

# THE TWO CULTURES: A DISCUSSION

---

Katrin Newger

Supervisor: Christoph Jansen M.Sc. and Dipl.-Math. Georg Schollmeyer

June 27, 2015

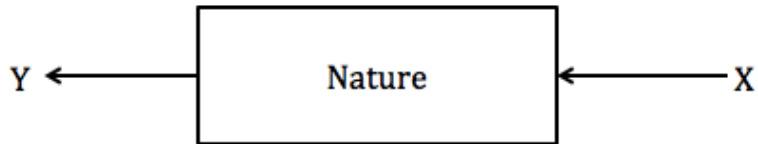
Department of Statistics, LMU Munich

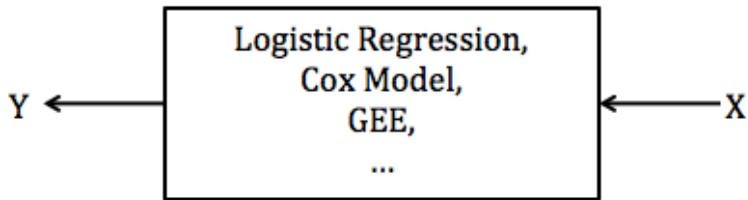
# TABLE OF CONTENTS

1. The Two Cultures
2. Breiman's Argument
3. Discussion
4. Personal Impressions and Conclusion

## THE TWO CULTURES

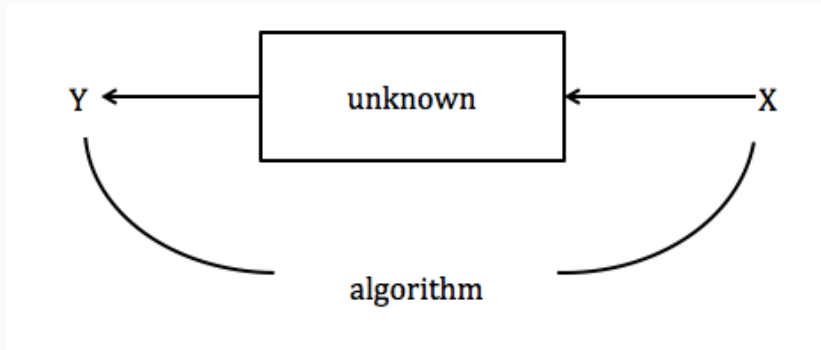
---





## Assumptions:

- Stochastic model
- Distribution of residuals
- Further model specific assumptions



**Goal:**

Function  $f(x)$  that minimizes loss  $L(Y, f(x))$

## Methods:

- Support vector machines
- Random forests
- Artificial neural networks
- ...

## BREIMAN'S ARGUMENT

---



- Critical model assumptions
- Conclusions about model, not about nature
- Wrong model → wrong conclusions about nature
- Algorithmic models only assume iid. variables

*“A few decades ago (...) the belief in data models was such that even simple precautions such as residual analysis or goodness-of-fit tests were not used” (Breiman 2001, p. 199)*

- Necessity of checking the model's fit
- Discussion of the fit is superficial
- Most popular: goodness-of-fit tests, residual analysis

## Goodness-of-Fit Tests

- Not useful if direction of alternative not precisely defined
- Extreme discrepancy to the data is needed

## Residual Analysis

- For more than four dimensions: interactions between variables  
→ manipulation of residual plots

Algorithmic modeling: cross-validation is standard procedure

- Different models → different assumptions  
→ different conclusions
- Neither model is able to trump
- Further problem: variable selection based on model
- **Algorithmic modeling: only iid. assumption**

- Common assumption:  $n \rightarrow \infty$  never fulfilled
- Testing on 5% level is arbitrary  
(*“suspect way to arrive at conclusions”, Breiman 2001, p. 203*)
- Algorithmic modeling: no inference

- Originally:  $n \gg p$   
↔ nowadays:  $p \gg n$
- Data models become too complex
- Common procedure: reducing dimensionality (e.g. principal component analysis) → loss of information
- **Algorithmic modeling: the more variables the more information**

- Prediction is more important than interpretation—always
- If prediction is bad, how can interpretation be good?
- Breiman's experience: algorithmic models are best predictors



- Everyone's choice which model is best

*“The best solution could be an algorithmic model, or maybe a data model, or maybe a combination” (Breiman 2001, p. 206)*

- Openness for new methods

## DISCUSSION

---

*“[The Bias] has to be lurking somewhere inside the theory” (Brad Efron, in Breiman 2001, p. 219)*

- In algorithmic modeling, small variance at cost of bias?
- Breiman avoids answer

- Does not concern prediction
- Just as well in algorithmic models
- Main difference between models: distribution
- Breiman manipulates reader

- Why not use known information (e.g. distribution)?
- Critical iid. assumption in data models and algorithmic models
- Alternatives if iid. assumption is violated?

- Rivaling abilities of models
- Often interpretation required
- Prediction sometimes indirectly related to data

*“The whole point of science is to open up black boxes, understand their insides, and build better boxes for the purposes of mankind” (Brad Efron, in Breiman 2001, p. 219)*

## PERSONAL IMPRESSIONS AND CON- CLUSION

---



Leo Breiman

**Statistical Modeling: The Two Cultures.**

*Statistical Science* 16(3), 2001: 199–231.



T. Hastie, R. Tibshirani and J. Friedman

***The Elements of Statistical Learning. Data Mining, Inference and Prediction.***

Heidelberg: Springer, 2009.



## QUESTIONS AND DISCUSSION