setequal, union

```
In [ ]: library(tidyr)
```

```
In [ ]: summary(df)
```

Athlete	V1	V2	V3
Min. : 1.00	Min. :-32767	Min. :-32767	Min. :-32767
1st Qu.: 7.75	1st Qu.: -15380	1st Qu.: -23405	1st Qu.: -24139
Median :14.50	Median : -5173	Median : -19200	Median : -20534
Mean :14.50	Mean : -5356	Mean : -17264	Mean : -17550
3rd Qu.:21.25	3rd Qu.: 2285	3rd Qu.: -13034	3rd Qu.: -14862
Max. :28.00	Max. : 32767	Max. : 32766	Max. : 32767

V4	V5	V6	V7
Min. :-32767	Min. :-32767	Min. :-32767	Min. :-32767
1st Qu.: 7361	1st Qu.: 528	1st Qu.: -24749	1st Qu.: 10705
Median : 15455	Median : 10668	Median : -21170	Median : 17112
Mean : 13566	Mean : 8668	Mean : -19045	Mean : 14653
3rd Qu.: 20631	3rd Qu.: 16947	3rd Qu.: -16230	3rd Qu.: 21092
Max. : 32766	Max. : 32766	Max. : 32767	Max. : 32766

V8	V9	V10	V11
Min. :-32767	Min. :-32767	Min. :-32767	Min. :-32767
1st Qu.: 10862	1st Qu.: 7689	1st Qu.: -13051	1st Qu.: -23050
Median : 15787	Median : 12329	Median : -3117	Median : -15726
Mean : 14637	Mean : 12131	Mean : -1425	Mean : -14123
3rd Qu.: 20044	3rd Qu.: 18981	3rd Qu.: 8959	3rd Qu.: -8102
Max. : 32766	Max. : 32766	Max. : 32766	Max. : 32766

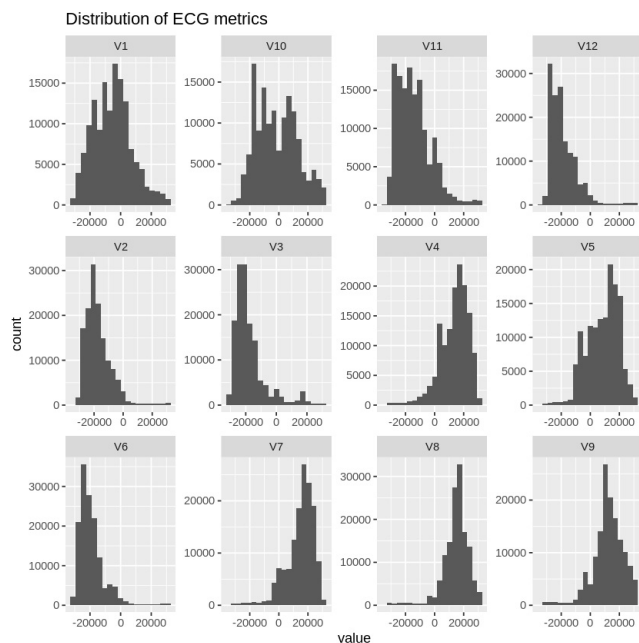
V12	Gender	Sport	Training2017
Min. :-32767	Length:140000	Length:140000	Min. : 321.8
1st Qu.: -25768	Class :character	Class :character	1st Qu.: 743.3
Median : -21094	Mode :character	Mode :character	Median : 821.8
Mean : -18844			Mean : 821.9
3rd Qu.: -14336			3rd Qu.: 900.7
Max. : 32766			Max. : 1328.7

Training2018	Training2019	Simulated Avg Training	Age
Min. : 304.8	Min. : 55.5	Min. : 453.9	Min. :20.00
1st Qu.: 744.8	1st Qu.: 683.2	1st Qu.: 760.6	1st Qu.:23.75
Median : 820.4	Median : 798.2	Median : 813.8	Median :25.00
Mean : 820.6	Mean : 797.8	Mean : 813.4	Mean :26.21
3rd Qu.: 896.7	3rd Qu.: 912.9	3rd Qu.: 866.3	3rd Qu.:31.00
Max. :1355.8	Max. :1563.9	Max. :1139.1	Max. :36.00

```
In [ ]: num_vars <- df %>%
  select(starts_with("V")) %>%
  gather(key = "variable", value = "value") %>%
  ggplot(aes(x = value)) +
  geom_histogram(bins = 20) +
  facet_wrap(~ variable, scales = "free") +
  labs(title = "Distribution of ECG metrics")
```

```
In [ ]: print(num_vars)
```



```
In [ ]: correlation_matrix <- cor(df %>%
  select(starts_with("V")))
```

```
In [ ]: install.packages("corrplot")
```

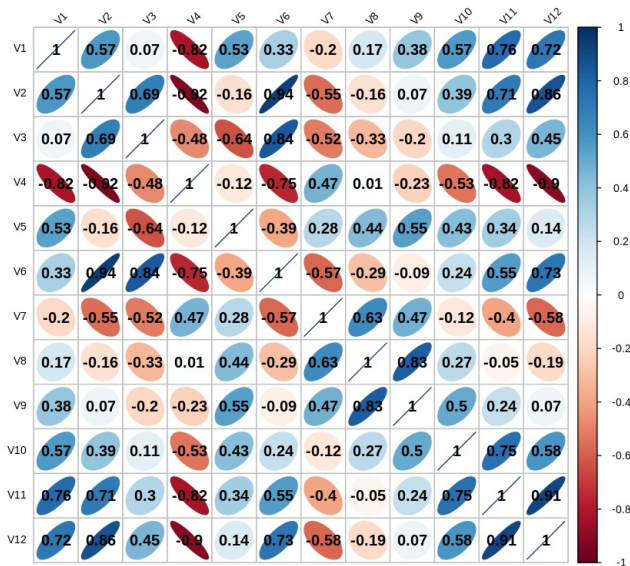
Installing package into '/usr/local/lib/R/site-library'  
(as 'lib' is unspecified)

```
In [ ]: # Using corrplot package for better visualization
library(corrplot)

corrplot(correlation_matrix, method = "ellipse",
  tl.col = "black", tl.srt = 45, tl.cex = 0.7,
```

```
addrect = 4, cl.cex = 0.7, addCoef.col = "black")
```

corrplot 0.92 loaded

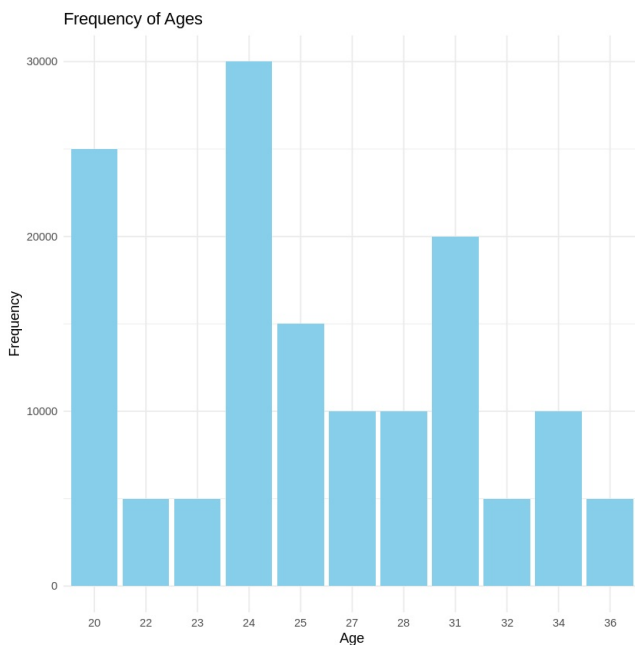


```
In [ ]: age_counts <- table(df$Age)

# Convert the result to a data frame
age_counts_df <- as.data.frame(age_counts)
names(age_counts_df) <- c("Age", "Frequency")

# Create the bar chart
bar_chart <- ggplot(age_counts_df, aes(x = Age, y = Frequency)) +
  geom_bar(stat = "identity", fill = "skyblue") +
  labs(title = "Frequency of Ages", x = "Age", y = "Frequency") +
  theme_minimal()
```

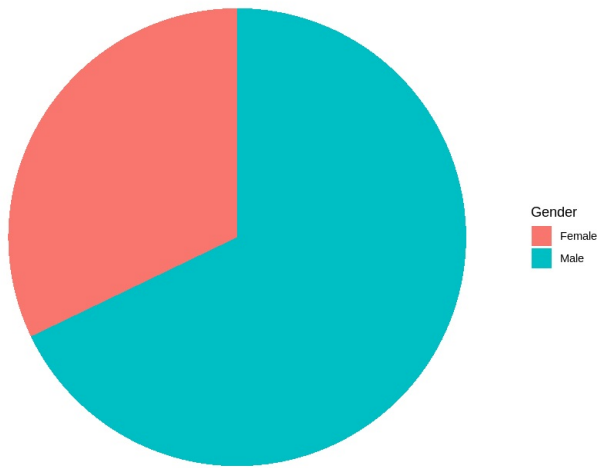
```
In [ ]: print(bar_chart)
```



```
In [ ]: gender_pie <- ggplot(df, aes(x = "", fill = Gender)) +
  geom_bar(width = 1) +
  coord_polar(theta = "y") +
  labs(title = "Gender Distribution") +
  theme_void()

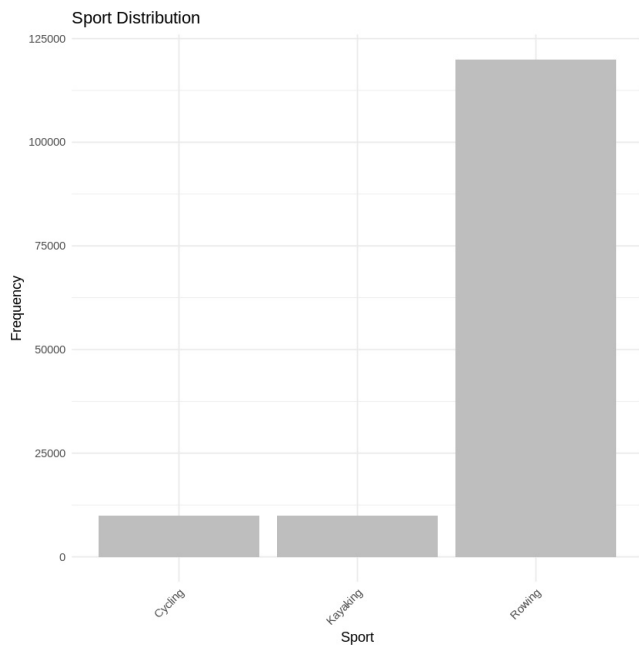
# Display the pie chart
print(gender_pie)
```

Gender Distribution



```
In [ ]: sport_bar <- ggplot(df, aes(x = Sport)) +
  geom_bar(fill = "grey") +
  labs(title = "Sport Distribution", x = "Sport", y = "Frequency") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

# Display the bar chart
print(sport_bar)
```

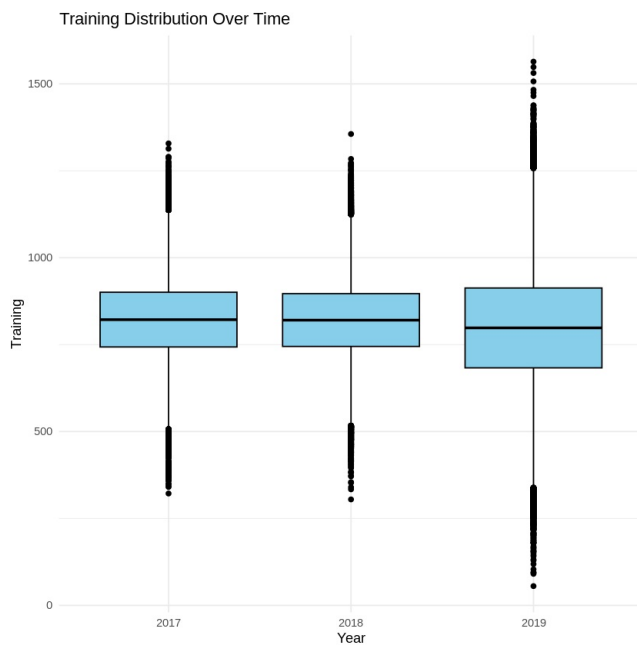


```
In [ ]: training_2017 <- df$Training2017
training_2018 <- df$Training2018
training_2019 <- df$Training2019

# Creating a data frame for boxplot
training_data <- data.frame(Year = factor(rep(c("2017", "2018", "2019"), each = nrow(df))),
  Training = c(training_2017, training_2018, training_2019))

# Creating boxplot for training distribution over time
training_boxplot <- ggplot(training_data, aes(x = Year, y = Training)) +
  geom_boxplot(fill = "skyblue", color = "black") +
  labs(title = "Training Distribution Over Time", x = "Year", y = "Training") +
  theme_minimal()

# Display the boxplot
print(training_boxplot)
```



## 't' TESTS

### 1. Gender Differences in Cardiac Electrical Activity Among Athletes: Insights from ECG Lead Measurements- 'T test'

**Lateral View (Leads I, aVL, V5, and V6):** This view is crucial due to its potential to reflect changes in the lateral wall of the left ventricle, which may undergo hypertrophy or exhibit enhanced cardiac function as a result of intense physical training. Such adaptations can differ by gender based on factors like training intensity, type of sport, and inherent physiological differences.

**Septal and Anterior Views (Leads V1 through V4):** These views are important for assessing the septal and anterior regions, where adaptations could also indicate enhanced cardiac performance or training-related changes. Differences in these areas could provide insights into how male and female athletes' hearts respond differently to their training regimens

```
In [ ]: # Lateral View: Leads I (V1), aVL (V5), V5 (V11), and V6 (V12) + gender
lateral_view <- df[, c('Gender', 'V1', 'V5', 'V11', 'V12')]

# Septal and Anterior Views: Leads V1 through V4 (V7 to V10) + gender
septal_anterior_view <- df[, c('Gender', 'V7', 'V8', 'V9', 'V10')]

head(septal_anterior_view)
```

A data.frame: 6 × 5

	Gender	V7	V8	V9	V10
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
1	Male	13580	11410	14721	16103
2	Male	13331	11096	14093	15416
3	Male	13829	10991	13644	14815
4	Male	13331	10467	13016	14128
5	Male	13331	10049	12477	13355
6	Male	12832	9630	11938	12926

#### Generalized Null and Alternative Hypotheses for the two views :

**Null Hypothesis ( $H_0$ ):** For views Lateral, Septal/Anterior there is no difference in the mean values of the ECG leads between genders. This hypothesis posits that gender does not influence the ECG lead measurements, suggesting that the physiological or electrical properties captured by these leads are similar across males and females.

**Alternative Hypothesis ( $H_1$ ):** For at least one of the two views, there is a difference in the mean values of the ECG leads between genders. This hypothesis suggests that gender may have an influence on the ECG lead measurements, indicating possible physiological or electrical differences in how males and females present ECG characteristics within at least one of the clinical views.

**Significance Level ( $\alpha$ ):**

The significance level, often denoted as  $\alpha$  (alpha), is the probability of rejecting the null hypothesis when it is actually true (Type I error).  $\alpha$  is equal 0.05

**Test Statistics:** Conduct independent T-tests for each lead within each view, comparing male and female groups. The aggregation of these tests across all views and leads provides a comprehensive examination of gender differences in ECG measurements

```
In [ ]: lateral_results <- list(
  'V1' = t.test(V1 ~ Gender, data=lateral_view),
  'V5' = t.test(V5 ~ Gender, data=lateral_view),
  'V11' = t.test(V11 ~ Gender, data=lateral_view),
  'V12' = t.test(V12 ~ Gender, data=lateral_view)
)

selected_leads <- c('V1', 'V5', 'V11', 'V12')
for (lead in selected_leads) {
  cat(sprintf("\nResults for %s:\n", lead))
  print(lateral_results[[lead]])
}
```

Results for V1:

Welch Two Sample t-test

data: V1 by Gender

t = -53.532, df = 111034, p-value < 2.2e-16

alternative hypothesis: true difference in means between group Female and group Male is not equal to 0

95 percent confidence interval:

-3588.013 -3334.553

sample estimates:

mean in group Female	mean in group Male
-7704.832	-4243.549

Results for V5:

Welch Two Sample t-test

data: V5 by Gender

t = -133.81, df = 89465, p-value < 2.2e-16

alternative hypothesis: true difference in means between group Female and group Male is not equal to 0

95 percent confidence interval:

-7929.877 -7700.930

sample estimates:

mean in group Female	mean in group Male
3364.951	11180.354

Results for V11:

Welch Two Sample t-test

data: V11 by Gender

t = -75.006, df = 109238, p-value < 2.2e-16

alternative hypothesis: true difference in means between group Female and group Male is not equal to 0

95 percent confidence interval:

-4573.856 -4340.905

sample estimates:

mean in group Female	mean in group Male
-17147.98	-12690.60

Results for V12:

Welch Two Sample t-test

data: V12 by Gender

t = -49.861, df = 99270, p-value < 2.2e-16

alternative hypothesis: true difference in means between group Female and group Male is not equal to 0

95 percent confidence interval:

-2673.448 -2471.216

sample estimates:

mean in group Female	mean in group Male
-20589.77	-18017.44

```
In [ ]: septal_anterior_results <- list(
  'V7' = t.test(V7 ~ Gender, data=septal_anterior_view),
  'V8' = t.test(V8 ~ Gender, data=septal_anterior_view),
  'V9' = t.test(V9 ~ Gender, data=septal_anterior_view),
  'V10' = t.test(V10 ~ Gender, data=septal_anterior_view)
)
```

```
In [ ]: selected_leads <- c('V7', 'V8', 'V9', 'V10')
for (lead in selected_leads) {
  cat(sprintf("\nResults for %s:\n", lead))
  print(septal_anterior_results[[lead]])
}
```

Results for V7:

Welch Two Sample t-test

```
data: V7 by Gender
t = -7.1282, df = 90863, p-value = 1.025e-12
alternative hypothesis: true difference in means between group Female and group Male is not equal to 0
95 percent confidence interval:
 -502.0964 -285.5277
sample estimates:
mean in group Female    mean in group Male
      14385.34             14779.15
```

Results for V8:

Welch Two Sample t-test

```
data: V8 by Gender
t = -14.598, df = 90680, p-value < 2.2e-16
alternative hypothesis: true difference in means between group Female and group Male is not equal to 0
95 percent confidence interval:
 -849.6893 -648.5349
sample estimates:
mean in group Female    mean in group Male
      14128.46             14877.57
```

Results for V9:

Welch Two Sample t-test

```
data: V9 by Gender
t = -70.952, df = 90440, p-value < 2.2e-16
alternative hypothesis: true difference in means between group Female and group Male is not equal to 0
95 percent confidence interval:
 -4261.499 -4032.388
sample estimates:
mean in group Female    mean in group Male
      9317.239             13464.182
```

Results for V10:

Welch Two Sample t-test

```
data: V10 by Gender
t = -11.699, df = 97764, p-value < 2.2e-16
alternative hypothesis: true difference in means between group Female and group Male is not equal to 0
95 percent confidence interval:
 -1050.8727 -749.2892
sample estimates:
mean in group Female    mean in group Male
      -2036.184             -1136.103
```

For all selected leads the p-values are extremely low, far below the significance level of 0.05. This indicates a statistically significant difference in the means of these variables between genders. Therefore, we reject the null hypothesis for each of these tests, concluding that there is a significant difference between males and females for the tested variables.

## Inference

The Welch Two Sample t-tests performed across several ECG leads (V1, V5, V11, V12) revealed statistically significant variations in mean values between genders among athletes, with p-values significantly lower than the customary threshold of 0.05. These findings strongly support the rejection of the null hypothesis, demonstrating that male and female athletes have different mean values in these ECG leads.

Similarly, leads V7 to V10 exhibit significant variances, indicating a similar pattern of gender inequality in cardiac electrical activity across the athletes tested. The amount and direction of these changes are further defined by confidence intervals and mean estimates, which supplement the statistical findings with measurable metrics of effect size.

## Clinical implications

The significant gender differences in ECG lead measurements indicate that male and female athletes may experience different cardiac adaptations in response to training, which could have implications for sports medicine, including training regimens, performance optimisation, and health monitoring strategies. Understanding these differences is critical for generating gender-specific suggestions that can help maximise athletic performance while lowering the risk of heart disease.

## Conclusion

By incorporating gender-specific data from these ECG leads into clinical and training procedures, the sports medicine community can improve the accuracy of athlete health care. This method not only improves sports performance but also greatly reduces the risk of long-



term unfavourable cardiac consequences. These findings enable a more detailed knowledge of athlete heart health, promoting a transition towards more personalised athlete care.

While the findings largely show disparities within the athlete population, they also contribute to a larger medical understanding of how athletic training affects heart function differently in men and women, hence enabling individualised medical and training methods.

\*\*\*\*\*END OF TEST  
1\*\*\*\*\*

## ANOVA and 'F' TESTS

### 2."Exploring the Impact of Sports played on ECG Measurements: An analysis Using ANOVA"

*Null Hypothesis (H0):* There is no difference in the mean ECG measurements across sports (rowing, kayaking, cycling) for each lead variable. Any observed differences are due to chance alone.

*Alternative Hypothesis (H1):* There is a significant difference in the mean ECG measurements across sports for at least one of the lead variables. The observed differences are not due to chance and reflect a true effect of sport type on ECG measurements.

*Significance Level ( $\alpha$ ):*

The significance level, often denoted as  $\alpha$  (alpha), is the probability of rejecting the null hypothesis when it is actually true (Type I error).  $\alpha$  is equal 0.05

This investigation used stratified sampling to overcome the large imbalance reported across sports groups. This strategy divides the population into discrete subgroups, or strata, based on specific characteristics—in this case, sport type. Stratified sampling ensures that each subgroup is proportionally represented in the sample, resulting in a more accurate and representative study. In essence, it reduces the influence of different sample sizes between groups, increasing the trustworthiness of statistical inferences generated from the data.

```
In [ ]: set.seed(123) # Ensure reproducibility

n_samples <- min(table(df$Sport))

# Perform stratified sampling
library(dplyr)
stratified_sample <- df %>%
  group_by(Sport) %>%
  sample_n(size = n_samples) %>%
  ungroup()

# Now, stratified_sample contains a balanced subset of your data
```

```
In [ ]: # Conduct ANOVA for each ECG lead on the stratified sample
results_anova_stratified <- list()
for(lead in c('V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10', 'V11', 'V12')) {
  formula <- as.formula(paste(lead, '~ Sport'))
  results_anova_stratified[[lead]] <- summary(aov(formula, data = stratified_sample))
}

# Loop through the list of ANOVA results and print each one
leads <- c('V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10', 'V11', 'V12')
for (lead in leads) {
  cat("\nANOVA Results for", lead, ":\n")
  print(results_anova_stratified[[lead]])
}
```

```

ANOVA Results for V1 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 8.442e+09 4.221e+09  36.26 <2e-16 ***
Residuals 29997 3.492e+12 1.164e+08
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V2 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 1.459e+11 7.295e+10  1127 <2e-16 ***
Residuals 29997 1.942e+12 6.475e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V3 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 4.841e+11 2.420e+11  3102 <2e-16 ***
Residuals 29997 2.340e+12 7.801e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V4 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 6.013e+10 3.007e+10  389.1 <2e-16 ***
Residuals 29997 2.318e+12 7.727e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V5 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 2.624e+11 1.312e+11  1313 <2e-16 ***
Residuals 29997 2.998e+12 9.995e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V6 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 2.164e+11 1.082e+11  1768 <2e-16 ***
Residuals 29997 1.836e+12 6.119e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V7 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 9.293e+11 4.647e+11  5268 <2e-16 ***
Residuals 29997 2.646e+12 8.821e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V8 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 5.224e+11 2.612e+11  3922 <2e-16 ***
Residuals 29997 1.998e+12 6.660e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V9 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 1.096e+10 5.479e+09   55.01 <2e-16 ***
Residuals 29997 2.988e+12 9.961e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V10 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 3.065e+10 1.532e+10   94.02 <2e-16 ***
Residuals 29997 4.889e+12 1.630e+08
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V11 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 1.782e+11 8.912e+10   977.8 <2e-16 ***
Residuals 29997 2.734e+12 9.114e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ANOVA Results for V12 :
      Df    Sum Sq   Mean Sq F value Pr(>F)
Sport    2 9.696e+10 4.848e+10   664.8 <2e-16 ***
Residuals 29997 2.188e+12 7.293e+07
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

## F test Results

The significant F-values found across all ECG leads in our ANOVA study clearly show that the variance between sports (rowing, kayaking, and cycling) is significantly bigger than the variance found within each sport. The huge difference in variances, as indicated by

the F-test, strongly supports the alternative hypothesis that mean ECG values differ statistically significantly among sports. We reject the null hypothesis as all leads' p-values are significantly lower than the specified significance level ( $\alpha$ ) of 0.05.

### Statistical Significance Across All Leads:

The "Sport" factor has a significant effect on the measurements for every ECG lead (V1-V12). This suggests that the type of sport an athlete participates in is connected with variances in ECG results across all investigated leads.

Beyond statistical significance, determining impact sizes (mean differences) and their practical implications in sports physiology and athlete health is critical. This includes examining how significant the disparities are and what they might entail for players in other sports.

The continuous statistical significance observed across all ECG leads emphasises the impact of sport type on cardiovascular measures. This results lends support to a further in-depth investigation into how and why different sports may result in varying ECG profiles across athletes, emphasising the significance of sport-specific factors in athlete health monitoring and research.

```
In [ ]: # Loop through all ECG leads and conduct Tukey HSD post-hoc test for each
leads <- c('V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10', 'V11', 'V12')

# Store the results in a list for easy access
tukey_results <- list()

for (lead in leads) {
  # Fit ANOVA model for the current lead
  anova_model <- aov(reformulate('Sport', response = lead), data = df)

  # Conduct Tukey HSD post-hoc test on the fitted ANOVA model
  tukey_post_hoc <- TukeyHSD(anova_model)

  # Store the post-hoc test results in the list
  tukey_results[[lead]] <- tukey_post_hoc

  # Print the results
  cat("\nTukey HSD Post-Hoc Test Results for", lead, ":\n")
  print(tukey_post_hoc)
}
```

Tukey HSD Post-Hoc Test Results for V1 :  
Tukey multiple comparisons of means  
95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

```
$Sport
      diff      lwr      upr  p adj
Kayaking-Cycling -572.6304 -985.1868 -160.074 0.003274
Rowing-Cycling    751.2669  447.6336 1054.900 0.000000
Rowing-Kayaking   1323.8973 1020.2640 1627.531 0.000000
```

Tukey HSD Post-Hoc Test Results for V2 :  
Tukey multiple comparisons of means  
95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

```
$Sport
      diff      lwr      upr p adj
Kayaking-Cycling -5401.248 -5695.996 -5106.500    0
Rowing-Cycling   -2676.700 -2893.628 -2459.771    0
Rowing-Kayaking   2724.548  2507.619  2941.476    0
```

Tukey HSD Post-Hoc Test Results for V3 :  
Tukey multiple comparisons of means  
95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

```
$Sport
      diff      lwr      upr p adj
Kayaking-Cycling -9037.142 -9382.7482 -8691.535    0
Rowing-Cycling   -7985.428 -8239.7874 -7731.068    0
Rowing-Kayaking   1051.714   797.3543  1306.073    0
```

Tukey HSD Post-Hoc Test Results for V4 :  
Tukey multiple comparisons of means  
95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

```
$Sport
      diff      lwr      upr p adj
Kayaking-Cycling 3303.7400 2971.1977 3636.2823    0
Rowing-Cycling   2571.1384 2326.3939 2815.8829    0
```

Rowing-Kayaking -732.6016 -977.3461 -487.8571 0

Tukey HSD Post-Hoc Test Results for V5 :

Tukey multiple comparisons of means

95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

\$Sport

	diff	lwr	upr	p adj
Kayaking-Cycling	-1052.438	-1406.010	-698.8656	0
Rowing-Cycling	5669.269	5409.047	5929.4906	0
Rowing-Kayaking	6721.706	6461.484	6981.9283	0

Tukey HSD Post-Hoc Test Results for V6 :

Tukey multiple comparisons of means

95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

\$Sport

	diff	lwr	upr	p adj
Kayaking-Cycling	-6405.420	-6686.060	-6124.781	0
Rowing-Cycling	-1984.741	-2191.286	-1778.196	0
Rowing-Kayaking	4420.679	4214.134	4627.224	0

Tukey HSD Post-Hoc Test Results for V7 :

Tukey multiple comparisons of means

95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

\$Sport

	diff	lwr	upr	p adj
Kayaking-Cycling	13570.845	13258.862	13882.827	0
Rowing-Cycling	5660.521	5430.909	5890.134	0
Rowing-Kayaking	-7910.324	-8139.936	-7680.711	0

Tukey HSD Post-Hoc Test Results for V8 :

Tukey multiple comparisons of means

95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

\$Sport

	diff	lwr	upr	p adj
Kayaking-Cycling	10192.351	9898.786	10485.915	0
Rowing-Cycling	4454.755	4238.697	4670.813	0
Rowing-Kayaking	-5737.596	-5953.654	-5521.538	0

Tukey HSD Post-Hoc Test Results for V9 :

Tukey multiple comparisons of means

95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

\$Sport

	diff	lwr	upr	p adj
Kayaking-Cycling	359.9783	12.62403	707.3326	0.040153
Rowing-Cycling	1405.4890	1149.84315	1661.1348	0.000000
Rowing-Kayaking	1045.5107	789.86485	1301.1565	0.000000

Tukey HSD Post-Hoc Test Results for V10 :

Tukey multiple comparisons of means

95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

\$Sport

	diff	lwr	upr	p adj
Kayaking-Cycling	2282.551	1819.11811	2745.9835	0.0000000
Rowing-Cycling	291.110	-49.96719	632.1873	0.1120972
Rowing-Kayaking	-1991.441	-2332.51799	-1650.3635	0.0000000

Tukey HSD Post-Hoc Test Results for V11 :

Tukey multiple comparisons of means

95% family-wise confidence level

Fit: aov(formula = reformulate("Sport", response = lead), data = df)

\$Sport

	diff	lwr	upr	p adj
--	------	-----	-----	-------

Kayaking-Cycling	-4298.664	-4675.612	-3921.717	0
Rowing-Cycling	1325.326	1047.900	1602.752	0
Rowing-Kayaking	5623.990	5346.564	5901.416	0

Tukey HSD Post-Hoc Test Results for V12 :  
 Tukey multiple comparisons of means  
 95% family-wise confidence level

```
Fit: aov(formula = reformulate("Sport", response = lead), data = df)
```

\$Sport		diff	lwr	upr	p adj
Kayaking-Cycling	-3971.9477	-4285.6278	-3658.2676	0.0000000	
Rowing-Cycling	-403.7036	-634.5659	-172.8413	0.0001232	
Rowing-Kayaking	3568.2441	3337.3818	3799.1064	0.0000000	

The Tukey HSD post-hoc test results confirm that there are substantial and statistically significant disparities in ECG measures between athletes from various sports. This detailed analysis lends support to your ANOVA's alternative hypothesis, revealing genuine effects of sport type on ECG data that beyond what could be expected by chance. These findings can help personalise athlete training regimens, potentially leading to increased athletic performance and better cardiovascular health management in sports medicine.

### Conclusion:

The study, which used ANOVA followed by Tukey HSD post-hoc testing, clearly shows that the type of sport has a substantial affect on ECG data across various leads. The ANOVA results, corroborated by substantial F-values, indicated that variances in ECG profiles exist across sports such as rowing, kayaking, and cycling, implying that each sport has a distinct effect on cardiac electrical activity. These findings are supported by the Tukey HSD tests, which revealed unique pairwise differences and emphasised the varying impact of various sports on the ECG.

This study demonstrates that athletic training in different sports results in distinct cardiac adaptations, emphasising the importance of sport-specific cardiovascular monitoring and personalised therapies. Our thorough statistical approach, particularly the successful use of the F-test in ANOVA, ensures that these findings are both scientifically robust and clinically useful, paving the way for optimal sports medicine procedures that can improve athletic performance while reducing health risks.

\*\*\*\*\*END OF TEST  
 2\*\*\*\*\*

## BAYESIAN NETWORKS

**How does the interaction between age group and average training intensity influence the ECG characteristics (V1-V12), and how do these influences differ when accounting for the mediating effects of gender and sport discipline?**

```
In [ ]: df$AgeGroup <- cut(df$Age, breaks=c(18, 25, 35, 45, 55, 65), labels=c("18-24", "25-34", "35-44", "45-54", "55-64"))
```

```
In [ ]: df$Simulated_Avg_Training_Group <- cut(df$Simulated_Avg_Training, breaks=quantile(df$Simulated_Avg_Training, probs=c(0.33, 0.66)), labels=c("Low", "Medium", "High"), include.lowest = TRUE)
```

```
In [ ]: df$Gender <- as.factor(df$Gender)
df$Sport <- as.factor(df$Sport)
df$Sport <- as.factor(df$Sport)
df$Simulated_Avg_Training_Group <- as.factor(df$Simulated_Avg_Training_Group)
```

```
In [ ]: install.packages("bnlearn")

Installing package into '/usr/local/lib/R/site-library'
(as 'lib' is unspecified)
```

```
In [ ]: library(bnlearn)
```

```
In [ ]: library(data.table)
```

Attaching package: 'data.table'

The following objects are masked from 'package:dplyr':

between, first, last

```
In [ ]: bn_data <- df[, c("Gender", "Sport", "AgeGroup", "V1", "V2", "V3", "V4", "V5", "V6", "V7", "V8", "V9", "V10", "V11", "V12")]
```

```
In [ ]: # Ensure all are factors
# Convert all categorical variables to factors if not already
```

```
bn_data[] <- lapply(bn_data, function(x) if(is.character(x) || is.integer(x)) factor(x) else x)
bn_data <- lapply(bn_data, factor)
```

```
In [ ]: # Assuming your subsetted data is stored in bn_data and it's a data.table
bn_data <- as.data.frame(bn_data)
```

```
In [ ]: head(bn_data,10)
```

A data.frame: 10 × 16

	Gender	Sport	AgeGroup	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	Simulated_Avg_Trainin
	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	<fct>	
1	Male	Rowing	25-34	10251	-1096	-10267	-3724	9391	-5395	13580	11410	14721	16103	6662	-3806	
2	Male	Rowing	25-34	8643	-2558	-10829	-1862	8973	-6448	13331	11096	14093	15416	5897	-4548	
3	Male	Rowing	25-34	5427	-3776	-9985	743	6469	-6711	13829	10991	13644	14815	5460	-5105	
4	Male	Rowing	25-34	5427	-4507	-10829	1116	7304	-7501	13331	10467	13016	14128	4696	-5848	
5	Male	Rowing	25-34	6231	-4751	-11673	1116	8138	-8027	13331	10049	12477	13355	3822	-6962	
6	Male	Rowing	25-34	5427	-5481	-11954	1860	8138	-8554	12832	9630	11938	12926	3385	-7704	
7	Male	Rowing	25-34	4623	-5969	-11954	2605	7721	-8817	12832	9316	11669	12410	2730	-8076	
8	Male	Rowing	25-34	4221	-6943	-12517	3722	7721	-9606	12583	9107	11400	11981	2074	-9004	
9	Male	Rowing	25-34	3416	-7187	-12235	4467	6886	-9606	12334	8897	11131	11723	1637	-9190	
10	Male	Rowing	25-34	4221	-7430	-13079	4094	8138	-10133	12085	8792	11041	11637	1419	-9375	

```
In [ ]: str(bn_data)
```

```
'data.frame': 140000 obs. of 16 variables:
 $ Gender      : Factor w/ 2 levels "Female","Male": 2 2 2 2 2 2 2 2 2 2 ...
 $ Sport       : Factor w/ 3 levels "Cycling","Kayaking",...: 3 3 3 3 3 3 3 3 3 3 ...
 $ AgeGroup    : Factor w/ 3 levels "18-24","25-34",...: 2 2 2 2 2 2 2 2 2 2 ...
 $ V1          : Factor w/ 4359 levels "-32767","-32766",...: 2911 2801 2580 2580 2635 2580 2524
2493 2434 2493 ...
 $ V2          : Factor w/ 5914 levels "-32767","-32645",...: 3729 3596 3477 3400 3375 3295 3244
3134 3109 3080 ...
 $ V3          : Factor w/ 4533 levels "-32767","-32650",...: 2038 1999 2058 1999 1943 1924 1924
1876 1898 1839 ...
 $ V4          : Factor w/ 5197 levels "-32767","-32766",...: 1837 1993 2207 2238 2238 2301 2368
2465 2530 2494 ...
 $ V5          : Factor w/ 3146 levels "-32767","-32766",...: 1882 1861 1737 1779 1825 1825 1803
1803 1756 1825 ...
 $ V6          : Factor w/ 5082 levels "-32767","-32643",...: 2812 2725 2702 2640 2591 2545 2518
2449 2449 2406 ...
 $ V7          : Factor w/ 4822 levels "-32767","-32766",...: 3130 3106 3154 3106 3106 3062 3062
3038 3012 2990 ...
 $ V8          : Factor w/ 8348 levels "-32767","-32766",...: 4490 4440 4425 4352 4289 4229 4181
4154 4120 4107 ...
 $ V9          : Factor w/ 9661 levels "-32767","-32682",...: 5894 5753 5657 5520 5399 5293 5234
5175 5118 5099 ...
 $ V10         : Factor w/ 9746 levels "-32767","-32656",...: 7334 7223 7125 7009 6871 6792 6693
6617 6568 6555 ...
 $ V11         : Factor w/ 8961 levels "-32767","-32710",...: 6567 6490 6441 6357 6249 6201 6114
6036 5985 5961 ...
 $ V12         : Factor w/ 6903 levels "-32767","-32696",...: 4164 4090 4028 3938 3801 3696 3649
3526 3499 3473 ...
 $ Simulated_Avg_Training_Group: Factor w/ 3 levels "Low","Medium",...: 3 3 1 3 1 3 1 2 2 2 ...
```

```
In [ ]: set.seed(123) # Set seed for reproducibility
bn_structure <- hc(bn_data, score = "bic")
```

```
In [ ]: if (!require("Rgraphviz")) install.packages("Rgraphviz", repos="http://bioconductor.org/packages/release/bioc")
```

```
Loading required package: Rgraphviz
```

```
Loading required package: graph
```

```
Loading required package: BiocGenerics
```

```
Attaching package: 'BiocGenerics'
```

```
The following object is masked from 'package:bnlearn':
```

```
score
```

```
The following objects are masked from 'package:dplyr':
```

```
combine, intersect, setdiff, union
```

```
The following objects are masked from 'package:stats':
```

```
IQR, mad, sd, var, xtabs
```

```
The following objects are masked from 'package:base':
```

```
anyDuplicated, aperm, append, as.data.frame, basename, cbind,  
colnames, dirname, do.call, duplicated, eval, evalq, Filter, Find,  
get, grep, grepl, intersect, is.unsorted, lapply, Map, mapply,  
match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,  
Position, rank, rbind, Reduce, rownames, sapply, setdiff, sort,  
table, tapply, union, unique, unsplit, which.max, which.min
```

```
Attaching package: 'graph'
```

```
The following objects are masked from 'package:bnlearn':
```

```
degree, nodes, nodes<-
```

```
Loading required package: grid
```

```
In [ ]: library(Rgraphviz)
```

```
In [ ]: fitted_bn <- bn.fit(bn_structure, data = bn_data)
```

```
In [ ]: print(bn_structure)
```

```
Bayesian network learned via Score-based methods
```

```
model:
```

```
[AgeGroup][Simulated_Avg_Training_Group][V5|AgeGroup][Gender|V5][Sport|V5]  
[V1|Gender:AgeGroup][V2|Gender][V3|Gender:AgeGroup][V4|Gender][V6|Gender]  
[V7|Gender][V8|Gender][V9|Gender][V10|Gender][V11|Gender][V12|Gender]
```

```
nodes: 16
```

```
arcs: 16
```

```
  undirected arcs: 0
```

```
  directed arcs: 16
```

```
average markov blanket size: 2.12
```

```
average neighbourhood size: 2.00
```

```
average branching factor: 1.00
```

```
learning algorithm: Hill-Climbing
```

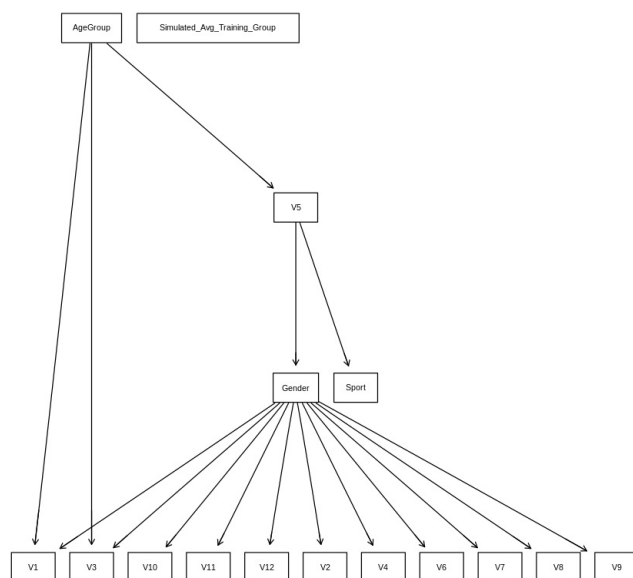
```
score: BIC (disc.)
```

```
penalization coefficient: 5.924699
```

```
tests used in the learning procedure: 360
```

```
optimized: TRUE
```

```
In [ ]: graphviz.plot(fitted_bn)
```



AgeGroup and Simulated\_Average\_Training\_Group are separate root nodes. V5 is a mediator variable that influences AgeGroup and Simulated\_Avg\_Training\_Group, which in turn affect Gender and Sport. Gender influences the remaining ECG characteristics (V1, V2, V3, V4, V6, V7, V8, V9, V10, V11, V12), and some are also affected by AgeGroup.

The structure suggests that AgeGroup and Simulated\_Avg\_Training\_Group may have a direct effect on one of the ECG characteristics (V5). This ECG characteristic (V5) then has downstream consequences on Gender and Sport, influencing all other ECG values.

### Inference

The inference proposes a model in which age and training influence ECG readings both directly and indirectly via their effects on gender and sport. In practice, this could imply that age and training intensity have both universal and specific effects on ECG readings, manifesting differently in males and females and across different sports disciplines.

Specific age groups or training intensities are linked to ECG readings indicating increased cardiovascular risk or athletic performance. Gender-specific pathways (mediated by V5) to various ECG readings reflect gender-specific cardiovascular responses or athletic performance characteristics.

Different sports have distinct ECG profiles that are influenced by underlying characteristics such as age and training intensity, presumably as a result of sport-specific physiological requirements. The actual nature of the inferences and subsequent study approaches would be determined by the variables' details (particularly the nature of the ECG features) and domain-specific information.

### Conclusion

*Age and Training:* Age group and average training intensity are key determinants in influencing ECG outcomes, demonstrating that both intrinsic (age) and extrinsic (training) elements are crucial for cardiovascular health and performance.

*Gender and Sport Mediation:* The characteristic V5 is not only directly affected by age and training, but it also acts as a mediator between these parameters and the variables gender and sport. This shows that therapies or training adjustments affecting V5 may have varying consequences depending on an individual's gender and sport.

*Impact on Other ECG features:* Gender and sport have a significant influence on a variety of other ECG features (V1, V2, V3, V4, V6, V7, V8, V9, V10, V11, V12), implying a complex interplay between physiological and demographic factors. This complexity suggests that personalised approaches to sports training and healthcare may be advantageous.

Potential applications include tailoring athletic training programmes or healthcare advice to individual age and training levels, which could improve cardiovascular results and athletic performance.

\*\*\*\*\*End of Analysis 3 \*\*\*\*\*

## CASUAL INFERENCE

### ECG Data Analysis Suggests Causal Inference:

The Bayesian network analysis reveals a causal relationship in which age and training intensity have a direct impact on the ECG characteristic V5. This attribute appears to influence other ECG parameters via the gender and sport discipline pathways.

*Simplified interpretation:*



The study suggests that the age and intensity with which we train may have a direct impact on a specific feature of our heart's electrical activity (V5). This element therefore appears to have a knock-on impact, influencing other heart activity metrics, possibly differently for men and women and across sports. While we notice these correlations in the data, verifying that one factor causes the other would necessitate a more thorough investigation. For the time being, this knowledge could lead to better, more personalised training recommendations that take into account a person's age and gender, resulting in healthier hearts and improved athletic performance.

\*\*\*\*\*End of Analysis 4\*\*\*\*\*

## **\*\*Cardiac Rhythms in Motion: Unraveling the ECG Patterns of Athletes Across Genders and other Disciplines\*\***

### **Abstract:**

This study looks into the multiple effects of age, training intensity, gender, and sport type on electrocardiogram (ECG) characteristics in a diverse group of athletes. Using data from the PhysioNet Norwegian Athlete ECG dataset, this study uses statistical methods to determine the impact of these factors on cardiac electrophysiology. Significant findings using comprehensive Welch Two Sample t-tests and ANOVA indicate the link between age and training intensity with ECG features, as well as gender-specific changes in ECG patterns. Furthermore, the effect of different sports disciplines on ECG features is investigated, providing insights into sport-specific cardiac adaptations. The findings seek to improve preventative sports medicine practices and athlete training programmes, resulting in better health outcomes and performance.

### **Introduction:**

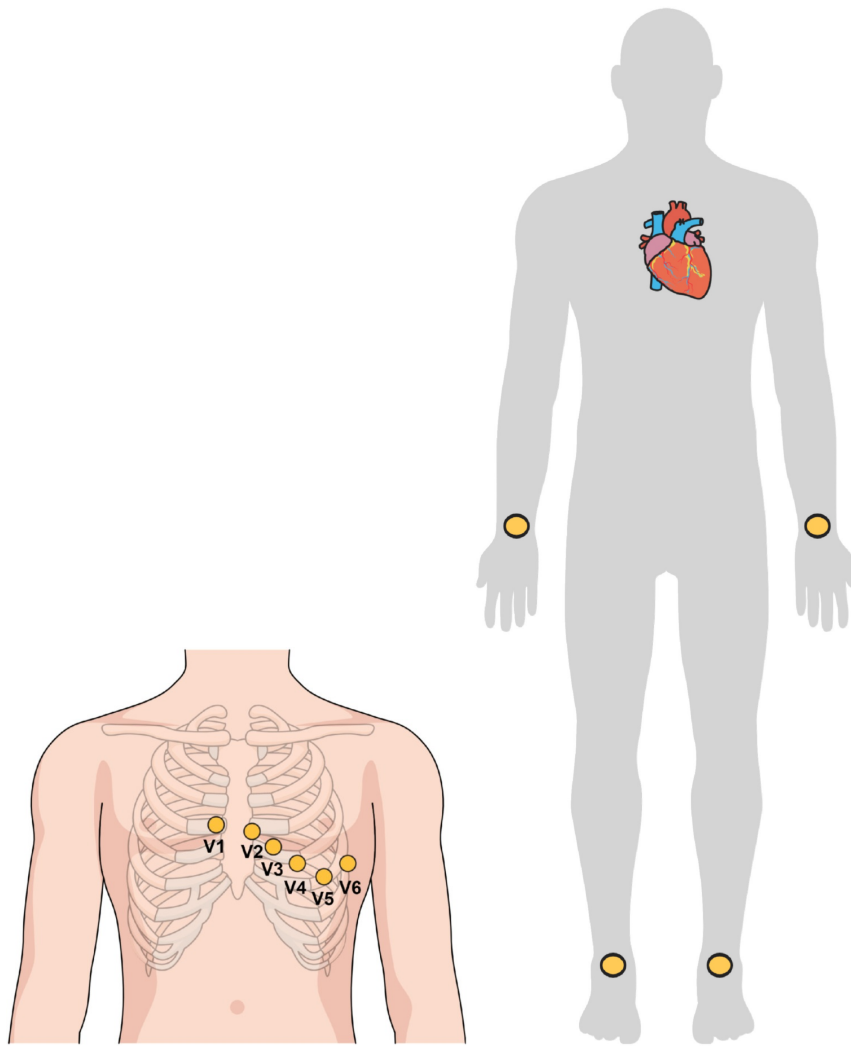
The human heart, a living and active organ, reacts intriguingly to the demands of athletic exercise. Investigating the delicate subtleties of how various elements influence heart function yields insights that are not only scientifically interesting, but also critical for athlete health and performance optimisation. An electrocardiogram (ECG) is a basic test that can reveal important information about the heart's condition. ECGs can identify adaptations or anomalies caused by rigorous physical exercise in athletes. Understanding these changes is crucial for improving performance and avoiding potential cardiac problems. Previous research has identified baseline variability in ECG readings among athletes, emphasising the effect of strenuous physical activity on heart function. There is a large corpus of research on gender-specific cardiac adaptations and the impact of various sports on heart shape and function.

### **Data Preparation**

The dataset for this study was obtained from PhysioNet's Norwegian Athlete ECG dataset, which is part of a bigger database designed to aid scientific research. The dataset for this observational study was created by sports cardiology experts, with the goal of providing thorough insights into Norwegian players' ECG characteristics. The file types were compatible with , each of the 28.dat files has a 12 x 5000 array, where 12 represents the number of leads and 5000 represents the number of samples within each lead.

Participants who volunteered their ECGs to this study were informed and provided written consent before data collection began; they also agreed to have their ECGs shared in an open database. The study protocol and permission form were approved by the Norwegian Centre for Research Data (application ID: 389013) and the University of Oslo. The ethical concerns were approved by the Regional Committees for Medical and Health Research Ethics (application ID: 51205).

This study presents a dataset of electrocardiograms from 28 competitive Norwegian endurance athletes. The electrocardiograms are typical 12-lead resting ECGs that were recorded for 10 seconds while the athletes were lying supine on a bench. The electrocardiograms were then evaluated by both an algorithm and a trained cardiologist.



The figure(i) shows how the precordial leads were placed on the test subjects

The figure(ii) shows how the limb leads were placed on the subjects in this study.

#### Questions of interest:

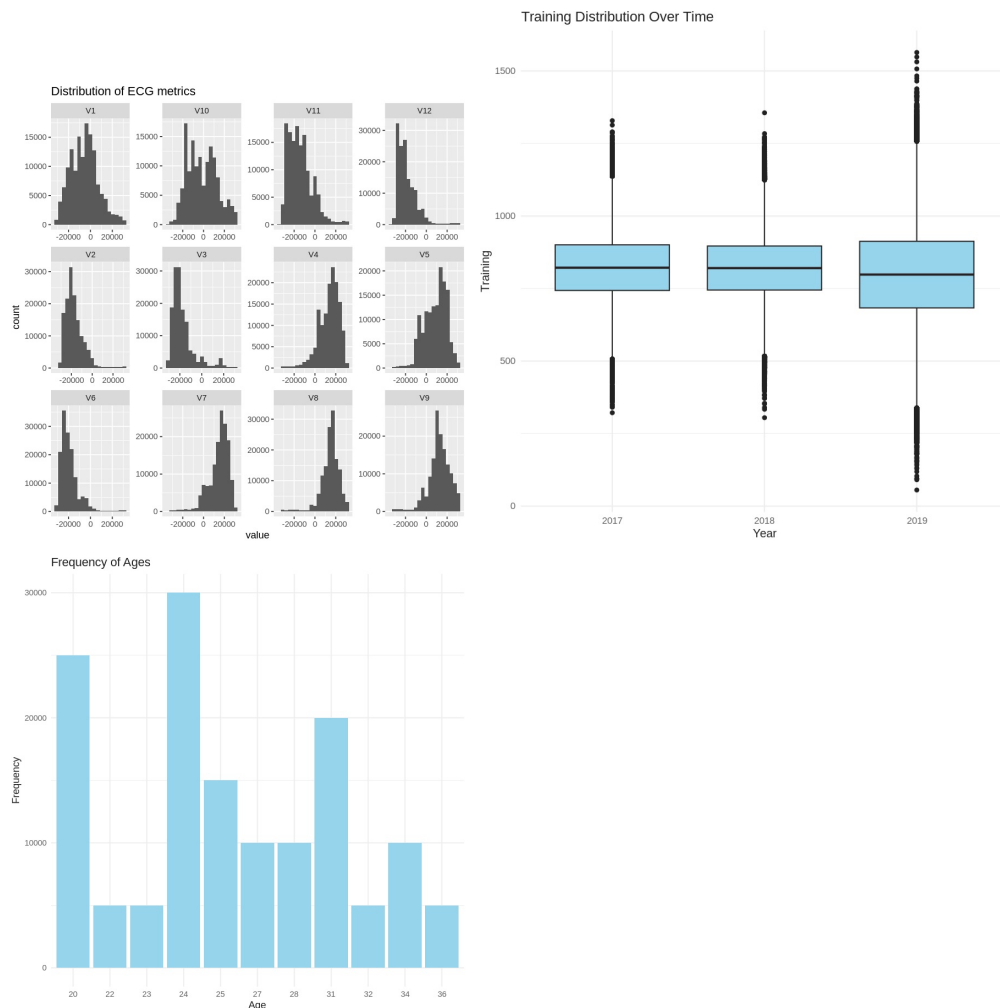
1. Do differences exist between male and female athletes concerning ECG parameters?
2. Does the type of sport an athlete is engaged in influence their ECG patterns?
3. How do age and training intensity impact ECG readings in athletes?

By delving into these issues, the study hopes to improve our understanding of cardiac electrophysiology in the context of sports medicine and athletic training.

#### Exploratory Data Analysis

The EDA used several visualisation techniques, such as showing the distribution of ECG data across multiple leads, to better comprehend the variability and detect any obvious patterns or outliers. Descriptive statistics summarised the central tendencies and dispersions for each lead. The dataset's multidimensionality was addressed by creating a correlation matrix, which was then visualised with heatmaps to reveal the correlations between distinct ECG leads. Gender differences were investigated using grouped histograms, and the possible impact of athletes' sports disciplines on their ECG parameters was assessed using box plots for each ECG lead.

This EDA laid a solid foundation for inferential statistical analysis, ensuring that any conclusions drawn were based on a thorough understanding of the dataset's underlying structure and properties.



## Methods/ Results(Experimental Design)

Our study sought to understand the subtle dynamics of ECG characteristics in athletes, including the connections between age, training intensity, gender, and sports discipline. We used Bayesian Networks, T-tests, and ANOVA to investigate the nuances of cardiac electrical activity as seen in ECG lead data. Electrocardiogram (ECG) leads provide a non-invasive view of cardiac electrical activity, which can be affected by a variety of factors such as age, training intensity, gender, and sport type. Understanding these effects is critical for athlete health and performance improvement.

### 1. Gender Differences in Cardiac Electrical Activity Among Athletes: Insights from ECG Lead Measurements- 'T test'

Independent T-tests were effective in revealing significant gender variations in ECG lead measurements, notably in the lateral and septal/anterior views of the heart. These statistical tests were based on several critical assumptions, including male and female data sets being independent, the normal distribution of ECG measures within each gender group for all leads, and variances being equal between the two groups. The results of these T-tests were illuminating, with significant differences between genders in lateral view leads I, aVL, V5, and V6, as well as septal and anterior views (Leads V1 through V4). This research suggests that male and female athletes may have different cardiac adaptations in response to physical exercise, presumably due to intrinsic physiological differences. These findings are crucial, as they provide statistical support for the concept that gender-specific characteristics must be considered when interpreting ECG readings in athletes.

### 2. "Exploring the Impact of Sports played on ECG Measurements: An analysis Using ANOVA"

The use of Analysis of Variance (ANOVA) provided a rigorous methodology for determining if the type of sport had an effect on ECG features across groups, followed by thorough pairwise comparisons using Tukey HSD post-hoc testing. This analysis was based on assumptions about the sports groups' independence, the normality of the distribution of ECG measures within each sporting category, and the homogeneity of variances between these categories. The ANOVA results revealed statistically significant differences in ECG measures among athletes from various sports disciplines, which were further supported by the Tukey HSD post-hoc tests, which identified the individual sports where these differences were most prominent. This comparative investigation demonstrated that not all sports have the same impact on cardiac electrical activity, emphasising the significance of customising cardiovascular monitoring and therapies to the unique demands of each sport.

### 3. How does the interaction between age group and average training intensity influence the ECG characteristics (V1-V12), and how do these influences differ when accounting for the mediating effects of gender and sport discipline?

(BAYESIAN NETWORKS)

The implementation of Bayesian Networks in this study provided remarkable insights into the complex interactions between age, training intensity, gender, and sport discipline as they influence ECG features. Using a probabilistic framework, these networks captured not just

the direct impacts but also the subtle interdependences between these factors. For example, the variable V5 appeared as a crucial mediator, implying that it may play an important role in the regulation of ECG outcomes in response to other variables. One critical assumption in this research was the notion of probabilistic rather than deterministic interactions, which recognises the inherent unpredictability and uncertainty in biological systems. It was also assumed that the data used were typical of the diverse athletic community, guaranteeing that the model's conclusions could be generalised. The discovery that age and training intensity have a substantial impact sheds light on how internal ageing processes and extrinsic demands of training regimens interact to affect cardiac electrophysiological characteristics.

### **Causal Relationship:**

In our analysis, the application of Bayesian networks revealed what appears to be a causal relationship, with age and training intensity having a direct influence on the ECG parameter V5. This conclusion is significant because V5 not only reflects individual cardiovascular responses, but it may also influence the link between other ECG parameters and characteristics including gender and sport discipline. In layman's terms, how old we are and how hard we train may have a specific effect on one aspect of our heart's electrical pattern, which may then ripple out to affect other parts of cardiac function. These effects may differ according to gender and sport type. Although these associations are correlative and additional research is needed to determine causation definitely, these findings pave the path for more personalised and nuanced training standards.

### **Inference and Findings**

The study's findings are based on a strong analytical strategy that includes Bayesian Networks, T-tests, and ANOVA, all of which target distinct aspects of athletes' cardiac health. Bayesian analysis revealed complex causal pathways in which age and training intensity directly influenced ECG features, particularly lead V5, which then mediated the impact of gender and sporting discipline. This demonstrates the diverse nature of cardiac electrophysiology and emphasises the importance of individualised athlete monitoring.

The T-tests revealed significant gender-based differences in ECG readings, particularly in the lateral and septal/anterior leads, implying gender-specific cardiac adaptations to physical training. These findings promote gender-specific training and medical interventions. Furthermore, ANOVA revealed significant differences in ECG patterns across different sports, which were supported by the Tukey HSD post-hoc analysis. This demonstrates that the type of sport has a significant impact on athletes' cardiac electrical activity, necessitating sport-specific cardiovascular examination and training regimes.

Overall, the findings from these studies support a complex approach to athlete health that takes into account the interdependence of age, training, gender, and sport discipline. These findings help to deepen our understanding of sports cardiology and pave the road for personalised, precision-based athlete treatment and training optimisation.

### **Conclusion**

Several key findings emerged from this thorough investigation, contributing to our understanding of cardiac electrophysiology in athletes. The exploratory data analysis (EDA) and thorough data cleaning provided a solid platform for the use of advanced statistical methods, ensuring the integrity and dependability of the results. Stratified sampling was critical in ensuring representative and unbiased data across all sports disciplines and genders, which was especially important given the diverse distribution of participants.

Subsequent analyses, including Bayesian Networks, independent t-tests, and ANOVA, revealed the complex correlations between age, training intensity, gender, and sport type on ECG parameters. The study highlighted how these parameters collectively influence cardiac electrical activity, emphasising the necessity of personalised treatments in sports medicine.

In essence, the initiative not only contributes to the area of sports cardiology by identifying crucial aspects that influence athletes' ECG readings, but it also pushes for personalised training and healthcare. It promotes the use of precision medicine in athletic training and healthcare, with the goal of achieving peak performance while protecting players' health. This endeavour has established a precedent for future research and practical applications, demonstrating the importance of data-driven approaches to improving athlete care.

### **Future Research:**

The next step of this research will focus on harnessing the amount of information contained within raw ECG signals by extracting comprehensive variables such as heart rate variability, QRS complex features, and ST-segment abnormalities. We hope to distil the raw ECG data into a set of potential indicators for cardiovascular health and athletic performance by using sophisticated signal processing techniques such as wavelet transforms and Fourier analysis.

Once these variables have been identified and confirmed, the dataset will be enhanced, making it ideal for the use of advanced machine learning techniques. Random forests, support vector machines, and neural networks will be used to identify patterns and relationships in the data that standard statistical methods may miss. The goal is to create prediction algorithms that can not only detect tiny variances throughout the athlete spectrum, but also identify probable cardiac problems, revolutionising preventative healthcare in sports medicine.

Furthermore, continuous monitoring and real-time data analysis enabled by wearable ECG technology will create new opportunities for in-the-moment assessments, allowing for quick feedback and adaptable training tactics. This strategy will ensure the seamless integration of data science and clinical expertise, paving the door for personalised and dynamic athlete management solutions.

### **Acknowledgment:**

*Norwegian Endurance Athlete ECG Database*

I will thank Professor Emeritus Knut Gjessdal for providing his medical expertise and interpreting all of the ECGs. This work was done at the University of Oslo and I will thank Professor Ørjan Grøttem Martinsen for providing appropriate facilities for ECG measurements.

## References:

- [1] Prior D. L. and Gerche A. L., "The athlete's heart," *Heart*, vol. 98, no. 12, pp. 947–955, Jun. 2012, doi: 10.1136/heartjnl-2011-301329. [PubMed] [CrossRef] [Google Scholar]
- [2] Drezner J. A. et al., "International criteria for electrocardiographic interpretation in athletes: Consensus statement," *Brit. J. Sports Med.*, vol. 51, no. 9, pp. 704–731, May 2017, doi: 10.1136/bjsports-2016-097331. [PubMed] [CrossRef] [Google Scholar]
- [3] Stokstad M. T., Berge H. M., and Gjesdal K., "Hjertescreening av unge idrettsutøvere," *Tidsskrift for Den norske legeforening*, vol. 133, no. 16, pp. 1722–1725, 2013, doi: 10.4045/tidsskr.13.0016. [PubMed] [CrossRef] [Google Scholar]
- [4] Berge H. M., Gjesdal K., Andersen T. E., Solberg E. E., and Steine K., "Prevalence of abnormal ECGs in male soccer players decreases with the Seattle criteria, but is still high," *Scand. J. Med. Sci. Sports*, vol. 25, no. 4, pp. 501–508, 2015, doi: 10.1111/sms.12274. [PubMed] [CrossRef] [Google Scholar]
- [5] Drezner J. A. et al., "Electrocardiographic interpretation in athletes: The 'Seattle criteria'," *Brit. J. Sports Med.*, vol. 47, no. 3, pp. 122–124, Feb. 2013, doi: 10.1136/bjsports-2012-092067. [PubMed] [CrossRef] [Google Scholar]
- [6] Abhimanyu U. et al., "Interpretation of the electrocardiogram of young athletes," *Circulation*, vol. 124, no. 6, pp. 746–757, Aug. 2011, doi: 10.1161/CIRCULATIONAHA.110.013078. [PubMed] [CrossRef] [Google Scholar]
- [7] Nabeel S. et al., "Comparison of electrocardiographic criteria for the detection of cardiac abnormalities in elite black and white athletes," *Circulation*, vol. 129, no. 16, pp. 1637–1649, Apr. 2014, doi: 10.1161/CIRCULATIONAHA.113.006179. [PubMed] [CrossRef] [Google Scholar]
- [8] Drezner J. A., "18 highlights from the International criteria for ECG interpretation in athletes," *Brit. J. Sports Med.*, vol. 54, no. 4, pp. 197–199, Feb. 2020, doi: 10.1136/bjsports-2019-101537. [PubMed] [CrossRef] [Google Scholar]
- [9] Bickerton M. and Pooler A., "Misplaced ECG electrodes and the need for continuing training," *Brit. J. Cardiac Nurs.*, vol. 14, no. 3, pp. 123–132, Mar. 2019, doi: 10.12968/bjca.2019.14.3.123. [CrossRef] [Google Scholar]
- [10] Berge H. M., Steine K., Andersen T. E., Solberg E. E., and Gjesdal K., "Visual or computer-based measurements: Important for interpretation of athletes' ECG," *Brit. J. Sports Med.*, vol. 48, no. 9, pp. 761–767, May 2014, doi: 10.1136/bjsports-2014-093412. [PubMed] [CrossRef] [Google Scholar]