

# A Comparison of Covariate-based Prediction Methods for FIFA World Cups

**A. Groll**

Faculty of Statistics,  
TU Dortmund University

(joint work with J. Abedieh, C. Ley, A. Mayr, T. Kneib, G. Schaubberger,  
G. Tutz & H. Van Eetvelde)

**Zurich R User Group Meetup**  
October 25<sup>th</sup> 2018, University of Zurich



# Who will celebrate?



Sources: youtube.com, EMAJ Magazine, youfrisky.com, Bailiwick Express

# Who will cry?



Sources: [youtube.com](https://www.youtube.com), [pinterest](https://www.pinterest.com), [BBC](https://www.bbc.com), [Daily Mail](https://www.dailymail.com)

# Theoretical Background

# Part I: Regression-based Methods

## Model for international soccer tournaments

$$y_{ijk} | \mathbf{x}_{ik}, \mathbf{x}_{jk} \sim \text{Pois}(\lambda_{ijk}) \quad i, j \in \{1, \dots, n\}, i \neq j$$

$$\lambda_{ijk} = \exp(\beta_0 + (\mathbf{x}_{ik} - \mathbf{x}_{jk})^\top \boldsymbol{\beta})$$

$n$ : Number of teams

$y_{ijk}$ : Number of goals scored by team  $i$  against opponent  $j$  at tournament  $k$

$\mathbf{x}_{ik}, \mathbf{x}_{jk}$ : Covariate vectors of team  $i$  and opponent  $j$  varying over tournaments

$\boldsymbol{\beta}$ : Parameter vector of covariate effects

## Regularized estimation

Maximize penalized log-likelihood

$$l_p(\beta_0, \beta) = l(\beta_0, \beta) - \lambda J(\beta)$$

## Regularized estimation

Maximize penalized log-likelihood

$$\begin{aligned}l_p(\beta_0, \boldsymbol{\beta}) &= l(\beta_0, \boldsymbol{\beta}) - \lambda J(\boldsymbol{\beta}) \\ &= l(\beta_0, \boldsymbol{\beta}) - \lambda \sum_{i=1}^p |\beta_i|,\end{aligned}$$

with lasso penalty term (Tibshirani, 1996):

$$J(\boldsymbol{\beta}) = \sum_{i=1}^p |\beta_i|.$$

The model can be estimated with the R-package `glmnet` (Friedman et al., 2010).



## Regularized estimation

Maximize penalized log-likelihood

$$\begin{aligned}l_p(\beta_0, \boldsymbol{\beta}) &= l(\beta_0, \boldsymbol{\beta}) - \lambda J(\boldsymbol{\beta}) \\ &= l(\beta_0, \boldsymbol{\beta}) - \lambda \sum_{i=1}^p |\beta_i|,\end{aligned}$$

with lasso penalty term (Tibshirani, 1996):

$$J(\boldsymbol{\beta}) = \sum_{i=1}^p |\beta_i|.$$

The model can be estimated with the R-package `glmnet` (Friedman et al., 2010).

Versions used for: EURO 2012 (Groll and Abedieh, 2013); World Cup 2014 (Groll et al., 2015); EURO 2016 (Groll et al., 2018)

# Part II: Ranking Methods

## Independent Poisson ranking model

$$Y_{ijm} \sim \text{Pois}(\lambda_{ijm}),$$

$$\lambda_{ijm} = \exp\left(\beta_0 + (r_i - r_j) + h \cdot \mathbb{I}(\text{team } i \text{ playing at home})\right)$$

$n$ : Number of teams

$M$ : Number of matches

$y_{ijm}$ : Number of goals scored by team  $i$  against opponent  $j$  in match  $m$

$r_i, r_j$ : strengths / ability parameters of team  $i$  and team  $j$

$h$ : home effect; added if team  $i$  plays at home

# Independent Poisson ranking model

**Likelihood function:**

$$L = \prod_{m=1}^M \left( \frac{\lambda_{ijm}^{y_{ijm}}}{y_{ijm}!} \exp(-\lambda_{ijm}) \cdot \frac{\lambda_{jim}^{y_{jim}}}{y_{jim}!} \exp(-\lambda_{jim}) \right)^{w_{type,m} \cdot w_{time,m}},$$

with weights

$$w_{time,m}(t_m) = \left(\frac{1}{2}\right)^{\overline{\text{Half period}}^{t_m}}$$

and

$$w_{type,m} \in \{1, 2, 3, 4\} \quad (\text{depending on type of match}).$$

# Independent Poisson ranking model

**Likelihood function:**

$$L = \prod_{m=1}^M \left( \frac{\lambda_{ijm}^{y_{ijm}}}{y_{ijm}!} \exp(-\lambda_{ijm}) \cdot \frac{\lambda_{jim}^{y_{jim}}}{y_{jim}!} \exp(-\lambda_{jim}) \right)^{w_{type,m} \cdot w_{time,m}},$$

with weights

$$w_{time,m}(t_m) = \left( \frac{1}{2} \right)^{\overline{\text{Half period}}^{t_m}}$$

and

$$w_{type,m} \in \{1, 2, 3, 4\} \quad (\text{depending on type of match}).$$

Different extensions, for example, **bivariate Poisson models**. Ley et al. (2018) show that bivariate Poisson with Half Period of 3 years is best for prediction.

# Part III: Random Forests

# Random Forests

- introduced by Breiman (2001)
- **principle**: aggregation of (large) number of **classification** / **regression trees**  
⇒ can be used both for classification & regression purposes

# Random Forests

- introduced by Breiman (2001)
- **principle**: aggregation of (large) number of **classification / regression trees**  
⇒ can be used both for classification & regression purposes
- **final predictions**: single tree predictions are aggregated, either by majority vote (classification) or by averaging (regression)



# Random Forests

- introduced by Breiman (2001)
- **principle**: aggregation of (large) number of **classification** / **regression trees**  
⇒ can be used both for classification & regression purposes
- **final predictions**: single tree predictions are aggregated, either by majority vote (classification) or by averaging (regression)
- feature space is partitioned recursively, each partition has its own prediction

# Random Forests

- introduced by Breiman (2001)
- **principle**: aggregation of (large) number of **classification / regression trees**  
⇒ can be used both for classification & regression purposes
- **final predictions**: single tree predictions are aggregated, either by majority vote (classification) or by averaging (regression)
- feature space is partitioned recursively, each partition has its own prediction
- find split with strongest difference between the two new partitions w.r.t. some criterion

# Random Forests

- introduced by Breiman (2001)
- **principle**: aggregation of (large) number of **classification / regression trees**  
⇒ can be used both for classification & regression purposes
- **final predictions**: single tree predictions are aggregated, either by majority vote (classification) or by averaging (regression)
- feature space is partitioned recursively, each partition has its own prediction
- find split with strongest difference between the two new partitions w.r.t. some criterion
- Observations within the same partition as similar as possible, observations from different partitions very different (w.r.t. response variable)

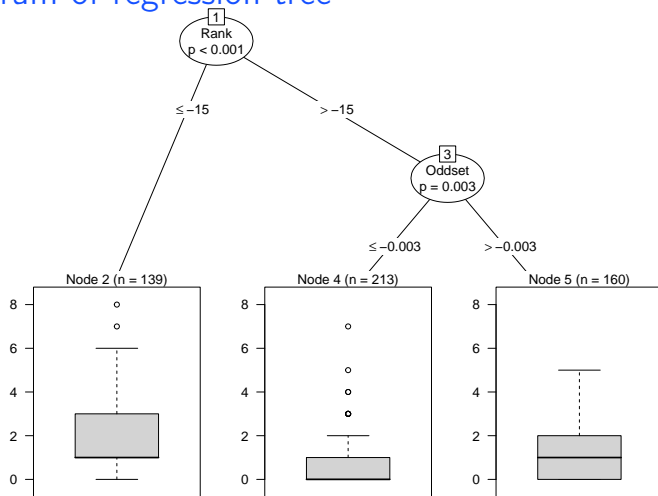
# Random Forests

- introduced by Breiman (2001)
- **principle**: aggregation of (large) number of **classification** / **regression trees**  
⇒ can be used both for classification & regression purposes
- **final predictions**: single tree predictions are aggregated, either by majority vote (classification) or by averaging (regression)
- feature space is partitioned recursively, each partition has its own prediction
- find split with strongest difference between the two new partitions w.r.t. some criterion
- Observations within the same partition as similar as possible, observations from different partitions very different (w.r.t. response variable)
- a single tree is usually pruned (lower variance but increases bias)

# Random Forests

- introduced by Breiman (2001)
- **principle**: aggregation of (large) number of **classification / regression trees**  
⇒ can be used both for classification & regression purposes
- **final predictions**: single tree predictions are aggregated, either by majority vote (classification) or by averaging (regression)
- feature space is partitioned recursively, each partition has its own prediction
- find split with strongest difference between the two new partitions w.r.t. some criterion
- Observations within the same partition as similar as possible, observations from different partitions very different (w.r.t. response variable)
- a single tree is usually pruned (lower variance but increases bias)
- visualized in dendrogram

## Dendrogram of regression tree



Exemplary regression tree for FIFA World Cup 2002 – 2014 data using the function `ctree` from the R-package `party` (Hothorn et al., 2006). **Response:** *Number of goals*; **predictors:** only *FIFA Rank* and *Oddset* are used.

# Random Forests

- repeatedly grow different regression trees
- main goal: decrease variance

# Random Forests

- repeatedly grow different regression trees
- main goal: decrease variance  $\implies$  decrease correlation between single trees.



# Random Forests

- repeatedly grow different regression trees
- main goal: decrease variance  $\implies$  decrease correlation between single trees.
- $\implies$  two different randomisation steps:
  - 1) trees are not applied to the original sample but to **bootstrap samples** or random subsamples of the data.
  - 2) at each node a **(random) subset of the predictors** is drawn that are used to find the best split.

# Random Forests

- repeatedly grow different regression trees
- main goal: decrease variance  $\implies$  decrease correlation between single trees.
- $\implies$  two different randomisation steps:
  - 1) trees are not applied to the original sample but to **bootstrap samples** or random subsamples of the data.
  - 2) at each node a **(random) subset of the predictors** is drawn that are used to find the best split.
- by de-correlating and combining many trees  $\implies$  predictions with **low bias and reduced variance**

# Random Forests for Soccer

- response: metric variable *Number of Goals*
- predefined number of trees  $B$  (e.g.,  $B = 5000$ ) is fitted based on (bootstrap samples of) the training data
- prediction of new observation: covariate values are dropped down each of the regression trees, resulting in  $B$  predictions  $\implies$  average
- use predicted expected value as event rate  $\hat{\lambda}$  of a Poisson distribution  $Po(\lambda)$

# Random Forests for Soccer

- response: metric variable *Number of Goals*
- predefined number of trees  $B$  (e.g.,  $B = 5000$ ) is fitted based on (bootstrap samples of) the training data
- prediction of new observation: covariate values are dropped down each of the regression trees, resulting in  $B$  predictions  $\implies$  average
- use predicted expected value as event rate  $\hat{\lambda}$  of a Poisson distribution  $Po(\lambda)$
- 2 slightly different variants:
  - 1) classical RF algorithm proposed by Breiman (2001) from the R-package *ranger* (Wright and Ziegler, 2017)
  - 2) RFs based conditional inference trees: *cforest* from the *party* package (Hothorn et al., 2006)

# Application to FIFA World Cups

# Covariates

**Data basis: World Cups 2002–2014**

# Covariates

## Data basis: World Cups 2002–2014

- **Economic Factors:**  
GDP per capita, population

# Covariates

## Data basis: World Cups 2002–2014

- **Economic Factors:**  
GDP per capita, population
  
- **Sportive Factors:**  
bookmaker's odds (Oddset), FIFA rank



# Covariates

## Data basis: World Cups 2002–2014

- **Economic Factors:**  
GDP per capita, population
- **Sportive Factors:**  
bookmaker's odds (Oddset), FIFA rank
- **Home advantage:**  
host of the world cup, same continent as host, continent

# Covariates

## Data basis: World Cups 2002–2014

- **Economic Factors:**  
GDP per capita, population
- **Sportive Factors:**  
bookmaker's odds (Oddset), FIFA rank
- **Home advantage:**  
host of the world cup, same continent as host, continent
- **Factors describing the team's structure**  
(Second) Maximum number of teammates, average age, number of Champions League & Europa League players, number of players abroad

# Covariates

## Data basis: World Cups 2002–2014

- **Economic Factors:**  
GDP per capita, population
- **Sportive Factors:**  
bookmaker's odds (Oddset), FIFA rank
- **Home advantage:**  
host of the world cup, same continent as host, continent
- **Factors describing the team's structure**  
(Second) Maximum number of teammates, average age, number of Champions League & Europa League players, number of players abroad
- **Factors describing the team's coach**  
age, nationality, tenure





# Covariates

## Data basis: World Cups 2002–2014

- **Economic Factors:**  
GDP per capita, population
- **Sportive Factors:**  
bookmaker's odds (Oddset), FIFA rank
- **Home advantage:**  
host of the world cup, same continent as host, continent
- **Factors describing the team's structure**  
(Second) Maximum number of teammates, average age, number of Champions League & Europa League players, number of players abroad
- **Factors describing the team's coach**  
age, nationality, tenure





All **variables** are incorporated as differences between the team whose goals are considered and its opponent!

## Extract of the design matrix

FRA  0:0  URU  
URU  1:2  DEN

Team	Age	Rank	Oddset	...
France	28.3	1	0.149	...
Uruguay	25.3	24	0.009	...
Denmark	27.4	20	0.012	...
⋮	⋮	⋮	⋮	⋮

## Extract of the design matrix

FRA  0:0  URU  
URU  1:2  DEN

Team	Age	Rank	Oddset	...
France	28.3	1	0.149	...
Uruguay	25.3	24	0.009	...
Denmark	27.4	20	0.012	...
⋮	⋮	⋮	⋮	⋮

Goals	Team	Opponent	Age	Rank	Oddset	...
0	France	Uruguay	3.00	-23	0.140	...
0	Uruguay	France	-3.00	23	-0.140	...
1	Uruguay	Denmark	-2.10	4	-0.003	...
2	Denmark	Uruguay	2.10	-4	0.003	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮

## Comparison of predictive performance: WC 2002-2014 data

1. Form a training data set containing 3 out of 4 World Cups.
2. Fit each of the methods to the training data.
3. Predict the left-out World Cup using each of the prediction methods.
4. Iterate steps 1-3 such that each World Cup is once the left-out one.
5. Compare predicted and real outcomes for all prediction methods.

## Comparison of predictive performance: WC 2002-2014 data

1. Form a training data set containing 3 out of 4 World Cups.
2. Fit each of the methods to the training data.
3. Predict the left-out World Cup using each of the prediction methods.
4. Iterate steps 1-3 such that each World Cup is once the left-out one.
5. Compare predicted and real outcomes for all prediction methods.

**We combine both the random forest and the LASSO with the ability estimates from the ranking method!**



## Prediction of match outcomes

- true ordinal match outcomes:  $\tilde{y}_1, \dots, \tilde{y}_N$  with  $\tilde{y}_i \in \{1, 2, 3\}$ , for all matches  $N$  from the 4 World Cups.
- predicted probabilities  $\hat{\pi}_{1i}, \hat{\pi}_{2i}, \hat{\pi}_{3i}$ ,  $i = 1, \dots, N$ ,
- Let  $G_{1i}$  and  $G_{2i}$  denote the goals scored by 2 competing teams in match  $i$   
 $\implies$  compute  $\hat{\pi}_{1i} = P(G_{1i} > G_{2i})$ ,  $\hat{\pi}_{2i} = P(G_{1i} = G_{2i})$  and  $\hat{\pi}_{3i} = P(G_{1i} < G_{2i})$   
based on the corresponding Poisson distributions  $G_{1i} \sim Po(\hat{\lambda}_{1i})$  and  $G_{2i} \sim Po(\hat{\lambda}_{2i})$  with estimates  $\hat{\lambda}_{1i}$  and  $\hat{\lambda}_{2i}$  (Skellam distribution)
- **benchmark: bookmakers**  $\implies$  compute the 3 quantities  $\tilde{\pi}_{ri} = 1/\text{odds}_r$ ,  $r \in \{1, 2, 3\}$ , normalize with  $c_i := \sum_{r=1}^3 \tilde{\pi}_{ri}$  (adjust for bookmakers' margins)  
 $\implies$  estimated probabilities  $\hat{\pi}_{ri} = \tilde{\pi}_{ri}/c_i$

# Prediction of match outcomes

## 3 Performance measures:

- (a) **multinomial likelihood** (probability of correct prediction): for single match defined as

$$\hat{\pi}_{1i}^{\delta_{1\tilde{y}_i}} \hat{\pi}_{2i}^{\delta_{2\tilde{y}_i}} \hat{\pi}_{3i}^{\delta_{3\tilde{y}_i}},$$

with  $\delta_{ri}$  denoting Kronecker's delta

- (b) **classification rate**: is match  $i$  correctly classified using the indicator function

$$\mathbb{I}(\tilde{y}_i = \arg \max_{r \in \{1,2,3\}} (\hat{\pi}_{ri}))$$

- (c) **rank probability score** (RPS; explicitly accounts for the ordinal structure):

$$\frac{1}{3-1} \sum_{r=1}^{3-1} \left( \sum_{l=1}^r \hat{\pi}_{li} - \delta_{l\tilde{y}_i} \right)^2$$

## Prediction of match outcomes

	Likelihood	Class. Rate	RPS
Hybrid Random Forest	0.419	0.556	0.187
Random Forest	0.410	0.548	0.192
Ranking	0.415	0.532	0.190
Lasso	0.419	0.524	0.198
Hybrid Lasso	0.429	0.540	0.194
Bookmakers	0.425	0.524	0.188

Comparison of different prediction methods for ordinal outcome based on multinomial likelihood, classification rate and ranked probability score (RPS)

## Prediction of exact numbers of goals

- let now  $y_{ijk}$ , for  $i, j = 1, \dots, n$  and  $k \in \{2002, 2006, 2010, 2014\}$ , denote the observed number of goals scored by team  $i$  against team  $j$  in tournament  $k$
- $\hat{y}_{ijk}$  the corresponding predicted value
- 2 quadratic errors:  $(y_{ijk} - \hat{y}_{ijk})^2$  and  $((y_{ijk} - y_{jik}) - (\hat{y}_{ijk} - \hat{y}_{jik}))^2$

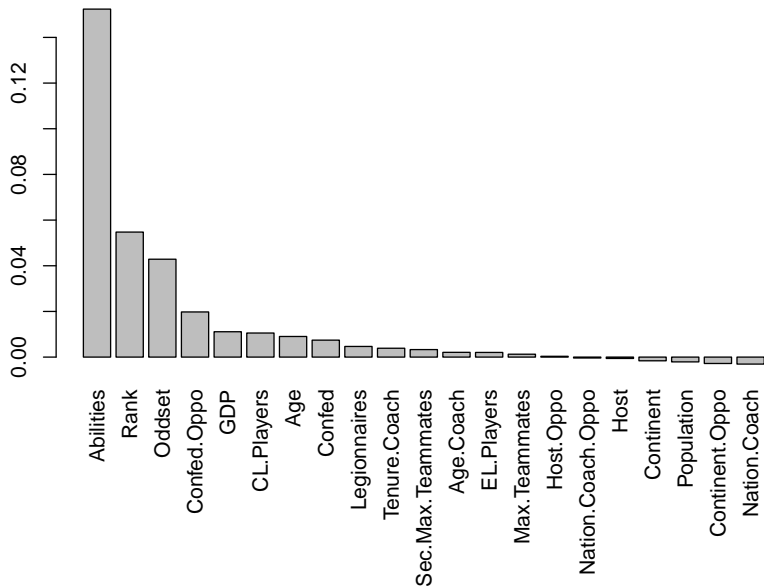
## Prediction of exact numbers of goals

	Goal Difference	Goals
Hybrid Random Forest	2.473	1.296
Random Forest	2.543	1.330
Ranking	2.560	1.349
Lasso	2.835	1.421
Hybrid Lasso	2.809	1.427













Comparison of different prediction methods for the exact number of goals and the goal difference based on MSE

# Prediction of FIFA World Cup 2018

## Variable importance



































## Winning probabilities

			Round of 16	Quarter finals	Semi finals	Final	World Champion	Oddset
1.		ESP	88.4	73.1	47.9	28.9	17.8	11.8
2.		GER	86.5	58.0	39.8	26.3	17.1	15.0
3.		BRA	83.5	51.6	34.1	21.9	12.3	15.0
4.		FRA	85.5	56.1	36.9	20.8	11.2	11.8
5.		BEL	86.3	64.5	35.7	20.4	10.4	8.3
6.		ARG	81.6	50.5	29.8	15.2	7.3	8.3
7.		ENG	79.8	57.0	29.8	15.6	7.1	4.6
8.		POR	67.5	46.1	19.8	7.3	2.5	3.8
9.		CRO	65.9	30.8	15.6	6.0	2.2	3.0
10.		SUI	58.9	30.6	13.1	5.6	2.2	1.0
11.		COL	79.2	33.1	14.0	5.7	2.1	1.8
12.		DEN	59.0	26.1	12.4	4.8	1.7	1.1
∴	∴	∴	∴	∴	∴	∴	∴	∴

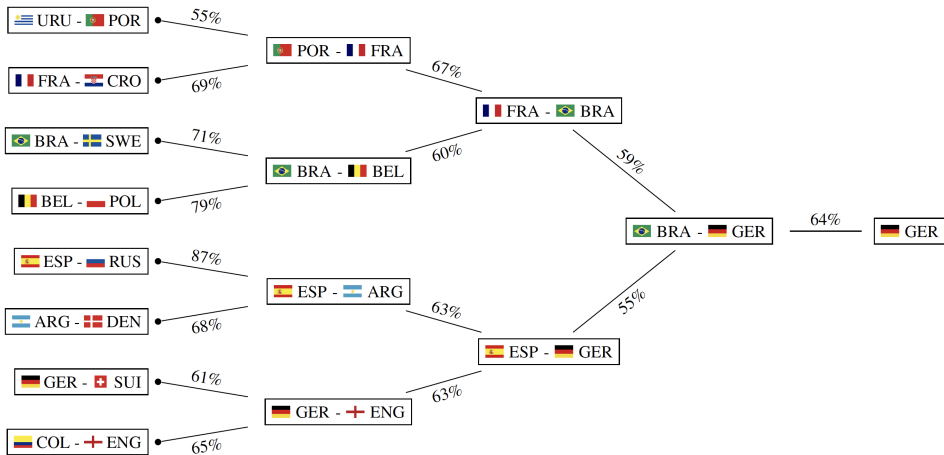


# Most probable group stage

Group A 28.7%	Group B 38.5%	Group C 31.5%	Group D 30.7%
1.  URU	1.  ESP	1.  FRA	1.  ARG
2.  RUS	2.  POR	2.  DEN	2.  CRO
 KSA	 MOR	 AUS	 ICE
 EGY	 IRN	 PER	 NGA

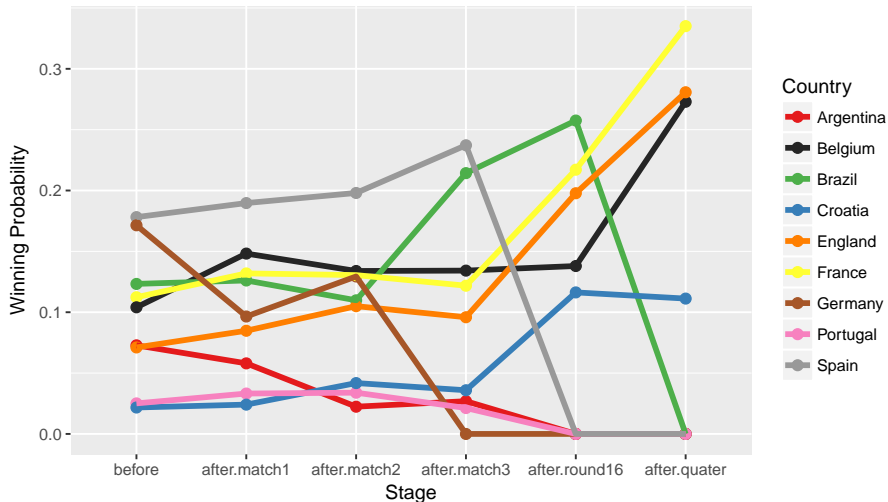
Group E 29.0%	Group F 29.9%	Group G 38.1%	Group H 26.5%
1.  BRA	1.  GER	1.  BEL	1.  COL
2.  SUI	2.  SWE	2.  ENG	2.  POL
 CRC	 MEX	 PAN	 SEN
 SRB	 KOR	 TUN	 JPN

# Most probable knockout stage



# Winning probabilities over time

Time course of the winning probabilities for the nine (originally) favored teams:



## Performance I

	Likelihood	Class. Rate	RPS
Hybrid Random Forest	0.440	0.609	0.188
Random Forest	0.433	0.609	0.191
Lasso	0.424	0.547	0.207
Hybrid Lasso	0.434	0.609	0.201
Ranking	0.423	0.578	0.197
Bookmakers	0.438	0.562	0.194

## Performance I

	Likelihood	Class. Rate	RPS
Hybrid Random Forest	0.440	0.609	0.188
Random Forest	0.433	0.609	0.191
Lasso	0.424	0.547	0.207
Hybrid Lasso	0.434	0.609	0.201
Ranking	0.423	0.578	0.197
Bookmakers	0.438	0.562	0.194

	Goal Difference	Goals
Hybrid Random Forest	1.181	2.113
Random Forest	1.209	2.177
Lasso	1.216	2.333
Hybrid Lasso	1.187	2.270
Ranking	1.253	2.171

# Performance II

Final standing in forecast competition [fifaexperts.com](http://fifaexperts.com) (> 500 participants):

Submit your forecasts

Check your results

Scoreboard

Your league

1. Esportes em Números: 4650 points
2. Andreas Groll: 4644 points
3. Danilo Lopes: 4634 points
4. Natanael Prata: 4634 points
5. Chance de Gol: 4611 points
6. Wilson Chaves: 4597 points
7. Sigma Benedek: 4589 points
8. Márcio Diniz: 4587 points
9. Francesco Beatrice: 4574 points
10. Alun Owen: 4565 points
11. Tolstói Tói: 4558 points
12. Magne Aldrin: 4557 points

# Performance III

Final standing in forecast competition **Kicktipp** (with colleagues):

## Gesamtübersicht

Spieltagspunkte ▼



Spieltage

Pos	Name	1	2	3	4	5	6	7	Ac	Vi	Ha	Fi	B	S	G
1	stats_model	14	13	14	9	12	10	19	13	7	4	4	28	2,50	147
2	Hendrik	20	14	9	9	11	5	8	12	9	4	0	28	1,83	129
3	Katharina	12	11	9	10	15	10	11	16	7	3	2	20	1,50	126
4	Katrin	12	14	8	6	12	4	15	18	7	4	2	24	0,83	126
5	Lukas	10	12	9	6	9	6	4	15	7	3	6	32	1,00	119
6	Jona	10	9	6	10	9	6	11	12	8	6	7	24	1,00	118
7	Hilsi	16	8	7	7	10	2	6	14	9	7	2	24	1,50	112
8	Borussenengel	13	10	10	11	14	2	5	14	5	4	2	16	1,00	106

## Performance IV

Final standing in **WC-forecast competition** from Prof. Claus Ekstrøm :

	log.loss
Groll, Ley, Schauburger, VanEetvelde	-11.69
Ekstrom (Skellam)	-11.72
Ekstrom (ELO)	-13.48
Random guessing	-14.56

And the winner is the prediction by **Groll, Ley, Schauburger, VanEetvelde** (although not by much). Well done! Time to prepare the prediction algorithms for the next tournament - and hopefully we can get more people to participate.



# Summary

## Regarded models & predictive performance:

- (Regularized) regression approaches vs. random forests vs. ranking methods
- random forests & ranking methods perform pretty good (almost as good as bookmakers)

# Summary

## Regarded models & predictive performance:

- (Regularized) regression approaches vs. random forests vs. ranking methods
- random forests & ranking methods perform pretty good (almost as good as bookmakers)
- $\implies$  combine random forests & ranking methods to **hybrid random forest**

# Summary

## Regarded models & predictive performance:

- (Regularized) regression approaches vs. random forests vs. ranking methods
- random forests & ranking methods perform pretty good (almost as good as bookmakers)
- $\implies$  combine random forests & ranking methods to **hybrid random forest**
- $\implies$  combination outperforms bookmakers (on FIFA WC 2002 – 2014 data)

# Summary

## Regarded models & predictive performance:

- (Regularized) regression approaches vs. random forests vs. ranking methods
- random forests & ranking methods perform pretty good (almost as good as bookmakers)
- $\implies$  combine random forests & ranking methods to **hybrid random forest**
- $\implies$  combination outperforms bookmakers (on FIFA WC 2002 – 2014 data)

## FIFA WC 2018 prediction:

- Spain favorite with 17.8%, closely follow by Germany (17.1%); then: Brazil, France, Belgium (**before the tournament start**)

# Summary

## Regarded models & predictive performance:

- (Regularized) regression approaches vs. random forests vs. ranking methods
- random forests & ranking methods perform pretty good (almost as good as bookmakers)
- $\implies$  combine random forests & ranking methods to **hybrid random forest**
- $\implies$  combination outperforms bookmakers (on FIFA WC 2002 – 2014 data)

## FIFA WC 2018 prediction:

- Spain favorite with 17.8%, closely follow by Germany (17.1%); then: Brazil, France, Belgium (**before the tournament start**)
- Performance: Germany & Spain already dropped out; **but**: very good performance **on average**

# Summary

## Regarded models & predictive performance:

- (Regularized) regression approaches vs. random forests vs. ranking methods
- random forests & ranking methods perform pretty good (almost as good as bookmakers)
- $\implies$  combine random forests & ranking methods to **hybrid random forest**
- $\implies$  combination outperforms bookmakers (on FIFA WC 2002 – 2014 data)

## FIFA WC 2018 prediction:

- Spain favorite with 17.8%, closely follow by Germany (17.1%); then: Brazil, France, Belgium (**before the tournament start**)
- Performance: Germany & Spain already dropped out; **but**: very good performance **on average**
- Conclusion: single match outcome / tournament winner almost impossible to predict, but in general very adequate model

# References

-  Breiman, L. (2001). Random forests. *Machine Learning* 45, 5–32.
-  Friedman, J., T. Hastie and R. Tibshirani (2010): Regularization paths for generalized linear models via coordinate descent, *Journal of Statistical Software*, 33, 1.
-  Groll, A. and J. Abedieh (2013). Spain retains its title and sets a new record - generalized linear mixed models on European football championships. *Journal of Quantitative Analysis in Sports* 9(1), 51–66.
-  Groll, A., G. Schauburger, and G. Tutz (2015). Prediction of major international soccer tournaments based on team-specific regularized Poisson regression: An application to the FIFA World Cup 2014. *Journal of Quantitative Analysis in Sports* 11(2), 97–115.
-  Groll, A., T. Kneib, A. Mayr, and G. Schauburger (2018). On the dependency of soccer scores – A sparse bivariate Poisson model for the UEFA European Football Championship 2016. *Journal of Quantitative Analysis in Sports* 14(2), 65–79.

## References II



Hothorn, T., K. Hornik, and A. Zeileis (2006). Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical Statistics* 15, 651–674.



Ley, C., T. Van de Wiele and H. Van Eetvelde (2018): Ranking soccer teams on basis of their current strength: a comparison of maximum likelihood approaches, *Statistical Modelling*, submitted.



Wright, M. N. and A. Ziegler (2017). ranger: A fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software* 77(1), 1–17.



Schauberger, G. and Groll, A. (2018). Predicting matches in international football tournaments with random forests. *Statistical Modelling*, to appear.



Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B* 58, 267–288.





Thank you for  
your attention!

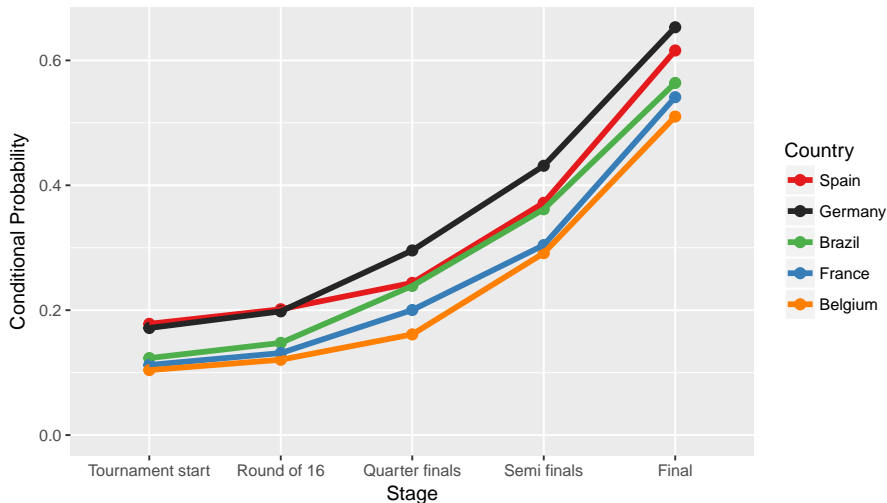
(Working paper on arXiv: <https://arxiv.org/pdf/1806.03208.pdf>)



















Sources: Forbes, JewishNews.com

## Conditional winning probabilities

Winning probabilities conditional on reaching the single stages of the tournament for the five favored teams:



## Winning probabilities after group stage

			Quarter finals	Semi finals	Final	World Champion
1.		ESP	88.2	61.1	42.2	23.7
2.		BRA	79.9	51.2	35.6	21.4
3.		BEL	85.1	40.9	24.1	13.4
4.		FRA	63.4	43.6	22.1	12.2
5.		ENG	71.6	45.4	20.1	9.6
6.		SUI	60.6	24.1	9.7	3.6
7.		CRO	56.1	20.8	10.2	3.6
8.		ARG	36.6	21.6	7.0	2.7
9.		DEN	43.9	15.2	6.8	2.4
10.		POR	55.1	19.0	5.5	2.1
11.		COL	28.4	15.9	5.2	1.8
12.		SWE	39.4	14.7	5.1	1.5
13.		URU	44.9	15.8	4.0	1.4
14.		MEX	20.1	4.7	1.2	0.3
15.		RUS	11.8	2.8	0.7	0.1
16.		JPN	14.9	3.1	0.6	0.1