# The two cultures of statistical modeling after Leo Breiman - 2001

Thomas Pietraho

# Sources

- Leo Breiman, *Statistical Modeling: The Two Cultures*, Statistical Science 2001. with comments by D. R. Cox, Brad Efron, Bruce Hoadley, Emanuel Parzen, and rejoinder by Breiman.
- Jelena Bradic, Yinchu Zhu, *Comments on Leo Breiman's paper: "Statistical Modeling: The Two Cultures"*, Observational Studies, 2021.
- Tyler H. McCormick, *The "given data" paradigm undermines both cultures*, Observational Studies, 2021.
- Andrew C. Miller, Nicholas J. Foti, Emily B. Fox, *Breiman's Two Cultures: You Don't Have to Choose Sides*, Observational Studies, 2021.

Nature functions to associate the input variables $\vec{x}$ with response variables $\vec{y}$ :



### Goals:

- *Prediction*: predict responses from future input variables, and
- *Information*: extract some information about how nature is associating response to input.

## Data modeling culture

Start by assuming a data model for the black box:

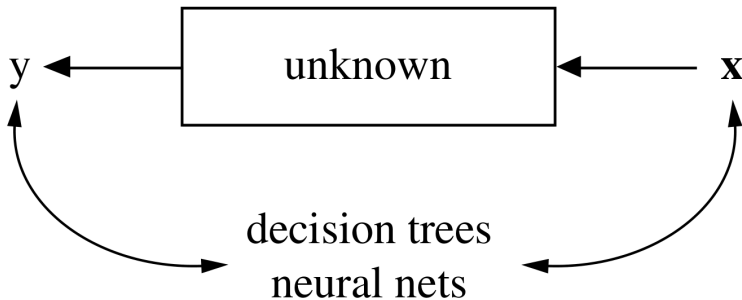*Example:* response variables = *f*(input variable, random noise, parameters)



*Workflow:*

- choose one of a finite number of models for data,
- estimate values of parameters,
- use model for prediction.

*Validation:* yes-no using goodness-of-fit tests.

## Algorithmic modelers

Assume the black box is too complex and unknown to model.



*Workflow:*

- choose a rich class of *surrogate functions*; e.g. universal approximators,
- find *f* within this class so that $f(\vec{x}) \approx \vec{y}$,
- use *f* for prediction.

*Validation:* Predictive accuracy.

## Data models: example

One of the oldest data science fields is cryptography.

| | |
|---|---|
| *input* plaintext: | **mathisoftenfun** |
| *output* ciphertext: | **xadsfowekdlsdf** |

*Workflow:*

- choose one of a finite number of models, i.e. encryption schemes,
- estimate values of parameters,
- reverse model for decryption.

*Validation:* goodness-of-fit is assessed by the plaintext recovered.

One of the oldest data science fields is cryptography.

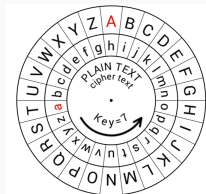| | |
|---|---|
| *input* plaintext: | **mathisoftenfun** |
| *output* ciphertext: | **xadsfowekdlsdf** |

*Workflow:*

- choose one of a finite number of models, i.e. encryption schemes,
- estimate values of parameters,
- reverse model for decryption.

*Validation:* goodness-of-fit is assessed by the plaintext recovered.

### Example

*Caesar cipher:* direct substitution of one character for another in plaintext yields ciphertext.

*Parameter*: recover permutation in $S_{26}$

## Data models: example

One of the oldest data science fields is cryptography.

| | |
|---|---|
| *input* plaintext: | **mathisoftenfun** |
| *output* ciphertext: | **xadsfowekdlsdf** |

*Workflow:*

- choose one of a finite number of models, i.e. encryption schemes,
- estimate values of parameters,
- reverse model for decryption.

*Validation:* goodness-of-fit is assessed by the plaintext recovered.

### Example

*Enigma machine:* electromechanical rotor mechanism that scrambles the 26 letters of the alphabet.

*Parameters:* rotors and their positions

## Data models: example

One of the oldest data science fields is cryptography.

*input* plaintext:     **mathisoftenfun**

*output* ciphertext:    **xadsfowekdlsdf**
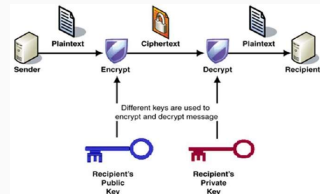
*Workflow:*

- choose one of a finite number of models, i.e. encryption schemes,
- estimate values of parameters,
- reverse model for decryption.

*Validation:* goodness-of-fit is assessed by the plaintext recovered.
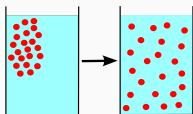
### Example

*RSA:* prime-based public-key cryptosystem that is widely used for secure data transmission

*Parameters:* two large primes

*Diffusion:* estimate concentration of particles at time *t* at position *x*.
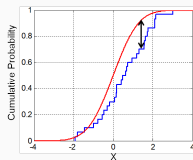


JrPol: wikipedia.com

*Workflow:*

- model from physics: $f(x) \propto e^{-\frac{(x-\mu)^2}{t}}$
- estimate $\mu$ and $t$ from measurements,
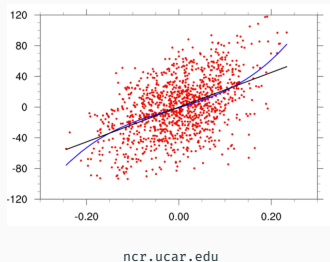- use $f$ for prediction.

*Validation:* Kolmogorov–Smirnov test goodness-of-fit test, for example



BScan: wikipedia.com

## Data models: example

*Regression:* $y = f(\vec{x}) = b_0 + \sum_i b_i x_i + \epsilon$ with $\epsilon \sim N(0, \sigma)$



ncr.ucar.edu

*Workflow:*

- assume a lineal regression model,
- estimate $b_i$ and $\sigma$
- use $f$ for prediction with error bounds for prediction

*Validation:* $R^2$ goodness-of-fit test, value in $[0, 1]$

$$R^2 = 1 - \frac{SS_{RES}}{SS_{TOT}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \overline{y})^2}$$

analyticsvidhya.com

There are tremendous advantages to data models:

- *Interpretability:* if a model is a good emulation of the data-generating mechanism, prediction and interpretation will be valuable. A scientist will know whether and why a scientific phenomenon has been observed. For example, regression models in econometrics:

$$y = f(\vec{x}) = b_0 + \sum_i b_i x_i + \epsilon$$

  Size of coefficients suggest interventions.

- *Derivative results*: A precise theoretical formulation of the models allows control of prediction error and derivative conclusions. For example, a model for equity prices allows one to predict prices for derivatives such as call options.

## Data models: summary

Breiman points out some flaws in data models.

The previous examples illustrate a continuum of models. From ones that clearly replicate the data-generating process to ones that are used only for the lack of anything better.

> *This enterprise has at its heart the belief that a statistician, by imagination and by looking at the data, can invent a reasonably good parametric class of models for a complex mechanism devised by nature. –Breiman*

If a model is a poor emulation of nature, its conclusions may be wrong. Three changes in perception:

- *Rashomon effect:* the multiplicity of good models,
- *Occam:* conflict between simplicity and accuracy, and
- *Bellman:* high dimensional data is a curse and a blessing.

*Rashomon* A 1950 Akira Kurosawa film known for a plot device that "involves various characters providing subjective, alternative and contradictory versions of the same incident."



PD-Japan-organization

Suppose two data scientists, each one using a different data model, fit different models to the same data set. Suppose that both pass goodness-of-fit tests.

**Question:** What have we leaned about the mechanism generating the data, i.e. *Nature*?

When data has more than a small number of dimensions, there *will* be a large number of models that pass goodness-of-fit tests.

## Rashomon effect and dimensionality

*Example (Breiman)*: Suppose that we have a data set in 30 variables and want to find the best 5 for linear regression. There are $\binom{30}{5} \approx 140,000$ choices of five-variable subsets. On a given set, Breiman found three that all passed a goodness-of-fit test and had RSS within 1% of each other:

$$f(\vec{x}) = y \approx \quad 2.1 + 3.8x_3 - 0.6x_8 + 83.2x_{12} - 2.1x_{17} + 3.2x_{27}$$
$$f(\vec{x}) = y \approx -8.9 + 4.6x_5 - 0.01x_6 + 12.0x_{15} - 17.5x_{21} + 0.2x_{22}$$
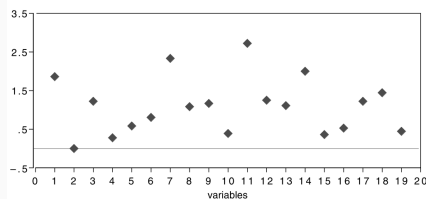$$f(\vec{x}) = y \approx -76.7 + 9.3x_2 + 22.0x_7 - 13.2x_8 + 3.4x_{11} + 7.2x_{28}$$

*Questions:*

- Which one should be use? Each one tells a different story about nature.
- If this is economic data and we would like to propose an intervention by directly affecting one of the variables, which one should we affect?

This is a common phenomenon in high-dimensional data. It is ubiquitous regardless of model or algorithm used.

## Rashomon effect and simplicity

*Example (Efron - Diaconis)*: Survival of 155 hepatitis patients, 19 variables. Which ones are important? Below are the coefficients of logistic regression:



Breiman

Variables 1, 7, 11, 14 seem most important. Their experiment: 500 *bootstrap* samples and estimate important coefficients.

**Conclusion:** Of the four variables originally selected, not one was selected in more than 60% of the samples. The variables identified cannot be taken too seriously.

*But:* Prediction error rate was 17.4% using all variables, and around 20.0% if only four were used.

*Question:* Should one give up accuracy of the model for (perhaps dubious) interpretability? On more modern data, this effect is more pronounced.

*Question:* Which variables are important in your data?

*Question:* Which variables are important in your data?

*Answers:*

- *Regression*: absolute value of regression coefficient, especially if same variable appears consistently when regression is performed on bootstrapped samples.

*Question:* Which variables are important in your data?

*Answers:*

- *Regression*: absolute value of regression coefficient, especially if same variable appears consistently when regression is performed on bootstrapped samples.
- *Neural networks and random forests*:
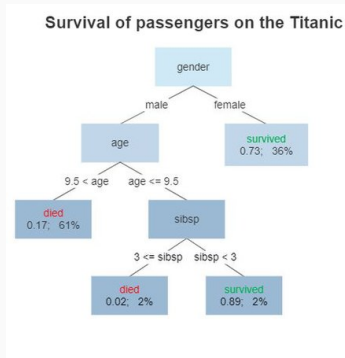  - Use $\frac{\partial N}{\partial x_i}$ for neural networks.

*Question:* Which variables are important in your data?

*Answers:*

- *Regression*: absolute value of regression coefficient, especially if same variable appears consistently when regression is performed on bootstrapped samples.
- *Neural networks and random forests*:
  - Use $\frac{\partial N}{\partial x_i}$ for neural networks.
  - No derivatives for random forests: permute *only* the values of the *i*th variable in the data. The variable with the most profound loss of accuracy is the most important.

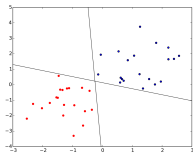Random forests are families of decision trees:



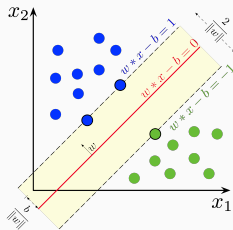Decision tree
Gilgoldm wikipedia.com

Each tree is found using a different random sample of the data. Output of random forest is a statistic of the individual trees.
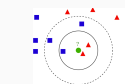
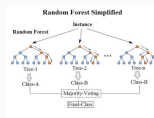# Classical algorithmic methods



Perceptron

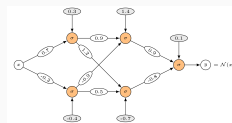Qwertyus wikipedia.com



kNN classifier

Antti Ajanki wikipedia.com



Random forest

Venkata Jagannath wikipedia.com



Support vector machine

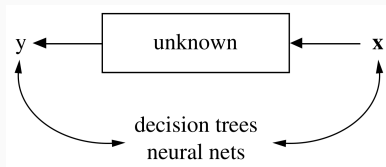Lahrmam wikipedia.com



Neural network

## Algorithmic modelers: again

Assume the black box is too complex and unknown to model.



In this way, one does not presume to understand what nature does, but hopes to mimic its actions.

*Workflow:*

- choose a rich class of *surrogate functions*; e.g. universal approximators,
- find $f$ within this class so that $f(\vec{x}) \approx \vec{y}$,
- use $f$ for prediction.

*Validation:* Predictive accuracy on hold-out data, i.e. *generalization error*

At the time of Breiman's writing, the split between data modelers and algorithmic modelers was estimated to be 98%/2%. It is much different today.

## Breiman's call to action

*Pro:* Data modelers use simple models with well-understood theoretical properties that provide the allure of interpretability

*Con:* Such beneficial properties have limited value if the underlying models are ill-suited to describe the data being analyzed. Practitioners often "clumsily try, often strenuously, to find the best of most appropriate model for a particular problem."

> *If our goal is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools. Nowhere is it written on a stone tablet what kind of model should be used to solve problems involving data; the goals in statistics are [simply] to use data to predict and to get information about the underlying data mechanism.*
> *–Breiman*

Breiman was as singular voice at the time. He was proved to be right. In many applications empirically-built models built solely to improve prediction accuracy have replaced more traditional data-models. They are sufficiently flexible to capture noisy, complex processes, and led to fundamental advances the fields of:

- vision and sound and
- natural language processing

## But there are problems

Algorithmic models can be uninterpretable, "making it difficult to explain, audit, or critique the predictions from a complex neural network or random forest."

- *Health applications:* a common goal is to predict the occurrence of rare events. But they are no well-represented in training data. How can we make accurate predictions in rare and critical situations?
- *Public policy:* this is a problem similar to the one above. For some populations, data is often scarce and of poor quality. How can we made good decisions based on algorithmic models in this environment?
- *Risk management:* without a strong theoretical framework, how do we assess the risk inherent in our model's predictions? Can we use financial data and assess the risk of ruin in an algorithmic trading strategy?

*Bradic and Zhu:* Prediction accuracy is no longer satisfactory and should not be the only measure of success. Some more desiderata:

- *Stability:* stability measures how a data result changes when the data and/or model are perturbed. Adversarial examples show that some neural network models are very unstable.
- *Reproducibility:* experiments must be sufficiently documented so other experts can reproduce them. This is a problem in many scientific fields.
- *Inference:* How do we test a hypothesis such as the effectiveness of a treatment. "The impact of machine-learning methods on [inferential] tasks has yet to be unlocked."

**Setting:** input vectors $\vec{x}_i \in \mathbb{R}^n$, model prediction $D_i \in \{0, 1\}$, and outcome $Y_i$.

**Question:** What was the average effect of the treatment? Did the intervention improve outcomes?

Bradic and Zhu (2021): Used a random forest to predict whether a treatment should be used and tried to estimate:

$$\mathbb{E}(Y_i(1) - Y_i(0))$$

or in English, the average difference in outcome for the population which received the treatment versus the population that did not. They were not successful in estimating the true effect. The theory underlying the estimation of average treatment effect is not robust to cover random forest predictions. And certainly not ones made by neural networks.

# Toward machine-learning systems

*Successful technological fields have a moment when they become pervasive, important, and noticed. They are deployed into the world and, inevitably, something goes wrong. A badly designed interface leads to an aircraft disaster. A buggy controller delivers a lethal dose of radiation to a cancer patient. The field must then choose to mature and take responsibility for avoiding the harms associated with what it is producing. Machine learning has reached this moment. [T]he community needs to adopt systematic approaches for creating robust artifacts that contribute to larger systems that impact the real human world. [We need to move] beyond narrow machine-learning algorithms to complete machine-learning systems. –Charles Isbell*