

Marathon Race Time Prediction

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.1      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.1
v purrr      1.0.2
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(fs)
library(here)
```

here() starts at /Users/ubc/marathon-analysis

```
library(tidymodels)
```

```
-- Attaching packages ----- tidymodels 1.2.0 --
v broom       1.0.6      v rsample     1.2.1
v dials       1.3.0      v tune        1.2.1
v infer       1.0.7      v workflows   1.1.4
v modeldata   1.4.0      v workflowsets 1.1.0
v parsnip     1.2.1      v yardstick   1.3.1
v recipes     1.1.0
-- Conflicts ----- tidymodels_conflicts() --
x scales::discard() masks purrr::discard()
x dplyr::filter()   masks stats::filter()
```

```
x recipes::fixed() masks stringr::fixed()
x dplyr::lag()      masks stats::lag()
x yardstick::spec() masks readr::spec()
x recipes::step()   masks stats::step()
* Learn how to get started at https://www.tidymodels.org/start/
```

```
library(testthat)
```

Attaching package: 'testthat'

The following object is masked from 'package:rsample':

matches

The following object is masked from 'package:dplyr':

matches

The following object is masked from 'package:purrr':

is_null

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

edition_get, local_edition

The following object is masked from 'package:tidyr':

matches

Summary

Over the last ten years, marathon running has become a popular physical activity around the world. This study aims to investigate the relationship between marathon runners' maximum distance ran per week during race training and their race time. We modeled a simple linear regression to predict marathon times based on their training patterns and tested on the model with the test dataset. By having this analysis, it helps marathon lovers to gain insights into how training volume influences race performance and to better prepare the race.

Introduction

Over the last ten years, marathon running has become a popular physical activity around the world (Zoladz & Nieckarz, 2021). It is commonly known as a high-level endurance exercise that requires the runners to have dedicated training (Kaufmann et al., 2020). The maximum distance ran per week during race training is a key metric that marathon lovers would care about and is a key reference to predict which athletes will perform better than others. Thus, this study will investigate how the maximum distance ran per week during race training will predict the time it takes a runner to finish the race. Specifically, we want to answer the question: What predicts which athletes will perform better than others? How the maximum distance ran per week (in miles) during race training predicts the time it takes a runner to finish the race? This study uses the dataset from a public dataset on GitHub, containing 13 variables about runners, such as age, bmi, maximum training distance per week (max, in miles) and actual marathon race time (time_hrs, in hours) and so on.

```
# download and save data

if (!fs::dir_exists(here::here("data"))) {
  fs::dir_create(here::here("data"))
}

if (!fs::file_exists(here::here("data/marathon.csv"))) {
  url <- "https://raw.githubusercontent.com/UBC-DSCI/dsci-100-student/refs/heads/master/marathon.csv"
  marathon <- readr::read_csv(url)
  readr::write_csv(marathon, here::here("data/marathon.csv"))
}
```

```
# read saved data
marathon <- readr::read_csv(here::here("data/marathon.csv"))
```

Rows: 929 Columns: 13

-- Column specification -----

Delimiter: ","

dbl (13): age, bmi, female, footwear, group, injury, mf_d, mf_di, mf_ti, max...

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.

Data Validation Check

Correct Data File Format

```
if (!inherits(marathon, "data.frame")) {  
  stop("The data file is not in the correct format. Expected a CSV to be read as a data frame")  
}
```

No Empty Observations

```
empty_rows <- apply(marathon, 1, function(row) all(is.na(row)))  
if (any(empty_rows)) {  
  warning("There are completely empty observations. Consider removing these rows.")  
  marathon <- marathon[!empty_rows, ]  
}
```

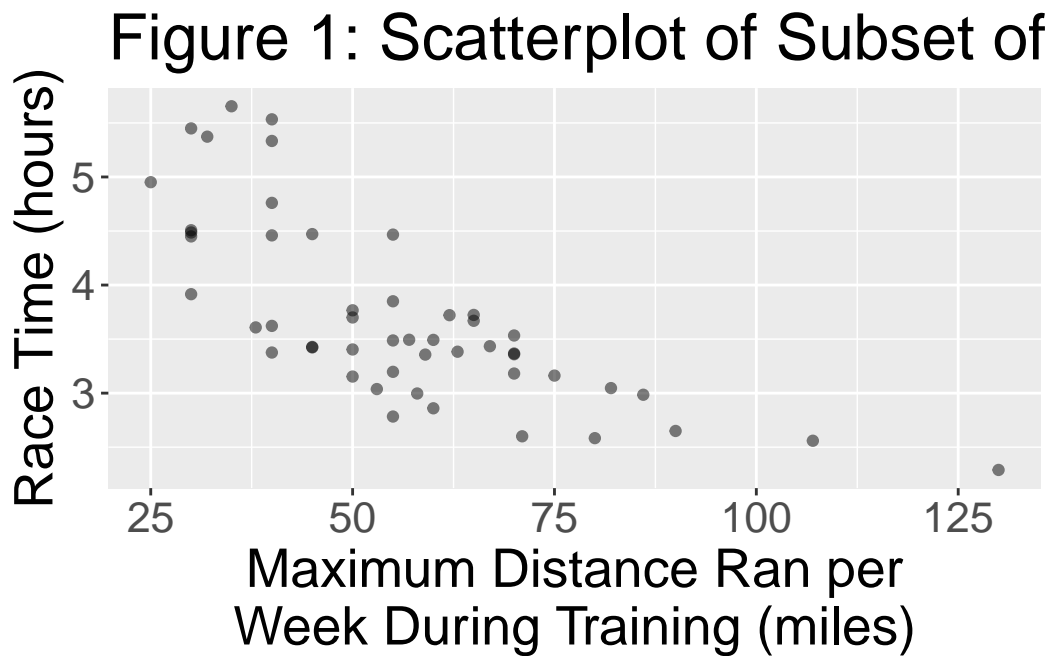
No Duplicate Observations

```
if (any(duplicated(marathon))) {  
  warning("There are duplicate observations. Removing them now.")  
  marathon <- marathon[!duplicated(marathon), ]  
}
```

We want to predict race time (in hours) (time_hrs) given a particular value of maximum distance ran per week (in miles) during race training (max). With this subset, we can plot a scatterplot to assess the relationship between these two variables.

```
set.seed(2000)  
  
marathon_50 <- marathon |>  
  dplyr::sample_n(50)  
  
marathon_50 |>  
  ggplot(aes(x = max, y = time_hrs)) +  
  geom_point(alpha = 0.5) +  
  xlab("Maximum Distance Ran per\nWeek During Training (miles)") +  
  ylab("Race Time (hours)") +
```

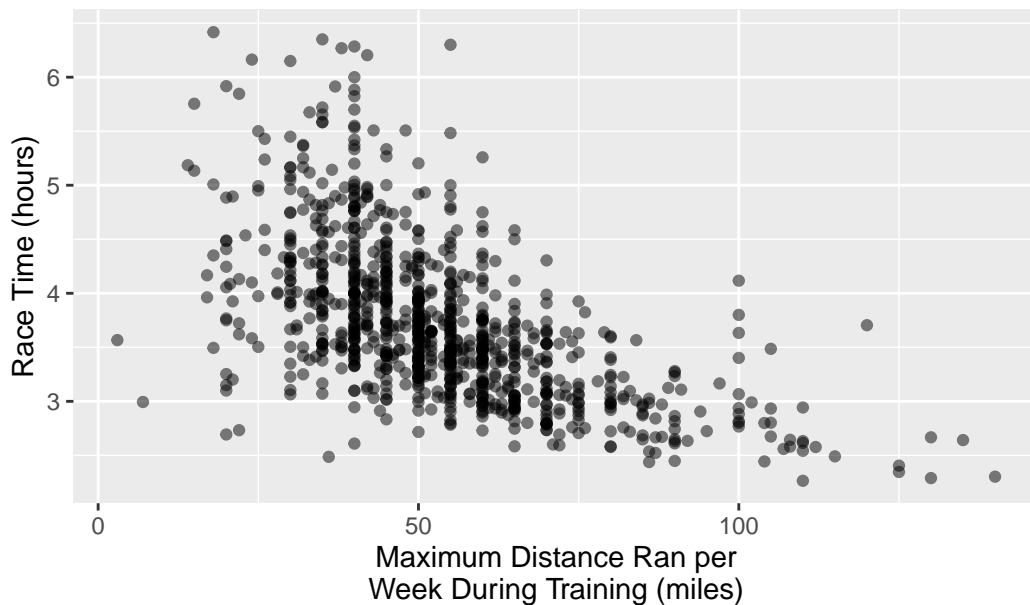
```
theme(text = element_text(size = 20)) +
ggtitle("Figure 1: Scatterplot of Subset of Maximum Distance Ran per Week vs. Race Time")
```



Here's the full dataset of values.

```
ggplot(marathon, aes(x = max, y = time_hrs)) +
  geom_point(alpha = 0.5) +
  xlab("Maximum Distance Ran per\nWeek During Training (miles)") +
  ylab("Race Time (hours)") +
  theme(text = element_text(size = )) +
  ggtitle("Figure 2: Scatterplot of Full Dataset Maximum Distance Ran per Week vs. Race Time")
```

Figure 2: Scatterplot of Full Dataset Maximum Distance Ran per



Analysis

We will analyze the data using simple linear regression to assess the relationship between maximum weekly distance and marathon race time. We will first split the dataset into the training and testing datasets, using 75% of the original data as the training data. The training set was used to fit the model, while the test set was used for performance evaluation. In the `strata` argument of the `initial_split` function, we will use the variable we are trying to predict.

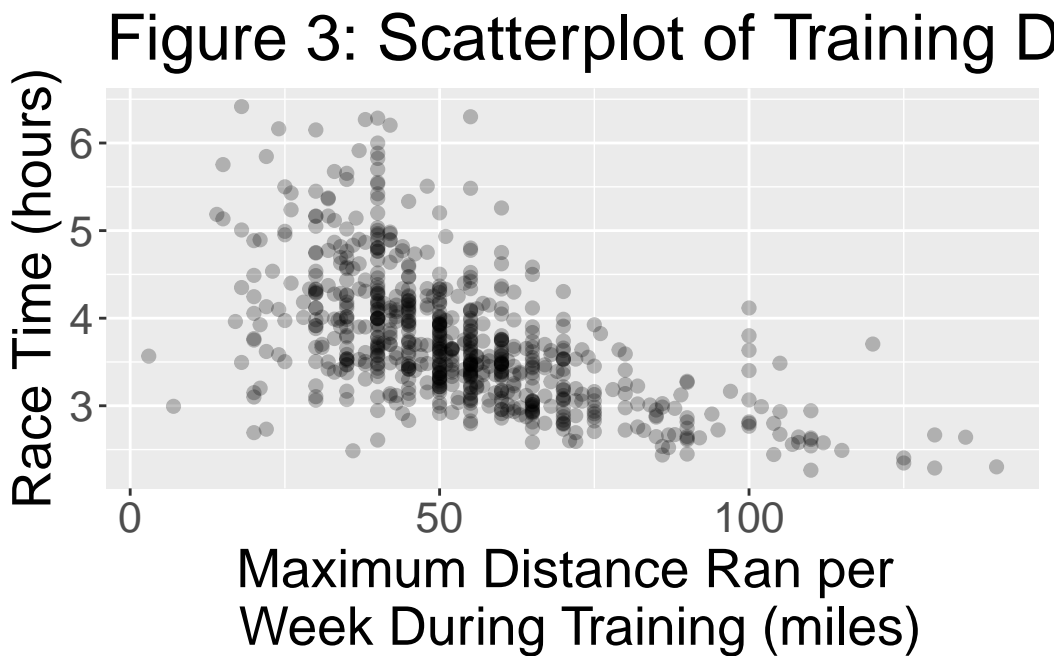
Training Data

```
set.seed(2000)

marathon_split <- rsample::initial_split(marathon, prop = 0.75, strata = time_hrs)
marathon_training <- rsample::training(marathon_split)
marathon_testing <- rsample::testing(marathon_split)
```

We can look at a scatterplot to assess the relationship between race time (`time_hrs`) and maximum distance ran per week during training (`max`) using only the observations in the training dataset.

```
marathon_training |>
  ggplot(aes(x = max, y = time_hrs)) +
  geom_point(alpha = 0.25, size = 2) +
  xlab("Maximum Distance Ran per \n Week During Training (miles)") +
  ylab("Race Time (hours)") +
  theme(text = element_text(size = 20)) +
  ggtitle("Figure 3: Scatterplot of Training Dataset Maximum Distance Ran per Week vs. Race Time")
```



Linear Regression

Now that we have our training data, the next step is to build a linear regression model specification.

```
lm_spec <- parsnip::linear_reg() |>
  parsnip::set_engine("lm") |>
  parsnip::set_mode("regression")

lm_spec
```

Linear Regression Model Specification (regression)

Computational engine: lm

After we have created our linear regression model specification, the next step is to create a recipe, establish a workflow analysis and fit our simple linear regression model.

```
lm_recipe <- recipes::recipe(time_hrs ~ max, data = marathon_training)

lm_fit <- workflows::workflow() |>
  workflows::add_recipe(lm_recipe) |>
  workflows::add_model(lm_spec) |>
  parsnip::fit(data = marathon_training)
```

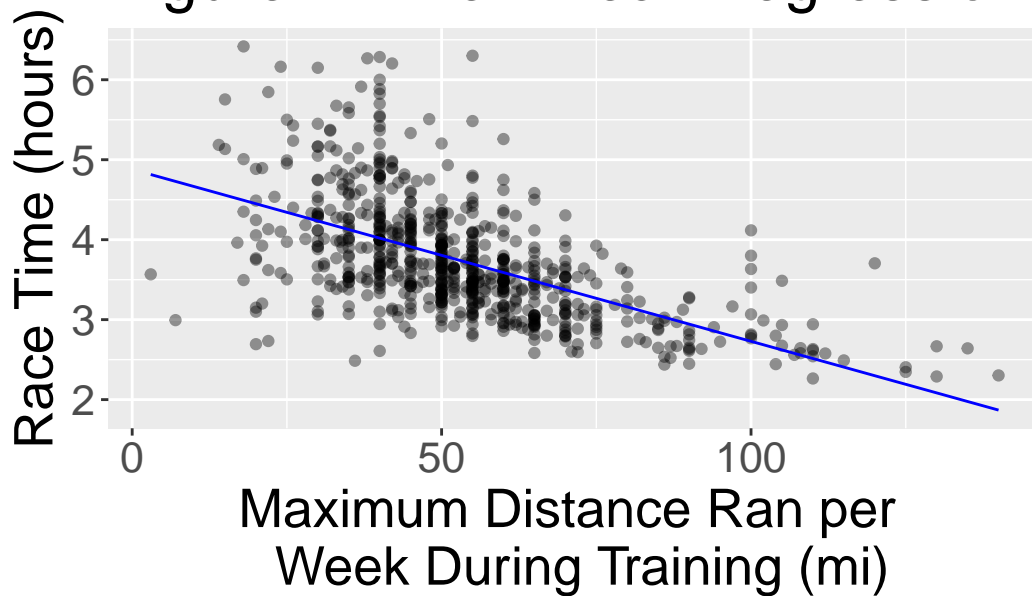
Results

Now, let's visualize the model predictions as a straight line overlaid on the training data.

```
marathon_preds <- lm_fit |>
  predict(marathon_training) |>
  dplyr::bind_cols(marathon_training)
```

```
marathon_preds |>
  ggplot(aes(x = max, y = time_hrs)) +
  geom_point(alpha = 0.4) +
  geom_line(
    mapping = aes(x = max, y = .pred),
    color = "blue") +
  xlab("Maximum Distance Ran per \n Week During Training (mi)") +
  ylab("Race Time (hours)") +
  theme(text = element_text(size = 20)) +
  ggtitle("Figure 4: The Linear Regression of Maximum Distance Ran per Week and Race Time")
```


Figure 4: The Linear Regression



Model Performance

We can look at our model performance by looking at the RMSPE on the test data.

```
lm_test_results <- lm_fit |>
  predict(marathon_testing) |>
  dplyr::bind_cols(marathon_testing) |>
  yardstick::metrics(truth = time_hrs, estimate = .pred)

lm_rmspe <- lm_test_results |>
  dplyr::filter(.metric == 'rmse') |>
  dplyr::select(.estimate) |>
  dplyr::pull()

lm_rmspe
```

```
[1] 0.5504829
```

The RMSPE stands for the Root Mean Squared Prediction Error, which is a scoring metric that reflects how right/wrong each prediction is. It measures the distance of the prediction from the actual values on the test data. Lower RMSPE values indicate better predictive performance,

meaning that the model's predicted values are close to the actual observed values (Taraji et al., 2017). In our case, the RMSPE represents the average error in predicting marathon race time (in hours) for the runners. From the result of 0.5504829, This means that, our simple linear regression model's predictions deviate from the actual race times by approximately 0.55 hours on average.

Prediction on test data

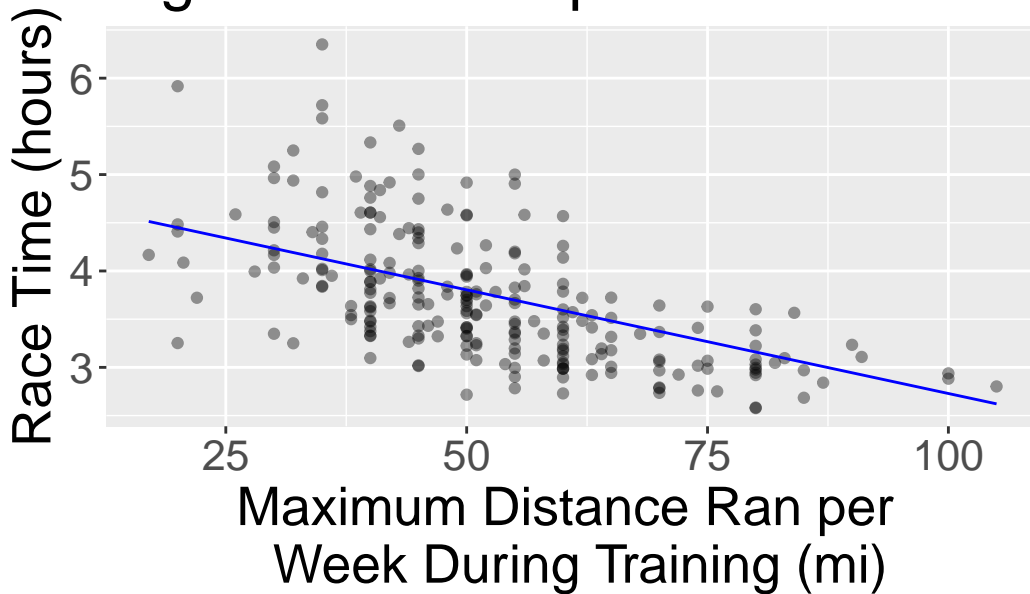
Visualize the model predictions as a straight line overlaid on the test data

```
test_preds <- lm_fit |>
  predict(marathon_testing) |>
  dplyr::bind_cols(marathon_testing)

lm_predictions_test <- test_preds |>
  ggplot(aes(x = max, y = time_hrs)) +
  geom_point(alpha = 0.4) +
  geom_line(
    mapping = aes(x = max, y = .pred),
    color = "blue") +
  ggtitle("Figure 5: Model predictions over test data") +
  xlab("Maximum Distance Ran per \n Week During Training (mi)") +
  ylab("Race Time (hours)") +
  theme(text = element_text(size = 20))

lm_predictions_test
```

Figure 5: Model predictions over t



```
lm_coefficients <- broom::tidy(lm_fit)
lm_coefficients
```

A tibble: 2 x 5

term	estimate	std.error	statistic	p.value
<chr>	<dbl>	<dbl>	<dbl>	<dbl>
1 (Intercept)	4.88	0.0651	74.9	0
2 max	-0.0215	0.00114	-18.8	6.03e-64

Interpretation

By having the result, the slope coefficient is -0.0215, which shows a negative relationship between the maximum distance ran per week and the race time. As the maximum distance ran increases by every additional mile per week, the race time will decrease by 0.0215 hours. It shows that adding more training distance per week is beneficial for improving marathon performance. The intercept is 4.88, which represents that when the maximum distance ran per week (max) is zero miles, the predicted race time is 4.88 hours. Although this scenario is unrealistic, it serves as the baseline value from which changes in race time are predicted based on the training distance. These coefficients with p-value of 6.03e-64 and 0 (less than 0.05) are statistically significant, meaning there is strong evidence that increasing the training mileage leads to faster marathon times.

If we want to manually calculate the marathon time, the formula will be like $times_{hrs} = 4.88 - 0.0215 \times max$, where $times_{hrs}$ is the race time and the max is the maximum distance (in miles) ran per week during training.

Discussion

However, there are several limitations of the model. As we are trying to find what predicts which athletes will perform better than others, other features about the runners may also be crucial to play roles in affecting the race time, such as their age, bmi. Therefore, we may need to add more features to the model to better predict the race time. Also, these features may not in a linear relationship with the race time, so using a Random Forest model to handle the complex, non-linear interactions between features may be useful.

To know whether the new model would be better than the simple linear regression or not, we can use those scoring metrics such as Mean squared error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE) and R-squared (R^2) on a test set or through cross-validation. A lower MSE, RMSE, MAE, MAPE would indicate that the new model makes more accurate predictions compared to the simple linear regression model. A higher R-squared would suggest that the model is explaining more variance in the target variable, implying better performance (Kolhatkar & Östblom, 2023).

Citation

Zoladz, J. A., & Nieckarz, Z. (2021). Marathon race performance increases the amount of particulate matter deposited in the respiratory system of runners: an incentive for “clean air marathon runs”. *PeerJ*, 9, e11562. <https://doi.org/10.7717/peerj.11562>

Kaufmann, C. C., Wegberger, C., Tscharre, M., Haller, P. M., Piackova, E., Vujasin, I., Kassem, M., Tentzeris, I., & Freynhofer, M. K. (2020). Effect of marathon and ultra-marathon on inflammation and iron homeostasis. *Scandinavian Journal of Medicine & Science in Sports*. <https://doi.org/10.1111/sms.13869>

Taraji, M., Haddad, P. R., Amos, R. I. J., Talebi, M., Szucs, R., Dolan, J. W., & Pohl, C. A. (2017). Error measures in quantitative structure-retention relationships studies. *Journal of Chromatography A*, 1524, 298-302. <https://doi.org/10.1016/j.chroma.2017.09.050>

Kolhatkar, V., & Östblom, J. (2023). Regression metrics [Lecture notes]. UBC GitHub Pages. https://pages.github.ubc.ca/mds-2024-25/DSCI_573_feat-model-select_students/lectures/02_regression-metrics.html