

# BDA - Project

Jan Nyberg, Carl-Victor Schauman

4/12/2022

## Contents

Introduction	2
Description of the data	2
Description of models and priors used	2
Rstan code	3
Running of stan model	4
Convergence diagnostics	7
Posterior predictive checks	9
Predictive performance assessment	11
Sensitivity analysis	11
Model comparison	12
Discussion of issues and potential improvements	15
Conclusion what was learned from the data analysis	15
Self-reflection	15

## Introduction

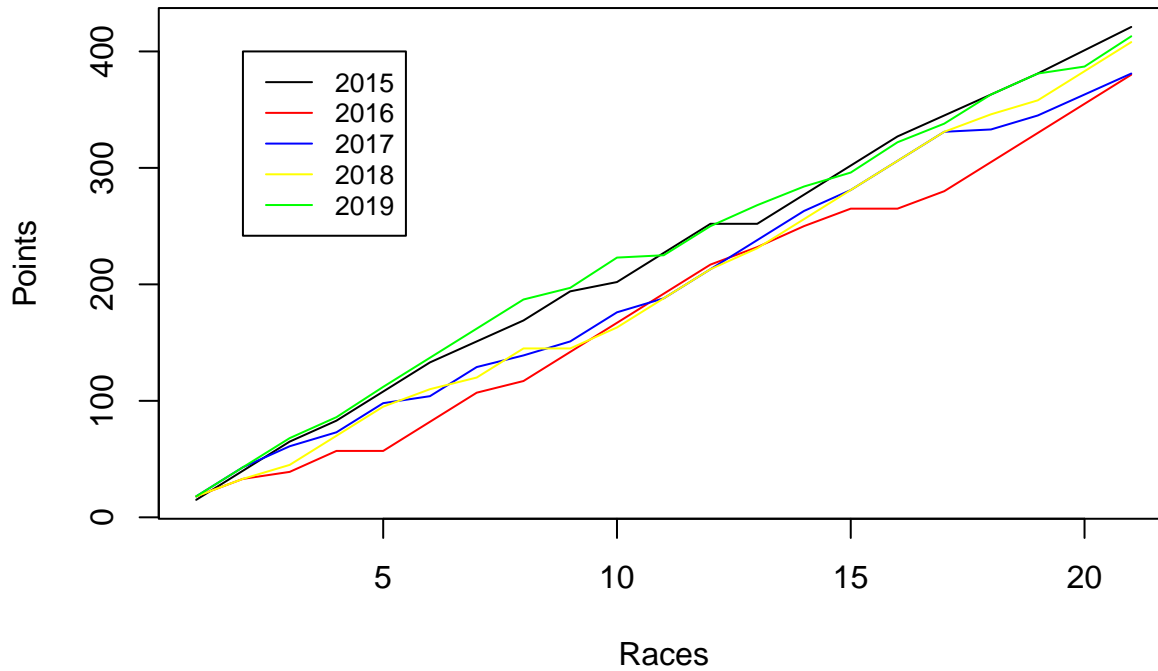
For this project, we are trying to predict Lewis Hamilton's average score in a season based on previous years scores. This is mostly just due to our curiosity if we are able to use this to somehow predict the score.

We are modeling his scores from five years, and trying to build a model using it. We want to see what kind of distribution the answer will be and how well it is able to estimate the following year. There are many factors we don't take into consideration, but we hope to see relatively good results and predictions.

## Description of the data

The data we use is from Kaggle and can be found [here](#). We took Lewis Hamilton out of the data and chose the years 2015-2019. We selected his scores from all the races from those years. One thing to note with the data is that every year doesn't have an equal amount of races. To account for this we chose to fill in the missing races with the median for the year. This makes the data a bit inaccurate, however, it shouldn't have too big an effect on the data. We couldn't find any similar analysis done before us.

### Hamilton cumulative points



## Description of models and priors used

The models we used, with priors are:

Hierarchical Model:

$$y_{ij} \sim \mathcal{N}(\mu_j, \sigma_j)$$

$$\mu_j \sim \mathcal{N}(\mu_{\text{hyperprior}}, \sigma_{\text{hyperprior}})$$

$$\sigma_j = \text{Gamma}(1, 1)$$

$$\mu_{\text{hyperprior}} = \mathcal{N}(0, \mu_{\text{prior}})$$

$$\tau = \text{Inv} - \chi^2(\sigma_{\text{prior}})$$

$$\sigma_{\text{prior}} = 7 \quad \mu_{\text{prior}} = 20$$

Pooled Model:

$$\begin{aligned}y_{ij} &\sim \mathcal{N}(\mu_j, \sigma_j) \\ \mu_j &\sim \mathcal{N}(0, \mu_{prior}) \\ \sigma_j &= Inv - \chi^2(\sigma_{prior}) \\ \sigma_{prior} &= 7 \\ \mu_{prior} &= 20\end{aligned}$$

The priors were chosen mostly based on looking at the data and, after some testing, seeing what fit well.

## Rstan code

### Hierarchical stan model

Below is the code for the hierarchical model for Hamilton's points.

```
data {
  int<lower=0> N;
  int<lower=0> J;
  vector[J] y[N];
  real<lower=0> mu_s;
  real<lower=0> sigma_prior;
}

parameters {
  real<lower=0> muh;
  real<lower=0> sigma;
  real<lower=0> tau;
  vector[J] mu;
}

model {
  //Hyperpriors
  tau ~ inv_chi_square(sigma_prior);
  muh ~ normal(0, mu_s);

  //priors
  sigma ~ gamma(1,1);
  mu ~ normal(muh, tau);

  for (j in 1:J)
    y[,j] ~ normal(mu[j], sigma);
}

generated quantities {
  vector[J] log_lik[N];
  real ypred;
  real ypred_6;
  ypred = normal_rng(mu[5], sigma);
  ypred_6 = normal_rng(muh, sigma);
  for (j in 1:J){
    for (n in 1:N){
      log_lik[n,j] = normal_lpdf(y[n,j] | mu[j], sigma);
    }
  }
}
```

```
}
```

## Pooled stan model

```
data {  
  int<lower=0> N;  
  int<lower=0> J;  
  vector[N*J] y;  
  real mean_mu;  
  real<lower=0> mean_sigma;  
  
}  
parameters {  
  real mu;  
  real<lower=0> sigma;  
}  
model {  
  // prior  
  mu ~ normal(0, mean_mu);  
  sigma ~ inv_chi_square(mean_sigma);  
  // likelihood  
  y ~ normal(mu, sigma);  
}  
generated quantities {  
  real ypred;  
  vector[N*J] log_lik;  
  
  // Distribution based on all seasons  
  ypred = normal_rng(mu, sigma);  
  for (jn in 1:J*N){  
    log_lik[jn] = normal_lpdf(y[jn] | mu, sigma);  
  }  
}
```

## Running of stan model

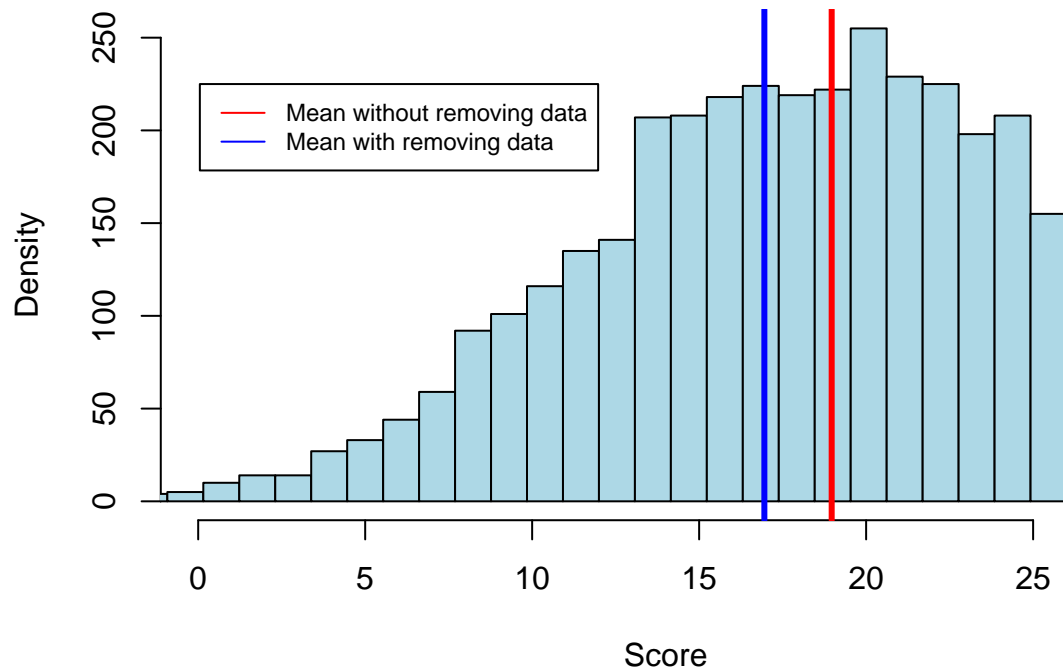
### Hierarchical model

Below is the hierarchical model run with the corresponding histogram with the data.

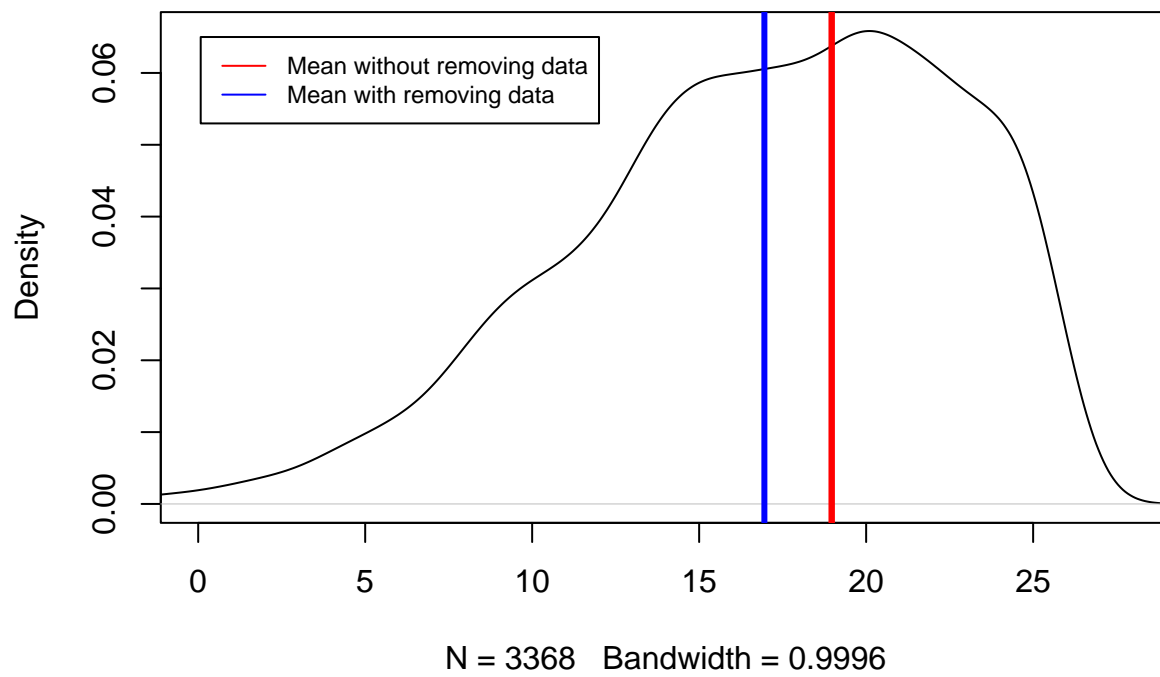
```
hier_data = list(  
  y = ham_data,  
  N = nrow(ham_data),  
  J = ncol(ham_data),  
  mu_s = 20,  
  sigma_prior = 7  
)  
  
hier_fit = sampling(  
  hier_ham,  
  data = hier_data,  
  chains = 4,  
  iter = 2000,
```

```
warmup = 1000,  
refresh = 0  
)
```

### Predictive distribution of the mean of Hamilton's points the next seas



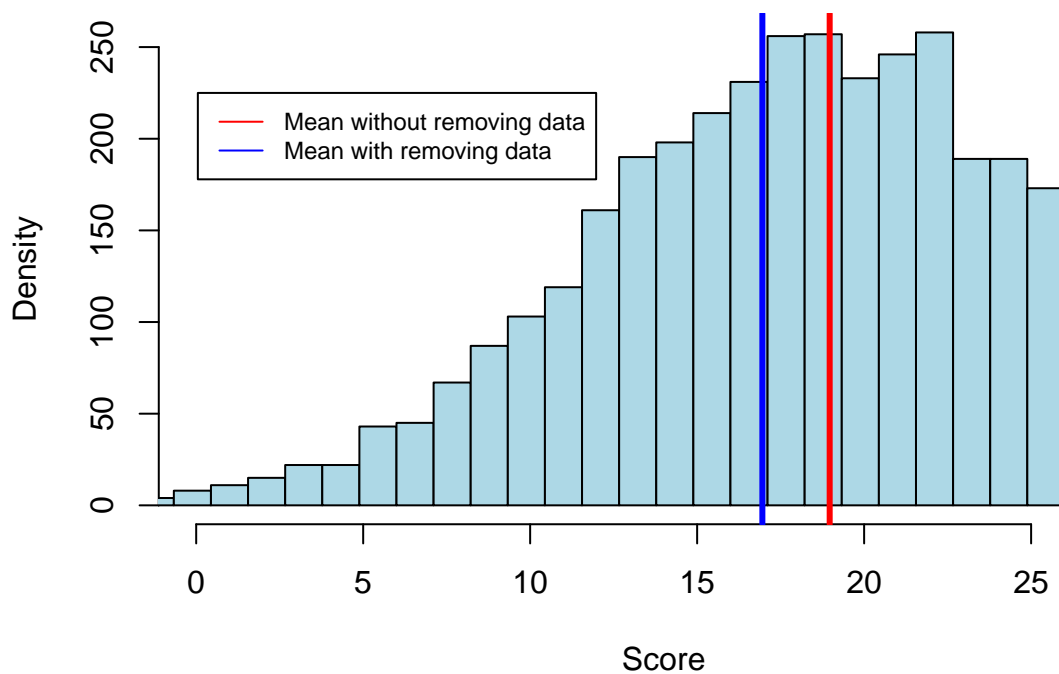
### Density plot of the mean of Hamilton's points the next season



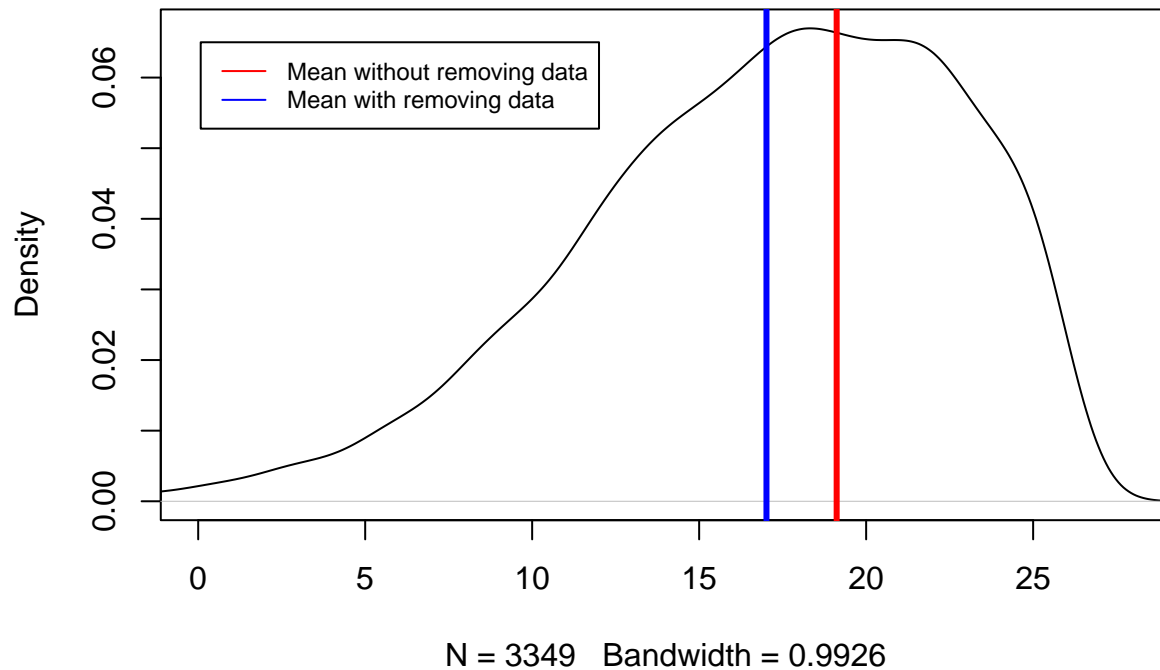
## Pooled model

```
pool_data = list(  
  y = unlist(ham_data),  
  N = nrow(ham_data),  
  J = ncol(ham_data),  
  mean_mu = 20,  
  mean_sigma = 7  
)  
  
pool_fit = sampling(  
  pool_ham,  
  data = pool_data,  
  chains = 4,  
  iter = 2000,  
  warmup = 1000,  
  refresh = 0  
)
```

## Predictive distribution of the mean of Hamilton's points the next seas



## Density plot of the mean of Hamilton's points the next season



Since the values of these normal distributions, go beyond the max points, i.e. 26, we have limited them a bit. We still plot the mean of both the limited and unlimited data. As can be seen, there isn't a lot of difference, however, over several races, this difference can be quite large. Below is also the histogram as a density plot.

## Convergence diagnostics

The  $\hat{R}$  for our fits are as follows:

- Hierarchical model: 1.01
- Pooled model: 1

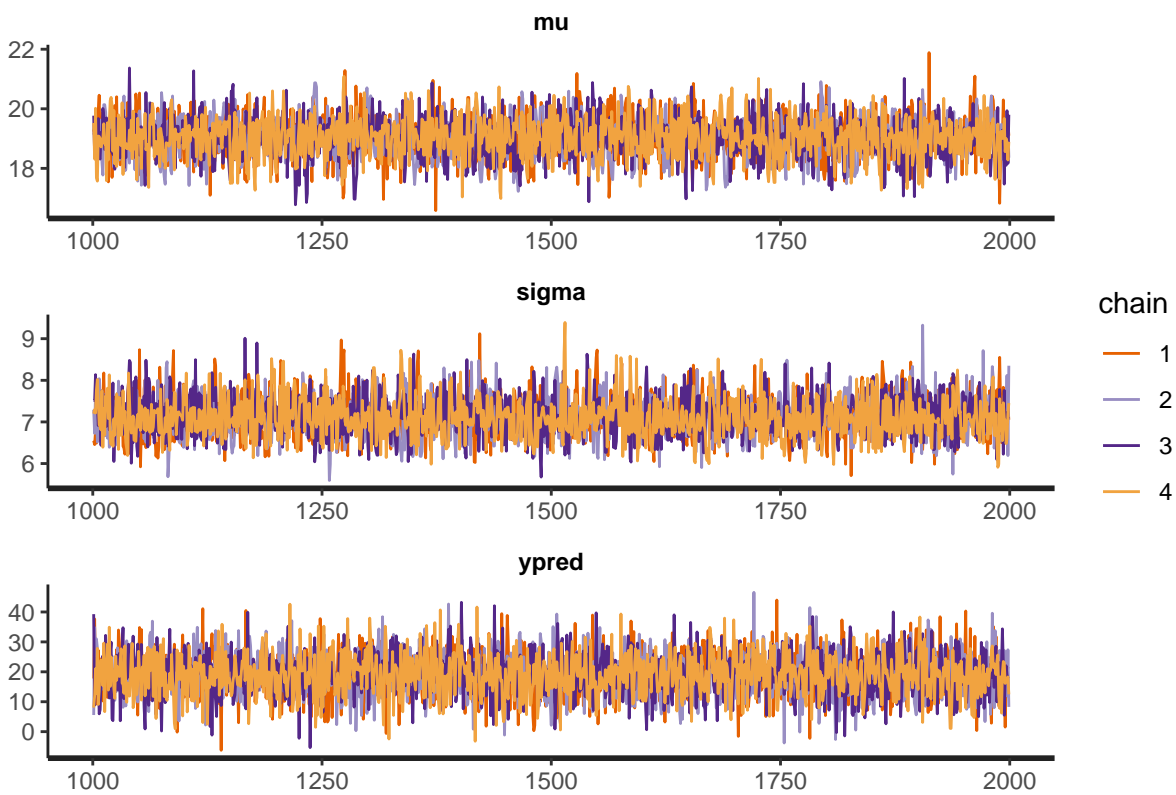
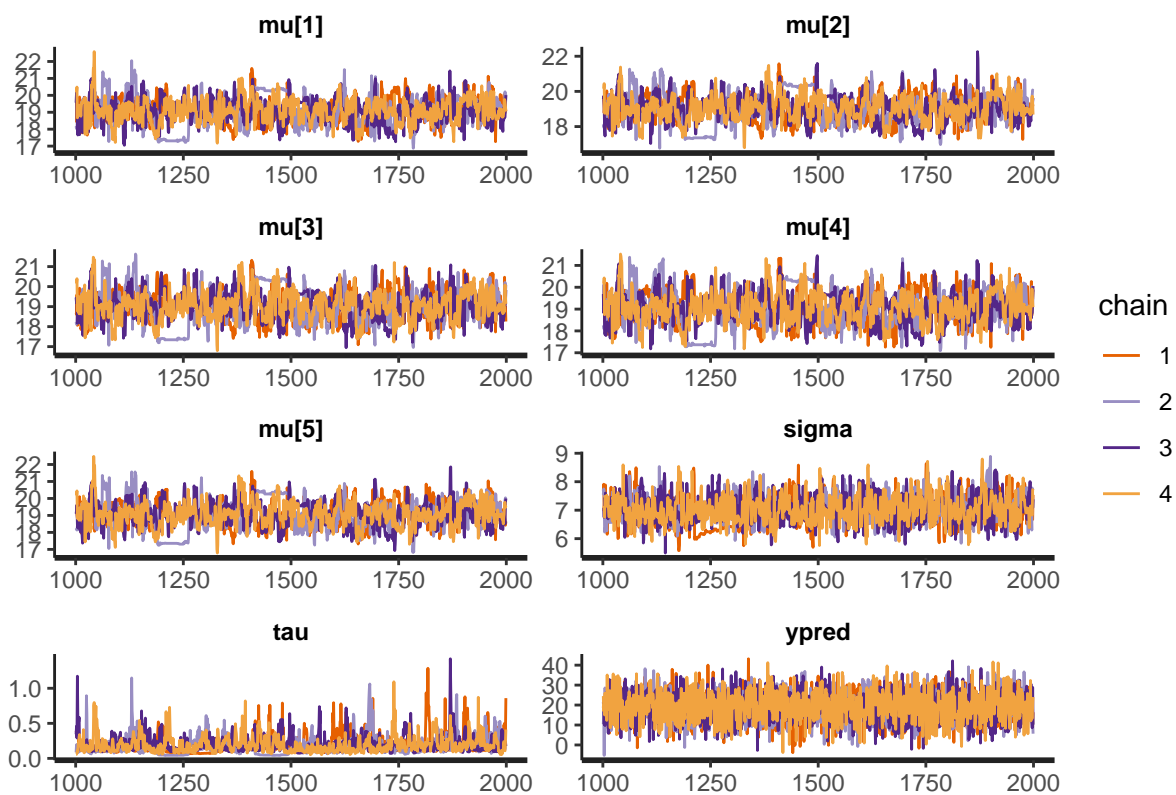
Since these  $\hat{R}$  values are under 1.05, the chains have most likely mixed well.

Another convergence diagnostic we can look at is the ESS value we get out of the fits.

- Bulk ESS of the hierarchical model: 264
- Tail ESS of the hierarchical model: 175.75
- Bulk ESS of the pooled model: 3697.33
- Tail ESS of the pooled model: 3214.33

These ESS values measure the crude effective sample size for the bulk and tail quantities. A value over 100 is good and all of our values are over it.

Below are the traceplots for the chains. As can be seen, they seem to converge well.





## Posterior predictive checks

To assess the performance of the models we simulate 1000 seasons with the help of our models and compare the outcomes with the real world data.

Year	Driver	Points
2015	Hamilton	413
2016	Hamilton	408
2017	Hamilton	363
2018	Hamilton	380
2019	Hamilton	381

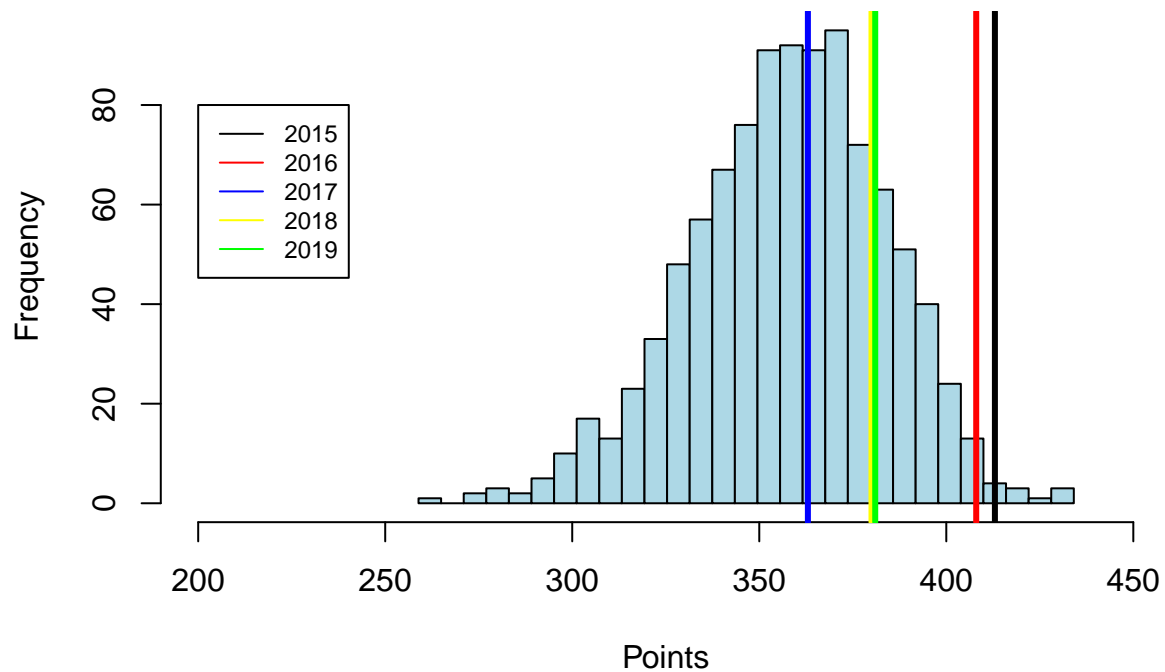
*# Code for simulating 1000 seasons*

```
pool_season_predictions = c()
hier_season_predictions = c()

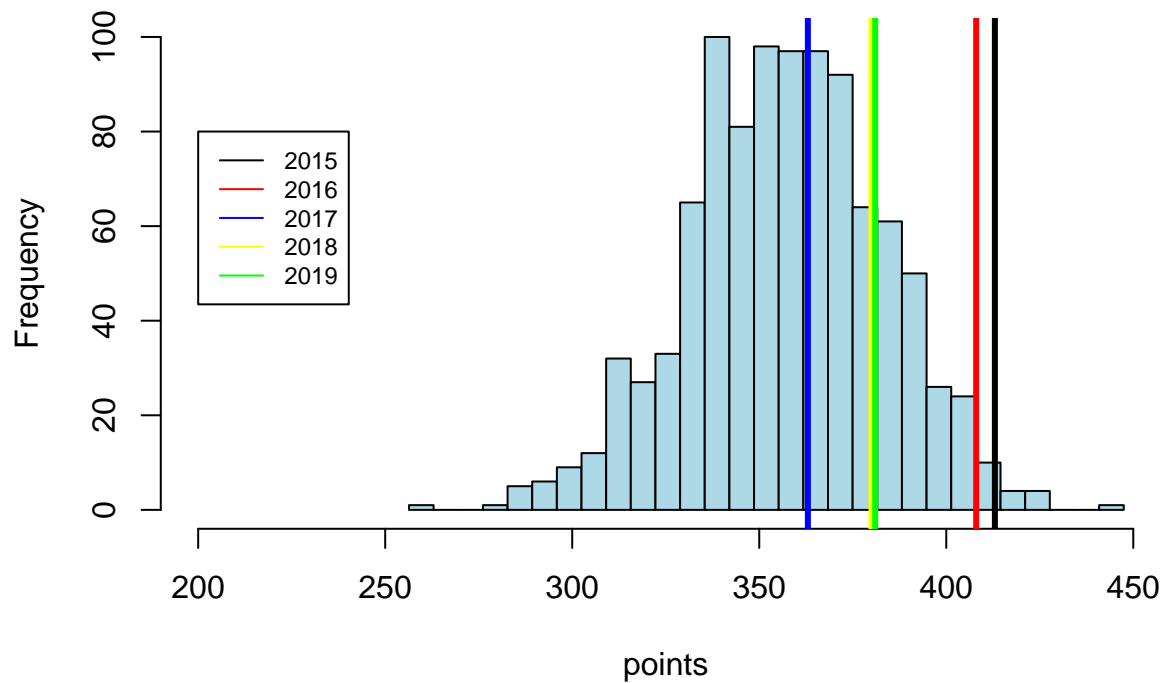
for(i in 1:1000) {
  dens = ham_pool_extracted_density
  N <- 21
  newx <- sample(x = dens$x, N,
    prob = dens$y, replace=TRUE)
  + rnorm(N, 0, dens$bw)
  pool_season_predictions <- append(pool_season_predictions, sum(newx))
}

for(i in 1:1000) {
  dens = ham_extracted_density
  N <- 21
  newx <- sample(x = dens$x, N,
    prob = dens$y, replace=TRUE)
  + rnorm(N, 0, dens$bw)
  hier_season_predictions <- append(hier_season_predictions, sum(newx))
}
```

### Simulated seasons – Pooled model



### Simulated seasons – Hierarchical model



When comparing the simulated data with the actual real world data we can see that it is in the right ball park, but slightly of.

## Predictive performance assessment

To assess the predictive performance of our model we look at the results from 2020-2022.

Year	Driver	Points
2020	Hamilton	347
2021	Hamilton	387.5
2022	Hamilton	240

When comparing the above results to the simulated ones we can see that our simulated points are well within our estimations. However the last season's result is not even in our simulated distribution. This speaks loudly about how unpredictable F1 is.

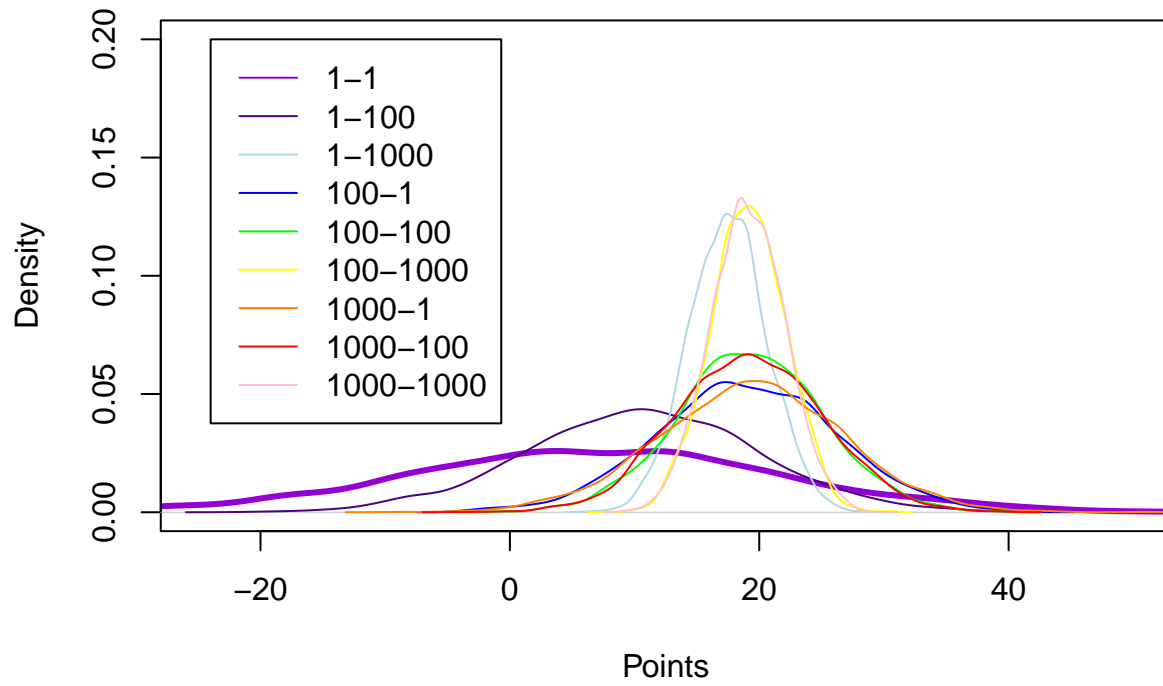
## Sensitivity analysis

To test the sensitivity of the model we use different combinations of small and large values of  $\mu$  and  $\sigma$ . The chosen values are

1, 100, 1000

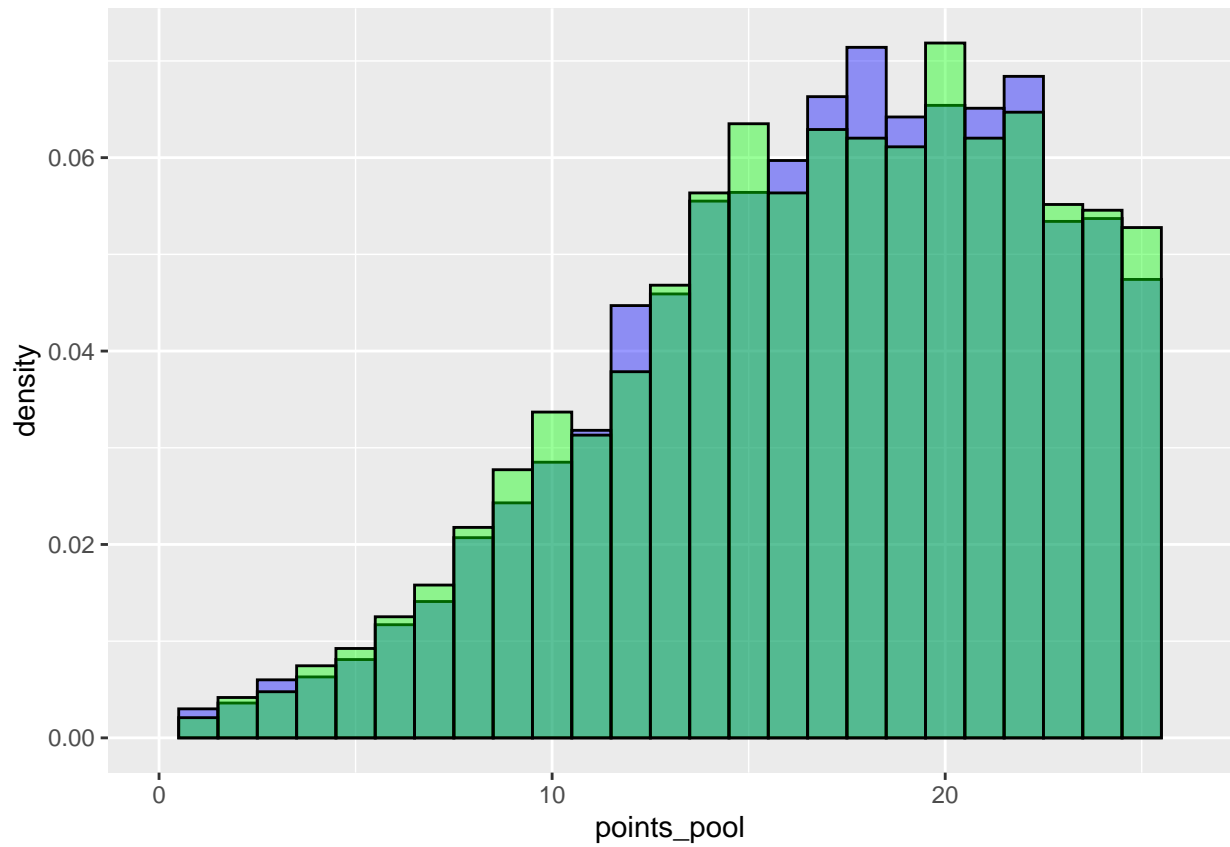
```
.  
sensitivity_mu = c(1, 100, 1000)  
sensitivity_sigma = c(1, 100, 1000)  
  
fits = list()  
  
for (i in sensitivity_mu) {  
  for (j in sensitivity_sigma) {  
    sensitivity_data = list(  
      y = unlist(ham_data),  
      N = nrow(ham_data),  
      J = ncol(ham_data),  
      mean_mu = i,  
      mean_sigma = j  
    )  
  
    sensitivity_fit = sampling(  
      pool_ham,  
      data = sensitivity_data,  
      chains = 4,  
      iter = 2000,  
      warmup = 1000,  
      refresh = 0  
    )  
  
    fits <- append(fits, list(sensitivity_fit, paste0(i, "-", j)))  
  }  
}
```

## Sensitivity plot



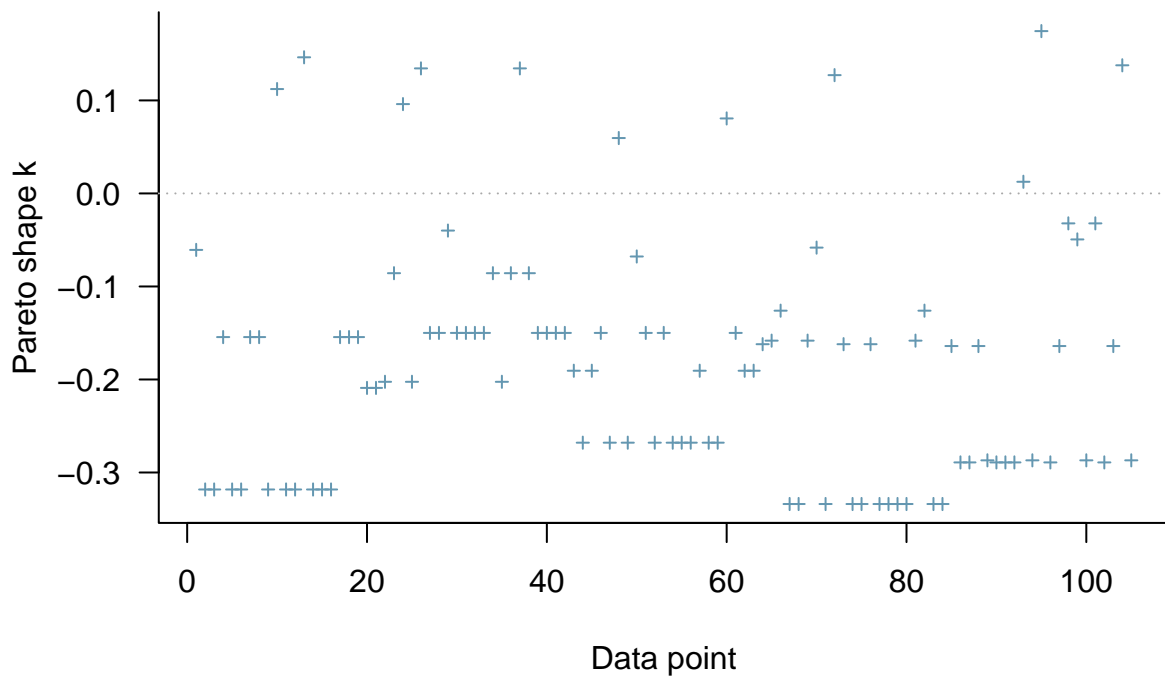
## Model comparison

Our models compared to each other in one graph looks like this. As we can see, they are relatively similar. To check performance and comparing them, we will also use PSIS-LOO values.



Hierarchical analysis

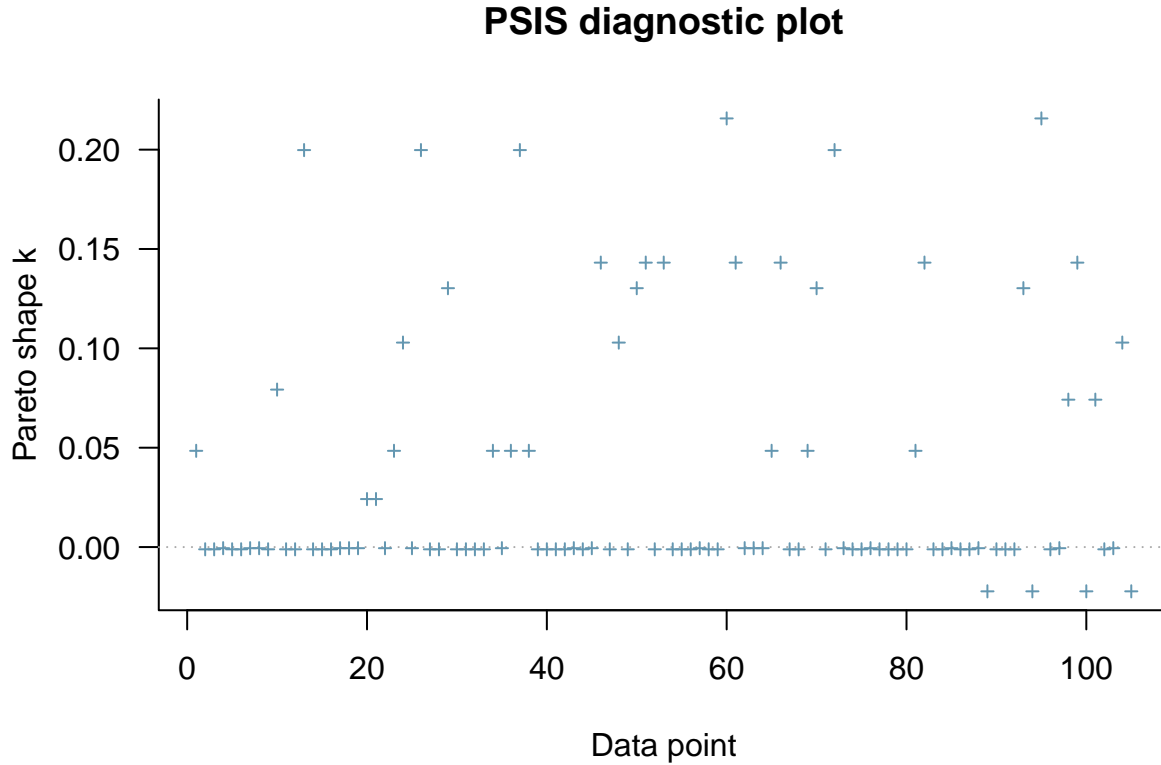
### PSIS diagnostic plot



	Estimate	SE
elpd_loo	-358.053749	8.5521592
p_loo	2.510075	0.5768406
looic	716.107499	17.1043184

As can be seen, all of our k values are under 0.5, which means they are all good, and our PSIS-LOO estimate can be considered reliable.

## Nonhierarchical analysis



	Estimate	SE
elpd_loo	-357.773141	8.2947834
p_loo	2.269913	0.5461605
looic	715.546282	16.5895668

Like in the previous one, all of our k values are under 0.5, which means they are all good, and our PSIS-LOO estimate can be considered reliable.

	elpd_diff	se_diff	elpd_loo	se_elpd_loo	p_loo	se_p_loo	looic	se_looic
model1	0.0000000	0.0000000	-357.7731	8.294783	2.269913	0.5461605	715.5463	16.58957
model2	-0.2806085	0.2584615	-358.0537	8.552159	2.510075	0.5768406	716.1075	17.10432

Here we can once again see that the values are very close to each other, and that both models perform quite similarly to each other. This wasn't a surprise considering the first graph in this section, which shows that both are very similar.

## Discussion of issues and potential improvements

Our problem to begin with is almost impossible to model correctly since F1 is dependent on so many different parameters. Who you are driving for, the current regulations and the engineering of the car are all important factors that the drivers can not simply affect in any way. The choice to model the data based on a normal distribution might also have been a bad idea since our rstan model is also predicting out values that are negative and larger than 26 (max points from one race). Lastly the data used had minor inconsistencies with number of races each season. 2015 had 2 less races and 2017 1 less. To make it easier for us we filled in the missing values with the mean score of the season.

In our project we decided to look at the results of an individual driver, but it might have also been a good idea to look at results of teams instead. However by looking at team results we just get different inconsistencies such as who is driving for the team and what their current form is. An other idea for improvements is to look at more parameters than just the points scored e.g. positions in qualifying and point difference between teammates.

## Conclusion what was learned from the data analysis

We found that we can relatively well describe Lewis Hamilton's score by using a model with his previous results. It is not perfect, but with his consistency, and his recent results, they all fall well within the expected values in our model. One thing we could see was that his 2022 performance was something that even the model thought was unlikely for the simulations, but it still happened. If we could get him to drive for, say 100 more years with relatively even consistency, we could probably plot the results close to our plot.

## Self-reflection

During this project, we learned that it is relatively easy to get some kind of model going. We didn't want to choose some really easy data set, which lead us to have to put in quite some work into getting rolling with the project. We aren't totally sure about the results, however, they seem to quite well simulate the real world even with the limited amount of parameters we used. We also learned while choosing a data set that it can be quite hard to model some things.