

UEFA EURO 2016 Tournament Predictive Modeling

Atanas Mirchev Georgi Dikov

Statistical Modeling and Machine Learning
8th July 2016

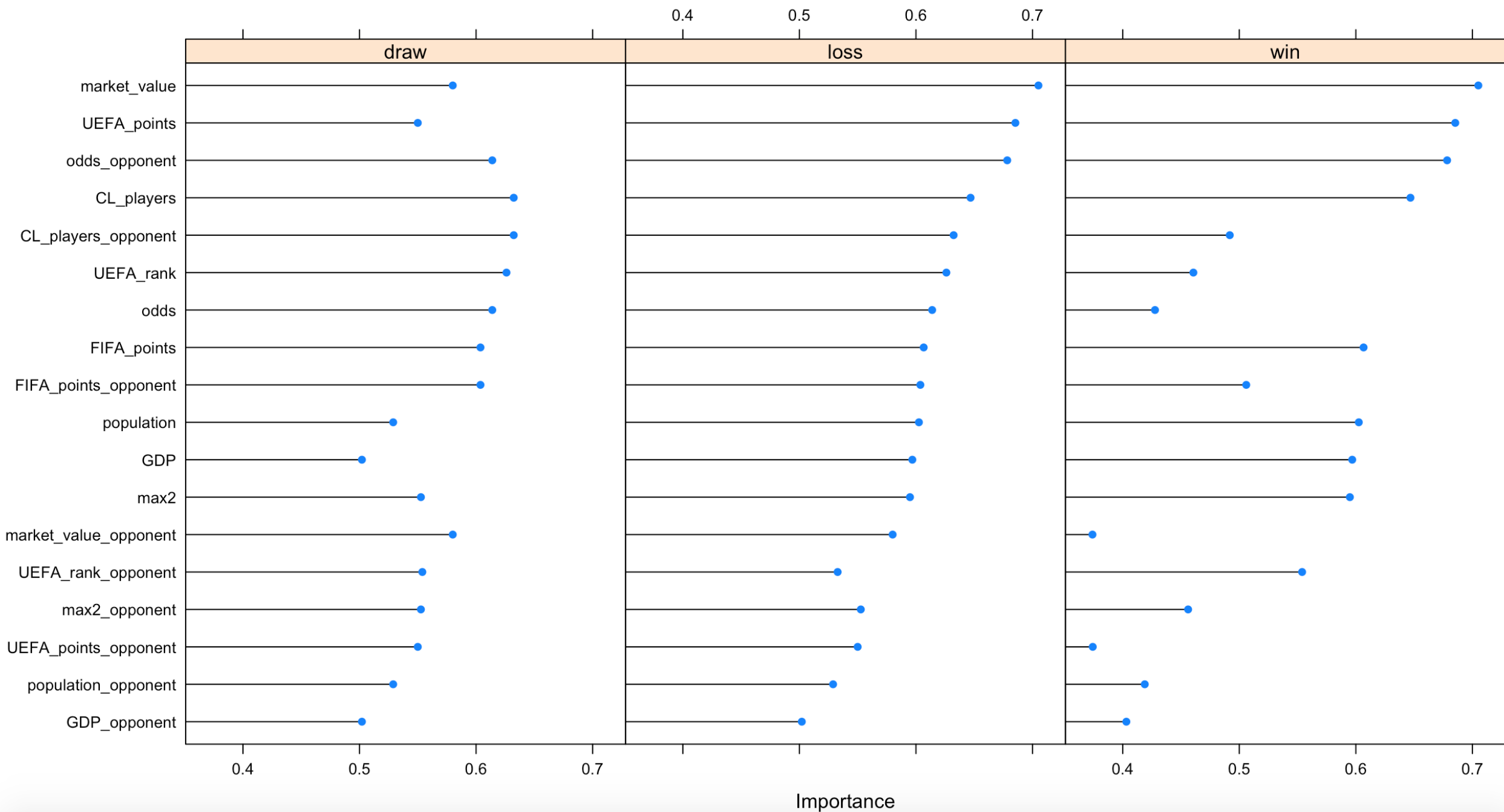
Modeling Concept

- Random Forest as a baseline
- Multivariate Logistic Regression with categorical dependent variable
 - Win/Draw/Loss classification + confidence
 - Interpretation of the co-variate's importance
 - Simple to implement, successful
- LogRegression with hard “draw probability” cut-off
 - Predict draw or not-draw
 - If not-draw then predict loss/win
- Other models for comparison: SVM, Poisson Regression

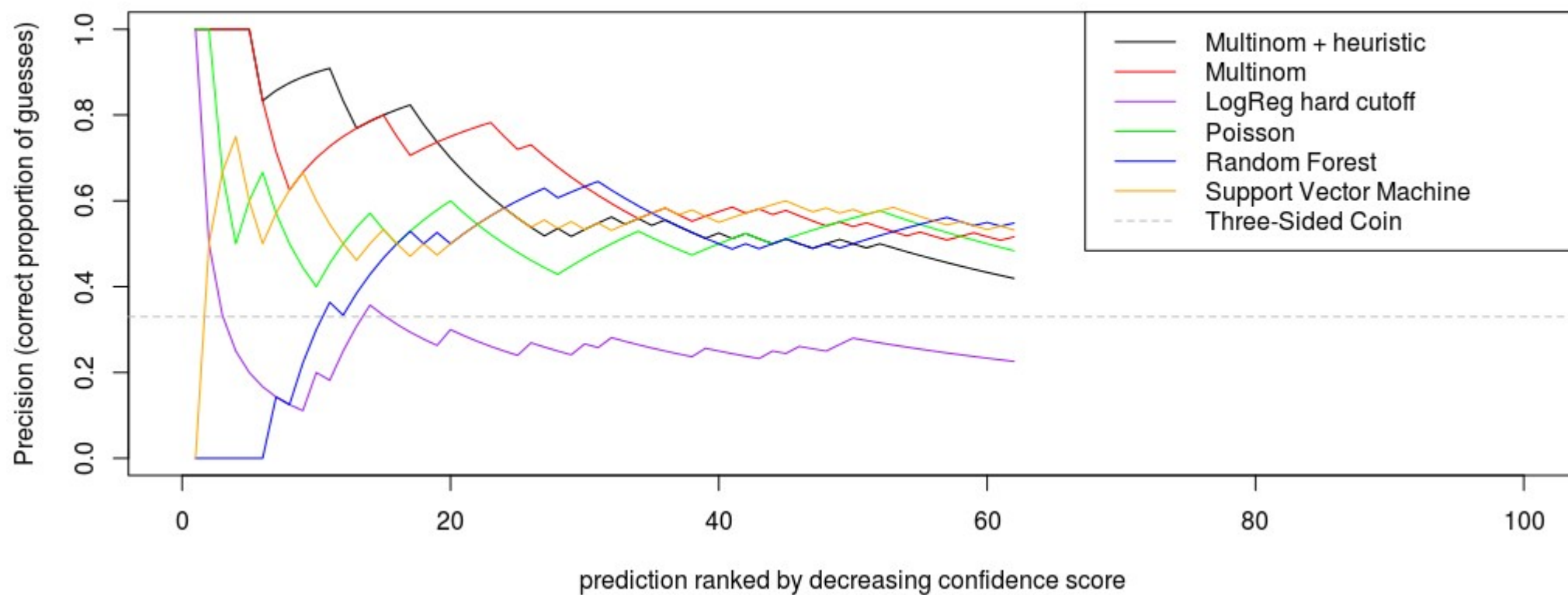
Data Preparation

Formatting, data clean-up and correction

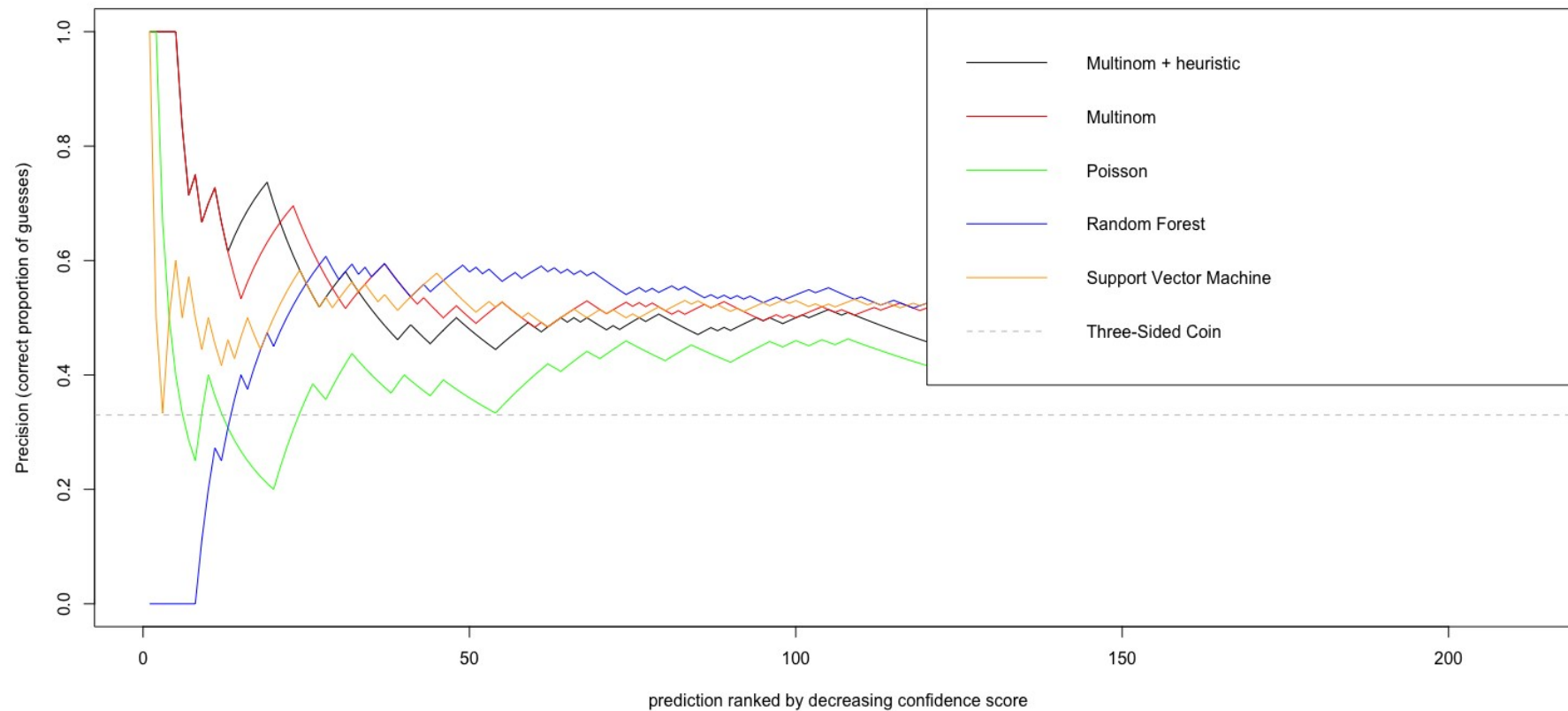
- Extending with new features
 - Nation's GDP, population
 - Weighted historical scoring in previous editions of the tournament
- Feature selection (LS-SVM)
- Data visualization (t-SNE)



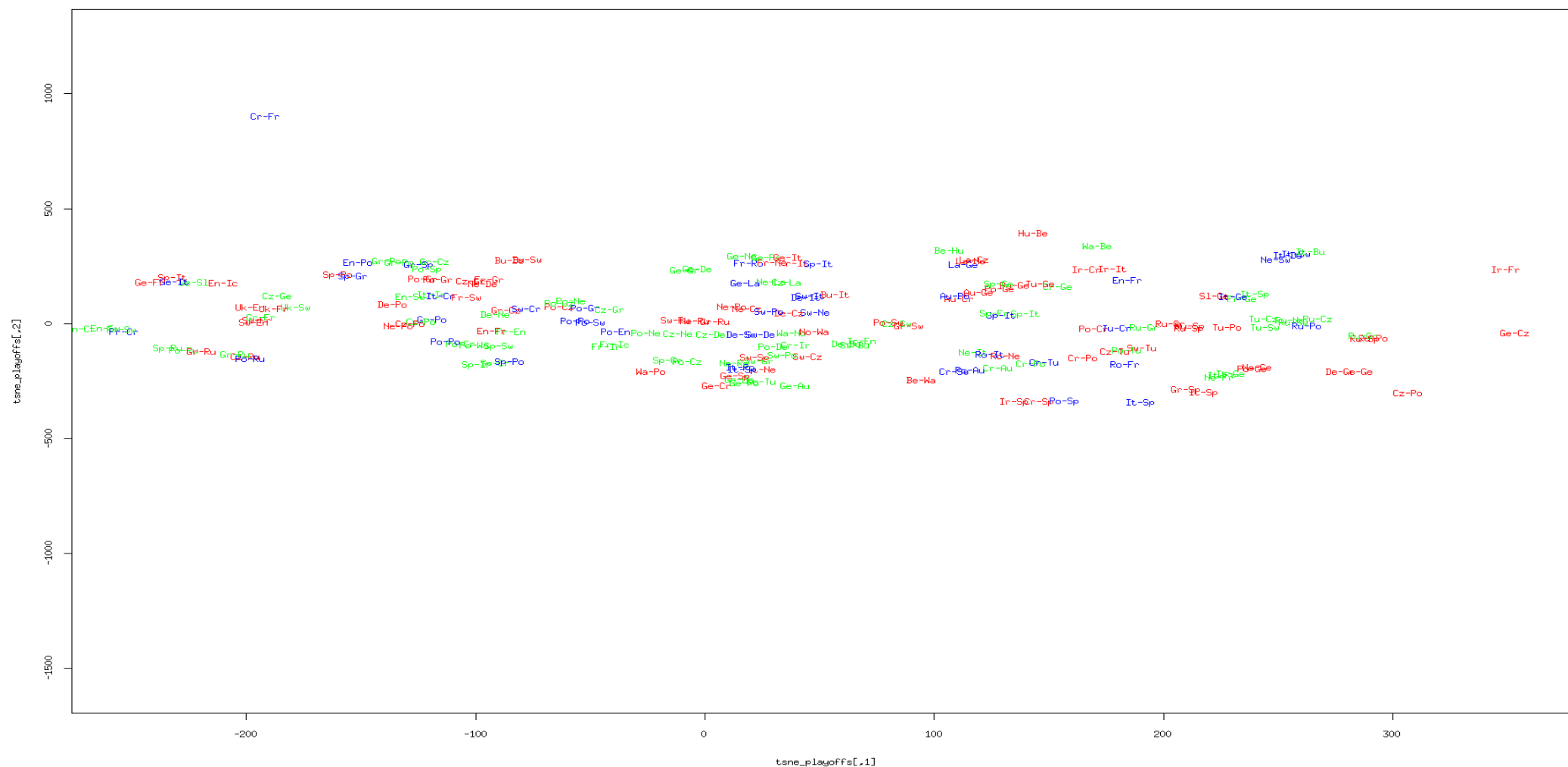
Evaluation and Model Selection



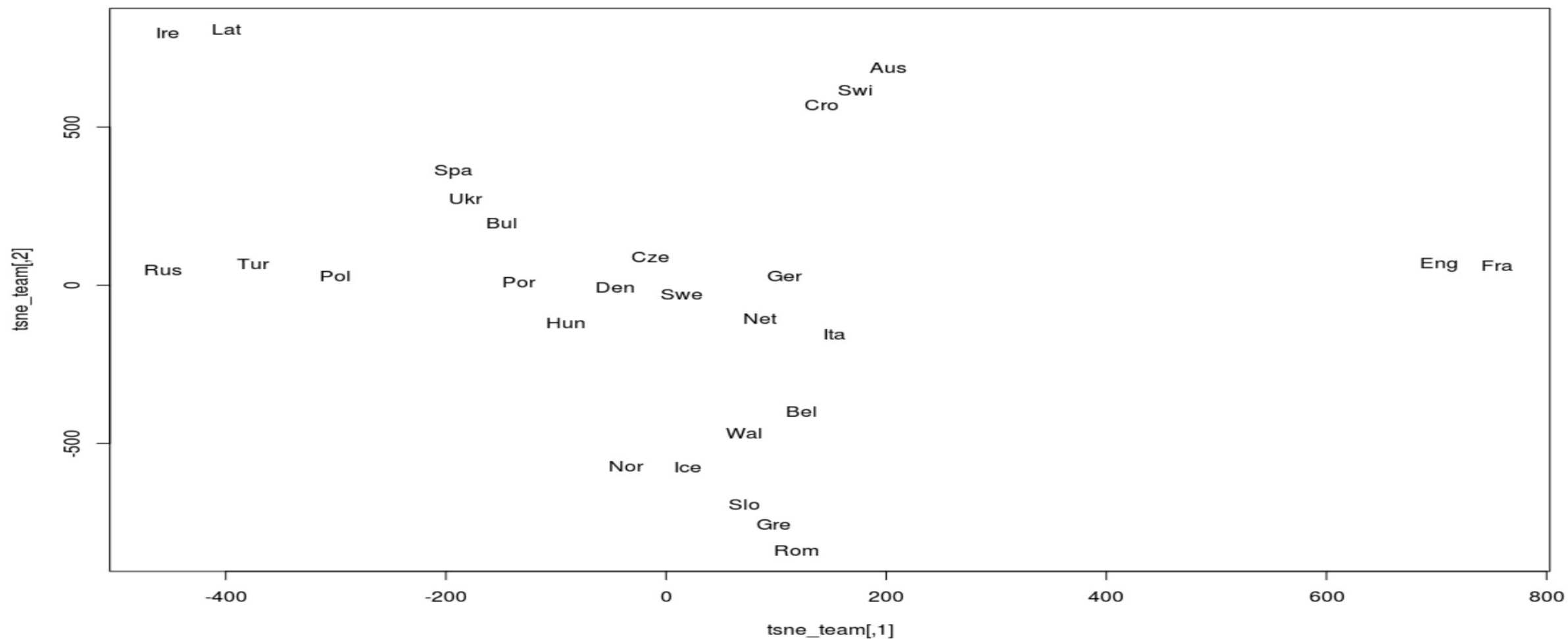
Evaluation and Model Selection



t-SNE on all games



t-SNE on team data only



Conclusions

- Probabilistic or not, models are similar in performance
- Modeling with Poisson, without additional features, throws away way to much information
- The lack of data decreases the confidence in our models' performance