



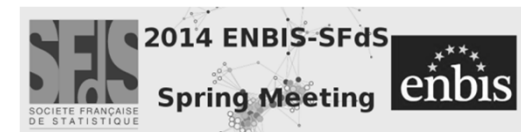
# DISCOVERING THE DRIVERS OF FOOTBALL MATCH OUTCOMES WITH DATA MINING

Maurizio Carpita, Marco Sandri,  
Anna Simonetto, Paola Zuccolotto



DMS StatLab – University of Brescia

Paris - 2014, April 11th





## **OUTLINE OF THE TALK**

- The case study: data and goals
- Football Mining: the Data Mining process
- Focus: variable selection
- Main outcomes
- Concluding remarks





# THE CASE STUDY: DATA AND GOALS

The top Italian professional football league «serie A»



20 teams  
38 matches/season



4 seasons



2008/2009 2009/2010 2010/2011 2011/2012

DATA



**PANINI DIGITAL**

*DigitalScout*



**1,300 recorded events** for each match (e.g. free kicks and shots, action type, fouls, crosses, recovered balls, goal assists, average time of ball possession, saves, goals on free kicks, etc.)

**482 variables**





# THE CASE STUDY: DATA AND GOALS

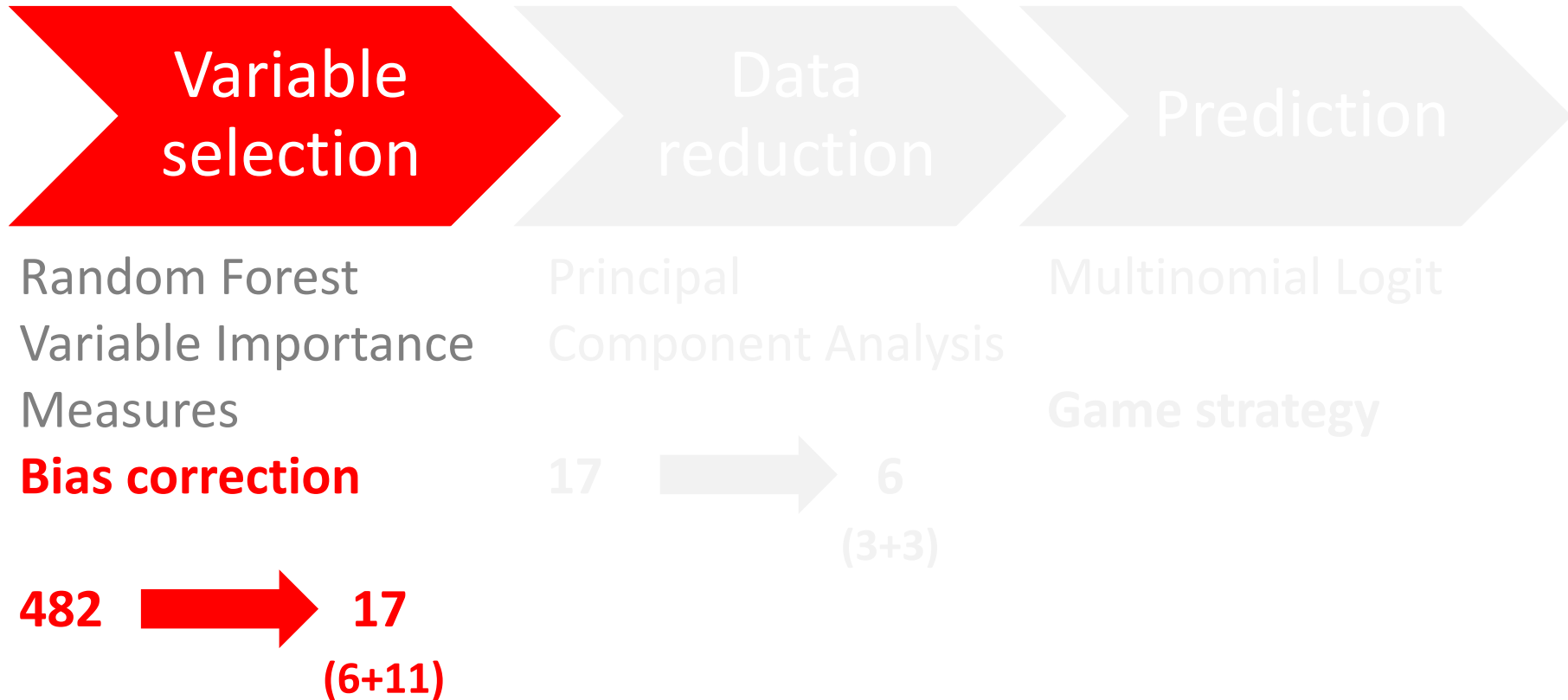


Identify the factors which mostly affect the probability of winning the match



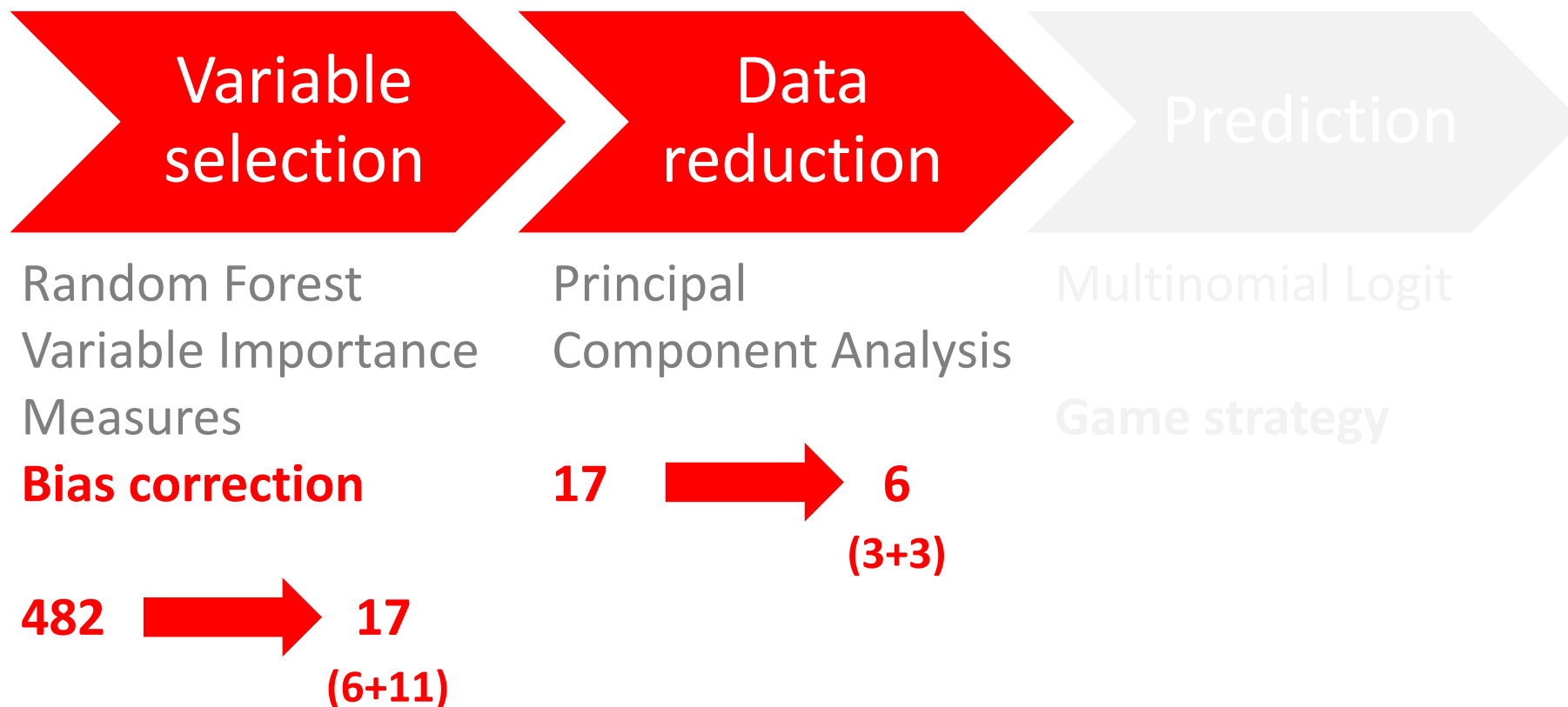


# FOOTBALL MINING: THE DATA MINING PROCESS



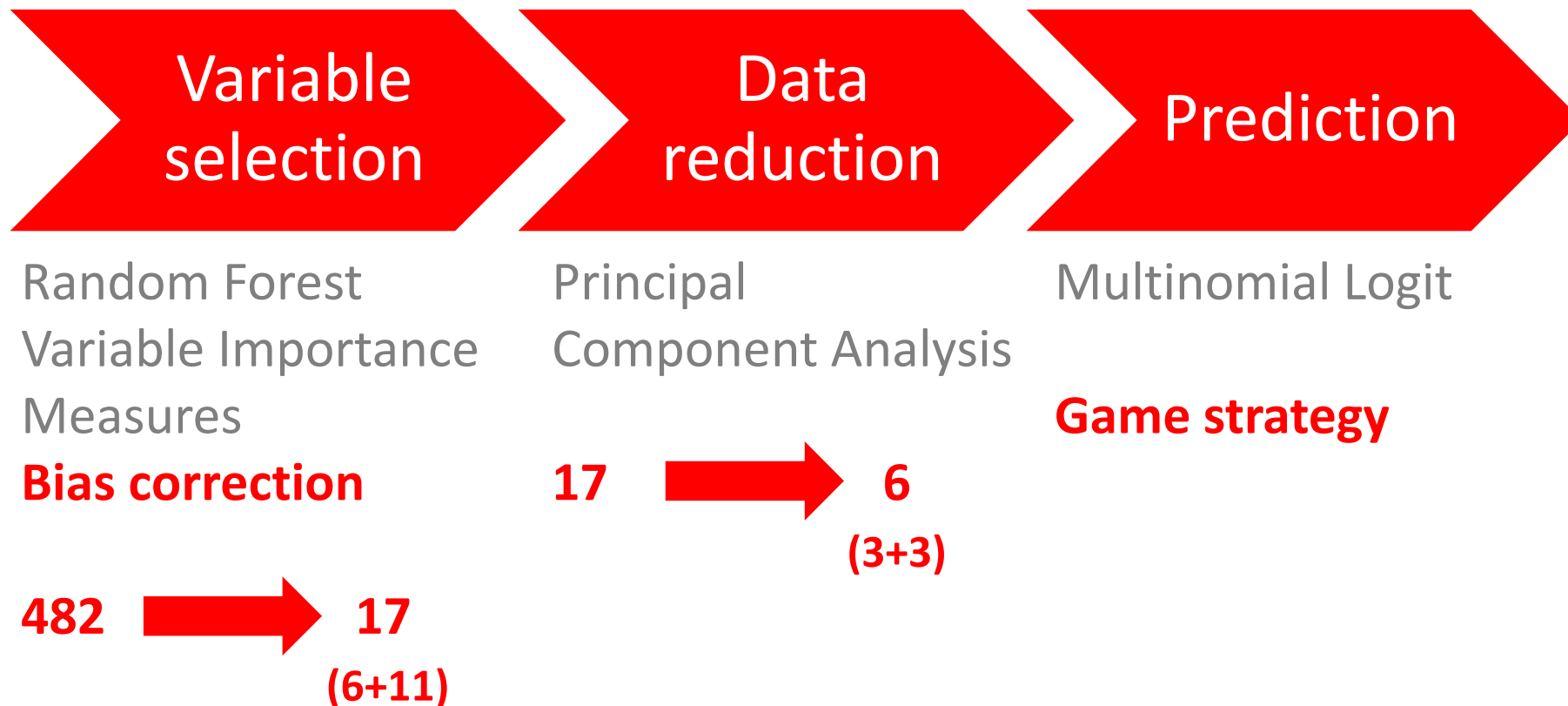


# FOOTBALL MINING: THE DATA MINING PROCESS



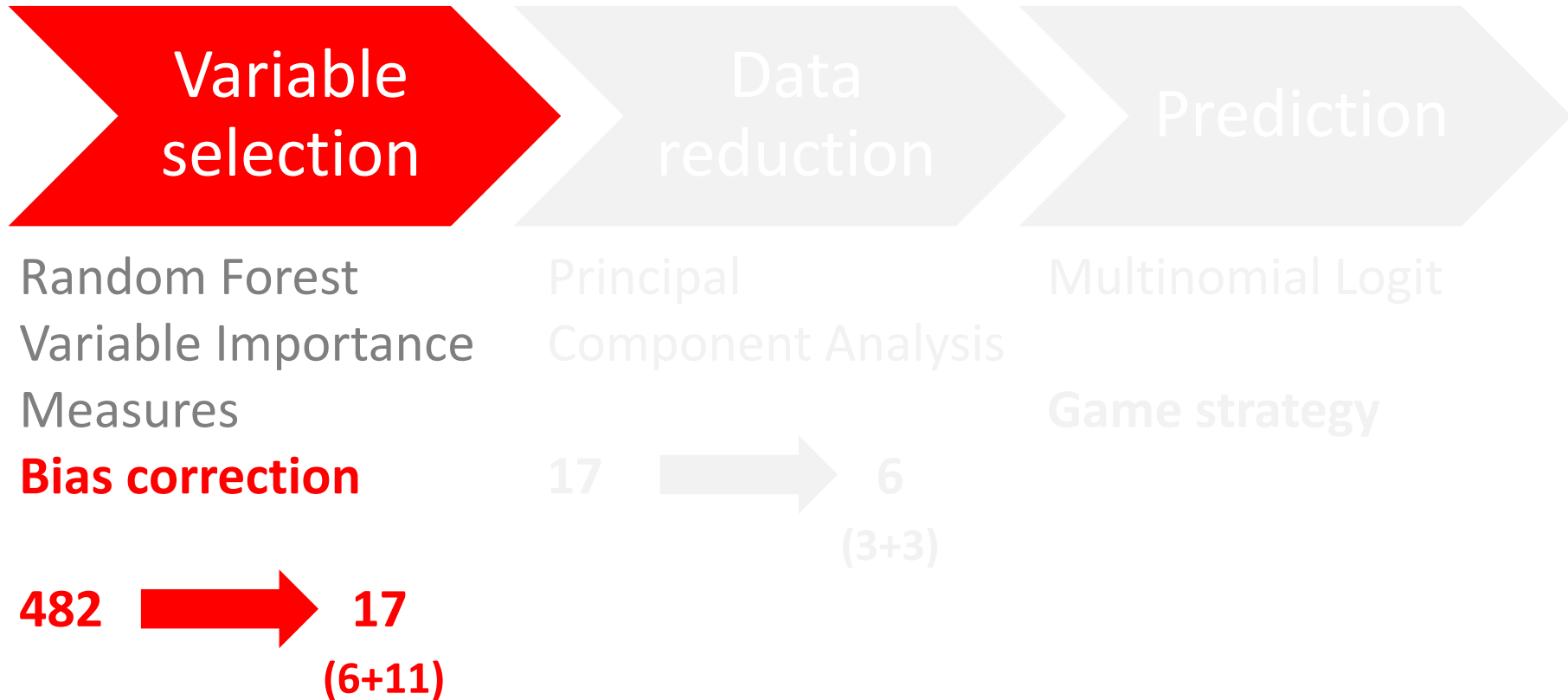


# FOOTBALL MINING: THE DATA MINING PROCESS





# FOCUS: VARIABLE SELECTION

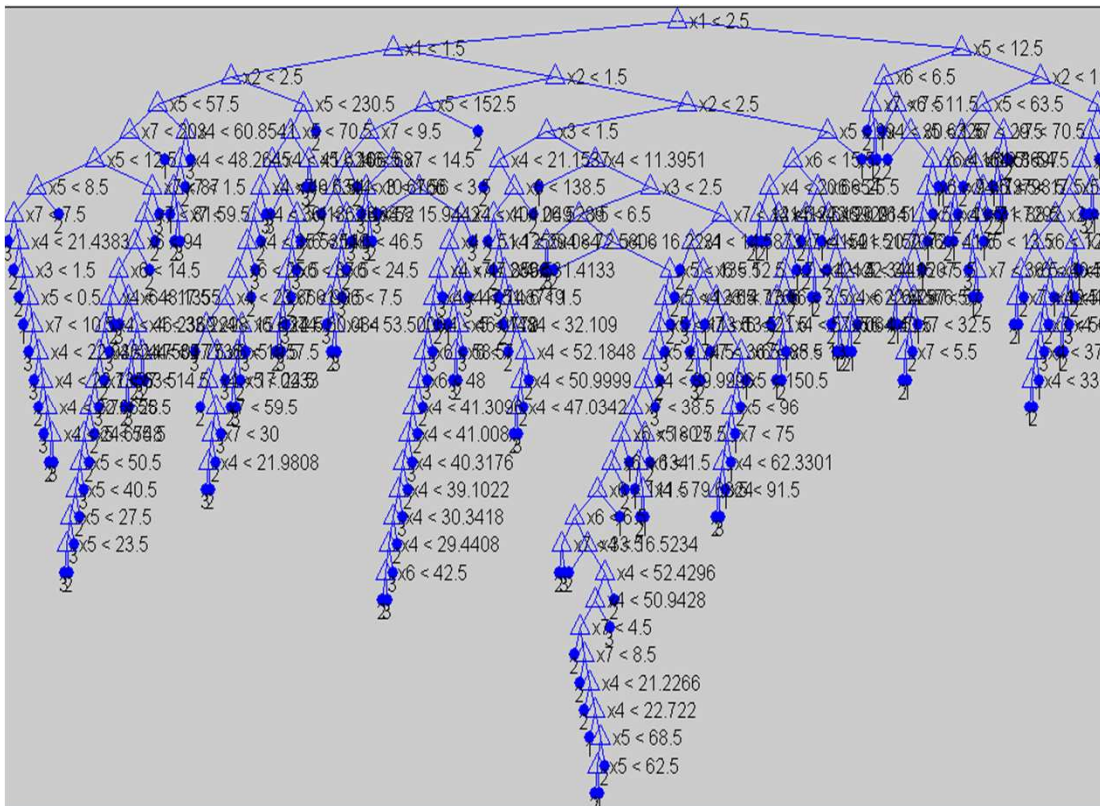






# FOCUS: VARIABLE SELECTION

## Random Forest Variable Importance Measures (TDNI VIMs)



single tree

$$\widehat{VI}_i(t) = \sum_{j \in J} \hat{d}_{ij} \cdot I_{ij}$$

RF

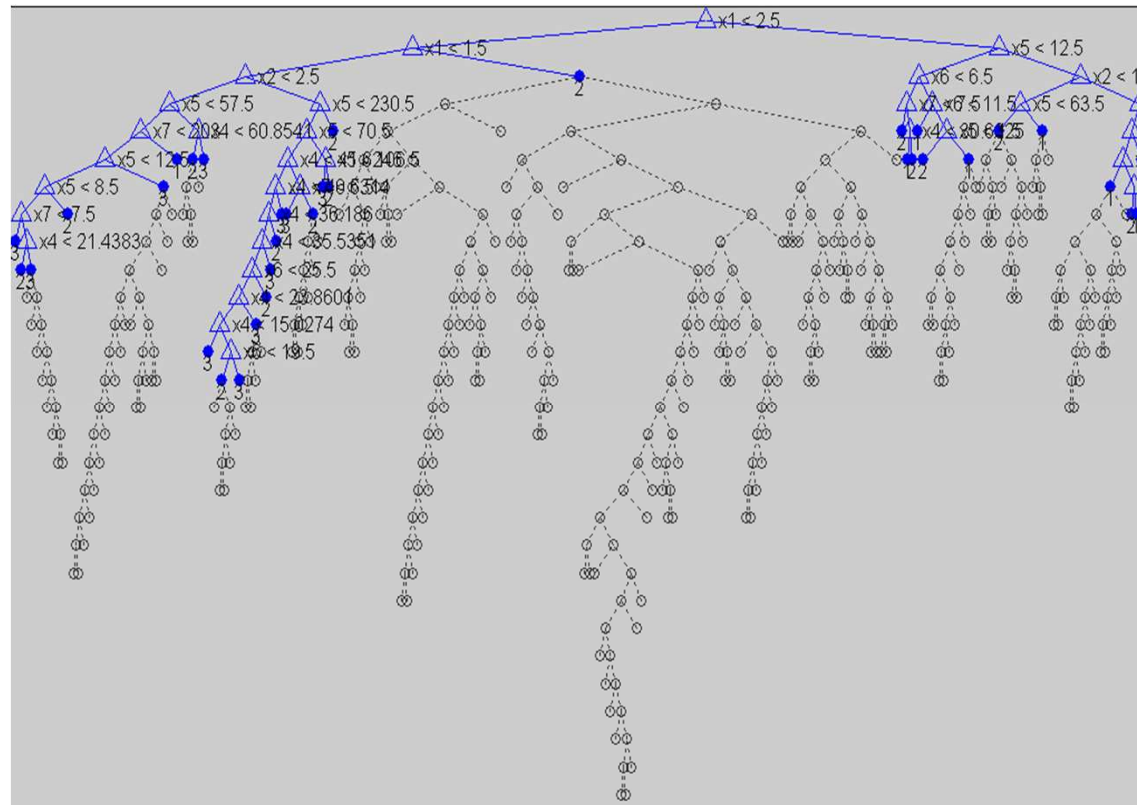
$$\widehat{VI}_i = \frac{1}{T} \sum_{t=1}^T \widehat{VI}_i(t)$$





# FOCUS: VARIABLE SELECTION

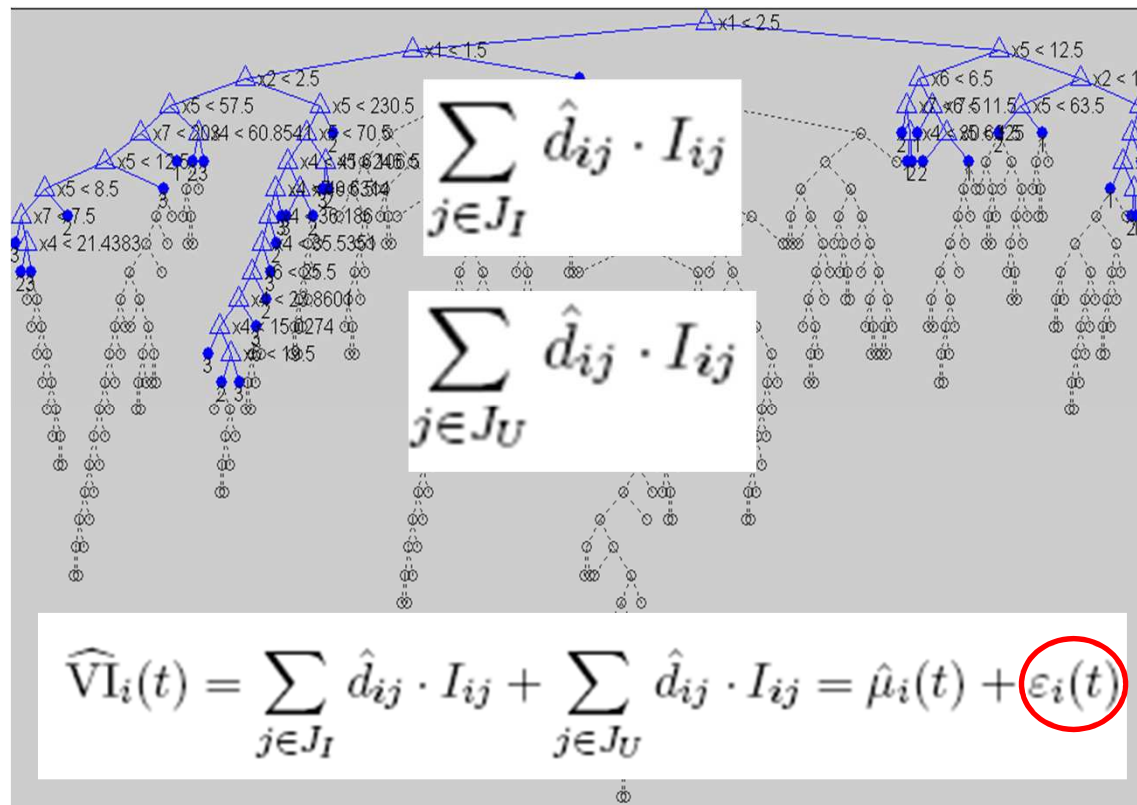
The source of **BIAS**: informative and uninformative splits  
(Sandri and Zuccolotto, 2008, 2010)





# FOCUS: VARIABLE SELECTION

The source of **BIAS**: informative and uninformative splits  
 (Sandri and Zuccolotto, 2008, 2010)



**BIAS**





# FOCUS: VARIABLE SELECTION

Correction of **BIAS**: PSEUDO-COVARIATES METHOD  
(Sandri and Zuccolotto, 2008, 2010)

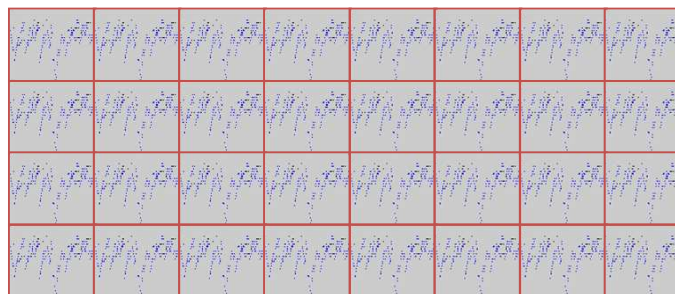
Original dataset

Soggetto	Y	X <sub>1</sub>	...	X <sub>s</sub>
1	y <sub>1</sub>	x <sub>11</sub>	...	x <sub>1s</sub>
2	y <sub>2</sub>	x <sub>21</sub>	...	x <sub>2s</sub>
...	...	...	...	...
N	y <sub>N</sub>	x <sub>N1</sub>	...	x <sub>Ns</sub>

Pseudo-covariates

Z <sub>1</sub>	...	Z <sub>s</sub>
z <sub>11</sub>	...	z <sub>1s</sub>
z <sub>21</sub>	...	z <sub>2s</sub>
...	...	...
z <sub>N1</sub>	...	z <sub>Ns</sub>

Soggetto	Y	X <sub>1</sub>	...	X <sub>s</sub>	Z <sub>1</sub>	...	Z <sub>s</sub>
1	y <sub>1</sub>	x <sub>11</sub>	...	x <sub>1s</sub>	z <sub>11</sub>	...	z <sub>1s</sub>
2	y <sub>2</sub>	x <sub>21</sub>	...	x <sub>2s</sub>	z <sub>21</sub>	...	z <sub>2s</sub>
...	...	...	...	...	...	...	...
N	y <sub>N</sub>	x <sub>N1</sub>	...	x <sub>Ns</sub>	z <sub>N1</sub>	...	z <sub>Ns</sub>



VIMs X  
VIMs Z

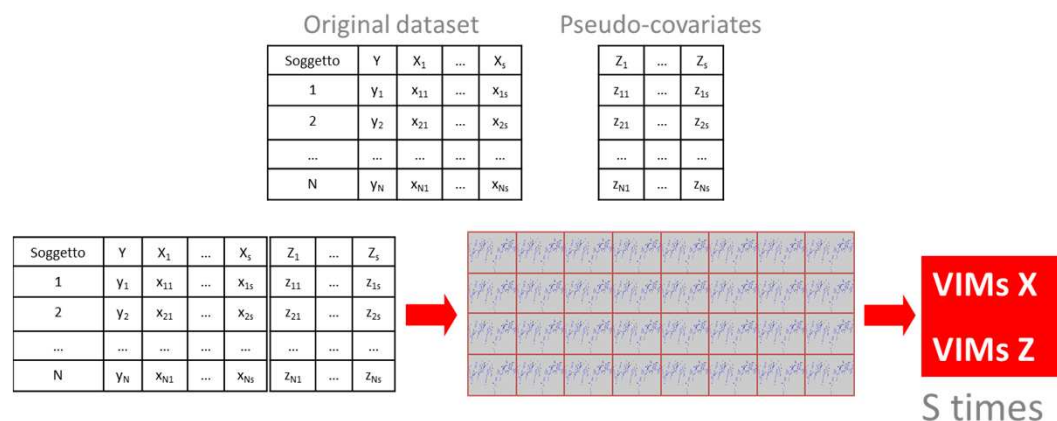
S times





# FOCUS: VARIABLE SELECTION

Correction of **BIAS**: PSEUDO-COVARIATES METHOD  
(Sandri and Zuccolotto, 2008, 2010)



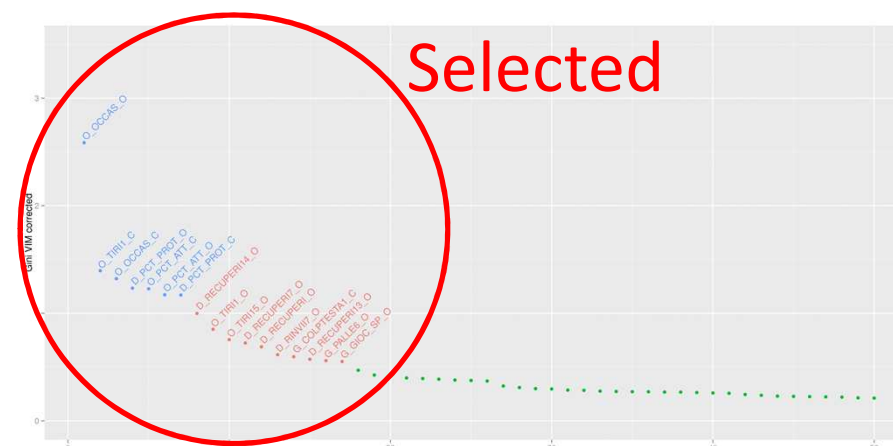
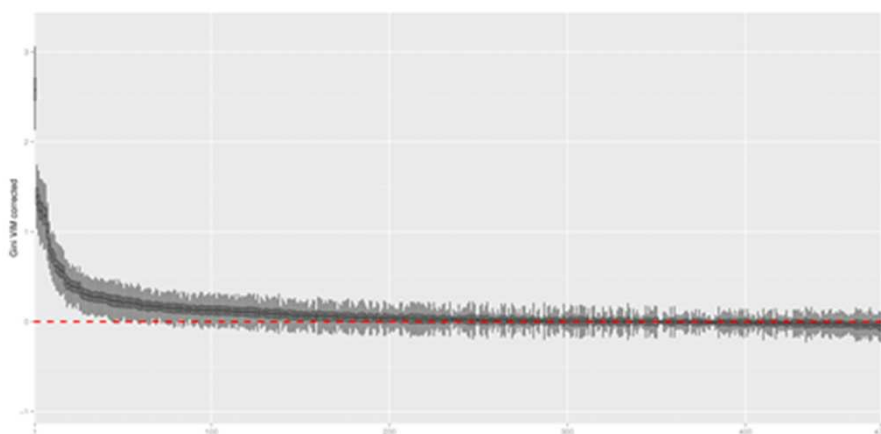
$$\overline{VI}_i = \frac{1}{S} \sum_{s=1}^S (\widehat{VI}_{X_i}^s - \widehat{VI}_{Z_i}^s)$$





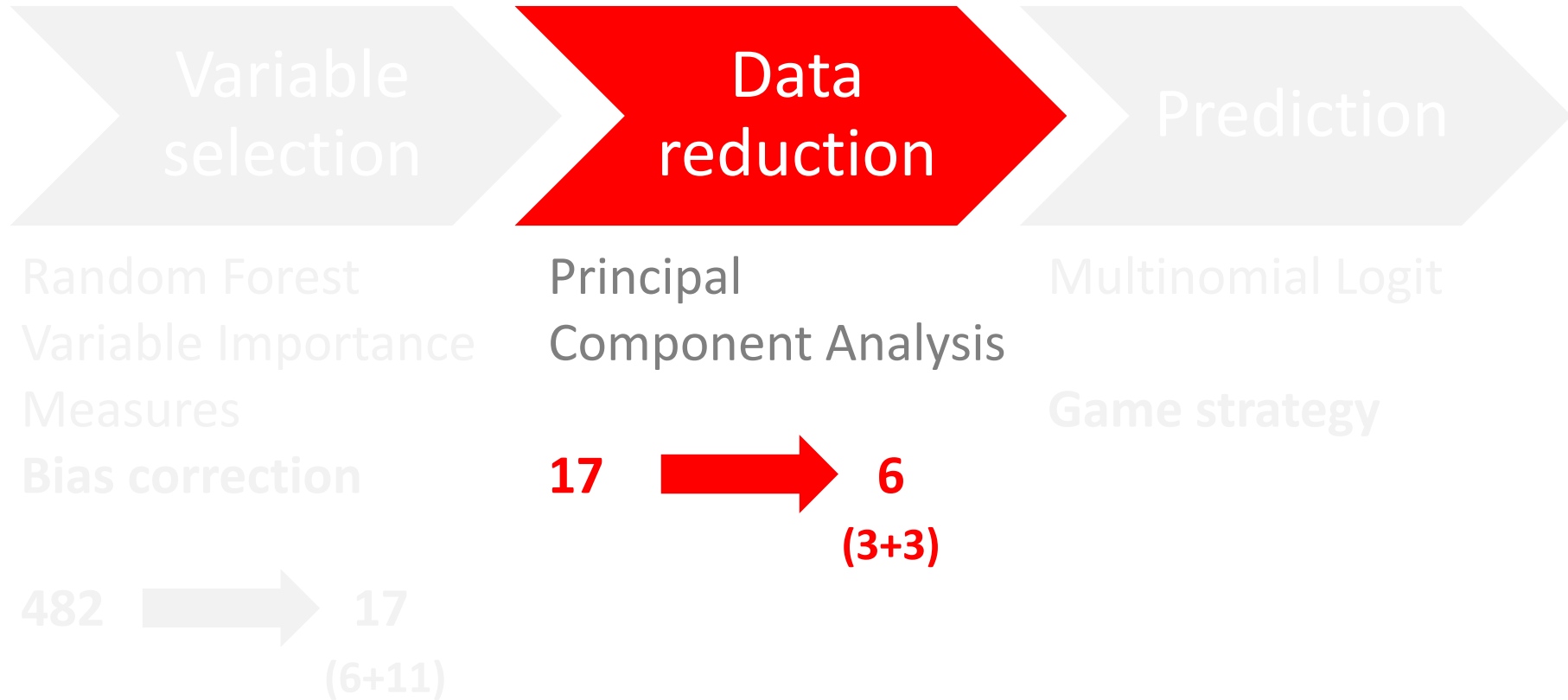
# FOCUS: VARIABLE SELECTION

Correction of **BIAS**: PSEUDO-COVARIATES METHOD  
(Sandri and Zuccolotto, 2008, 2010)





# MAIN OUTCOMES





# MAIN OUTCOMES

## HOME TEAM

- **shot.attack.home**

ability to create opportunities to make shots on goal

- **aerial.attack.home**

aerial abilities (crosses and heading) when the team is on the attack

- **defense.home**

general defense abilities







# MAIN OUTCOMES

## AWAY TEAM

- **midfield-defense.counterattack.away**

general defense abilities, long-range kicks and sudden counterattacks, with specific reference to actions in the midfield

- **shot.attack.away**

ability to create opportunities to make shots on goal

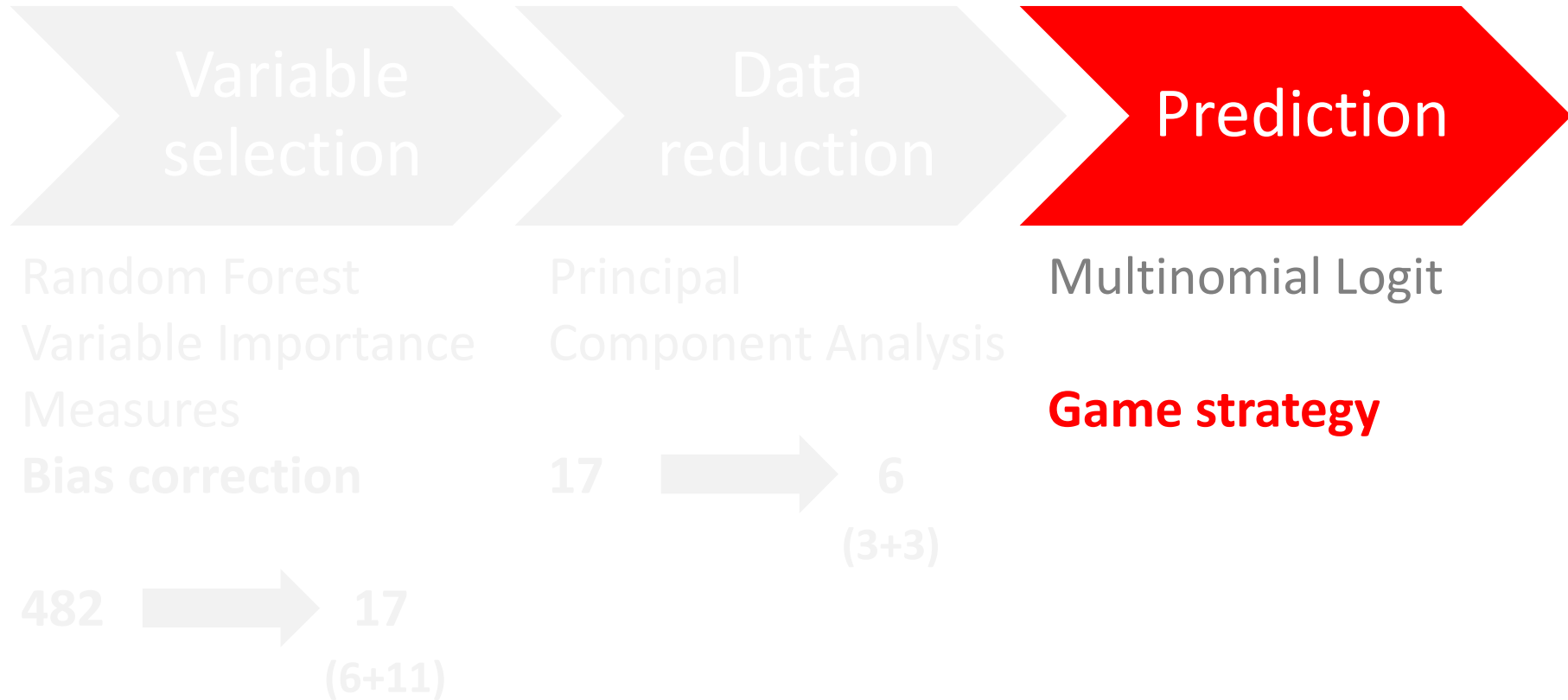
- **area-defense.away**

attitude to condense defense in the crucial penalty area








# MAIN OUTCOMES








# MAIN OUTCOMES

- shot.attack.home 
- aerial.attack.home 
- defense.home 



Identify the factors which mostly affect the probability of winning the match

- midfield-defense.counterattack.away 
- shot.attack.away 
- area-defense.away 





## CONCLUDING REMARKS

- Data Mining techniques offer interesting insights into sport strategies
- An effective variable selection technique is the starting point, as sport data are often big data
- The results of our analysis remain stable along the 4 examined seasons





# MAIN REFERENCES

1. **Agresti A.** (2003). Logit Models for Multinomial Responses. In: Categorical Data Analysis, 2nd Edition, John Wiley & Sons, Inc., Hoboken, NJ, USA.
2. **Albert J., Koning R.H.** (2008). Statistical Thinking in Sports, Chapman & Hall, Boca Raton.
3. **Breiman L.** (2001a). Random forests, Machine Learning, 45(1), 5-32.
4. **Carpita M., Sandri M., Simonetto A., Zuccolotto P.** (2014). Football Mining with R. In: Data Mining Applications with R (Edited by Y. Zhao, Y. Cen), Chapter 14. Elsevier.
5. **Hastie T., Tibshirani R., and Friedman J.H.** (2001). The elements of statistical learning: data mining, inference, and prediction. Springer, New York.
6. **Jolliffe I.T.** (2002). Principal Component Analysis. Springer Verlag, New York.
7. **Sandri M., Zuccolotto P.** (2008). A Bias Correction Algorithm for the Gini Variable Importance Measure in Classification Trees. Journal of Computational and Graphical Statistics, 17(3), 611-628.
8. **Sandri M., Zuccolotto P.** (2010). Analysis and correction of bias in Total Decrease in Node Impurity measures for tree-based algorithms. Statistics and Computing, 20, 393-407.

