

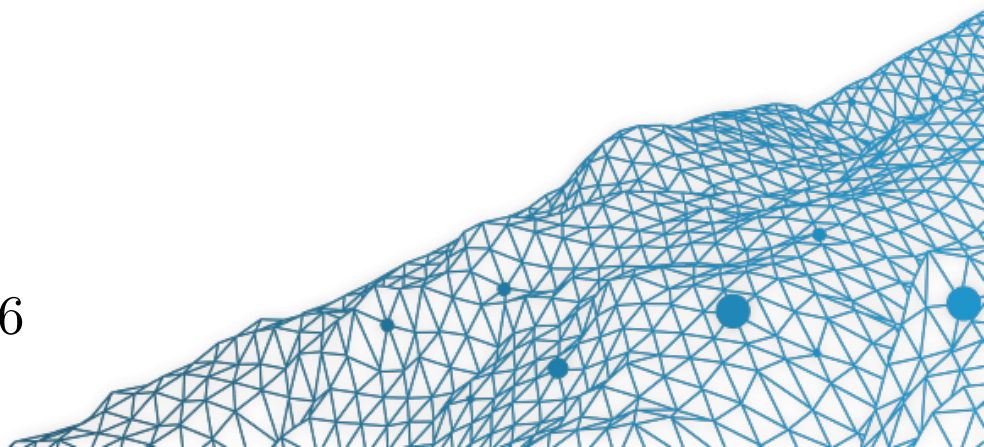
# Econometrics: Learning from 'Statistical Learning' Techniques

Arthur Charpentier (Université de Rennes 1 & UQàM)

Centre for Central Banking Studies

Bank of England, London, UK, May 2016

<http://freakonometrics.hypotheses.org>



# Econometrics: Learning from 'Statistical Learning' Techniques

Arthur Charpentier (Université de Rennes 1 & UQàM)

Professor, Economics Department, Univ. Rennes 1

In charge of Data Science for Actuaries program, IA

Research Chair *actinfo* (Institut Louis Bachelier)

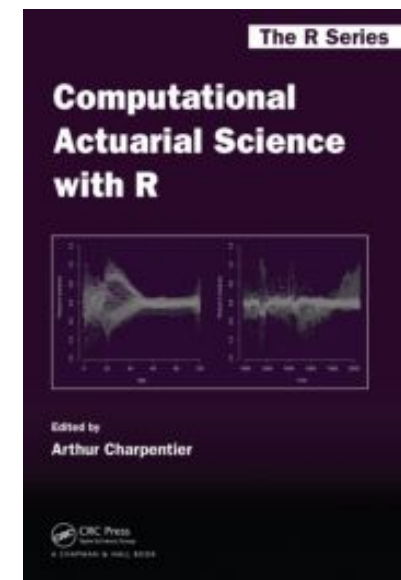
(previously Actuarial Sciences at UQàM & ENSAE Paristech  
actuary in Hong Kong, IT & Stats FFSA)

PhD in Statistics (KU Leuven), Fellow Institute of Actuaries

MSc in Financial Mathematics (Paris Dauphine) & ENSAE

Editor of the [freakonometrics.hypotheses.org](http://freakonometrics.hypotheses.org)'s blog

Editor of Computational Actuarial Science, CRC



## Agenda

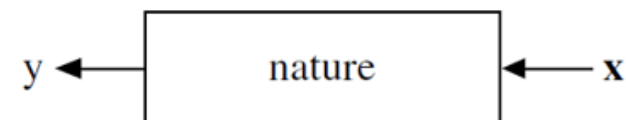
*“the numbers have no way of speaking for themselves. We speak for them. [...] Before we demand more of our data, we need to demand more of ourselves”* from Silver (2012).

- (big) data
- econometrics & probabilistic modeling
- algorithmics & statistical learning
- different perspectives on classification
- bootstrapping, PCA & variable selection

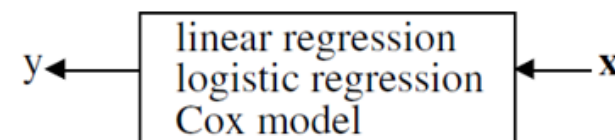
see Berk (2008), Hastie, Tibshirani & Friedman (2009), but also Breiman (2001)

*Statistical Science*  
2001, Vol. 16, No. 3, 199–231

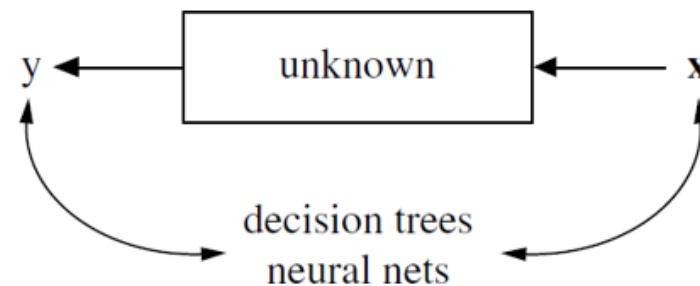
## Statistical Modeling: The Two Cultures



### The Data Modeling Culture



### The Algorithmic Modeling Culture



## Data and Models

From  $\{(y_i, \mathbf{x}_i)\}$ , there are different stories behind, see [Freedman \(2005\)](#)

- the **causal story** :  $x_{j,i}$  is usually considered as independent of the other covariates  $x_{k,i}$ . For all possible  $\mathbf{x}$ , that value is mapped to  $m(\mathbf{x})$  and a noise is attached,  $\varepsilon$ . The goal is to recover  $m(\cdot)$ , and the residuals are just the difference between the response value and  $m(\mathbf{x})$ .
- the **conditional distribution story** : for a linear model, we usually say that  $Y$  given  $\mathbf{X} = \mathbf{x}$  is a  $\mathcal{N}(m(\mathbf{x}), \sigma^2)$  distribution.  $m(\mathbf{x})$  is then the conditional mean. Here  $m(\cdot)$  is assumed to really exist, but no causal assumption is made, only a conditional one.
- the **explanatory data story** : there is no model, just data. We simply want to summarize information contained in  $\mathbf{x}$ 's to get an accurate summary, close to the response (i.e.  $\min\{\ell(\mathbf{y}, m(\mathbf{x}))\}$ ) for some loss function  $\ell$ .

See also [Varian \(2014\)](#)

## Data, Models & Causal Inference

We cannot differentiate data and model that easily..

After an operation, should I stay at hospital, or go back home ?

as in Angrist & Pischke (2008),

$$(\text{health} \mid \text{hospital}) - (\text{health} \mid \text{stayed home}) \quad [\text{observed}]$$

should be written

$$(\text{health} \mid \text{hospital}) - (\text{health} \mid \textit{had} \text{ stayed home}) \quad [\text{treatment effect}]$$

$$+ (\text{health} \mid \textit{had} \text{ stayed home}) - (\text{health} \mid \text{stayed home}) \quad [\text{selection bias}]$$

Need randomization to solve selection bias.

## Econometric Modeling

Data  $\{(y_i, \mathbf{x}_i)\}$ , for  $i = 1, \dots, n$ , with  $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^p$  and  $y_i \in \mathcal{Y}$ .

A model is a  $m : \mathcal{X} \mapsto \mathcal{Y}$  mapping

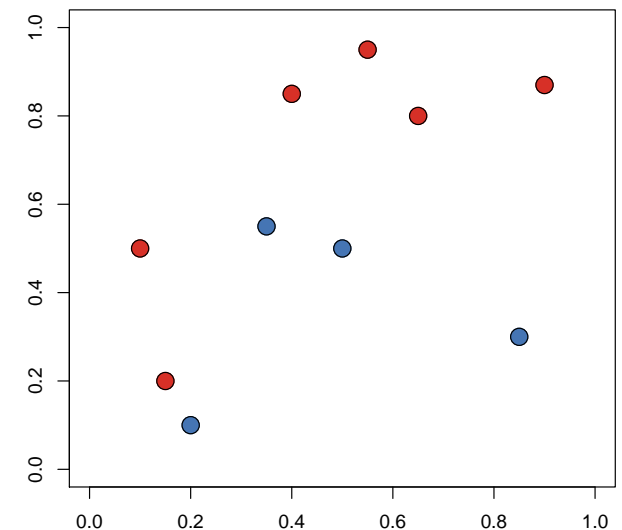
- regression,  $\mathcal{Y} = \mathbb{R}$  (but also  $\mathcal{Y} = \mathbb{N}$ )
- classification,  $\mathcal{Y} = \{0, 1\}, \{-1, +1\}, \{\bullet, \bullet\}$   
(binary, or more)

Classification models are based on two steps,

- **score** function,  $s(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) \in [0, 1]$



- **classifier**  $s(\mathbf{x}) \rightarrow \hat{y} \in \{0, 1\}$ .



## High Dimensional Data (not to say 'Big Data')

See Bühlmann & van de Geer (2011) or Koch (2013),  $\mathbf{X}$  is a  $n \times p$  matrix

Portnoy (1988) proved that maximum likelihood estimators are asymptotically normal when  $p^2/n \rightarrow 0$  as  $n, p \rightarrow \infty$ . Hence, **massive data**, when  $p > \sqrt{n}$ .

More interesting is the **sparcity** concept, based not on  $p$ , but on the effective size. Hence one can have  $p > n$  and convergent estimators.

High dimension might be scary because of **curse of dimensionality**, see Bellman (1957). The volume of the unit sphere in  $\mathbb{R}^p$  tends to 0 as  $p \rightarrow \infty$ , i.e. space is sparse.

## Computational & Nonparametric Econometrics

Linear Econometrics: estimate  $g : \mathbf{x} \mapsto \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$  by a linear function.

Nonlinear Econometrics: consider the approximation for some **functional basis**

$$g(\mathbf{x}) = \sum_{j=0}^{\infty} \omega_j g_j(\mathbf{x}) \text{ and } \hat{g}(\mathbf{x}) = \sum_{j=0}^h \omega_j g_j(\mathbf{x})$$

or consider a **local model**, on the neighborhood of  $\mathbf{x}$ ,

$$\hat{g}(x) = \frac{1}{n_{\mathbf{x}}} \sum_{i \in \mathcal{I}_{\mathbf{x}}} y_i, \text{ with } \mathcal{I}_{\mathbf{x}} = \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}_i - \mathbf{x}\| \leq h\},$$

see [Nadaraya \(1964\)](#) and [Watson \(1964\)](#).

Here  $h$  is some **tunning parameter**: not estimated, but chosen (optimally).



## Econometrics & Probabilistic Model

The primary goal in a regression analysis is to understand, as far as possible with the available data, how the conditional distribution of the response  $y$  varies across subpopulations determined by the possible values of the predictor or predictors. Since this is the central idea, it will be helpful to have a convey from [Cook & Weisberg \(1999\)](#), see also [Haavelmo \(1965\)](#).

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

**Linear Model:**  $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta}$

**Homoscedasticity:**  $\text{Var}[Y|\mathbf{X} = \mathbf{x}] = \sigma^2.$

## Conditional Distribution and Likelihood

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

The log-likelihood is

$$\log \mathcal{L}(\beta_0, \boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{x}) = -\frac{n}{2} \log[2\pi\sigma^2] - \frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2}_{\text{}}.$$

Set

$$(\hat{\beta}_0, \hat{\boldsymbol{\beta}}, \hat{\sigma}^2) = \operatorname{argmax} \{ \log \mathcal{L}(\beta_0, \boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{x}) \}.$$

First order condition  $\mathbf{X}^\top [\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}] = \mathbf{0}$ . If matrix  $\mathbf{X}$  is a full rank matrix

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}.$$

Asymptotic properties of  $\hat{\boldsymbol{\beta}}$ ,

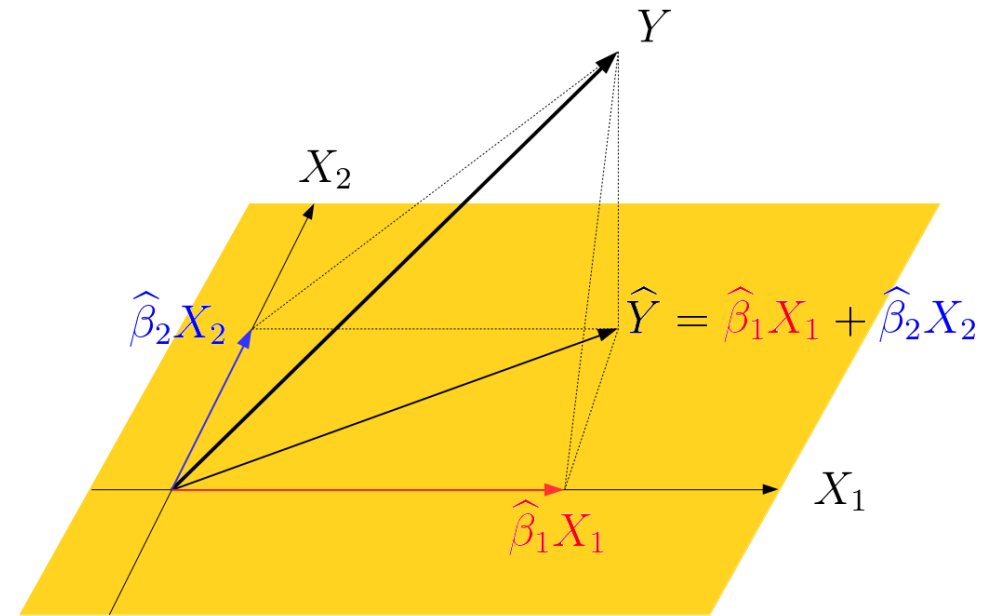
$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}) \text{ as } n \rightarrow \infty$$

## Geometric Perspective

Define the orthogonal projection on  $\mathcal{X}$ ,

$$\Pi_{\mathbf{X}} = \mathbf{X}[\mathbf{X}^{\top} \mathbf{X}]^{-1} \mathbf{X}^{\top}$$

$$\hat{\mathbf{y}} = \underbrace{\mathbf{X}[\mathbf{X}^{\top} \mathbf{X}]^{-1} \mathbf{X}^{\top}}_{\Pi_{\mathbf{X}}} \mathbf{y} = \Pi_{\mathbf{X}} \mathbf{y}.$$



**Pythagoras' theorem** can be written

$$\|\mathbf{y}\|^2 = \|\Pi_{\mathbf{X}} \mathbf{y}\|^2 + \|\Pi_{\mathbf{X}^{\perp}} \mathbf{y}\|^2 = \|\Pi_{\mathbf{X}} \mathbf{y}\|^2 + \|\mathbf{y} - \Pi_{\mathbf{X}} \mathbf{y}\|^2$$

which can be expressed as

$$\underbrace{\sum_{i=1}^n y_i^2}_{n \times \text{total variance}} = \underbrace{\sum_{i=1}^n \hat{y}_i^2}_{n \times \text{explained variance}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{n \times \text{residual variance}}$$

## Geometric Perspective

Define the angle  $\theta$  between  $\mathbf{y}$  and  $\Pi_{\mathcal{X}}\mathbf{y}$ ,

$$R^2 = \frac{\|\Pi_{\mathcal{X}}\mathbf{y}\|^2}{\|\mathbf{y}\|^2} = 1 - \frac{\|\Pi_{\mathcal{X}^\perp}\mathbf{y}\|^2}{\|\mathbf{y}\|^2} = \cos^2(\theta)$$

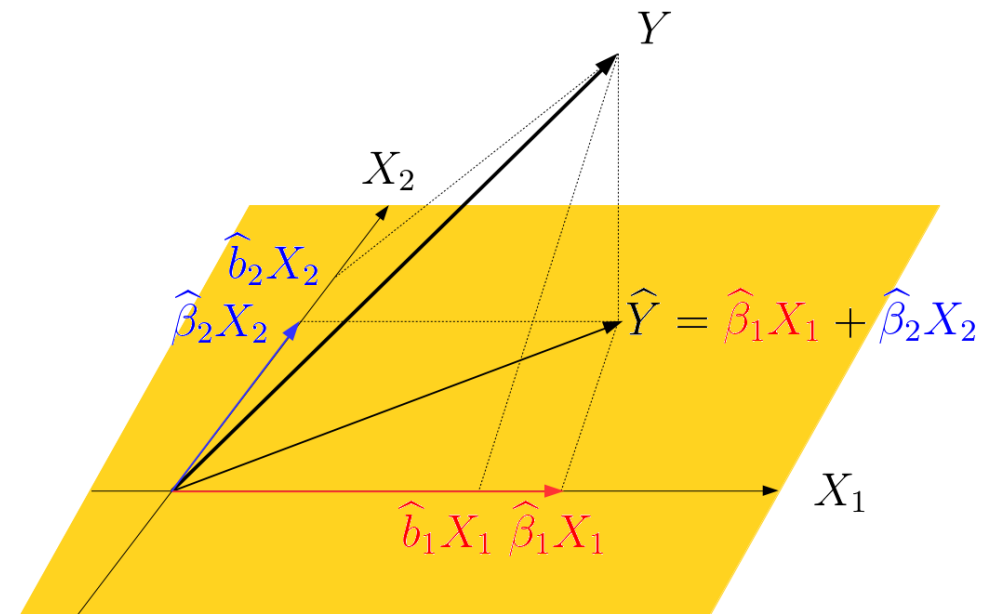
see Davidson & MacKinnon (2003)

$$\mathbf{y} = \beta_0 + \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \varepsilon$$

If  $\mathbf{y}_2^* = \Pi_{\mathcal{X}_1^\perp}\mathbf{y}$  and  $\mathbf{X}_2^* = \Pi_{\mathcal{X}_1^\perp}\mathbf{X}_2$ , then

$$\hat{\beta}_2 = [\mathbf{X}_2^{*\top}\mathbf{X}_2^*]^{-1}\mathbf{X}_2^{*\top}\mathbf{y}_2^*$$

$\mathbf{X}_2^* = \mathbf{X}_2$  if  $\mathbf{X}_1 \perp \mathbf{X}_2$ ,  
Frisch-Waugh theorem.



## From Linear to Non-Linear

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \underbrace{\mathbf{X}[\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top}_{\mathbf{H}} \mathbf{y} \text{ i.e. } \hat{y}_i = \mathbf{h}_{\mathbf{x}_i}^\top \mathbf{y},$$

with  $\mathbf{h}_{\mathbf{x}}$  for the linear regression -  $\mathbf{h}_{\mathbf{x}} = \mathbf{X}[\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{x}$ .

One can consider some smoothed regression, see [Nadaraya \(1964\)](#) and [Watson \(1964\)](#), with some smoothing matrix  $\mathbf{S}$

$$\hat{m}_h(x) = \mathbf{s}_x^\top \mathbf{y} = \sum_{i=1}^n s_{x,i} y_i \text{ with } s_{x,i} = \frac{K_h(x - x_i)}{K_h(x - x_1) + \dots + K_h(x - x_n)}$$

for some kernel  $K(\cdot)$  and some bandwidth  $h > 0$ .

## From Linear to Non-Linear

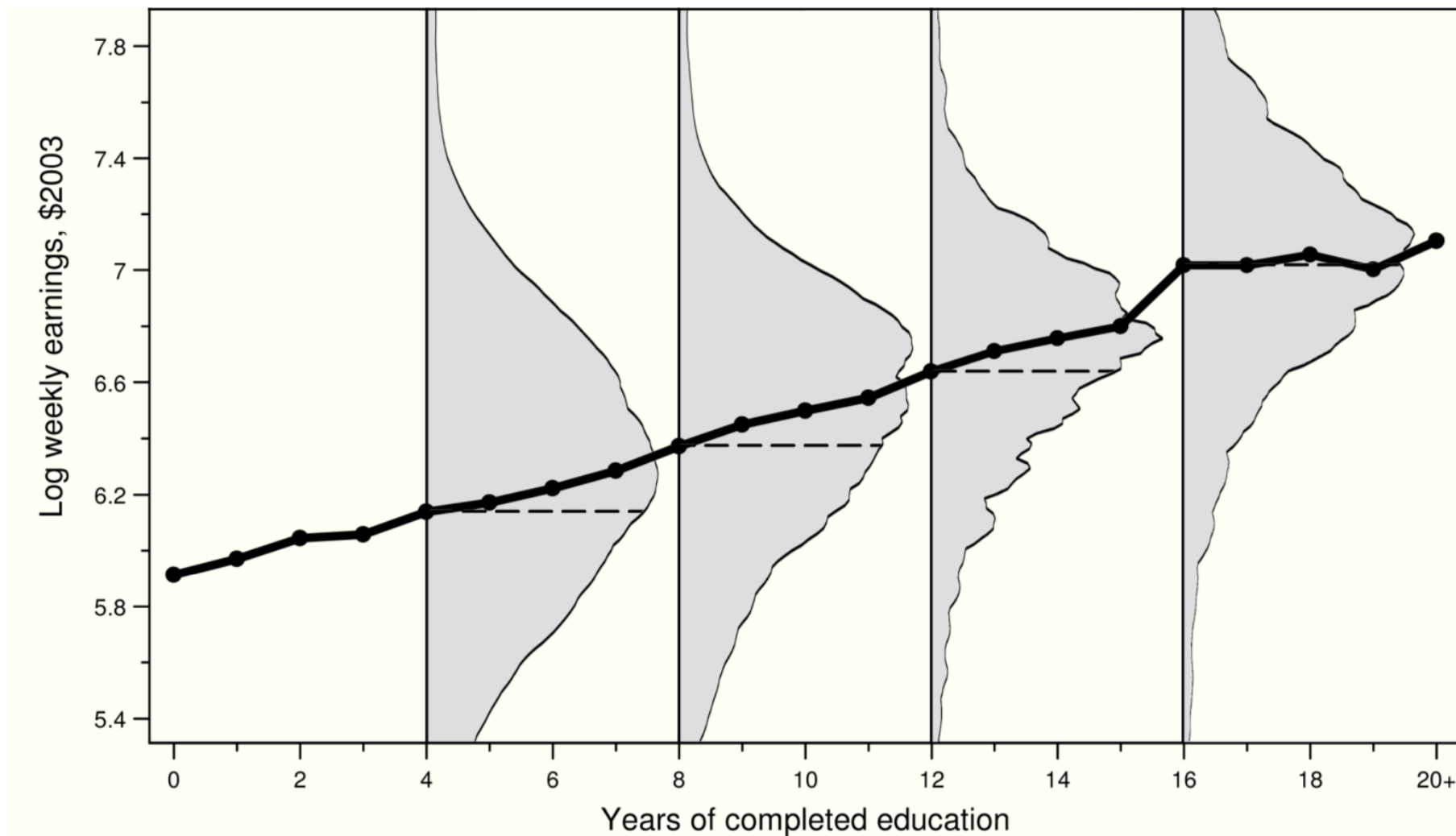
$$T = \frac{\|\mathbf{S}\mathbf{y} - \mathbf{H}\mathbf{y}\|}{\text{trace}([\mathbf{S} - \mathbf{H}]^{\top}[\mathbf{S} - \mathbf{H}])}$$

can be used to test for linearity, [Simonoff \(1996\)](#).  $\text{trace}(\mathbf{S})$  is the equivalent number of parameters, and  $n - \text{trace}(\mathbf{S})$  the degrees of freedom, [Ruppert et al. \(2003\)](#).

Nonlinear Model, but Homoscedastic - Gaussian

- $(Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{N}(\mu(\mathbf{x}), \sigma^2)$
- $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = \mu(\mathbf{x})$

## Conditional Expectation



from Angrist & Pischke (2008),  $x \mapsto \mathbb{E}[Y|X = x]$ .

## Exponential Distributions and Linear Models

$$f(y_i|\theta_i, \phi) = \exp\left(\frac{y_i\theta_i - b(\theta_i)}{a(\phi)} + c(y_i, \phi)\right) \text{ with } \theta_i = h(\mathbf{x}_i^\top \boldsymbol{\beta})$$

Log likelihood is expressed as

$$\log \mathcal{L}(\boldsymbol{\theta}, \phi|\mathbf{y}) = \sum_{i=1}^n \log f(y_i|\theta_i, \phi) = \frac{\sum_{i=1}^n y_i\theta_i - \sum_{i=1}^n b(\theta_i)}{a(\phi)} + \sum_{i=1}^n c(y_i, \phi)$$

and first order conditions

$$\frac{\partial \log \mathcal{L}(\boldsymbol{\theta}, \phi|\mathbf{y})}{\partial \boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{W}^{-1}[\mathbf{y} - \boldsymbol{\mu}] = \mathbf{0}$$

as in Müller (2001), where  $\mathbf{W}$  is a weight matrix, function of  $\boldsymbol{\beta}$ .

We usually specify the **link** function  $g(\cdot)$  defined as

$$\hat{y} = m(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = g^{-1}(\mathbf{x}^\top \boldsymbol{\beta}).$$



## Exponential Distributions and Linear Models

Note that  $\mathbf{W} = \text{diag}(\nabla g(\hat{\mathbf{y}}) \cdot \text{Var}[\mathbf{y}])$ , and set

$$\mathbf{z} = g(\hat{\mathbf{y}}) + (\mathbf{y} - \hat{\mathbf{y}}) \cdot \nabla g(\hat{\mathbf{y}})$$

the the maximum likelihood estimator is obtained iteratively

$$\hat{\boldsymbol{\beta}}_{k+1} = [\mathbf{X}^\top \mathbf{W}_k^{-1} \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{W}_k^{-1} \mathbf{z}_k$$

Set  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_\infty$ , so that

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, I(\boldsymbol{\beta})^{-1})$$

with  $I(\boldsymbol{\beta}) = \phi \cdot [\mathbf{X}^\top \mathbf{W}_\infty^{-1} \mathbf{X}]$ .

Note that  $[\mathbf{X}^\top \mathbf{W}_k^{-1} \mathbf{X}]$  is a  $p \times p$  matrix.

## Exponential Distributions and Linear Models

Generalized Linear Model:

- $(Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{L}(\theta_{\mathbf{x}}, \varphi)$
- $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = h^{-1}(\theta_{\mathbf{x}}) = g^{-1}(\mathbf{x}^{\top}\boldsymbol{\beta})$

e.g.  $(Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{P}(\exp[\mathbf{x}^{\top}\boldsymbol{\beta}])$ .

Use of maximum likelihood techniques for inference.

Actually, more a moment condition than a distribution assumption.

## Goodness of Fit & Model Choice

From the variance decomposition

$$\underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}_{\text{total variance}} = \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{residual variance}} + \underbrace{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{explained variance}}$$

and define

$$R^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

More generally

$$\text{Deviance}(\beta) = -2 \log[\mathcal{L}] = 2 \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \text{Deviance}(\hat{\mathbf{y}})$$

The null deviance is obtained using  $\hat{y}_i = \bar{y}$ , so that

$$R^2 = \frac{\text{Deviance}(\bar{y}) - \text{Deviance}(\hat{\mathbf{y}})}{\text{Deviance}(\bar{y})} = 1 - \frac{\text{Deviance}(\hat{\mathbf{y}})}{\text{Deviance}(\bar{y})} = 1 - \frac{D}{D_0}$$

## Goodness of Fit & Model Choice

One usually prefers a penalized version

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p} = R^2 - \underbrace{(1 - R^2) \frac{p - 1}{n - p}}_{\text{penalty}}$$

See also **Akaike** criteria  $AIC = \text{Deviance} + 2 \cdot p$

or **Schwarz**,  $BIC = \text{Deviance} + \log(n) \cdot p$

In high dimension, consider a corrected version

$$AIC_c = \text{Deviance} + 2 \cdot p \cdot \frac{n}{n - p - 1}$$

## Stepwise Procedures

### Forward algorithm

1. set  $j_1^* = \operatorname{argmin}_{j \in \{\emptyset, 1, \dots, n\}} \{AIC(\{j\})\}$
2. set  $j_2^* = \operatorname{argmin}_{j \in \{\emptyset, 1, \dots, n\} \setminus \{j_1^*\}} \{AIC(\{j_1^*, j\})\}$
3. ... until  $j^* = \emptyset$

### Backward algorithm

1. set  $j_1^* = \operatorname{argmin}_{j \in \{\emptyset, 1, \dots, n\}} \{AIC(\{1, \dots, n\} \setminus \{j\})\}$
2. set  $j_2^* = \operatorname{argmin}_{j \in \{\emptyset, 1, \dots, n\} \setminus \{j_1^*\}} \{AIC(\{1, \dots, n\} \setminus \{j_1^*, j\})\}$
3. ... until  $j^* = \emptyset$

## Econometrics & Statistical Testing

Standard test for  $H_0 : \beta_k = 0$  against  $H_1 : \beta_k \neq 0$  is **Student- $t$**  est  $t_k = \hat{\beta}_k / \text{se}_{\hat{\beta}_k}$ ,

Use the  **$p$ -value**  $\mathbb{P}[|T| > |t_k|]$  with  $T \sim t_\nu$  (and  $\nu = \text{trace}(\mathbf{H})$ ).

In high dimension, consider the FDR (False Discovery Ratio).

With  $\alpha = 5\%$ , 5% variables are wrongly significant.

If  $p = 100$  with only 5 significant variables, one should expect also 5 false positive, i.e. 50% FDR, see [Benjamini & Hochberg \(1995\)](#) and Andrew Gelman's talk.

## Under & Over-Identification

**Under-identification** is obtained when the true model is

$$y = \beta_0 + \mathbf{x}_1^\top \boldsymbol{\beta}_1 + \mathbf{x}_2^\top \boldsymbol{\beta}_2 + \varepsilon, \text{ but we estimate } y = \beta_0 + \mathbf{x}_1^\top \mathbf{b}_1 + \eta.$$

Maximum likelihood estimator for  $\mathbf{b}_1$  is

$$\begin{aligned} \hat{\mathbf{b}}_1 &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top [\mathbf{X}_{1,i} \boldsymbol{\beta}_1 + \mathbf{X}_{2,i} \boldsymbol{\beta}_2 + \varepsilon] \\ &= \boldsymbol{\beta}_1 + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \boldsymbol{\beta}_2}_{\boldsymbol{\beta}_{12}} + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon}_{\nu_i} \end{aligned}$$

so that  $\mathbb{E}[\hat{\mathbf{b}}_1] = \boldsymbol{\beta}_1 + \boldsymbol{\beta}_{12}$ , and the bias is null when  $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$  i.e.  $\mathbf{X}_1 \perp \mathbf{X}_2$ , see Frisch-Waugh).

**Over-identification** is obtained when the true model is  $y = \beta_0 + \mathbf{x}_1^\top \boldsymbol{\beta}_1 + \varepsilon$ , but we fit  $y = \beta_0 + \mathbf{x}_1^\top \mathbf{b}_1 + \mathbf{x}_2^\top \mathbf{b}_2 + \eta$ .

Inference is unbiased since  $\mathbb{E}(\mathbf{b}_1) = \boldsymbol{\beta}_1$  but the estimator is not efficient.

## Statistical Learning & Loss Function

Here, no probabilistic model, but a **loss function**,  $\ell$ . For some set of functions  $\mathcal{M}$ ,  $\mathcal{X} \rightarrow \mathcal{Y}$ , define

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

**Quadratic loss** functions are interesting since

$$\bar{y} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \frac{1}{n} [y_i - m]^2 \right\}$$

which can be written, with some underlying probabilistic model

$$\mathbb{E}(Y) = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \|Y - m\|_{\ell_2}^2 \right\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \mathbb{E}([Y - m]^2) \right\}$$

For  $\tau \in (0, 1)$ , we obtain the **quantile regression** (see Koenker (2005))

$$m^* = \operatorname{argmin}_{m \in \mathcal{M}_0} \left\{ \sum_{i=1}^n \ell_\tau(y_i, m(\mathbf{x}_i)) \right\} \text{ avec } \ell_\tau(x, y) = |(x - y)(\tau - \mathbf{1}_{x \leq y})|$$



## Boosting & Weak Learning

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

is hard to solve for some very large and general space  $\mathcal{M}$  of  $\mathcal{X} \rightarrow \mathcal{Y}$  functions.

Consider some iterative procedure, where we learn from the errors,

$$m^{(k)}(\cdot) = \underbrace{m_1(\cdot)}_{\sim \mathbf{y}} + \underbrace{m_2(\cdot)}_{\sim \boldsymbol{\varepsilon}_1} + \underbrace{m_3(\cdot)}_{\sim \boldsymbol{\varepsilon}_2} + \cdots + \underbrace{m_k(\cdot)}_{\sim \boldsymbol{\varepsilon}_{k-1}} = m^{(k-1)}(\cdot) + m_k(\cdot).$$

Formely  $\boldsymbol{\varepsilon}$  can be seen as  $\nabla \ell$ , the gradient of the loss.

## Boosting & Weak Learning

It is possible to see this algorithm as a gradient descent. Not

$$\underbrace{f(\mathbf{x}_k)}_{\langle f, \mathbf{x}_k \rangle} \sim \underbrace{f(\mathbf{x}_{k-1})}_{\langle f, \mathbf{x}_{k-1} \rangle} + \underbrace{(\mathbf{x}_k - \mathbf{x}_{k-1})}_{\alpha_k} \underbrace{\nabla f(\mathbf{x}_{k-1})}_{\langle \nabla f, \mathbf{x}_{k-1} \rangle}$$

but some kind of dual version

$$\underbrace{f_k(\mathbf{x})}_{\langle f_k, \mathbf{x} \rangle} \sim \underbrace{f_{k-1}(\mathbf{x})}_{\langle f_{k-1}, \mathbf{x} \rangle} + \underbrace{(f_k - f_{k-1})}_{a_k} \underbrace{\star}_{\langle f_{k-1}, \nabla \mathbf{x} \rangle}$$

where  $\star$  is a gradient in some functional space.

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

for some simple space  $\mathcal{F}$  so that we define some **weak learner**, e.g. step functions (so called stumps)

## Boosting & Weak Learning

Standard set  $\mathcal{F}$  are stumps functions but one can also consider splines (with non-fixed knots).

One might add a **shrinkage** parameter to learn even more weakly, i.e. set  $\varepsilon_1 = y - \alpha \cdot m_1(\mathbf{x})$  with  $\alpha \in (0, 1)$ , etc.

## Big Data & Linear Model

Consider some linear model  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$  for all  $i = 1, \dots, n$ .

Assume that  $\varepsilon_i$  are i.i.d. with  $\mathbb{E}(\varepsilon) = 0$  (and finite variance). Write

$$\underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}, n \times 1} = \underbrace{\begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,p} \end{pmatrix}}_{\mathbf{X}, n \times (p+1)} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}}_{\boldsymbol{\beta}, (p+1) \times 1} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\boldsymbol{\varepsilon}, n \times 1}.$$

Assuming  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I})$ , the maximum likelihood estimator of  $\boldsymbol{\beta}$  is

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}^\top \boldsymbol{\beta}\|_{\ell_2}\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

... under the assumption that  $\mathbf{X}^\top \mathbf{X}$  is a full-rank matrix.

What if  $\mathbf{X}^\top \mathbf{X}$  cannot be inverted? Then  $\hat{\boldsymbol{\beta}} = [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{y}$  does not exist, but  $\hat{\boldsymbol{\beta}}_\lambda = [\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^\top \mathbf{y}$  always exist if  $\lambda > 0$ .

## Ridge Regression & Regularization

The estimator  $\hat{\beta} = [\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^\top \mathbf{y}$  is the **Ridge** estimate obtained as solution of

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n [y_i - \beta_0 - \mathbf{x}_i^\top \beta]^2 + \lambda \underbrace{\|\beta\|_{\ell_2}^2}_{\mathbf{1}^\top \beta^2} \right\}$$

for some tuning parameter  $\lambda$ . One can also write

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_2} \leq s}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2} \}$$

There is a Bayesian interpretation of that regularization, when  $\beta$  has some prior  $\mathcal{N}(\beta_0, \tau \mathbb{I})$ .

## Over-Fitting & Penalization

Solve here, for some norm  $\|\cdot\|$ ,

$$\min \left\{ \sum_{i=1}^n \ell(y_i, \beta_0 + \mathbf{x}^\top \boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\| \right\} = \min \left\{ \text{objective}(\boldsymbol{\beta}) + \text{penalty}(\boldsymbol{\beta}) \right\}.$$

Estimators are **no longer unbiased**, but might have a smaller mse.

Consider some i.i.d. sample  $\{y_1, \dots, y_n\}$  from  $\mathcal{N}(\theta, \sigma^2)$ , and consider some estimator proportional to  $\bar{y}$ , i.e.  $\hat{\theta} = \alpha \bar{y}$ .  $\alpha = 1$  is the maximum likelihood estimator.

Note that

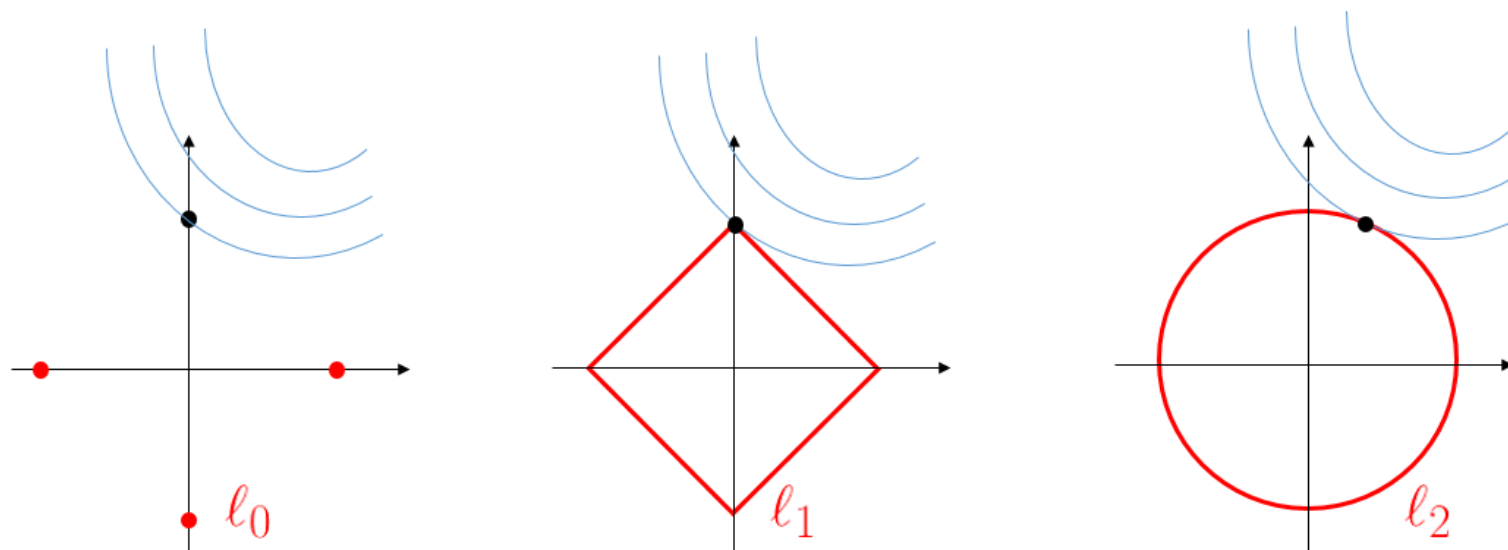
$$\text{mse}[\hat{\theta}] = \underbrace{(\alpha - 1)^2 \mu^2}_{\text{bias}[\hat{\theta}]^2} + \underbrace{\frac{\alpha^2 \sigma^2}{n}}_{\text{Var}[\hat{\theta}]}$$

$$\text{and } \alpha^* = \mu^2 \cdot \left( \mu^2 + \frac{\sigma^2}{n} \right)^{-1} < 1.$$

$$(\hat{\beta}_0, \hat{\beta}) = \operatorname{argmin} \left\{ \sum_{i=1}^n \ell(y_i, \beta_0 + \mathbf{x}^\top \beta) + \lambda \|\beta\| \right\},$$

can be seen as a **Lagrangian** minimization problem

$$(\hat{\beta}_0, \hat{\beta}) = \operatorname{argmin}_{\beta; \|\beta\| \leq s} \left\{ \sum_{i=1}^n \ell(y_i, \beta_0 + \mathbf{x}^\top \beta) \right\}$$



## LASSO & Sparsity

In several applications,  $p$  can be (very) large, but a lot of features are just noise:  $\beta_j = 0$  for many  $j$ 's. Let  $s$  denote the number of **relevant features**, with  $s \ll p$ , cf [Hastie, Tibshirani & Wainwright \(2015\)](#),

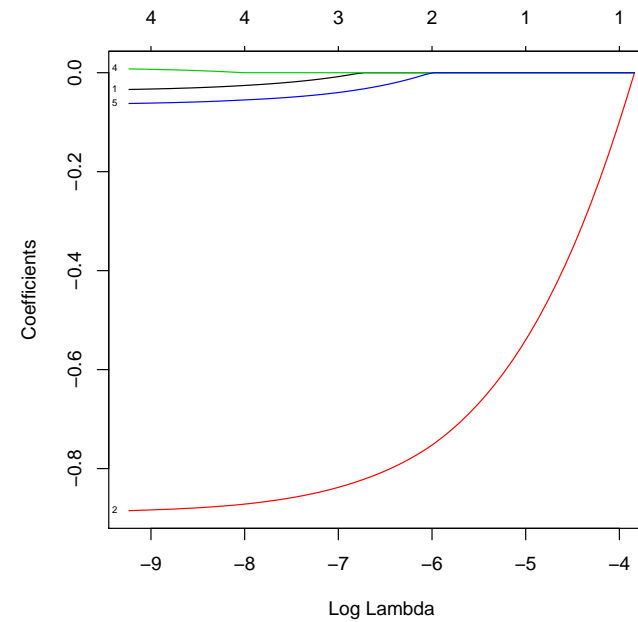
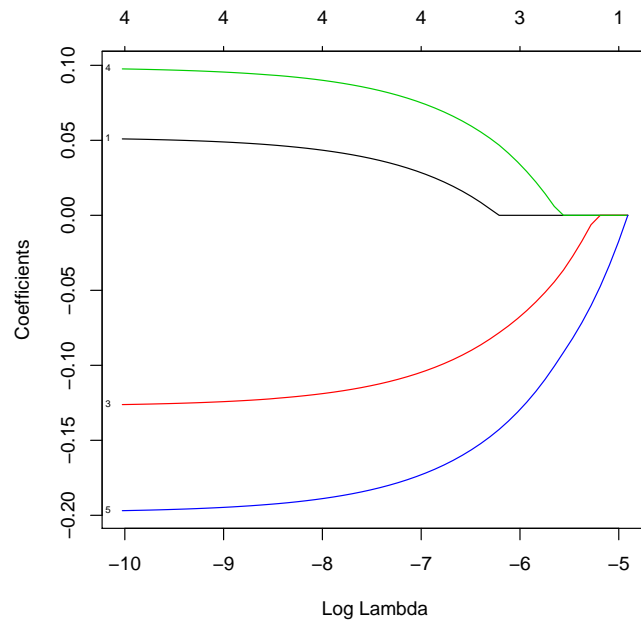
$$s = \text{card}\{\mathcal{S}\} \text{ where } \mathcal{S} = \{j; \beta_j \neq 0\}$$

The true model is now  $y = \mathbf{X}_{\mathcal{S}}^{\top} \boldsymbol{\beta}_{\mathcal{S}} + \varepsilon$ , where  $\mathbf{X}_{\mathcal{S}}^{\top} \mathbf{X}_{\mathcal{S}}$  is a full rank matrix.



## LASSO & Sparsity

Evolution of  $\hat{\beta}_\lambda$  as a function of  $\log \lambda$  in various applications



## In-Sample & Out-Sample

Write  $\hat{\beta} = \hat{\beta}((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ . Then (for the linear model)

$$\text{Deviance}_{\text{IS}}(\hat{\beta}) = \sum_{i=1}^n [y_i - \mathbf{x}_i^{\top} \hat{\beta}((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))]^2$$

With this “in-sample” deviance, we cannot use the central limit theorem

$$\frac{\text{Deviance}_{\text{IS}}(\hat{\beta})}{n} \not\rightarrow \mathbb{E} \left( [Y - \mathbf{X}^{\top} \beta] \right)$$

Hence, we can compute some “out-of-sample” deviance

$$\text{Deviance}_{\text{OS}}(\hat{\beta}) = \sum_{i=n+1}^{m+n} [y_i - \mathbf{x}_i^{\top} \hat{\beta}((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))]^2$$

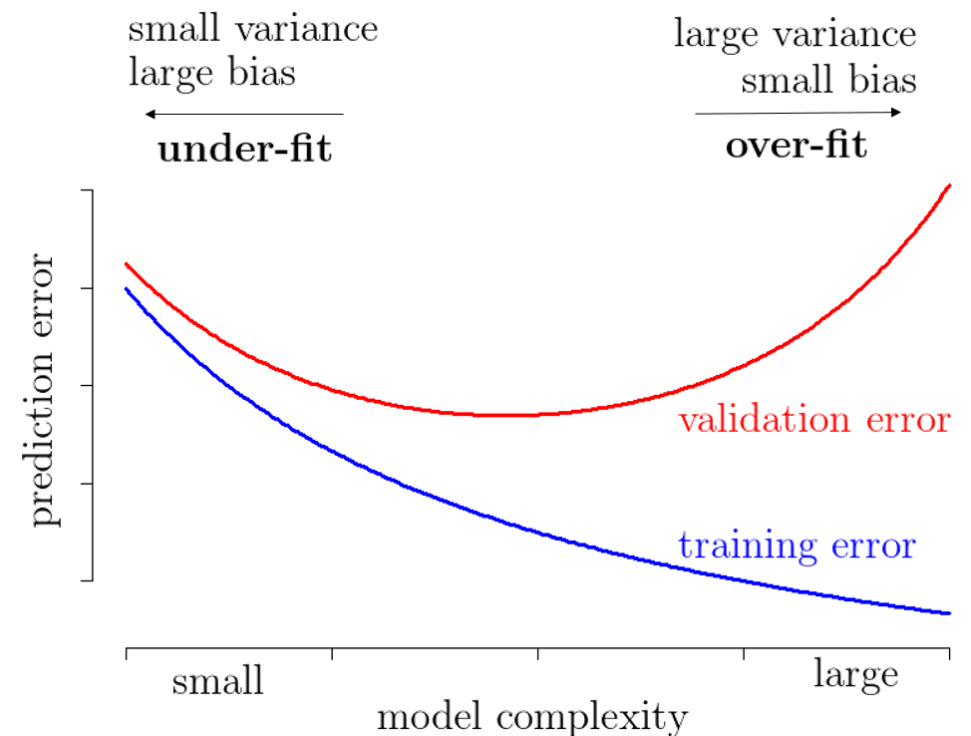
## In-Sample & Out-Sample

Observe that there are connexions with Akaike penalty function

$$\text{Deviance}_{IS}(\hat{\beta}) - \text{Deviance}_{OS}(\hat{\beta}) \approx 2 \cdot \text{degrees of freedom}$$

From [Stone \(1977\)](#), minimizing AIC is closed to cross validation,

From [Shao \(1997\)](#) minimizing BIC is closed to  $k$ -fold cross validation with  $k = n / \log n$ .



## Overfit, Generalization & Model Complexity

Complexity of the model is the degree of the polynomial function

## Cross-Validation

See Jackknife technique [Quenouille \(1956\)](#) or [Tukey \(1958\)](#) to reduce the bias.

If  $\{y_1, \dots, y_n\}$  is an i.i.d. sample from  $F_\theta$ , with estimator  $T_n(\mathbf{y}) = T_n(y_1, \dots, y_n)$ , such that  $\mathbb{E}[T_n(\mathbf{Y})] = \theta + O(n^{-1})$ , consider

$$\tilde{T}_n(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n T_{n-1}(\mathbf{y}_{(i)}) \text{ avec } \mathbf{y}_{(i)} = (y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n).$$

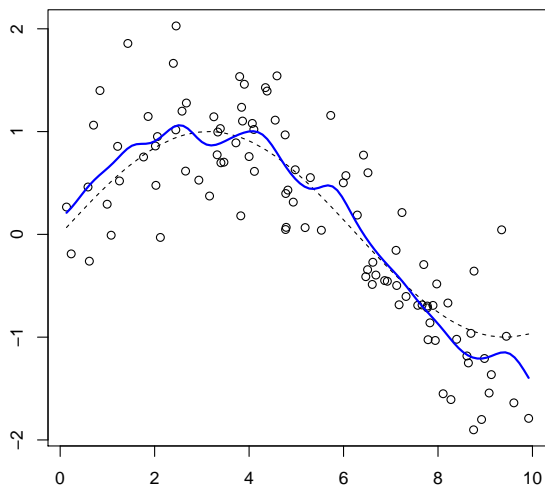
Then  $\mathbb{E}[\tilde{T}_n(\mathbf{Y})] = \theta + O(n^{-2})$ .

Similar idea in [leave-one-out cross validation](#)

$$\text{Risk} = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}_{(i)}(\mathbf{x}_i))$$

## Rule of Thumb vs. Cross Validation

$$\widehat{m}^{[h^*]}(x) = \widehat{\beta}_0^{[x]} + \widehat{\beta}_1^{[x]}x \text{ with } (\widehat{\beta}_0^{[x]}, \widehat{\beta}_1^{[x]}) = \underset{(\beta_0, \beta_1)}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \omega_{h^*}^{[x]} [y_i - (\beta_0 + \beta_1 x_i)]^2 \right\}$$



$$\text{set } h^* = \underset{h}{\operatorname{argmin}} \{ \text{mse}(h) \} \text{ with } \text{mse}(h) = \frac{1}{n} \sum_{i=1}^n [y_i - \widehat{m}_{(i)}^{[h]}(x_i)]^2$$

## Exponential Smoothing for Time Series

Consider some exponential smoothing filter, on a time series  $(x_t)$ ,  $\hat{y}_{t+1} = \alpha\hat{y}_t + (1 - \alpha)y_t$ , then consider

$$\alpha^* = \operatorname{argmin} \left\{ \sum_{t=2}^T \ell(\hat{y}_t, y_t) \right\},$$

see Hyndman *et al.* (2003).

## Cross-Validation

Consider a partition of  $\{1, \dots, n\}$  in  $k$  groups with the same size,  $\mathcal{I}_1, \dots, \mathcal{I}_k$ , and set  $\mathcal{I}_{\bar{j}} = \{1, \dots, n\} \setminus \mathcal{I}_j$ . Fit  $\hat{m}_{(j)}$  on  $\mathcal{I}_{\bar{j}}$ , and

$$\text{Risk} = \frac{1}{k} \sum_{j=1}^k \text{Risk}_j \text{ where } \text{Risk}_j = \frac{k}{n} \sum_{i \in \mathcal{I}_j} \ell(y_i, \hat{m}_{(j)}(\mathbf{x}_i))$$



## Randomization is too important to be left to chance!

Consider some **bootstrapped** sample,  $\mathcal{I}_b = \{i_{1,b}, \dots, i_{n,b}\}$ , with  $i_{k,b} \in \{1, \dots, n\}$

Set  $n_i = \mathbf{1}_{i \notin \mathcal{I}_1} + \dots + \mathbf{1}_{i \notin \mathcal{I}_B}$ , and fit  $\hat{m}_b$  on  $\mathcal{I}_b$

$$\text{Risk} = \frac{1}{n} \sum_{i=1}^n \frac{1}{n_i} \sum_{b: i \notin \mathcal{I}_b} \ell(y_i, \hat{m}_b(\mathbf{x}_i))$$

Probability that  $i$ th obs. is not selection  $(1 - n^{-1})^n \rightarrow e^{-1} \sim 36.8\%$ ,

see training / validation samples (2/3-1/3).

## Bootstrap

From Efron (1987), generate samples from  $(\Omega, \mathcal{F}, \mathbb{P}_n)$

$$\hat{F}_n(y) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(y_i \leq y) \text{ and } \hat{F}_n(y_i) = \frac{\text{rank}(y_i)}{n}.$$

If  $U \sim \mathcal{U}([0, 1])$ ,  $F^{-1}(U) \sim F$

If  $U \sim \mathcal{U}([0, 1])$ ,  $\hat{F}_n^{-1}(U)$  is uniform

on  $\left\{ \frac{1}{n}, \dots, \frac{n-1}{n}, 1 \right\}$ .

Consider some **bootstraped sample**,

- either  $(y_{i_k}, \mathbf{x}_{i_k}), i_k \in \{1, \dots, n\}$
- or  $(\hat{y}_k + \hat{\varepsilon}_{i_k}, \mathbf{x}_k), i_k \in \{1, \dots, n\}$

## Classification & Logistic Regression

Generalized Linear Model when  $Y$  has a **Bernoulli distribution**,  $y_i \in \{0, 1\}$ ,

$$m(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}] = \frac{e^{\beta_0 + \mathbf{x}^\top \boldsymbol{\beta}}}{1 + e^{\beta_0 + \mathbf{x}^\top \boldsymbol{\beta}}} = H(\beta_0 + \mathbf{x}^\top \boldsymbol{\beta})$$

Estimate  $(\beta_0, \boldsymbol{\beta})$  using maximum likelihood techniques

$$\mathcal{L} = \prod_{i=1}^n \left( \frac{e^{\mathbf{x}_i^\top \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^\top \boldsymbol{\beta}}} \right)^{y_i} \left( \frac{1}{1 + e^{\mathbf{x}_i^\top \boldsymbol{\beta}}} \right)^{1-y_i}$$

$$\text{Deviance} \propto \sum_{i=1}^n \left[ \log(1 + e^{\mathbf{x}_i^\top \boldsymbol{\beta}}) - y_i \mathbf{x}_i^\top \boldsymbol{\beta} \right]$$

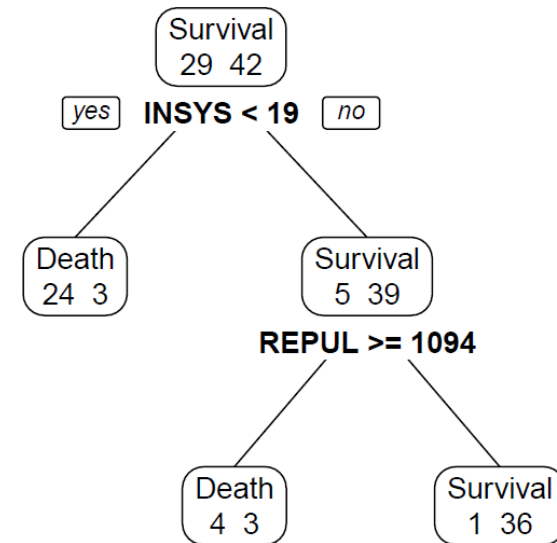
Observe that

$$D_0 \propto \sum_{i=1}^n [y_i \log(\bar{y}) + (1 - y_i) \log(1 - \bar{y})]$$

## Classification Trees

To split  $\{N\}$  into two  $\{N_L, N_R\}$ , consider

$$\mathcal{I}(N_L, N_R) = \sum_{x \in \{L, R\}} \frac{n_x}{n} \mathcal{I}(N_x)$$



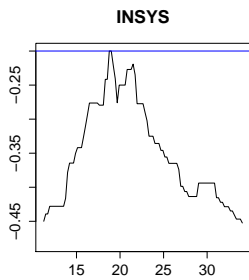
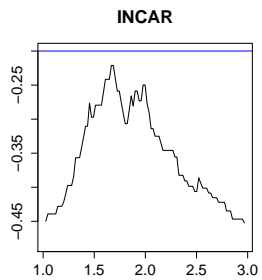
e.g. **Gini index** (used originally in CART, see [Breiman et al. \(1984\)](#))

$$\text{gini}(N_L, N_R) = - \sum_{x \in \{L, R\}} \frac{n_x}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_x} \left( 1 - \frac{n_{x,y}}{n_x} \right)$$

and the **cross-entropy** (used in C4.5 and C5.0)

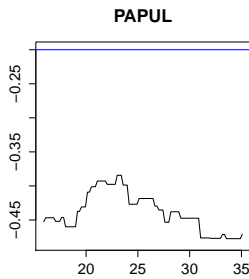
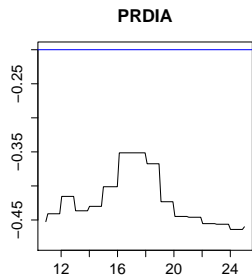
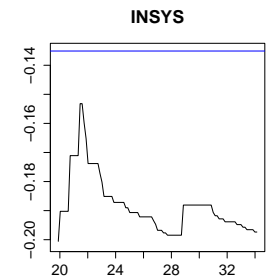
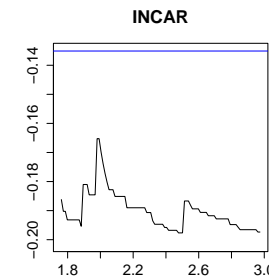
$$\text{entropy}(N_L, N_R) = - \sum_{x \in \{L, R\}} \frac{n_x}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_x} \log \left( \frac{n_{x,y}}{n_x} \right)$$

# Classification Trees

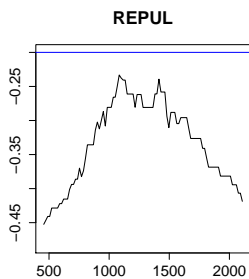
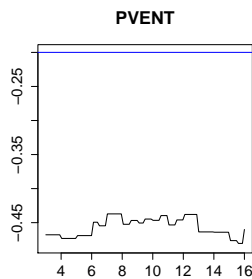
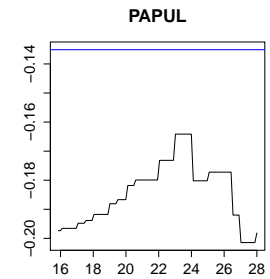
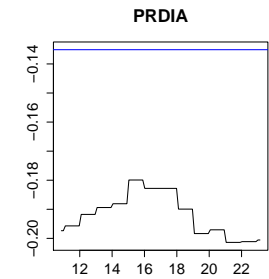


$$N_L: \{x_{i,j} \leq s\} \quad N_R: \{x_{i,j} > s\}$$

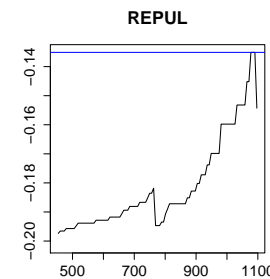
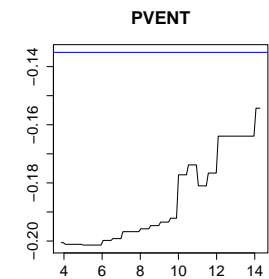
$$\text{solve } \max_{j \in \{1, \dots, k\}, s} \{\mathcal{I}(N_L, N_R)\}$$



← first split



second split →



## Trees & Forests

Bootstrap can be used to define the concept of **margin**,

$$\text{margin}_i = \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i^{(b)} = y_i) - \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i^{(b)} \neq y_i)$$

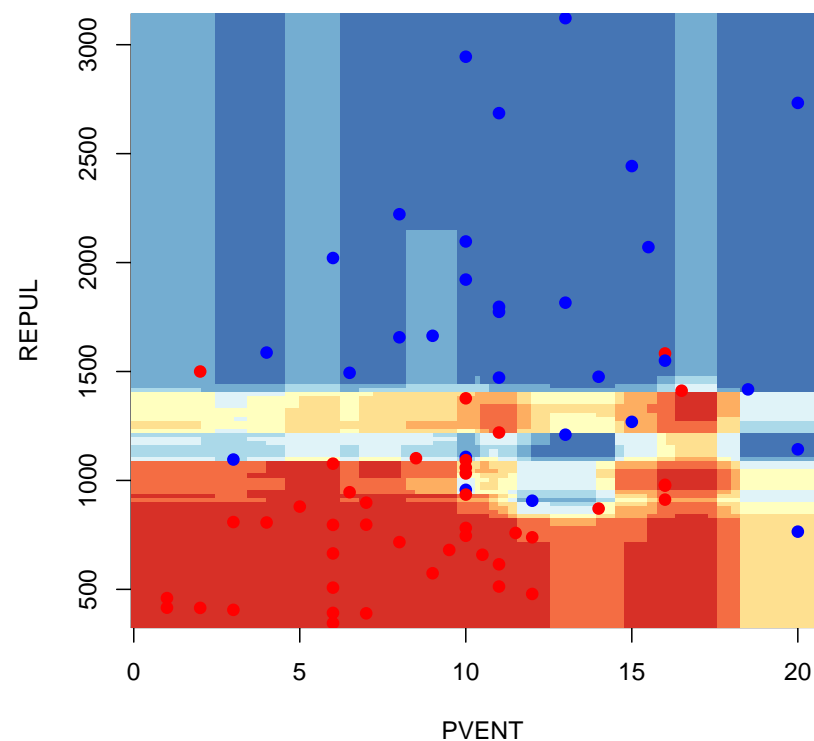
Subsampling of variable, at each knot (e.g.  $\sqrt{k}$  out of  $k$ )

Concept of **variable importance**: given some random forest with  $M$  trees,

$$\text{importance of variable } k \quad I(X_k) = \frac{1}{M} \sum_m \sum_t \frac{N_t}{N} \Delta \mathcal{I}(t)$$

where the first sum is over all trees, and the second one is over all nodes where the split is done based on variable  $X_k$ .

## Trees & Forests



See also discriminant analysis, SVM, neural networks, etc.

## Model Selection & ROC Curves

Given a scoring function  $m(\cdot)$ , with  $m(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$ , and a threshold  $s \in (0, 1)$ , set

$$\widehat{Y}^{(s)} = \mathbf{1}[m(\mathbf{x}) > s] = \begin{cases} 1 & \text{if } m(\mathbf{x}) > s \\ 0 & \text{if } m(\mathbf{x}) \leq s \end{cases}$$

Define the confusion matrix as  $\mathbf{N} = [N_{u,v}]$

$$N_{u,v}^{(s)} = \sum_{i=1}^n \mathbf{1}(\widehat{y}_i^{(s)} = u, y_i = v) \text{ for } (u, v) \in \{0, 1\}.$$

	$Y = 0$	$Y = 1$	
$\widehat{Y}_s = 0$	$\text{TN}_s$	$\text{FN}_s$	$\text{TN}_s + \text{FN}_s$
$\widehat{Y}_s = 1$	$\text{FP}_s$	$\text{TP}_s$	$\text{FP}_s + \text{TP}_s$
	$\text{TN}_s + \text{FP}_s$	$\text{FN}_s + \text{TP}_s$	$n$



## Model Selection & ROC Curves

ROC curve is

$$\text{ROC}_s = \left( \frac{\text{FP}_s}{\text{FP}_s + \text{TN}_s}, \frac{\text{TP}_s}{\text{TP}_s + \text{FN}_s} \right) \text{ with } s \in (0, 1)$$

## Model Selection & ROC Curves

In machine learning, the most popular measure is  $\kappa$ , see [Landis & Koch \(1977\)](#).

Define  $\mathbf{N}^\perp$  from  $\mathbf{N}$  as in the chi-square independence test. Set

$$\text{total accuracy} = \frac{\text{TP} + \text{TN}}{n}$$

$$\text{random accuracy} = \frac{\text{TP}^\perp + \text{TN}^\perp}{n} = \frac{[\text{TN} + \text{FP}] \cdot [\text{TP} + \text{FN}] + [\text{TP} + \text{FP}] \cdot [\text{TN} + \text{FN}]}{n^2}$$

and

$$\kappa = \frac{\text{total accuracy} - \text{random accuracy}}{1 - \text{random accuracy}}.$$

See Kaggle competitions.

## Reducing Dimension with PCA

Use **principal components** to reduce dimension (on centered and scaled variables): we want  $d$  vectors  $z_1, \dots, z_d$  such that

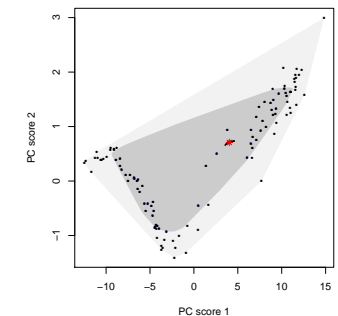
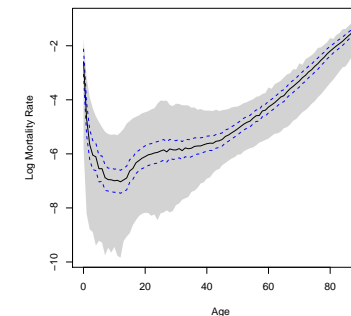
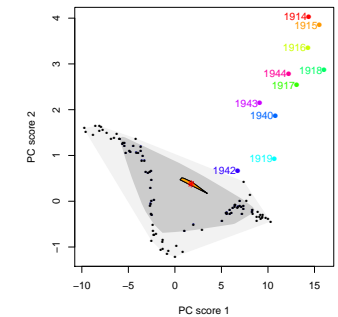
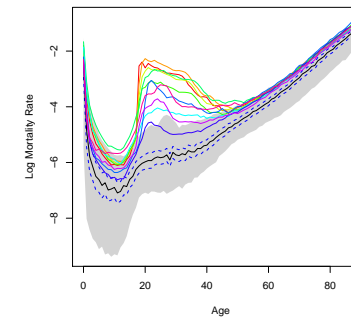
First Component is  $z_1 = X\omega_1$  where

$$\omega_1 = \operatorname{argmax}_{\|\omega\|=1} \left\{ \|\mathbf{X} \cdot \omega\|^2 \right\} = \operatorname{argmax}_{\|\omega\|=1} \left\{ \omega^\top \mathbf{X}^\top \mathbf{X} \omega \right\}$$

Second Component is  $z_2 = X\omega_2$  where

$$\omega_2 = \operatorname{argmax}_{\|\omega\|=1} \left\{ \|\widetilde{\mathbf{X}}^{(1)} \cdot \omega\|^2 \right\}$$

$$\text{with } \widetilde{\mathbf{X}}^{(1)} = \mathbf{X} - \underbrace{\mathbf{X}\omega_1\omega_1^\top}_{z_1}$$



## Reducing Dimension with PCA

A regression on (the  $d$ ) principal components,  $y = \mathbf{z}^\top \mathbf{b} + \boldsymbol{\eta}$  could be an interesting idea, unfortunately, principal components have no reason to be correlated with  $y$ . First component was  $\mathbf{z}_1 = \mathbf{X}\boldsymbol{\omega}_1$  where

$$\boldsymbol{\omega}_1 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \{ \|\mathbf{X} \cdot \boldsymbol{\omega}\|^2 \} = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \{ \boldsymbol{\omega}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\omega} \}$$

It is a non-supervised technique.

Instead, use **partial least squares**, introduced in [Wold \(1966\)](#). First component is  $\mathbf{z}_1 = \mathbf{X}\boldsymbol{\omega}_1$  where

$$\boldsymbol{\omega}_1 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \{ \langle \mathbf{y}, \mathbf{X} \cdot \boldsymbol{\omega} \rangle \} = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \{ \boldsymbol{\omega}^\top \mathbf{X}^\top \mathbf{y} \mathbf{y}^\top \mathbf{X} \boldsymbol{\omega} \}$$

(etc.)

## Instrumental Variables

Consider some instrumental variable model,  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$  such that

$$\mathbb{E}[Y_i | \mathbf{Z}] = \mathbb{E}[\mathbf{X}_i | \mathbf{Z}]^\top \boldsymbol{\beta} + \mathbb{E}[\varepsilon_i | \mathbf{Z}]$$

The estimator of  $\boldsymbol{\beta}$  is

$$\hat{\boldsymbol{\beta}}_{\text{IV}} = [\mathbf{Z}^\top \mathbf{X}]^{-1} \mathbf{Z}^\top \mathbf{y}$$

If  $\dim(\mathbf{Z}) > \dim(\mathbf{X})$  use the Generalized Method of Moments,

$$\hat{\boldsymbol{\beta}}_{\text{GMM}} = [\mathbf{X}^\top \boldsymbol{\Pi}_{\mathbf{Z}} \mathbf{X}]^{-1} \mathbf{X}^\top \boldsymbol{\Pi}_{\mathbf{Z}} \mathbf{y} \text{ with } \boldsymbol{\Pi}_{\mathbf{Z}} = \mathbf{Z}[\mathbf{Z}^\top \mathbf{Z}]^{-1} \mathbf{Z}^\top$$

## Instrumental Variables

Consider a standard two step procedure

- 1) regress columns of  $\mathbf{X}$  on  $\mathbf{Z}$ ,  $\mathbf{X} = \mathbf{Z}\boldsymbol{\alpha} + \boldsymbol{\eta}$ , and derive predictions  $\widehat{\mathbf{X}} = \boldsymbol{\Pi}_{\mathbf{Z}}\mathbf{X}$
- 2) regress  $Y$  on  $\widehat{\mathbf{X}}$ ,  $y_i = \widehat{\mathbf{x}}_i^\top \boldsymbol{\beta} + \varepsilon_i$ , i.e.

$$\widehat{\boldsymbol{\beta}}_{\text{IV}} = [\mathbf{Z}^\top \mathbf{X}]^{-1} \mathbf{Z}^\top \mathbf{y}$$

See [Angrist & Krueger \(1991\)](#) with 3 up to 1530 instruments : 12 instruments seem to contain all necessary information.

Use LASSO to select necessary instruments, see [Belloni, Chernozhukov & Hansen \(2010\)](#)

## Take Away Conclusion

### Big data mythology

- $n \rightarrow \infty$ : 0/1 law, everything is simplified (either true or false)
- $p \rightarrow \infty$ : higher algorithmic complexity, need variable selection tools

### Econometrics vs. Machine Learning

- **probabilistic interpretation** of econometric models  
(unfortunately sometimes misleading, e.g.  $p$ -value)  
can deal with non-i.id data (time series, panel, etc)
- machine learning is about **predictive modeling** and generalization  
**algorithmic** tools, based on **bootstrap** (sampling and sub-sampling),  
**cross-validation**, **variable selection**, **nonlinearities**, **cross effects**, etc

Importance of **visualization** techniques (forgotten in econometrics publications)