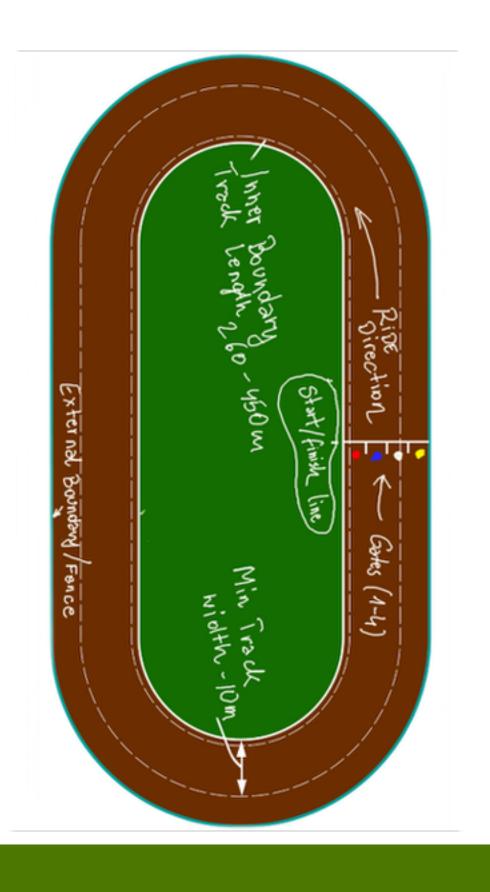# Using online update algorithms to predict speedway results.

Dawid Kałędkowski @ClickMeeting

# Speedway is a motorcycle sport.

- **Motorbike** racing on **dirt** circuit.

- No gears and **no breaks**.

- 4 riders competing over four laps.

- Start indicated by **tape rise** and light signal.

- Exclusions for falling, engine failure, retiring, touching tape etc.

- Sliding scale for scoring (known as the **3-2-1-0 method**)

# More than 1000 events annualy.

- 1000 riders from 31 countries

- 100 different competitions.

- Individual competitions - Grand-Prix

- League matches - PGE Ekstraliga,
  Premiership, Elitserien, Bundesliga etc.

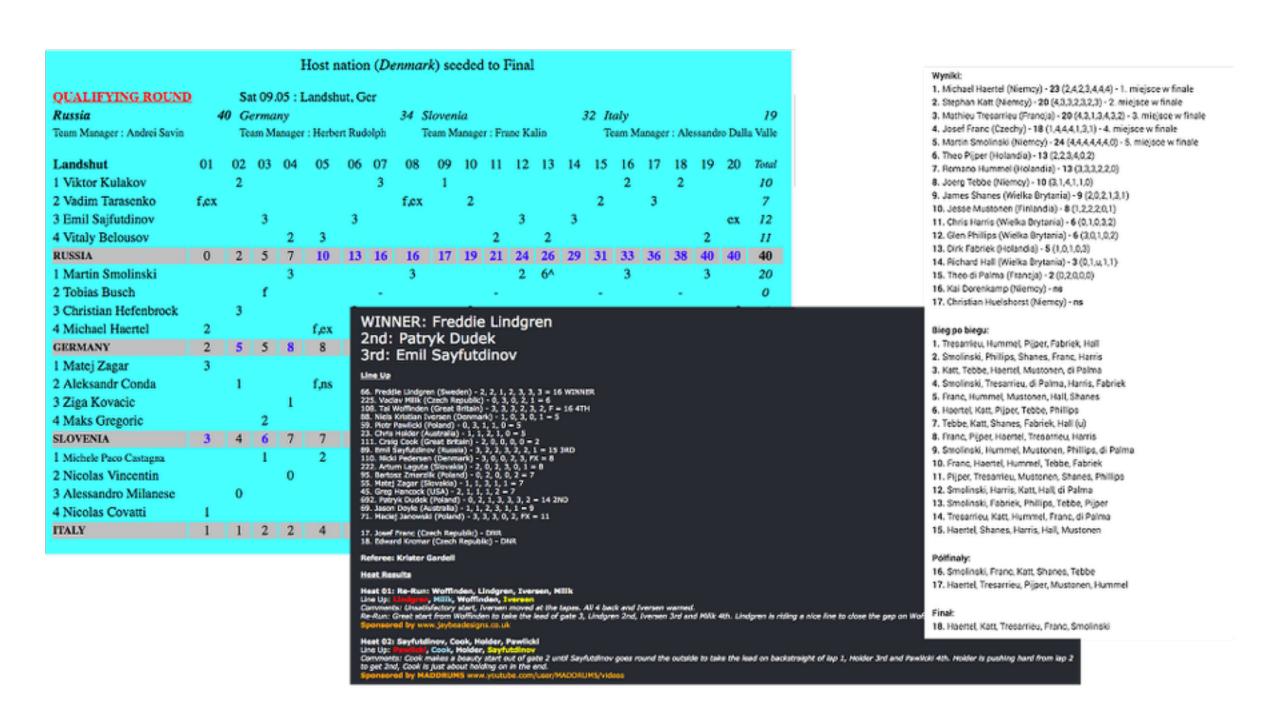- Team competitions - SWC/SON

- Pair competitions - SBP

# Statistics matter.

# Approach to data is old-fashioned



source: sportowefakty.pl, speedwayresults.com, speedwayupdates.proboards.com

# Data gathering partially automated.

- Multiple scrapers for different pages
  (`rvest`, `RSelenium` and a lot of REGEXP)

- PDF parsers (`tabulizer`)

- Finding appropriate rider-heat pattern

- Additional verifications

- Approximate string matching for language
  specific letters (`stringdist:amatch`)

- Apps to input and browse data
  (`shiny`,`rhandsontable`,`RMySQL`)

# Actual speedway database

- 10528 events in 148 competitions.

- Almost 700k individual performances.

- Unified rider names, competitions and
  places.

- 241 Speedway stadiums with coordinates

# Analytical challange

- Output is a **ranking**.

- Need for **continuous updates**.

- Some **commonly known effects** need to be examined:

  - Riders form changes in time.

  - Many interactions (gate*heat*stadium).

  - Field advantage.

  - Winter break

# Online update algorithms.

- Ranking modelled as **Rank Ordered Logit/BT Model.**

- Models are estimated using **Bayesian Approximation Method.**

$$R_i' \leftarrow R_i + K * (Y_i - \hat{Y}_i)$$

- Update doesn't require previous data.

- Computationaly efficient.

# Elo rating system

- Rating calculated using formula specified by FIDE.

- Points gain depends on rating difference

- No rating deviation in Elo formula

| D Rtg Dif | PD H  L | D Rtg Dif | PD H  L | D Rtg Dif | PD H  L | D Rtg Dif | PD H  L |
|---|---|---|---|---|---|---|---|
| 0-3 | .50 .50 | 92-98 | .63 .37 | 198-206 | .76 .24 | 345-357 | .89 .11 |
| 4-10 | .51 .49 | 99-106 | .64 .36 | 207-215 | .77 .23 | 358-374 | .90 .10 |
| 11-17 | .52 .48 | 107-113 | .65 .35 | 216-225 | .78 .22 | 375-391 | .91 .09 |
| 18-25 | .53 .47 | 114-121 | .66 .34 | 226-235 | .79 .21 | 392-411 | .92 .08 |
| 26-32 | .54 .46 | 122-129 | .67 .33 | 236-245 | .80 .20 | 412-432 | .93 .07 |
| 33-39 | .55 .45 | 130-137 | .68 .32 | 246-256 | .81 .19 | 433-456 | .94 .06 |
| 40-46 | .56 .44 | 138-145 | .69 .31 | 257-267 | .82 .18 | 457-484 | .95 .05 |
| 47-53 | .57 .43 | 146-153 | .70 .30 | 268-278 | .83 .17 | 485-517 | .96 .04 |
| 54-61 | .58 .42 | 154-162 | .71 .29 | 279-290 | .84 .16 | 518-559 | .97 .03 |
| 62-68 | .59 .41 | 163-170 | .72 .28 | 291-302 | .85 .15 | 560-619 | .98 .02 |
| 69-76 | .60 .40 | 171-179 | .73 .27 | 303-315 | .86 .14 | 620-735 | .99 .01 |
| 77-83 | .61 .39 | 180-188 | .74 .26 | 316-328 | .87 .13 | > 735 | 1.0 .00 |
| 84-91 | .62 .38 | 189-197 | .75 .25 | 329-344 | .88 .12 | | |

# Glicko rating system (`sport::glicko`)

- First bayesian rating system.

- Rating change depends on ratings (R) and ratings deviation (RD).

$$\hat{Y}_i = P(X_i > X_q) = \frac{1}{1 + 10^{-g(RD_{iq})*(R_i - R_q)/400}}$$

$$R'_i = R_i + \frac{1}{\frac{1}{RD_i^2} + \frac{1}{d_i^2}} * \sum_j g(RD_j) * (Y_{ij} - \hat{Y}_{ij})$$

$$RD'_i = \sqrt{(\frac{1}{RD_i^2} + \frac{1}{d_i^2})^{-1}}$$

# Glicko2 rating system (`sport::glicko2`)

- Volatile parameter σ added. Measures expected fluctuations.

- Updated rating deviation based on the 'Illinois Algorithm'.

$$\hat{Y}_{ij} = \frac{1}{1 + e^{-g(\phi_{ij})*(\mu_i - \mu_j)}}$$

$$\phi'_i = \frac{1}{\sqrt{\frac{1}{\phi_i^2 + \sigma_i'^2} + \frac{1}{v}}}$$

$$\mu'_i = \mu_i + \phi'_i * \sum_j g(\phi_j) * (Y_{ij} - \hat{Y}_{ij})$$

# Bayesian Bradley Terry (`sport::bbt`)

- Extends algorithms to multi-player teams.

- Team rating/variance is a sum of team players ratings/variances.

$$\hat{Y}_{ij} = P(X_i > X_j) = \frac{e^{R_i/c_{i_j}}}{e^{R_i/c_{ij}} + e^{R_j/c_{ij}}}$$

$$R'_i = R_i + \sum_j \frac{RD_i^2}{c_{ij}} * (Y_{ij} - \hat{Y}_{ij})$$

$$RD'_i = RD_i * [1 - \frac{RD_{ij}^2}{RD_i^2} \sum_j \gamma_j * (\frac{RD_i}{c_{ij}})^2 * \hat{Y}_{ij}\hat{Y}_{ji}]$$

# Bayesian Dynamic Logit (`sport::bdl`)

- EKF with logistic function as measurement equation.

- Teams/players are treated as alternative in discrete choice models.

$$Y_t = \frac{e^z}{1 + e^{z_t}}$$

$$z_t = \beta_{it}^T x_{it} - \beta_{jt}^T x_{jt}$$

$$\hat{w}_t = \hat{w}_{t-1} + \Sigma_t x_t (z_t - y_t)$$

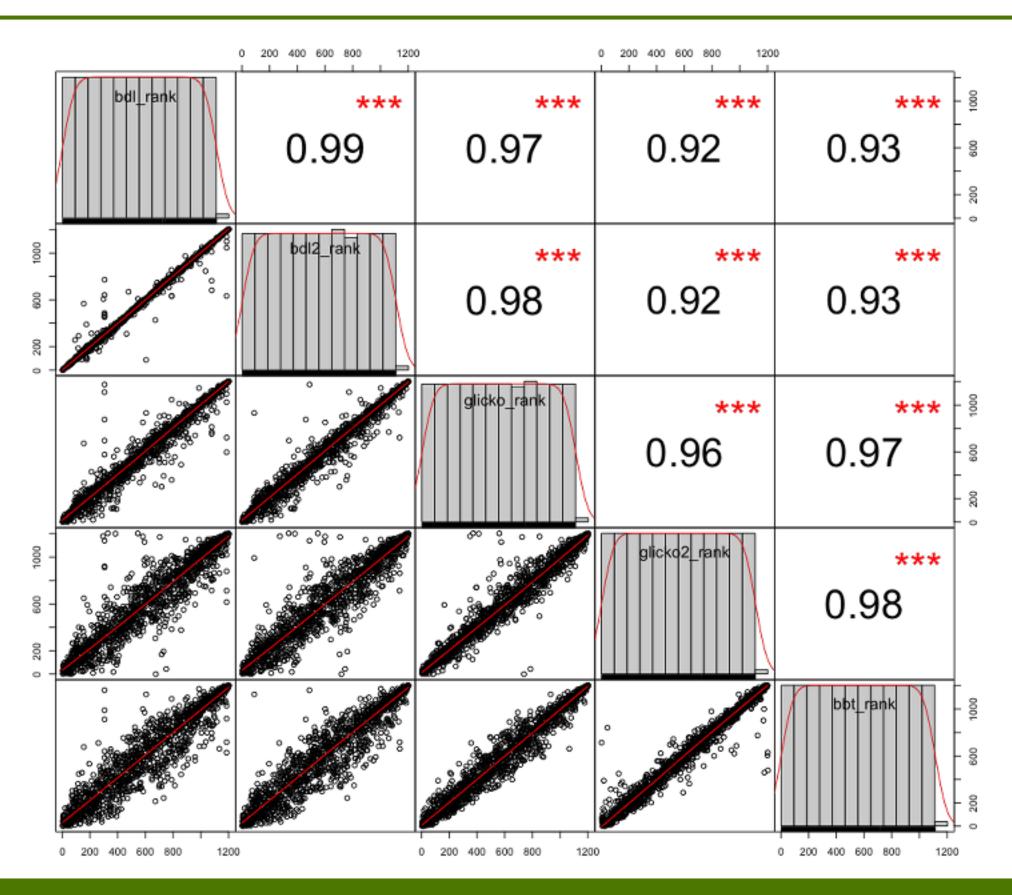$$s_t^2 = x_t^T (\Sigma_{t-1} + q_t I) x_t$$

# `sport` - package for sport analytics.

```r
# devtools::install_github("gogonzo/sport")
library(sport)

list_glicko  <- glicko_run(  formula = rank|id ~ rider_name, data = gpheats )
list_glicko2 <- glicko2_run( formula = rank|id ~ rider_name, data = gpheats )
list_bbt     <- bbt_run(     formula = rank|id ~ rider_name, data = gpheats )
list_bdl     <- bdl_run(     formula = rank|id ~ rider_name, data = gpheats )
> names(list_glicko)
[1] "r"        "pairs"    "final_r"  "final_rd"
> head(list_glicko$r)
  id        names        r       rd
1  1   Tomasz Gollob 1586.327 203.5029
2  1   Gary Havelock 1241.019 203.5029
3  1     Chris Louis 1758.981 203.5029
4  1 Tony Rickardsson 1413.673 203.5029
5  2   Sam Ermolenko 1758.981 203.5029
6  2  Jan Staechmann 1241.019 203.5029

> tail(list_glicko$pairs)
        id          team1            team2          P Y
61043 5063   Tai Woffinden Fredrik Lindgren 0.6333719 0
61044 5063   Tai Woffinden    Patryk Dudek 0.4877610 0
61045 5063   Tai Woffinden Emil Sajfutdinow 0.5177363 0
61046 5063 Emil Sajfutdinow Fredrik Lindgren 0.6167259 0
61047 5063 Emil Sajfutdinow    Patryk Dudek 0.4700839 0
61048 5063 Emil Sajfutdinow   Tai Woffinden 0.4822637 1
```

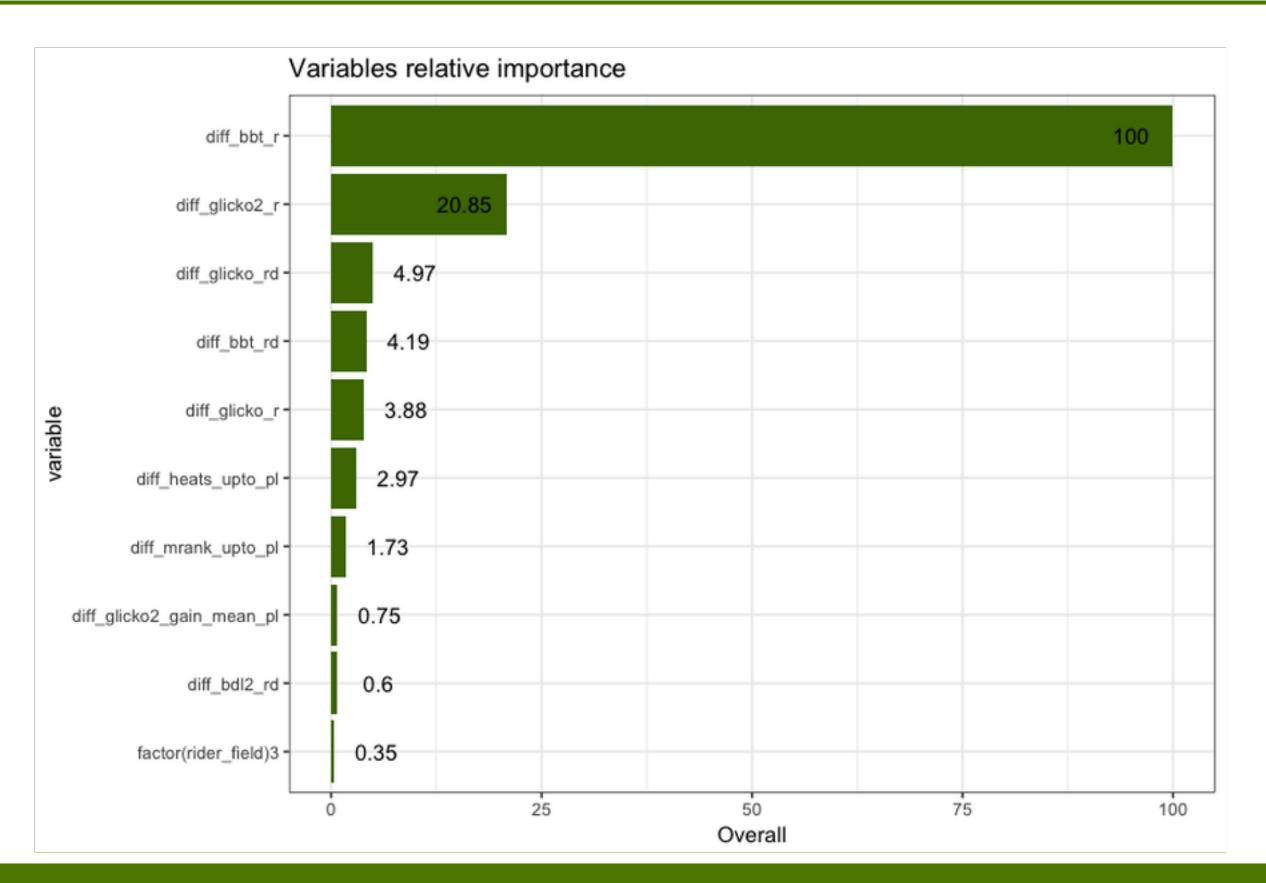# Methods perform similarly.

# ML methods don't improve model.

| Algorithm | Accuracy |
|-----------|----------|
| Glicko | 69.75% |
| Glicko2 | 70.20% |
| BBT | 71.06% |
| BDL | 69.09% |
| BDL2 | 68.81% |
| xgb1 | 71.01% |
| xgb2 | 71.00% |
| xgb3 | 71.02% |
| xgb4 | 71.67% |
| treeb1 | 70.16% |
| treeb2 | 69.65% |
| treeb3 | 69.97% |
| treeb4 | 69.03% |
| rf | 70.25% |

- xgBoost, Random Forests nor Boosted trees didn't improve accuracy.

- Additional variables have no additional predictive abilities.

# Non-rating vars contribution < 5%



Variables relative importance

| variable | Overall |
|---|---|
| diff_bbt_r | 100 |
| diff_glicko2_r | 20.85 |
| diff_glicko_rd | 4.97 |
| diff_bbt_rd | 4.19 |
| diff_glicko_r | 3.88 |
| diff_heats_upto_pl | 2.97 |
| diff_mrank_upto_pl | 1.73 |
| diff_glicko2_gain_mean_pl | 0.75 |
| diff_bdl2_rd | 0.6 |
| factor(rider_field)3 | 0.35 |

# The Best speedway riders.



Actual ratings

Ranking vs BBT Rating

| Rider | |
|---|---|
| Nicki Pedersen (41) | |
| Sebastian Tresarrieu (32) | |
| Artiom Łaguta (27) | |
| Tai Woffinden (27) | |
| Bartosz Zmarzlik (23) | |
| Greg Hancock (48) | |
| Fredrik Lindgren (32) | |
| Billy Janniro (38) | |
| Leon Madsen (29) | |
| Janusz Kołodziej (34) | |
| Maciej Janowski (26) | |
| Jarosław Hampel (36) | |
| Jason Doyle (32) | |
| Grigorij Łaguta (34) | |
| Emil Sajfutdinow (28) | |
| Patryk Dudek (26) | |
| Piotr Pawlicki (23) | |
| Michael Jepsen Jensen (26) | |
| Matej Žagar (35) | |
| Damian Sperz (28) | |

country_iso3
- AUS
- DNK
- FRA
- GBR
- POL
- RUS
- SVN
- SWE
- USA

# What comes next

- Searching for perfect predictive model.

- Improve `sport` package.

- Simulate events results.

- Promote data-science in speedway.

# References

- Mark E. Glickman (1999): Parameter estimation in large dynamic paired comparison experiments. Applied Statistics, 48:377-394.
URL http://www.glicko.net/research/glicko.pdf

- Mark E. Glickman (2001): Dynamic paired comparison models with stochastic variances, Journal of Applied Statistics, 28:673-689.
URL http://www.glicko.net/research/dpcmsv.pdf

- Mark E. Glickman (1995): A Comprehensive guide to chess ratings. American Chess Journal, 3, pp. 59--102.
URL http://www.glicko.net/research/acjpaper.pdf

- Ruby C. Weng and Chih-Jen Lin (2011): A Bayesian Approximation Method for Online Ranking. Journal of Machine Learning Research,12:267-300.
URL http://jmlr.csail.mit.edu/papers/volume12/weng11a/weng11a.pdf

- William D. Penny and Stephen J. Roberts (1999): Dynamic Logistic Regression, Departament of Electrical and Electronic Engineering, Imperial College

# Thank you

github.com/gogonzo
dawid.kaledkowski@clickmeeting.com
linkedin.com/in/dawidkaledkowski