# Métodos cuantitativos y estadísticos en las ciencias del lenguaje:
## panorama, predicciones y recomendaciones

**Damián E. Blasi**

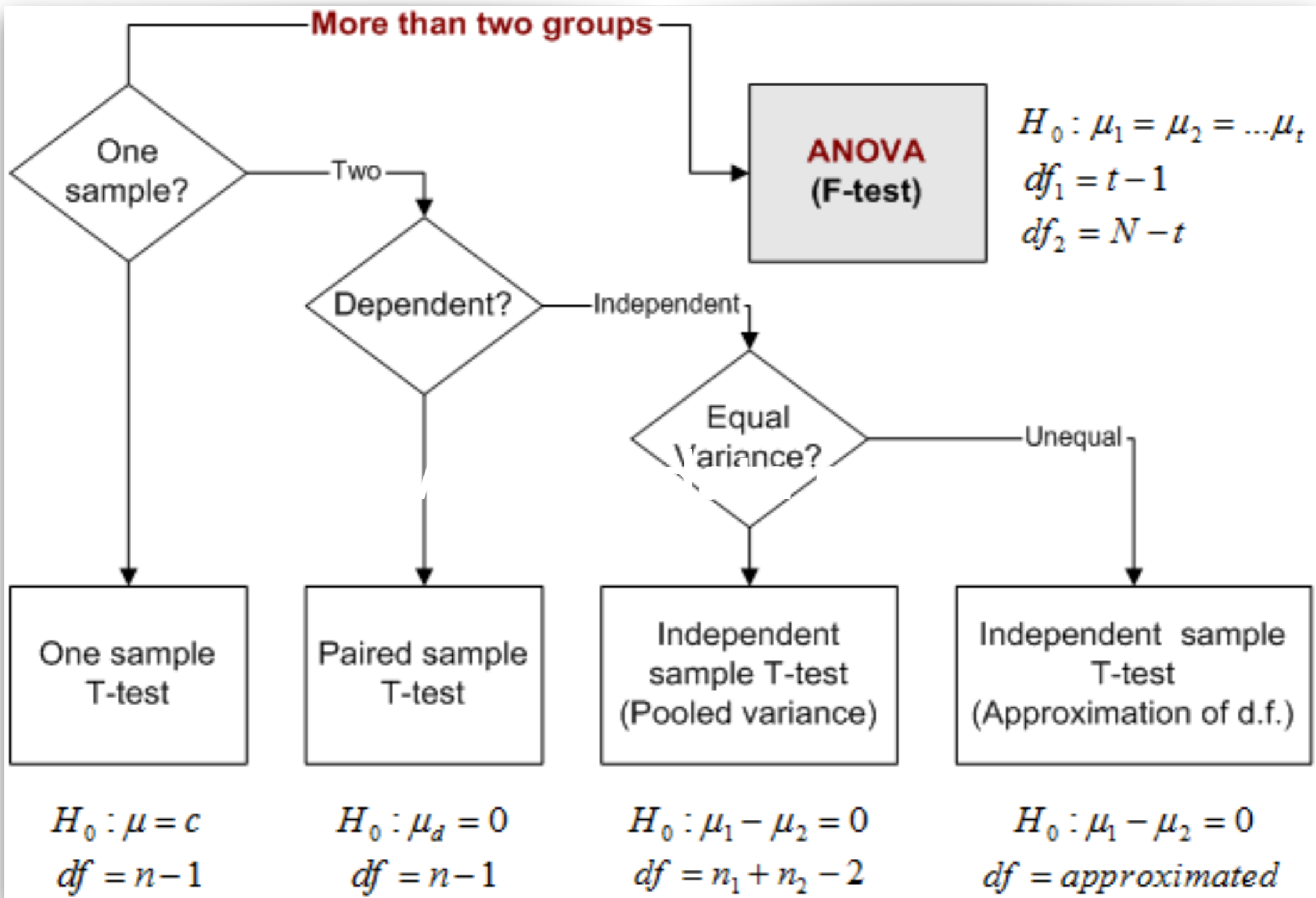Department of Comparative Linguistics,
**University of Zürich**

Department of Linguistic and Cultural Evolution,
**Max Planck Institute for the Science of Human History**

Human Relation Area Files,
**Yale University**

# Cancer and Smoking

THE curious associations with lung cancer found in relation to smoking habits do not, in the minds of some of us, lend themselves easily to the simple conclusion that the products of combustion reaching the surface of the bronchus induce, though after a long interval, the development of a cancer. If, for example, it were possible to infer that smoking cigarettes is a cause of this disease, it would equally be possible to infer on exactly similar grounds that inhaling cigarette smoke was a practice of considerable prophylactic value in preventing the disease, for the practice of inhaling is rarer among patients with cancer of the lung than with others.
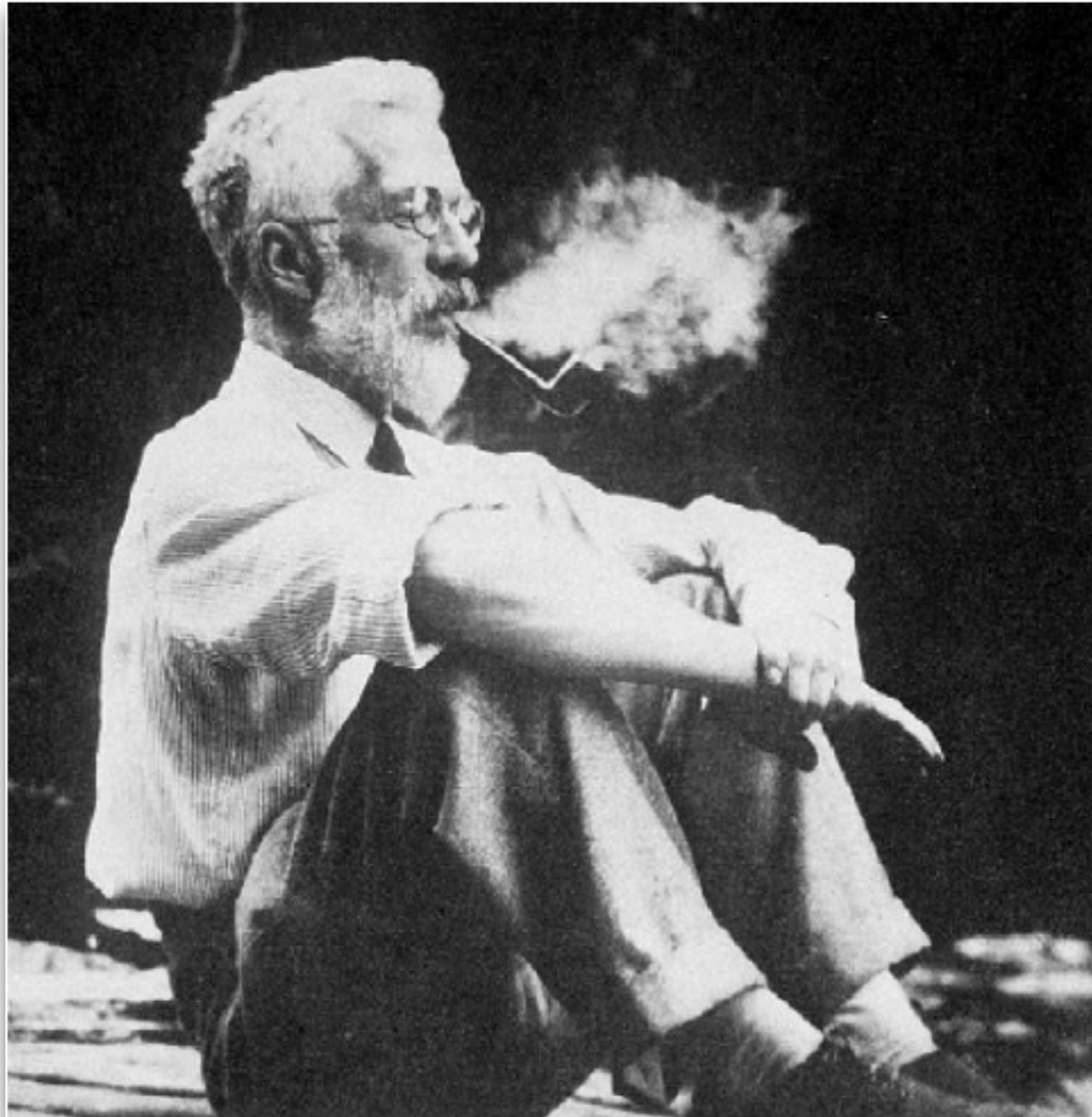
Ronald Fisher

# What is data science?

Sometimes the assumptions and the statements of statistical methods look *way too complex* to have anything to do with the real world.

E.g. assumptions behind a basic linear regression:
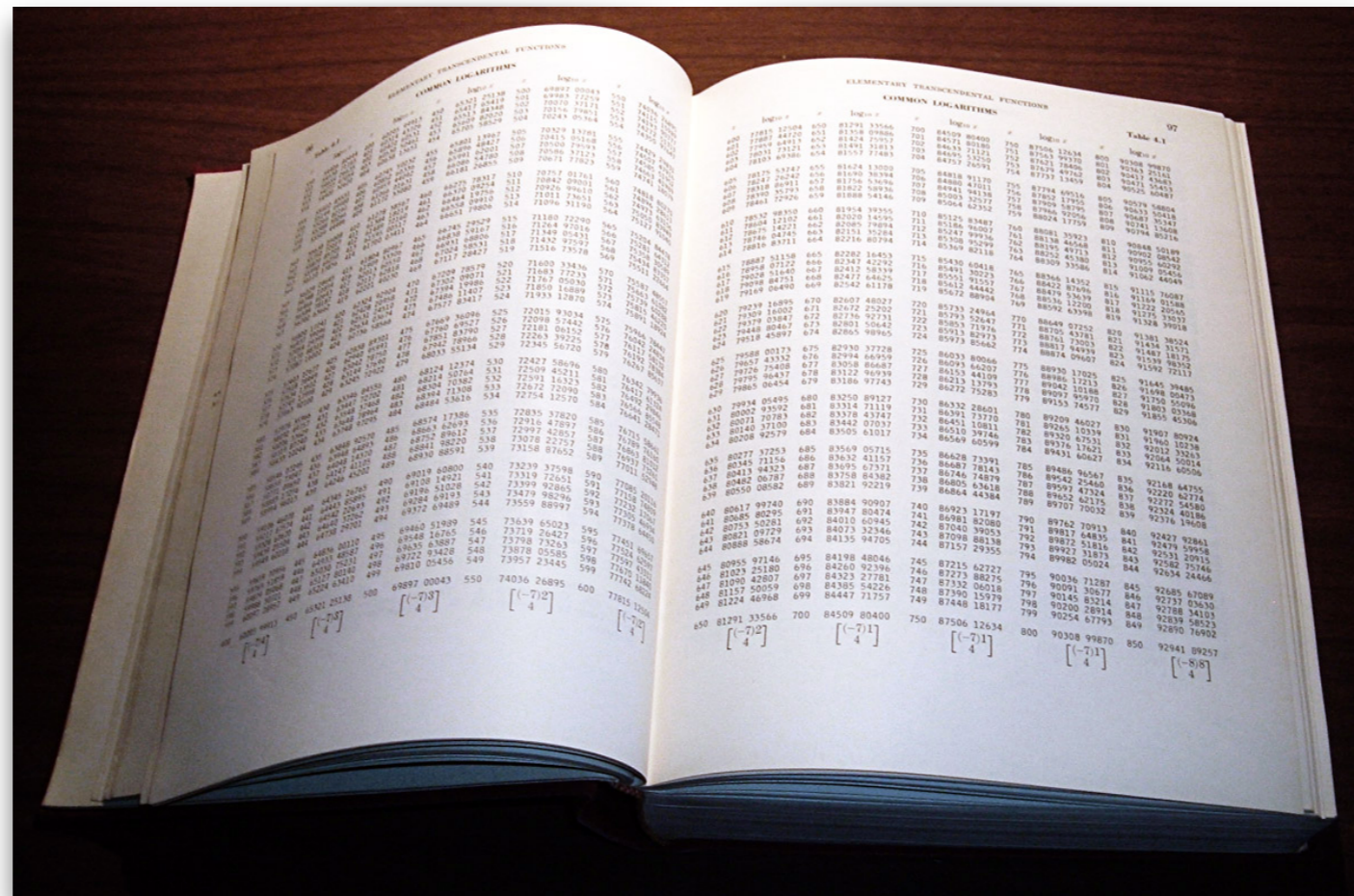**homoskedastic,
independent,
normally distributed data**

*No data I know in linguistics/ anthropology/cognitive sciences satisfy this!*

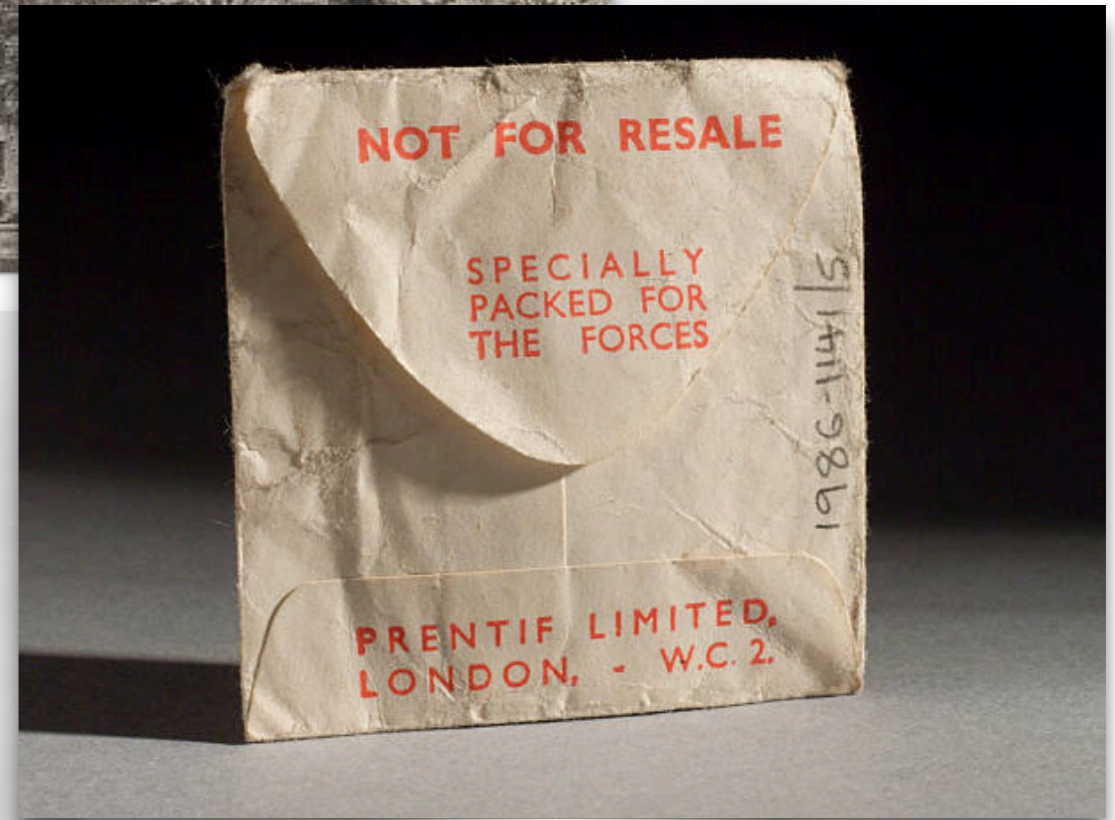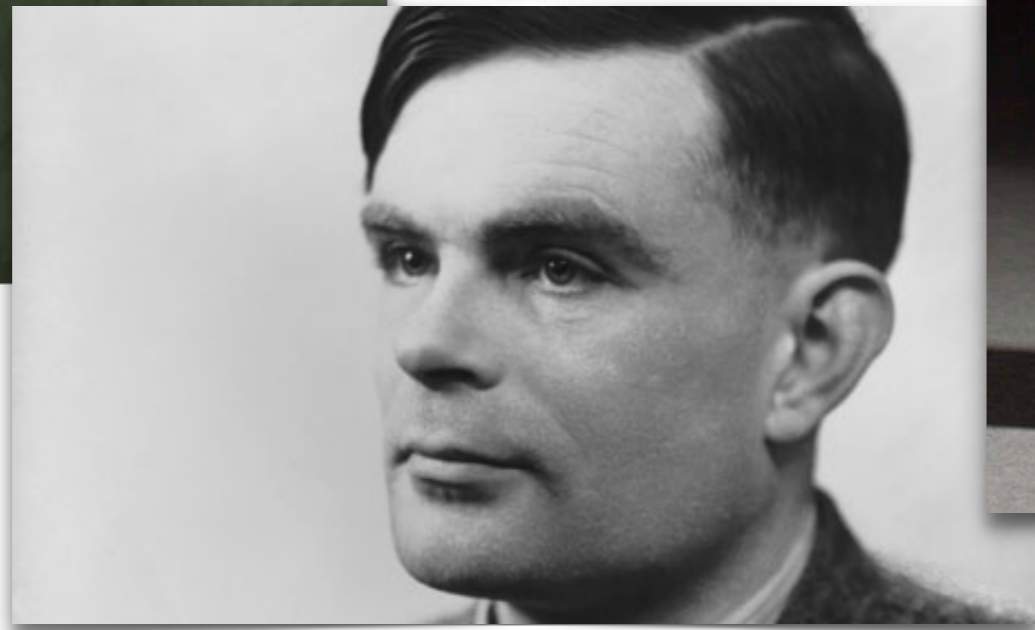# These assumptions were motivated by sheer convenience

*That explains the quirky math, but what about the recipe-like flavor of stats?*

George Barnard

The development of early methods was **strongly** practically oriented

NOT FOR RESALE

SPECIALLY
PACKED FOR
THE FORCES

PRENTIF LIMITED.
LONDON, - W.C. 2.

The usual statistics course in the social and health sciences results from the need of making standardized decisions

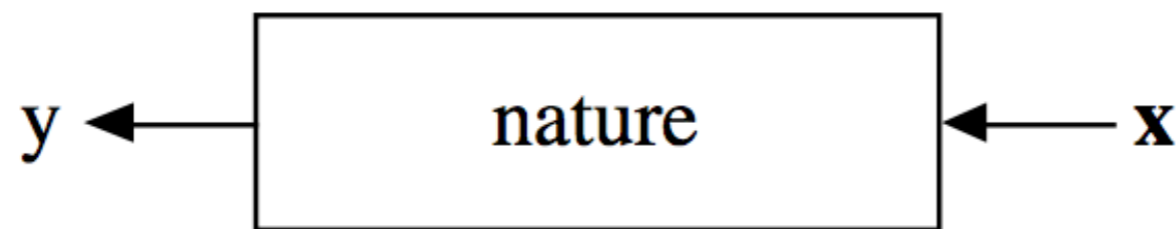Computers have partially solved for us the first issue - we don't require Procrustean assumptions anymore
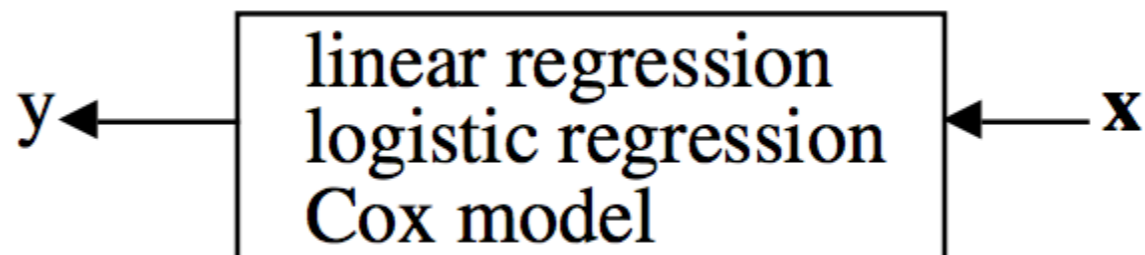
# Statistical Modeling: The Two Cultures

**Leo Breiman**

Statistics starts with data. Think of the data as being generated by a black box in which a vector of input variables $\mathbf{x}$ (independent variables) go in one side, and on the other side the response variables $\mathbf{y}$ come out. Inside the black box, nature functions to associate the predictor variables with the response variables, so the picture is like this:

$$\mathbf{y} \longleftarrow \boxed{\text{nature}} \longleftarrow \mathbf{x}$$

The values of the parameters are estimated from the data and the model then used for information and/or prediction. Thus the black box is filled in like this:

$$\mathbf{y} \longleftarrow \boxed{\begin{array}{l}\text{linear regression}\\\text{logistic regression}\\\text{Cox model}\end{array}} \longleftarrow \mathbf{x}$$
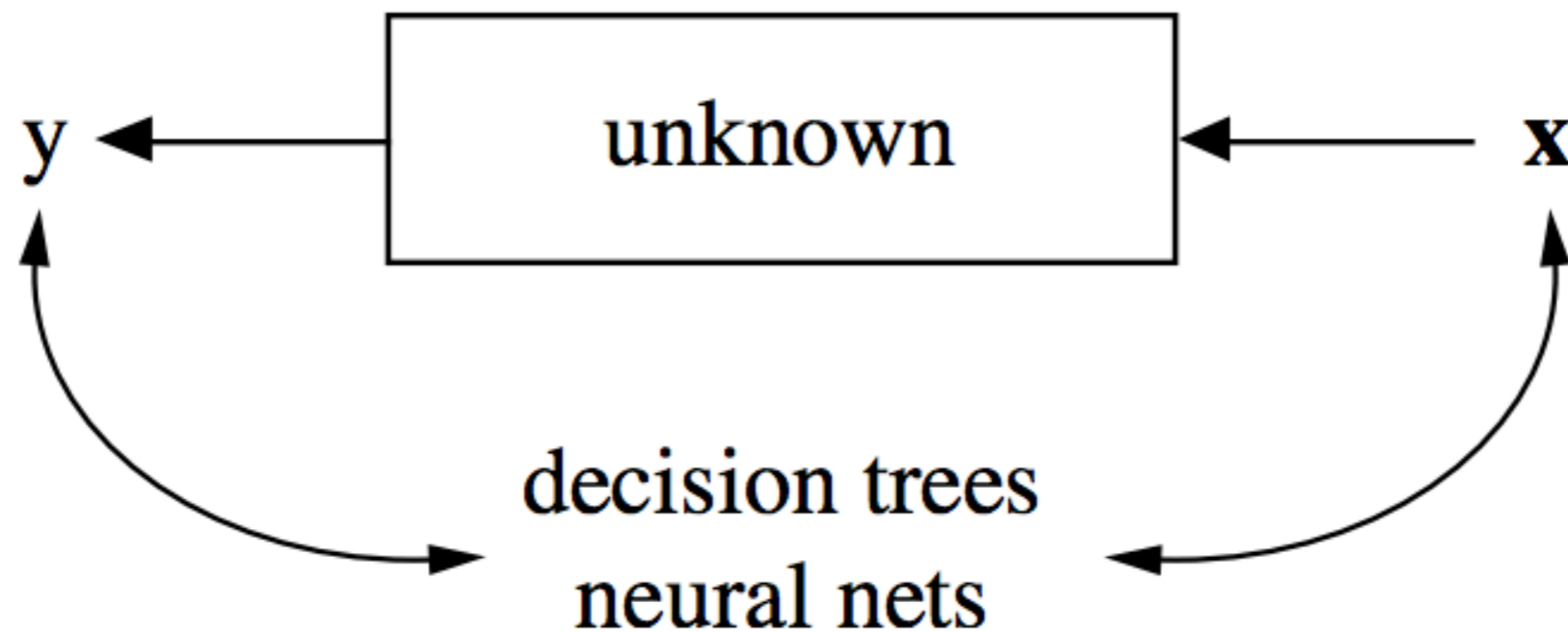
- The conclusions are about the model's mechanism, and not about nature's mechanism.

It follows that:

- If the model is a poor emulation of nature, the conclusions may be wrong.
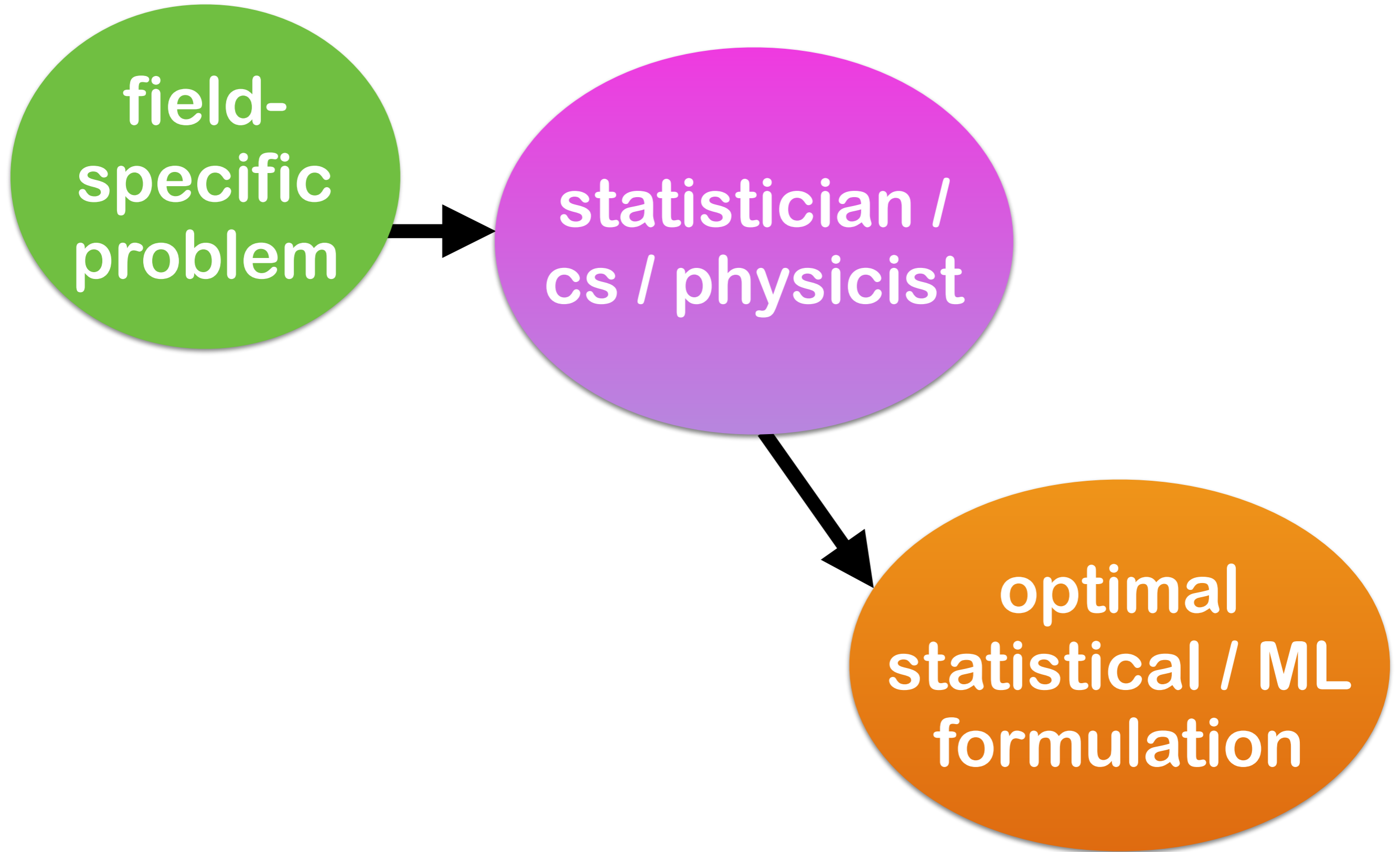
The analysis in this culture considers the inside of the box complex and unknown. Their approach is to find a function $f(\mathbf{x})$—an algorithm that operates on $\mathbf{x}$ to predict the responses $\mathbf{y}$. Their black box looks like this:

So now we have very powerful models that are able to describe (and predict) subtle associations in data without yielding an explicit mechanistic pathway

# Modelling the dynamics of language death

**Figure 1** The dynamics of language death. Symbols show the proportions of speakers over time of: **a**, Scottish Gaelic in Sutherland, Scotland[9]; **b**, Quechua in Huanuco, Peru; **c**, Welsh in Monmouthshire, Wales[10]; **d**, Welsh in all of Wales, from historical data[10] (blue) and a single modern census[11] (red). Fitted curves show solutions of the model in equation (1), with parameters $c$, $s$, $a$ and $x(0)$ estimated by least absolute-values regression. Where possible, data were obtained from several population censuses collected over a long timespan; otherwise, a single recent census with age-structured data was used (although errors are introduced, the size of which are reflected in the differing fits in **d**). Using the fraction of Catholic masses offered in Quechua in Peru as an indicator, we reconstructed an approximate history of the language's decline.

Of the remaining parameters, status, $s$, is the most relevant linguistically; it could serve as a useful measure of the threat to a given language. Quechua, for example, still has many speakers in Huanuco, Peru, but its low status is driving a rapid shift to Spanish, which leads to an unfortunate situation in which a child cannot communicate with his or her grandparents.

Contrary to the model's stark prediction, bilingual societies do, in fact, exist. But the

Contrary to the model's stark prediction, bilingual societies do, in fact, exist.
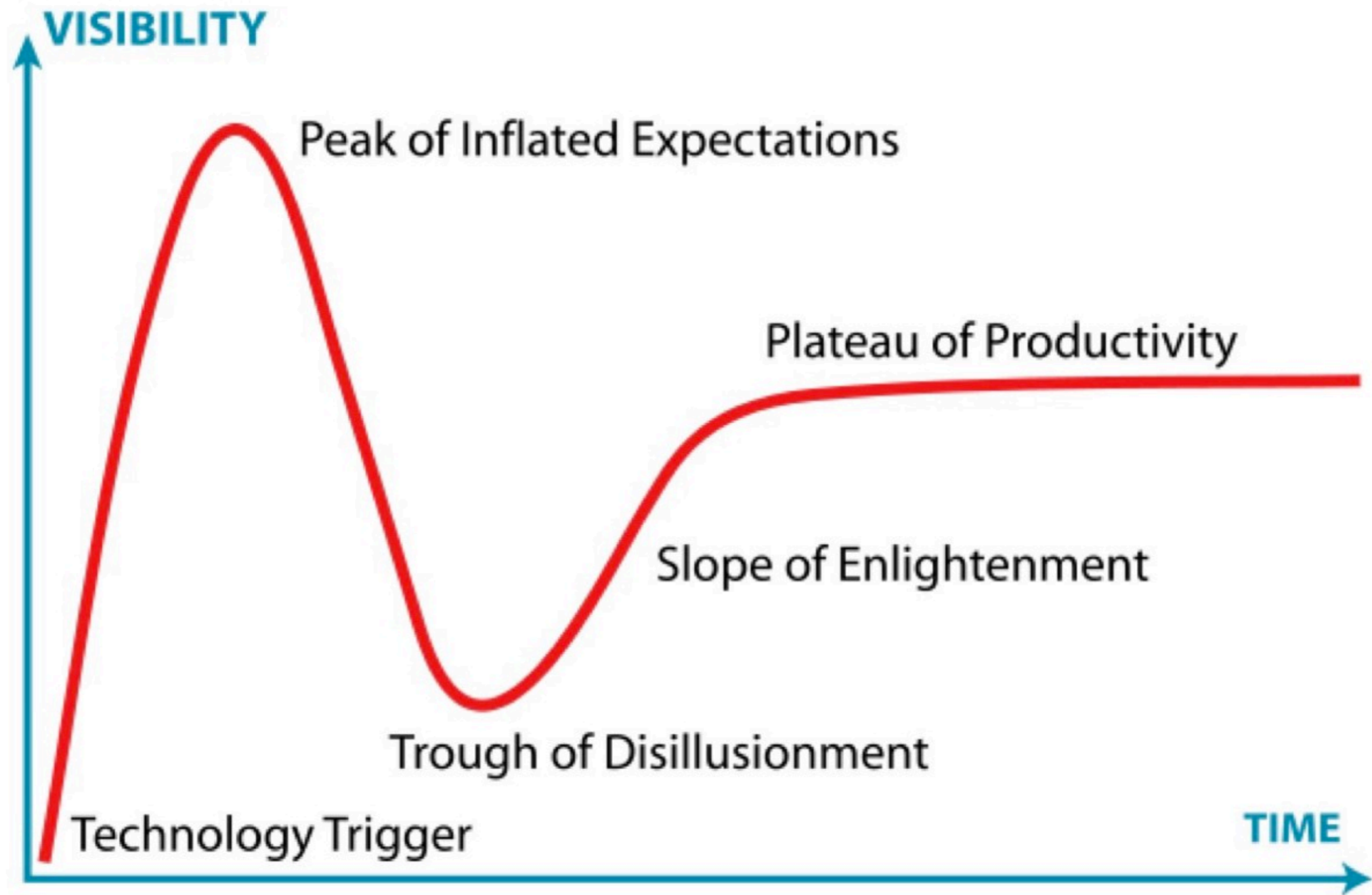
There is a blossoming data science of language that is *not* restricted to traditionally quantitative areas (such as e.g. psycholx, computational linguistics, phonetics, etc)

Importantly, more and more methods are informed by the special structure of the problem at hand
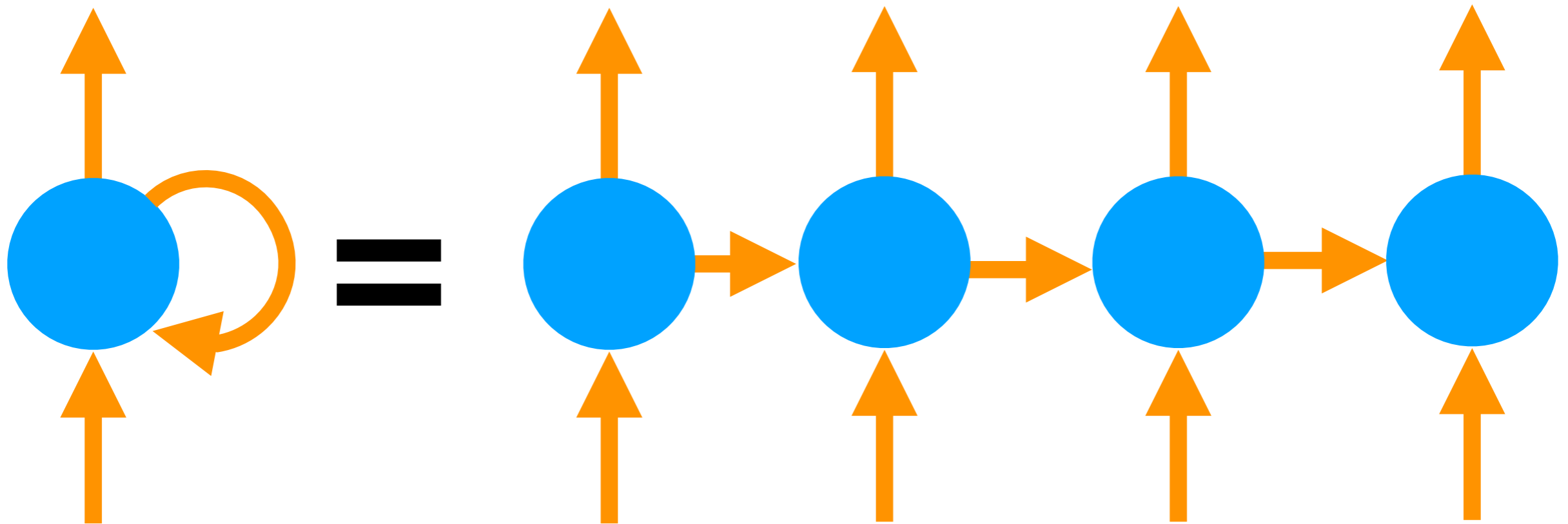
*So, where are we *not* going from here?*

**Explaining Character-Aware Neural Networks for Word-Level Prediction: Do They Discover Linguistic Rules?**

Frédéric Godin, Kris Demuynck, Joni Dambre, Wesley De Neve and Thomas Demeester

IDLab, Ghent University -
firstname.lastr

**Under the Hood: Using Diagnostic Classifiers to Investigate and Improve how Language Models Track Agreement Information**

**Mario Giulianelli**
University of Amsterdam

**Jack Harding**
University of Amsterdam

**Florian Mohnert**
University of Amsterdam

{mario.giulianelli, jack.harding, florian.mohnert}@student.uva.nl

pkes
Amsterdam
va.nl

**Willem Zuidema**
ILLC, University of Amsterdam
w.h.zuidema@uva.nl

# Representation of Linguistic Form and Function in Recurrent Neural Networks

Ákos Kádár , Grzegorz Chrupała

**Revisiting the poverty of the stimulus: hierarchical generalization without a hierarchical bias in recurrent neural networks**

R. Thomas McCoy (tom.mccoy@jhu.edu)
Department of Cognitive Science, Johns Hopkins University

Robert Frank (bob.frank@yale.edu)
Department of Linguistics, Yale University

Tal Linzen (tal.linzen@jhu.edu)
Department of Cognitive Science, Johns Hopkins University

*Fin*