

GRUPPO: CORES

STATISTICA COMPUTAZIONALE – REPORT FINALE

Maccianti Federico Rapacioli Nicola Riva Pietro
(909656) (915439) (908813)

1 Introduzione

Il presente report analizza un dataset ottenuto tramite il repository [TracingInsights](#)¹ Lo studio si concentra sui dati di telemetria relativi alle sessioni di qualifica della stagione di Formula 1 2025, selezionando per ciascun pilota il singolo giro migliore.

Il dataset originale include le seguenti variabili:

Variabile	Unità di misura	Supporto	Tipo
Gran premio	–	{Nomi GP}	character
Pilota	–	{Sigle Piloti}	character
Tempo dal via	s	\mathbb{R}^+	numeric
Distanza percorsa	m	\mathbb{R}^+	numeric
Distanza relativa	–	$[0, 1]$	numeric
Velocità	km/h	\mathbb{N}	numeric
Regime motore	RPM	\mathbb{N}	numeric
Marcia	–	$\{1, \dots, 8\}$	integer
Freno	–	$\{0, 1\}$	factor
Acceleratore	%	$[0, 100]$	numeric
DRS	–	$\{0, 1\}$	factor
Accelerazione laterale e longitudinale	g	\mathbb{R}	numeric
Coordinate spaziali x,y,z	m	\mathbb{R}	numeric
Tempo Giro	s	\mathbb{R}^+	numeric
Gomma	–	{Tipi di gomma}	character
Vita della gomma	–	\mathbb{N}	integer
Strategia	–	$\{1, 12, 21\}$	integer

Il dataset presenta le seguenti codifiche: le feature binarie assumono valore 0 in assenza dell’evento e 1 in sua presenza. La telemetria dell’acceleratore misura l’intensità dell’input del pilota, mentre per le distanze indicano la posizione progressiva rispetto allo start. La ‘vita della gomma’ è definita come il numero di giri già percorsi dallo pneumatico in uso. La variabile ‘strategia’ identifica il pattern di utilizzo delle gomme: 1 per l’uso esclusivo della Soft, 12 per la transizione Soft \rightarrow Medium e 21 per la sequenza inversa Medium \rightarrow Soft.

L’obiettivo del presente report è analizzare tramite il dataset descritto le qualifiche del 2025, al fine di caratterizzare il comportamento dei piloti durante il giro di qualifica attraverso le principali variabili telemetriche disponibili.

2 Analisi Esplorativa

2.1 Considerazioni sulle variabili

Poiché lo stile di guida non è riconducibile a variabili di tipo posizionale, le coordinate spaziali e le misure di distanza, sia assolute sia relative, vengono escluse dall’analisi.

In questa fase preliminare, lo stile di guida viene descritto attraverso variabili dinamiche quali l’utilizzo dell’acceleratore e del freno e le accelerazioni, longitudinale e laterale, che consentono di caratterizzare rispettivamente le modalità di decelerazione in ingresso curva, l’intensità con cui la curva viene affrontata comprendendo anche altri tipi di comportamenti riferibili allo stile di guida. A

¹Codice di estrazione dati nell’appendice A.

sostegno delle ipotesi sopra citate, si riporta la Figura 1 che confronta nel Gran Premio degli USA le accelerazioni per i piloti: Charles Leclerc, Lando Norris e Max Verstappen. Notando infatti come nelle variazioni repentine si contraddistinguono meglio i piloti.

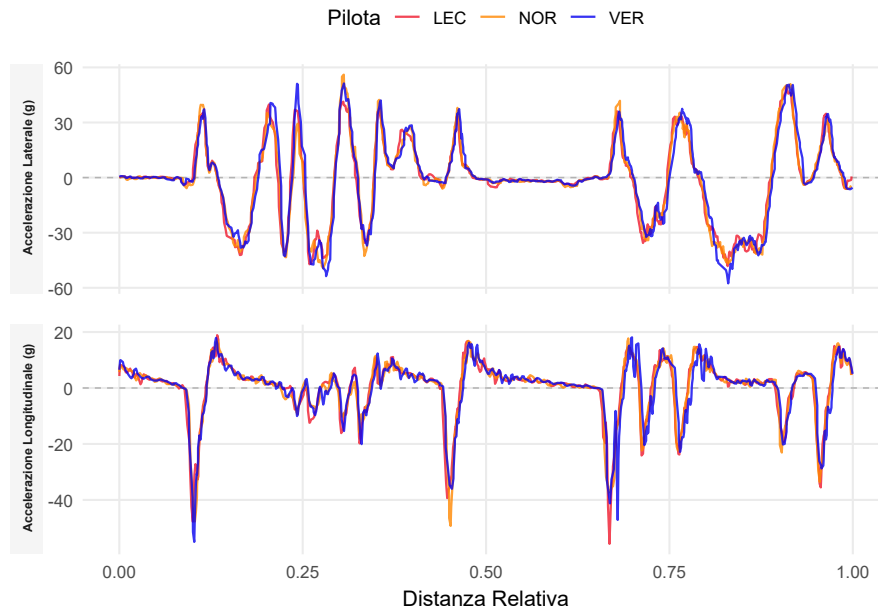


Figura 1: Accelerazioni VER, LEC, NOR.

Dalle prime analisi descrittive sul dataset originale², viene osservato come per il pilota Russell al Gran Premio di Miami si producano degli NA per quanto riguarda la variabile distanza relativa. Ricontrollando i dati grezzi ci si accorge di come probabilmente i sensori telemetrici abbiano avuto un'avaria, in quanto i record sono la maggior risulta pari a zero e le accelerazioni costanti durante l'intero giro. Per i motivi elencati sopra, si procede dunque ad eliminare il record.

Un ultimo accorgimento riguarda la variabile Tempo Giro, per alcuni piloti risulta infatti **None**, osservando i dati², si verifica per i piloti Tsunoda, Bearman e Hadjar rispettivamente nei Gran premi di Emilia Romagna, Australia e Stati Uniti.

Analizzando la realtà si scopre che il motivo è che non hanno terminato il giro di qualifica a causa di un incidente. I dati che si hanno a disposizione sono infatti il riscaldamento delle gomme come si nota nella figura 2. Anche se non è presente un tempo di giro si è registrata comunque la telemetria. Al fine di eliminare rumore nell'analisi si rimuovono queste osservazioni.



Figura 2: Confronto Emilia Romagna Tsunoda

²Estratto codice R consultabile nell'appendice B.

2.1.1 Trasformazione e creazione di nuove variabili

Al fine di rendere confrontabili i diversi tracciati dei Gran Premi, la variabile di accelerazione laterale viene riscalata³, per ciascun Gran Premio e Pilota, nell'intervallo $[-1, 1]$. Successivamente, viene trasformata tramite valore assoluto, così da ottenere una misura dell'intensità complessiva (magnitudo), il supporto quindi risulta compreso in $[0, 1]$.

Per l'accelerazione longitudinale vengono create due nuove variabili³.

Accelerazione: essa assimila tutti i valori positivi della variabile originaria, viene successivamente normalizzata nell'intervallo $[0, 1]$ per ciascun Gran Premio e Pilota, per riuscire a normalizzare i tracciati.

Decelerazione : essa assimila al contrario della precedente tutti i valori negativi della variabile originaria, successivamente viene applicato il valore atteso e normalizzata nell'intervallo $[0, 1]$. Per le restanti variabili telemetriche si mantiene invece la forma grezza.

Si creano tre nuove variabili di variazione percentuali⁴ avendo notato dalla figura 1 che quando ci sono variazioni repentine che si contraddistinguono lo stile di guida dei piloti:

- `_lag1` : variazione percentuale rispetto all'osservazione precedente
- `_lag3` : variazione percentuale rispetto a tre osservazioni precedenti
- `_lag5` : variazione percentuale rispetto a cinque osservazioni precedenti

Vengono calcolate indicando con x_t l'unità statistica al tempo t viene calcolata tramite la seguente formula:

$$\Delta = \frac{x_t - x_{t-k}}{x_{t-k}}$$

con $k = 1, 3, 5$ numero di lag per le misure di accelerazione, nel caso in cui $x_{t-k} = 0$ e $x_t = 0$ non si applica la formula e la variabile `_lag` assume 0 per pre-costruzione, successivamente quando si osserva solamente $x_{t-k} = 0$ procede con la sostituzione $x_{t-k} = 0.01$ per riuscire a mantenere la validità del calcolo e continuità delle dinamiche telemetriche.

Queste variabili consentono di catturare la dinamica delle grandezze nel tempo e sono utili per caratterizzare l'evoluzione del comportamento di guida tra un istante e l'altro.

2.2 Statistiche riassuntive

Si procede al calcolo delle statistiche descrittive⁵ al fine di valutare in che modo esse possano essere associate allo stile di guida, considerando separatamente ciascun Gran Premio e Pilota.

Per le misure di accelerazione laterale, longitudinale e decelerazione longitudinale vengono calcolate media e deviazione standard.

La variabile relativa alla frenata non viene invece inclusa nell'analisi descrittiva, in quanto la sua natura binaria e la forte dipendenza dalle caratteristiche specifiche del singolo Gran Premio rendono poco informative misure sintetiche come media e deviazione standard. Tali indicatori risulterebbero infatti più rappresentativi del tracciato e delle condizioni di gara che dello stile di guida del pilota.

Per le variabili di tipo `_lagk` vengono calcolate medie e deviazioni standard distinguendo tra variazioni positive e negative, così da evidenziare eventuali asimmetrie nel comportamento dinamico del pilota.

Il nuovo dataset pertanto contiene 58 variabili, che forniscono informazione sullo stile di guida del pilota nella specifica gara.

Si procede dunque con una analisi delle componenti principali per ridurre la dimensionalità, e comprendere quali siano le variabili più significative nel fornire l'informazione.

³Estratto codice R consultabile nell'appendice ??.

⁴Estratto codice R consultabile nell'appendice D

⁵Estratto codice R consultabile nell'appendice E.

2.3 Riduzione della dimensionalità

2.3.1 Analisi correlazioni

Per limitare i casi di multicollinearità ed eliminare il rumore, avendo correlazioni molto elevate ($cor(x, y) > 0.9$), attraverso la funzione `findCorrelation` implementata nella libreria : `{caret}`⁶, si cerca di eliminare la ridondanza delle variabili.

In particolare la funzione analizza le coppie di variabili con correlazione eccessiva e rimuove quella che, in media, risulta più correlata con tutte le altre. Questo processo di scrematura ha permesso di ridurre la dimensionalità a 27 variabili.

2.3.2 Analisi delle componenti principali

Seppur la dimensionalità sia diminuita, 27 dimensioni sono eccessive sia ai fini di interpretabilità che per sforzo computazionale.

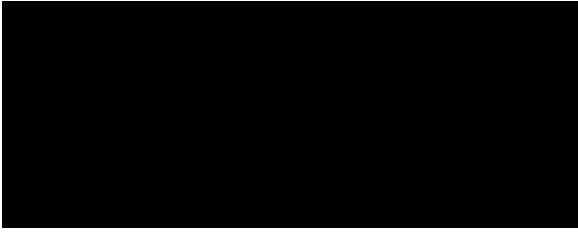
Si passa dunque da un'analisi delle componenti principali (PCA) ai fini di ridurre la dimensionalità. Per l'applicazione della PCA si opta di inserire come input la matrice di correlazione, in quanto, le variabili hanno unità di misura diverse⁷.

Dall'analisi risulta che le prime 8 componenti spiegano l'80% della varianza, i pesi associati ad ogni componente sono invece

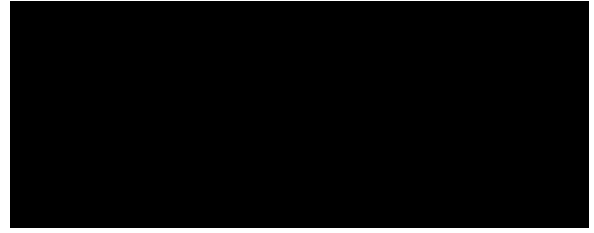
⁶Estratto del codice R consultabile nell'appendice F

⁷Estratto del codice R consultabile nell'appendice G

Come si evince dalla Figura 3, la distribuzione ...



(a) Distribuzione Variabile X



(b) Scatterplot X vs Y

Figura 3: Analisi esplorativa iniziale delle variabili principali.

3 Modellazione

Abbiamo applicato un modello di regressione lineare:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \epsilon_i \tag{1}$$

A Codice R Commentato Estrazione dei dati

```
1 library(jsonlite)
2 library(tidyverse)
3 rm(list=ls())
4 #Cartella origine
5 cartella <- "C:/Users/feder/Documents/datasets/Computazionale/2025"
6
7 #filtro sui file con i giri e poi qualifiche
8 file_giri <- list.files(path = cartella,
9 pattern = "laptimes\\.json$",
10 full.names = TRUE,
11 recursive = TRUE)
12
13 files_qualifiche_lap <- file_giri[grepl("Qualifying", file_giri) & !grepl("Sprint",
14 file_giri)]
15 lista_dati <- list()
16
17 #Ciclo di estrazione
18 for(i in 1:length(files_qualifiche_lap)) {
19
20   cartella_giro <- files_qualifiche_lap[i]
21
22   laptimes_data <- fromJSON(cartella_giro)
23
24   giro <- as.numeric(laptimes_data$lap[which.min(laptimes_data$time)])
25
26   cartella_pilota <- dirname(cartella_giro)
27   if(length(giro) > 0){
28     path_telemetry <- paste0(cartella_pilota, "/", giro, "_tel.json")
29   } else{
30     giro <- 1
31     path_telemetry <- paste0(cartella_pilota, "/", giro, "_tel.json") }
32
33   testo <- read_file(path_telemetry)
34   testo <- str_replace_all(testo, "NaN", "null")
35   json_data <- fromJSON(testo)[["tel"]]
36   json_data <- as_tibble(json_data)
37
38   parti_cartella <- str_split(cartella_giro, "/")[1]
39   n <- length(parti_cartella)
40
41   json_data$pilota <- parti_cartella[n-1]
42   json_data$GP <- parti_cartella[n-3]
43   json_data$gamma <- rep(laptimes_data$compound[giro],nrow(json_data))
44   json_data$lap_time <- rep(laptimes_data$time[giro],nrow(json_data))
45   json_data$strategy <- rep(laptimes_data$status[giro],nrow(json_data))
46   json_data$life <- rep(laptimes_data$life[giro],nrow(json_data))
47   json_data$life <- as.numeric(json_data$life)
48   lista_dati[[i]] <- json_data
49 }
50 #merge dei dati
51 dataset_completo <- bind_rows(lista_dati)
52 #pulizia
53 dataset_completo <- dataset_completo %>%
54 select(-dataKey)
55 setwd("C:/Users/feder/Documents/datasets/Computazionale/F1/data")
56 saveRDS(dataset_completo, file = "dataset_completo_best_tel.rds")
57
58
```

B Codice R Commentato Esplorazione iniziale

```

1  # A seguito del caricamento delle librerie e del dataset si inizia l'esplorazione
2  con le iniziali statistiche descrittive
3  summary(tel)
4  # Sistemazione delle variabili
5  tel$throttle <- ifelse(tel$throttle > 100, 100, tel$throttle)
6
7  #Trattamento degli NA e None
8
9  tel[which(is.na(tel$rel_distance)),]
10 tel <- tel %>% filter(pilota != "RUS" & GP != "Miami Grand Prix")
11
12 tel %>%
13 group_by(GP,pilota) %>%
14 filter(lap_time == "None") %>% summarize(n(),lap_time=max(lap_time))
15
16 tel <- tel %>% filter(lap_time != "None")

```

C Codice R Commentato Trattamento variabili

```

1  tel.guida <- tel %>%
2  group_by(GP,pilota) %>%
3  mutate(
4    dec_x = if_else(acc_x < 0, abs(acc_x), 0),
5    acc_x = if_else(acc_x > 0, acc_x, 0) %>%
6    mutate(acc_x=rescale(acc_x,to=c(0,1)),
7    dec_x=rescale(dec_x,to=c(0,1),
8    acc_y=abs(rescale(acc_y,to = c(-1, 1)))) %>%
9    ungroup()
10

```

D Codice R Commentato Creazione variabili

```

1  #Creazione delle variabili lag
2  tel.guida <- tel.guida %>%
3  select(GP,pilota,throttle,acc_x,acc_y,dec_x,rel_distance) %>%
4  arrange(GP,pilota,rel_distance) %>%
5  group_by(GP,pilota) %>%
6  mutate(
7    across(
8    c(throttle, acc_x, acc_y,dec_x),
9    list(
10     lag1 = ~round(
11     ifelse(
12     lag(.x, 1) == 0 & .x == 0,
13     0,
14     ifelse(
15     lag(.x, 1) == 0,
16     (.x-0.01)/0.01,
17     (.x - lag(.x, 1)) / lag(.x, 1))
18     ),
19     4),
20     ...
21     ),
22     .names = "{.col}_{.fn}"
23     )
24     ) %>%
25     ungroup() %>%
26     select(-rel_distance)
27

```

E Codice R Commentato Statistiche Riassuntive

```
1 tel.guida_summary <- tel.guida %>%
2 group_by(GP, pilota) %>%
3 {
4   lag_cols <- names(.) %>% .[str_detect(., "lag")]
5
6   summarise(.,
7     across(
8       c(throttle, acc_x, acc_y, dec_x),
9       list(
10        mean = ~round(mean(.x, na.rm = TRUE), 4),
11        sd = ~round(sd(.x, na.rm = TRUE), 4)
12      ),
13      .names = "{.col}_{.fn}"
14    ),
15
16    across(
17      all_of(lag_cols),
18      ~round(mean(.x[.x > 0], na.rm = TRUE), 4),
19      .names = "{.col}_mean_pos"
20    ),
21
22    across(
23      all_of(lag_cols),
24      ~round(mean(.x[.x <= 0], na.rm = TRUE), 4),
25      .names = "{.col}_mean_neg"
26    ),
27
28    across(
29      all_of(lag_cols),
30      ~round(sd(.x[.x > 0], na.rm = TRUE), 4),
31      .names = "{.col}_sd_pos"
32    ),
33
34    across(
35      all_of(lag_cols),
36      ~round(sd(.x[.x <= 0], na.rm = TRUE), 4),
37      .names = "{.col}_sd_neg"
38    ),
39
40    .groups = "drop"
41  )
42 }
43
```

F Codice R Commentato Analisi delle correlazioni

```
1 corr <- round(cor(tel.guida_summary %>% select(where(is.numeric))),4)
2
3 #Ricerca di variabili dipendenti
4
5 variabili_dipendenti <- findCorrelation(corr, cutoff = 0.9, names = TRUE, exact = T
6   ,verbose = T)
7
8 print(variabili_dipendenti)
9
10 tel.pca <- tel.guida_summary %>% select(-all_of(variabili_dipendenti))
11
```


G Codice R Commentato Analisi delle componenti principali

```
1 PCA <- princomp(tel.pca %>%
2   select(where(is.numeric)),cor=T)
3
4   summary(PCA)
5
6   PCA$loadings
7
8   fviz_contrib(PCA, choice = "var", axes = 1, top = 10)
9   fviz_contrib(PCA, choice = "var", axes = 2, top = 10)
10  fviz_contrib(PCA, choice = "var", axes = 3, top = 10)
11  fviz_contrib(PCA, choice = "var", axes = 4, top = 10)
12  fviz_contrib(PCA, choice = "var", axes = 5, top = 10)
13  fviz_contrib(PCA, choice = "var", axes = 6, top = 10)
14
```