

Econometrics & Machine Learning

A. Charpentier, E. Flachaire & A. Ly

<https://arxiv.org/abs/1708.06992>

Probabilistic Foundations of Econometrics

THE PROBABILITY APPROACH
IN ECONOMETRICS

By
TRYGVE HAAVELMO
RESEARCH ASSOCIATE
COWLES COMMISSION FOR
RESEARCH IN ECONOMICS

SUPPLEMENT TO ECONOMETRICA, VOLUME 12, JULY, 1944

THE ECONOMETRIC SOCIETY
THE UNIVERSITY OF CHICAGO
CHICAGO 37, ILLINOIS

see Haavelmo (1944)

THE FORMATION OF ECONOMETRICS

A Historical Perspective

QIN DUO

秦朵

1. The Probability Foundations of Econometrics	7
1.1 The Eve of the Probability Revolution	9
1.2 Introduction of Probability Theory	13
1.3 The Haavelmo Revolution	19
1.4 Alternative Approaches	26
1.5 An Incomplete Revolution	31

see Duo (1997)

Probabilistic Foundations of Econometrics

There is a probabilistic space $(\Omega, \mathcal{F}, \mathbb{P})$ such that data $\{y_i, \mathbf{x}_i\}$ can be seen as realizations of random variables $\{Y_i, \mathbf{X}_i\}$.

Consider a conditional distribution for $Y|\mathbf{X}$, e.g.

$$(Y|\mathbf{X} = \mathbf{x}) \stackrel{\mathcal{L}}{\sim} \mathcal{N}(\mu(\mathbf{x}), \sigma^2) \text{ where } \mu(\mathbf{x}) = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta}, \text{ and } \boldsymbol{\beta} \in \mathbb{R}^p.$$

for the linear model, with some extension when it is in the exponential family (GLM).

Then use maximum likelihood for inference, to derive confidence interval (quantification of uncertainty), etc.

Importance of **unbiased estimators** (Cramer Rao lower bound for the variance).

Loss Functions

Gneiting (2009) in a statistical context: a statistics T is said to be **ellicitable** if there is $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ such that

$$T(Y) = \operatorname{argmin}_{x \in \mathbb{R}} \left\{ \int_{\mathbb{R}} l(x, y) dF(y) \right\} = \operatorname{argmin}_{x \in \mathbb{R}} \left\{ \mathbb{E}[l(x, Y)] \text{ where } Y \stackrel{\mathcal{L}}{\sim} F \right\}$$

(e.g. $T(\mathbf{x}) = \bar{x}$ and $\ell = \ell_2$, or $T(\mathbf{x}) = \text{median}(\bar{x})$ and $\ell = \ell_1$).

In machine-learning, we want to solve

$$m^*(\mathbf{x}) = \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x})) \right\}$$

using optimization techniques, in a not-too-complex space \mathcal{M} .

Overfit ?

Underfit: true $y_i = \beta_0 + \mathbf{x}_1^\top \boldsymbol{\beta}_1 + \mathbf{x}_2^\top \boldsymbol{\beta}_2 + \varepsilon_i$ vs. fitted model $y_i = b_0 + \mathbf{x}_1^\top \mathbf{b}_1 + \eta_i$.

$$\hat{\mathbf{b}}_1 = (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} = \boldsymbol{\beta}_1 + \underbrace{(\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \boldsymbol{\beta}_2}_{\boldsymbol{\beta}_{12}} + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon}_{\nu_i}$$

i.e. $\mathbb{E}[\hat{\mathbf{b}}_1] = \boldsymbol{\beta}_1 + \boldsymbol{\beta}_{12} \neq \boldsymbol{\beta}_1$, unless $\mathbf{X}_1 \perp \mathbf{X}_2$ (see Frish-Waugh Theorem).

Overfit: true $y_i = \beta_0 + \mathbf{x}_1^\top \boldsymbol{\beta}_1 + \varepsilon_i$ vs. fitted model $y_i = b_0 + \mathbf{x}_1^\top \mathbf{b}_1 + \mathbf{x}_2^\top \mathbf{b}_2 + \eta_i$.

In that case $\mathbb{E}[\hat{\mathbf{b}}_1] = \boldsymbol{\beta}_1$ but no-longer efficient.

Occam's Razor and Parsimony

Importance of **penalty** in order to avoid a too complex model (overfit).

Consider some linear predictor, $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ where $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}$, or more generally $\hat{\mathbf{y}} = \mathbf{S}\mathbf{y}$ for some smoothing matrix \mathbf{S} .

In econometrics, consider some ex-post penalty for model selection,

$$AIC = -2 \log \mathcal{L} + 2p = \text{deviance} + 2p$$

where p is the dimension (i.e. $\text{trace}(\mathbf{S})$ in a more general setting).

In machine-learning, penalty is added in the objective function, see Ridge or LASSO regression

$$(\hat{\beta}_{0,\lambda}, \beta_\lambda) = \underset{(\beta_0, \hat{\beta})}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \ell(y_i, \beta_0 + \mathbf{x}_i^\top \beta) + \lambda \|\beta\| \right\}$$

Goodhart's law, “when a measure becomes a target, it ceases to be a good measure”

Boosting, or learning from previous errors

Construct a sequence of models

$$m^{(k)}(\mathbf{x}) = m^{(k-1)}(\mathbf{x}) + \alpha \cdot f^*(\mathbf{x})$$

where

$$f^* = \operatorname{argmin}_{f \in \mathcal{W}} \left\{ \sum_{i=1}^n \ell(y_i - m^{(k-1)}(\mathbf{x}_i), f(\mathbf{x}_i)) \right\}$$

for some set of **weak learner** \mathcal{W} .

Problem: where to stop to avoid overfit...