

BAYESIAN WORKFLOW ON TIME TO REACH GLOBAL TOP 20 IN TENNIS

Alireza Honarvar & Reza Soumi

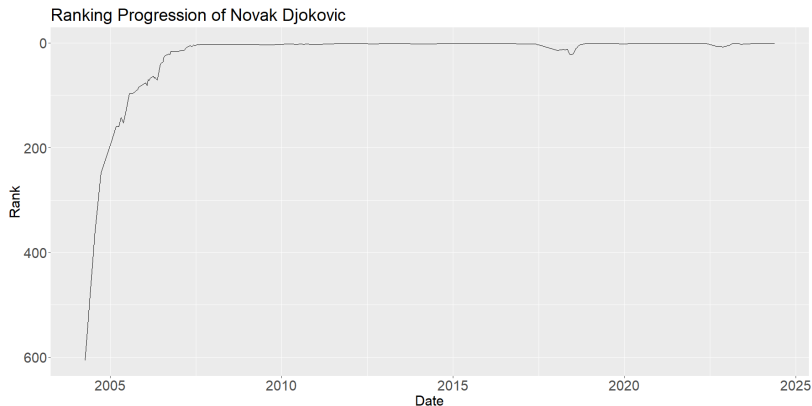
Aalto University
December 13, 2024

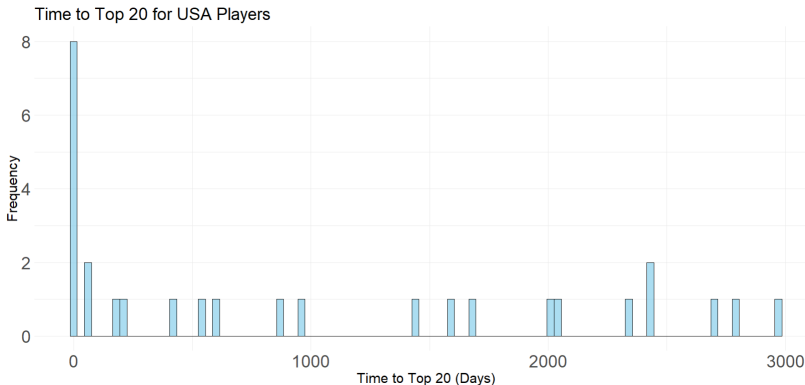
Overview

- Dataset: 106,716 rows containing 49 variables from Kaggle, such as:
 - **Match-Level:** Date, Surface, Duration
 - **Player Stats:** Player Name, Nationality, Age
 - **Rank Info:** Winner Rank, Loser Rank, Rank Points



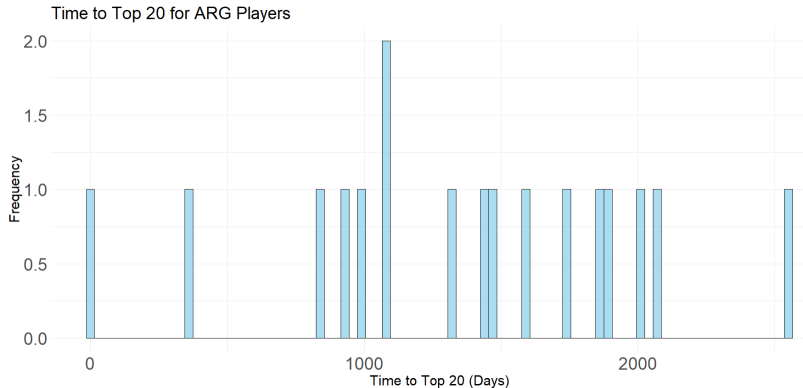
How long does it take players from different regions to reach the top 20?





$$y_i \sim \text{Gamma}(\lambda_i, \theta), \quad \log(\lambda_i \theta) = \alpha + u_{\text{region}[i]} + u_{\text{country}[i]}$$

$$u_{\text{region}[i]} \sim \mathcal{N}(0, \sigma_{\text{region}}), \quad u_{\text{country}[i]} \sim \mathcal{N}(0, \sigma_{\text{country}})$$



$$y_i \sim \mathcal{N}(\mu_i, \sigma), \quad \mu_i = \alpha + u_{\text{region}[i]} + u_{\text{country}[i]}$$
$$u_{\text{region}[i]} \sim \mathcal{N}(0, \sigma_{\text{region}}), \quad u_{\text{country}[i]} \sim \mathcal{N}(0, \sigma_{\text{country}})$$

- **Gaussian Model:**

- Intercept: $\text{Normal}(1000, 350)$
- Standard Deviations (Random Effects): $\text{Cauchy}(0, 2)$

- **Gamma Model:**

- Intercept: $\text{Normal}(7, 1)$ (log-scale; approx. mean time of 1094 days)
- Standard Deviations (Random Effects): $\text{Cauchy}(0, 0.5)$

Summary

- **Gaussian Model:**

- Stable convergence, no divergences
- Bulk ESS: 4959–8583, Tail ESS: 2059–2957
- $\hat{R} = 1.00$

- **Gamma Model:**

- Stable convergence, no divergences
- Smaller effective sample sizes (Bulk ESS: 1346–4907, Tail ESS: 1346–3216)
- $\hat{R} = 1.00$

Findings

- Gaussian model demonstrated more reliable and stable sampling

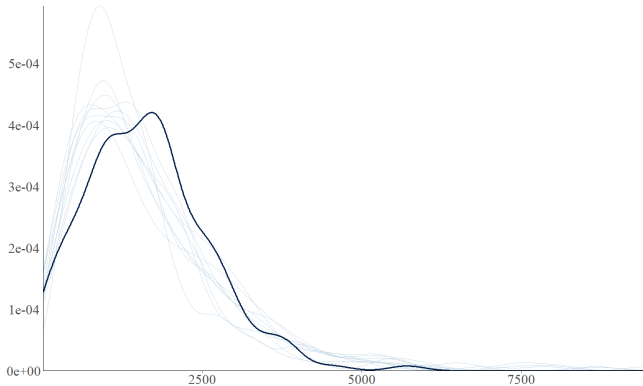


FIGURE: `pp_check` for Gamma

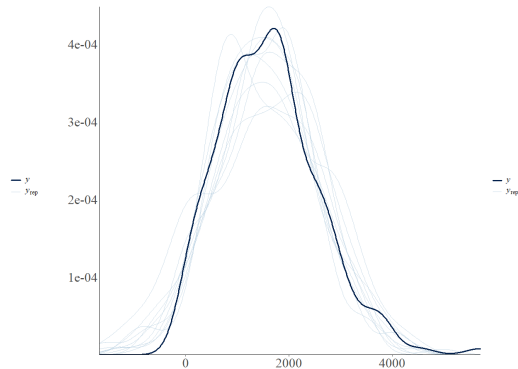


FIGURE: `pp_check` for Gaussian

Model Comparison (LOO-CV)

- **Gaussian Model:**

- $\text{elpd_loo} = -1724.7$, $\text{SE} = 12.0$

- **Gamma Model:**

- $\text{elpd_loo} = -1543.9$, $\text{SE} = 8.9$

Conclusion

- Gamma model performs better on fit and has better generalization.

- **Gaussian Model:**

- Shift $\mathcal{N}(1000, 1000) \rightarrow \mathcal{N}(10000, 10)$
- $\text{elpd_loo} = -1726.6$, $\text{SE} = 12.1$

- **Gamma Model:**

- Shift $\mathcal{N}(7, 1) \rightarrow \mathcal{N}(15, 1)$
- Shift $\text{Cauchy}(0, 0.5) \rightarrow \text{Cauchy}(0, 1)$
- $\text{elpd_loo} = -1546.0$, $\text{SE} = 8.9$

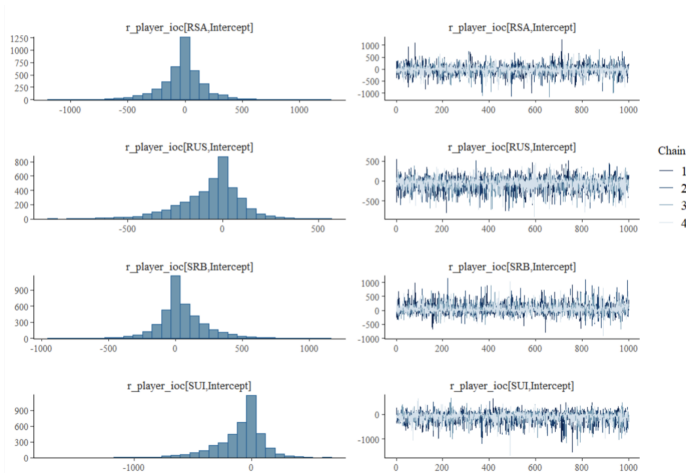


FIGURE: Player IOC and region effects on time taken to reach top 20

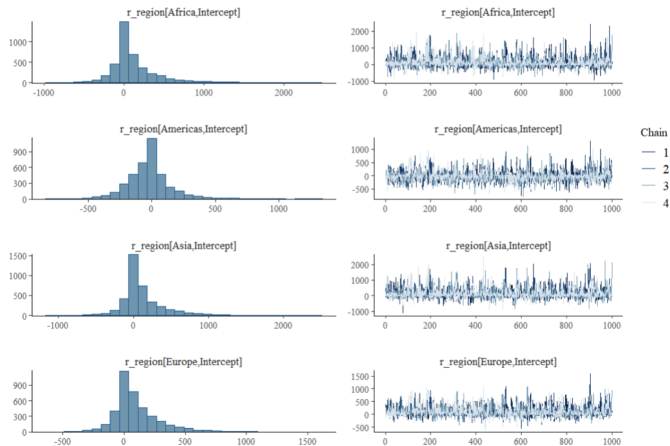


FIGURE: Player IOC effects on time taken to reach top 20

Problem Definition

- Does region affect performance on different surfaces?

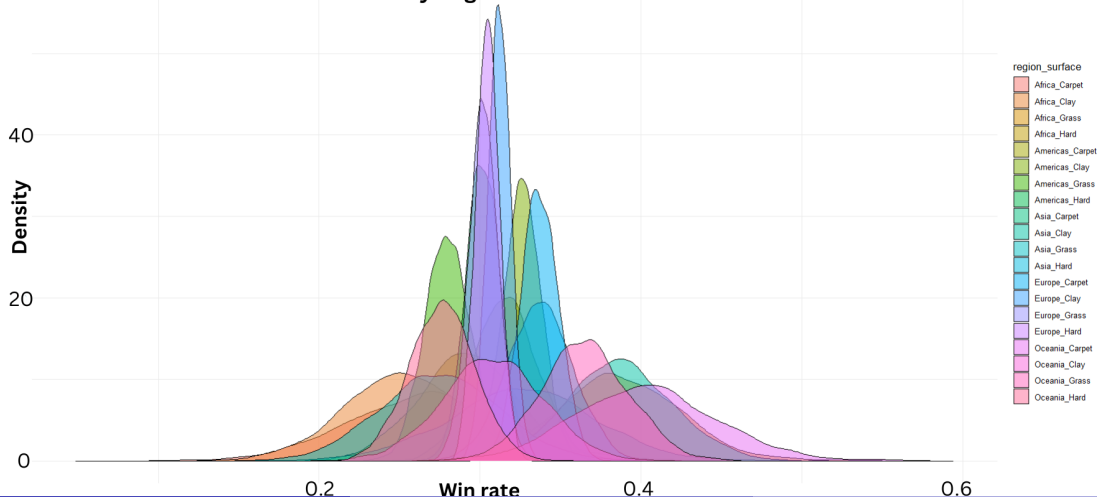
Model

$$y_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_{ij})$$

Priors

$$\mu_{ij} \sim \mathcal{N}(0.35, 0.35), \quad \sigma_{ij} \sim \text{Cauchy}(0, 2)$$

Posterior Distributions of Win Rates by Region and Surface



```
% GAUSSIAN
bayesian_model <- brm(
  time_to_top_20 ~ 1 + (1 | region) + (1 | player_ioc),
  data = time_to_top_20_regional_nonzero,
  family = gaussian(),
  prior = c(
    prior(normal(1000, 1000), class = "Intercept"),
    prior(cauchy(0, 3), class="sd", group="player_ioc"), # prior(normal(0, 1000)
      , class="sd", group="player_ioc")
    prior(cauchy(0, 3), class = "sd", group="region") # prior(normal(0, 1000),
      class = "sd", group="region")
  ),
  chains = 4,
  iter = 2000,
  warmup = 1000,
  cores = 4,
  control = list(adapt_delta = 0.99, max_treedepth = 15)
)
```

```
% GAMMA
bayesian_model_gamma <- brm(
  time_to_top_20 ~ 1 + (1 | region) + (1 | player_ioc),
  data = time_to_top_20_regional_nonzero,
  family = Gamma(link = "log"), # Gamma distribution with log-link
  prior = c(
    prior(normal(1000, 1000), class = "Intercept"), # Approx. log(mean time
      around 1000)
    prior(cauchy(0, 2), class = "sd", group = "player_ioc"), # Prior on random
      effect std. dev
    prior(cauchy(0, 2), class = "sd", group = "region") # Prior on random effect
      std. dev
  ),
  chains = 4,
  iter = 2000,
  warmup = 1000,
  cores = 4,
  control = list(adapt_delta = 0.99, max_treedepth = 15)
)
```



```
% EXTRA
bayesian_models <- list()
for (i in 1:nrow(bayesian_data)) {
  region <- bayesian_data$region[i]
  surface <- bayesian_data$surface[i]
  win_rates <- bayesian_data$win_rates[[i]]
  bayesian_models[[paste(region, surface, sep = "_")]] <- brm(
    win_rates ~ 1, # Single group model for win rates
    data = data.frame(win_rates = win_rates),
    family = gaussian(),
    prior = c(
      prior(normal(0.35, 0.35), class = "Intercept"), # Centered around 0.5,
      wide SD
      prior(cauchy(0, 0.3), class = "sigma") # Prior for standard deviation
    ),
    chains = 4,
    iter = 2000,
    cores = 4,
    seed = 123
  )
}
```