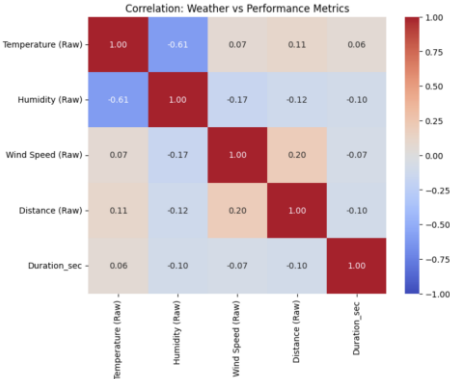# Data Analysis' Team (Week 1-6)

Project 3: ReflexionPro
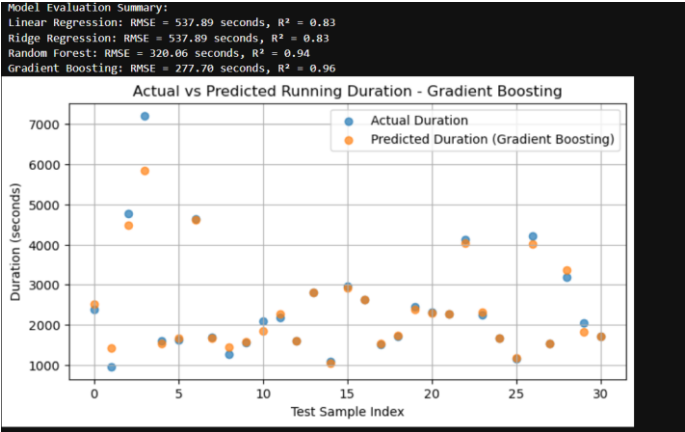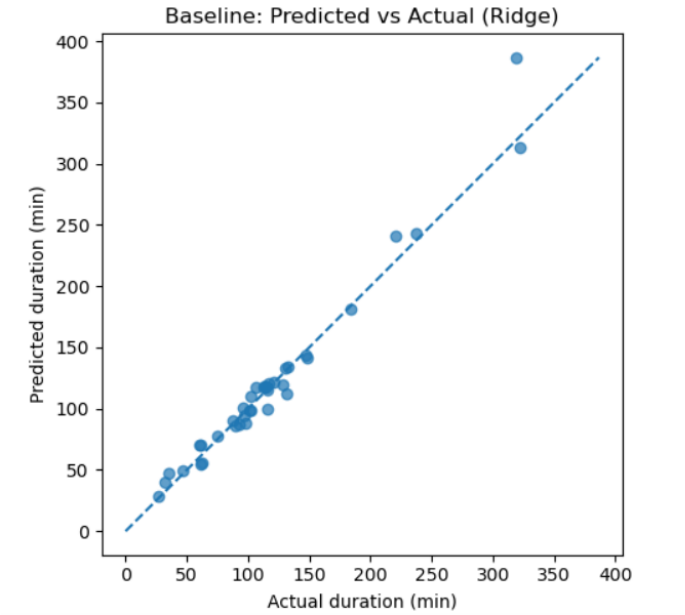
Mentor: Ben Dang

Co-leadership: Billy Thai & Shin Nguyen

| Name & Student | Task | Contribution & Evidence (Code & Output) |
|---|---|---|
| 1. Billy Thai (223339662) | - Expanded dataset & retrain model<br>- Develop model ( Fatigue prediction, injury risk assessment, recovery time)<br>- Visualise dashboard using Power BI ( ongoing) | **Expanded dataset**<br>- Transform the Garmin activity log into an analysis-ready dataset that captures workload and physiological context.<br><br>Link: <br><br>**Fatigue Prediction Model**<br>- Designed and implemented an end-to-end machine-learning pipeline that predicts next-day fatigue status for athletes by analysing their wearable-derived training data.<br>- Generated permutation importance plots, which highlighted acute load and chronic load as the two most influential predictors. |
| 2. Shin Nguyen | - I have successfully prepared, cleaned, and processed raw fitness tracker datasets (Running & Cycling) into an analysis-ready format, implementing robust preprocessing steps such as datetime conversion, duration transformation, numeric coercion, and missing value handling.<br>- I designed and implemented machine learning models (Random Forest Regressor) to predict running | <br><br>- The correlation heatmap was generated to explore relationships among the input variables and the target variable. The heatmap highlights both **positive and negative correlations**, with darker shades representing stronger relationships. This visualization helps to identify which features are most influential for predicting the outcome and to detect possible multicollinearity among predictors. |

| | | | |
|---|---|---|---|
| | duration based on distance, elevation, pace, and environmental data, achieving an average R² score of 0.7863, and created correlation heatmaps as well to analyse performance and environmental impact | ```
# Train and Evaluate Running Duration Model
# Initialize model
rf_model = RandomForestRegressor(random_state=42)

# Evaluate with our reusable CV function
evaluate_model(rf_model, X_run, y_run, scoring='r2')

Cross-validation r2 scores: [0.93727729 0.10664602 0.92689275 0.98167289 0.97914234]
Average r2: 0.7863
array([0.93727729, 0.10664602, 0.92689275, 0.98167289, 0.97914234])
``` | - The cross- |
| | | validation R² scores for the model were: **[0.937, 0.107, 0.927, 0.982, 0.979]**. These results indicate that the model performed well on most folds, with R² values close to 1.0, suggesting a strong predictive power.<br>- The **average R² score of 0.7863** demonstrates that the model generally explains a high proportion of the variance, but improvements may be needed to ensure more consistent performance across folds. | |
| 3. Benjamin Cole | | | |
| 4. ASMITHRA RAVICHANDRAN | - Correlate environmental factors (Wind speed, Humidity, Temperature) with athlete performance<br>- Collect and clean sample humidity data<br>- Link humidity values with athlete performance metrics (speed, stamina)<br>- Generate correlation matrix and visualisation | • Wrote a Python script to clean and merge humidity + performance dataset.<br>• Calculated correlation between humidity and athlete speed.<br>• Generated a scatter plot showing negative correlation (as humidity ↑, speed ↓).<br>• Correlation Matrix:<br><br>• Scatter plot:<br> | |
| 5. WIMANSA DESHANI | | | |
| 6. Chiranth Gowda | | | |
| 7. Dhivyaa Lakshimi (224744156) | - Correlate environmental factors (Wind speed, Humidity, | - Cleaned the Garmin dataset (`activities_cleaned.csv`) and prepared for analysis. | |

| | | |
|---|---|---|
| | Temperature) with athlete performance<br>- Completed **Wind Speed** analysis | - Focused on Wind Speed vs Athlete Avg Speed.<br>- Used Google Colab Notebook (`windspeed.ipynb`) to run analysis.<br>- Calculated correlation:<br>- **Pearson = 0.12**<br>- **Spearman = 0.09**<br>- Produced outputs:<br>- **Summary table** (`wind_summary.csv`)<br>- **Scatter plot** (Wind Speed vs Avg Speed with regression line)<br>- **Correlation heatmap** (Wind, AvgSpeed, Temp, Humidity)<br><br><br><br> |
| 8. Hrithik Joseph (223058217) | Build a predictive model for running duration | * Developed an Accurate Predictive Model for Running Duration<br>I successfully built and implemented a robust regression model that predicts running duration using key input variables such as distance, elevation gain, and pace<br><br>*Enhanced Model Reliability through Multi-Model Evaluation and Visualization<br>I evaluated multiple machine learning algorithms—including Linear Regression, Ridge Regression, Random Forest, and Gradient Boosting—to identify the best performing model based on RMSE and $R^2$ |

| | | |
|---|---|---|
| | | metrics. I also created clear visualizations comparing actual and predicted durations, enhancing interpretability and actionable insights for the coaching team.<br><br>Model Evaluation Summary:<br>Linear Regression: RMSE = 537.89 seconds, R² = 0.83<br>Ridge Regression: RMSE = 537.89 seconds, R² = 0.83<br>Random Forest: RMSE = 320.06 seconds, R² = 0.94<br>Gradient Boosting: RMSE = 277.70 seconds, R² = 0.96 |
| 9. Jeswin Joy<br>(224832817) | | I loaded the dataset, standardised key columns, checked missing values, explored distributions, and flagged the outliers.<br>I normalised the unit duration from hours to minutes and distance to kilometres, removed the duplicates and incomplete rows.<br>I trained two baselines to predict duration in minutes: Linear Regression and Ridge with cross-validated regularisation in a scaling pipeline.<br>I reported MAE, RMSE, and $R^2$ and plotted predicted vs actual values. |
| 10. Madhav Grover | | |
| 11. Nathan Riley | Progressed the predicting sleep efficiency using the sleep efficiency dataset | Sleep Efficiency Prediction<br>The sleep efficiency dataset tracks athletes and non-athletes sleep and rates the efficiency of that sleep |

| | | |
|---|---|---|
| | (ongoing) Srihari and I working on our assigned task regarding documentation | by analysing metrics such as no. of exercises performed, alcohol consumed, duration, etc. Using models, I'm attempting to analyse and predict these ratings so it can be used as a tool for athletes to achieve the highest rating possible so it can lead to peak performance by maximising recovery. I have incorporated a couple new models, and having determined RF as the optimal model for the prediction, I'm at the concluding stages where I'm inputting new data which can hopefully strengthen the model after comparing it to some previous entries. I'll begin looking into new things I can target this week. The task with Srihari is ongoing with discussions still going as well. |
| | |  |
| 12. Saumya Parasbhai | Perform cross validations on datasets | Cricket Performance Prediction developed and evaluated a number of prediction models that use batting characteristics like balls faced, strike rate, boundaries, and sixes to estimate the number of runs scored. Five-fold cross-validation was used for KNN, Random Forest, and Linear Regression. With the lowest mean absolute error and the highest R2 score, the results demonstrated that Linear Regression consistently performed better than the others. This demonstrated that the model was dependable for analyzing cricket performance since runs scored were well clarified by linear relationships between batting inputs and results. |

```
Linear Regression:
  MAE: 3.061183925144868
  R²: 0.976829429446577

Random Forest:
  MAE: 3.543190476190476
  R²: 0.9494601834358297

KNN:
  MAE: 4.62
  R²: 0.9313317047061667
```

**Model Comparison with 5-Fold Cross Validation**

| Average MAE (Lower is Better) | Average R² (Higher is Better) |
| --- | --- |



**VO$_2$ Max Prediction**

Using physiological and training data from Garmin, a cross-validated regression model was developed to forecast athletes' VO$_2$ Max values. With a high R2 and a low mean absolute error across folds, the Random Forest model showed good predictive power. Consistent model stability was validated by visualizing fold-wise results, suggesting that the model performs well when applied to athlete data that has not yet been seen. The possibility of using wearable data to track physical activity and provide data-driven performance insights is confirmed by these findings.

```
Average MAE: 0.8772196318402095
Average R²: 0.9241144974941177
```

**K-Fold Cross-Validation Results - VO$_2$ Max Prediction**



| 13. Scott Ailaiti (224348919) | Exploratory Data Analysis (EDA)<br><br>Data Cleaning & Preprocessing<br><br>Feature Engineering<br><br>Data Transformation (Normalization & Encoding) | I contributed to the data preprocessing pipeline by performing essential cleaning and transformation tasks on the provided running dataset. My responsibilities included loading and previewing the data, checking and confirming the absence of missing values, and removing duplicate rows and extreme outliers (e.g., heart rates above 250 bpm) to ensure data quality. I also implemented a time conversion function that transformed "Duration (h:m:s)" and "Moving Duration (h:m:s)" columns into total seconds to facilitate accurate numerical analysis. These foundational steps laid the |
| --- | --- | --- |

| | Dataset Splitting (Train/Test) Implement Initial Model Prototypes Experiment Tracking (Note capture & screenshots) | groundwork for future modeling by ensuring the dataset was clean, consistent, and ready for further exploration and machine learning tasks.  |
|---|---|---|
| 14. Siddh Parasbhai | Perform Cross Validation on the datasets and get insights. | An end-to-end machine learning pipeline was created to forecast the aerobic capacity ($VO_2max$) of cyclists based on power output data from 20-kilometer time trials. With an R2 of -0.52, the model showed low efficacy, failing to capture significant physiological correlations between $VO_2max$ and cycling performance. The limited physiological range of $VO_2max$ values (3.0-6.0 ml/kg/min) concealed basic model flaws, even with numerically minor mistakes (MAE: 0.76 ml/kg/min). The main failure point identified by permutation importance analysis was insufficient feature representation, underscoring the necessity of including extra physiological indicators such as heart rate variability, cadence patterns, and individual anthropometric data in order to establish reliable power-to-$VO_2max$ relationships for precise fitness evaluation. |

| | | |
|---|---|---|
| | |  Predicted vs. Actual VO₂max (ml/kg/min) |

created a predictive system that uses power output measurements to estimate the heart rate response in real time during cycling time trials. With clinically relevant error margins (MAE: 7.33 bpm) and moderate explanatory power (R2: 0.39), the model captured fundamental patterns of exertion but lacked individual precision. Power output fluctuations were found to be the main driver by feature importance analysis, which also revealed significant gaps in the accounting for environmental components and cumulative fatigue effects. The methodology has promise for monitoring coarse exertion in training scenarios; nevertheless, in order to attain actionable accuracy for performance optimization, temporal factors such as acute:chronic workload ratios, temperature acclimatization measurements, and individual recovery patterns must be integrated.

Prediction for Heart Rate (bpm):
Average R²: 0.39
Average MAE: 7.33
Average RMSE: 8.21


Predicted vs. Actual Heart Rate (bpm)

| | | |
|---|---|---|
| 15. Srihari Chandran | | |
| 16. Xiangning Wang | **Analyse personal activity metrics to identify fitness pattern** | 1. Total_minutes<br>2. Active_minutes_pct<br>3. Active_distiance_pct |

| | **Identify metrics related to personal activity within provided dataset**<br><br>**Data summary and exploration, including visualisations and data normalisation** | These 3 new variables can be considered as composite variables.<br><br>This is the basic summary for activity related variables data<br><br><br><br>Here is the activity intensity distribution by minutes portrayed in a boxplot:<br><br><br><br>Here is the correlation matrix between all related variable/metrics<br><br><br><br>Based on the summary statistics and distribution, we have created the following activity profiles: Highly Active, Moderately Active, Balanced, Lightly Active and Sedentary. Here is the adjusted distribution visualised through bar chart. |
|---|---|---|

| | | |
|---|---|---|
| | | Distribution of Activity Profiles<br><br> |
| | | |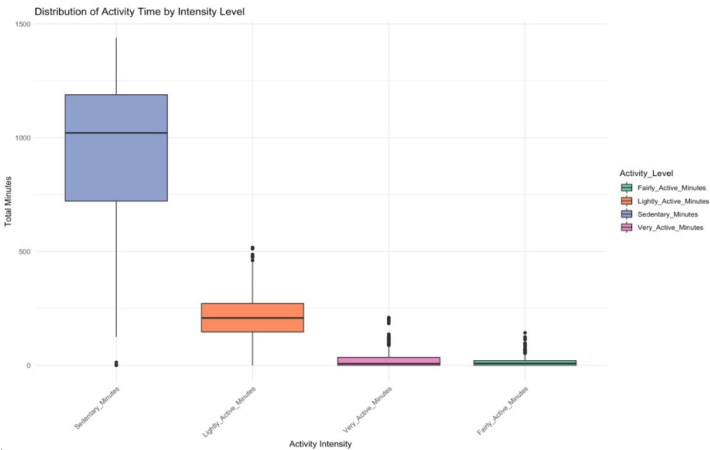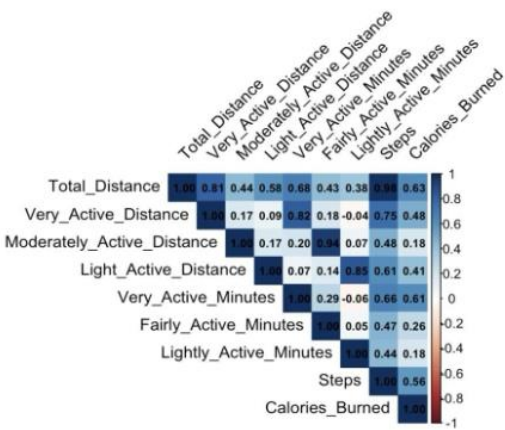