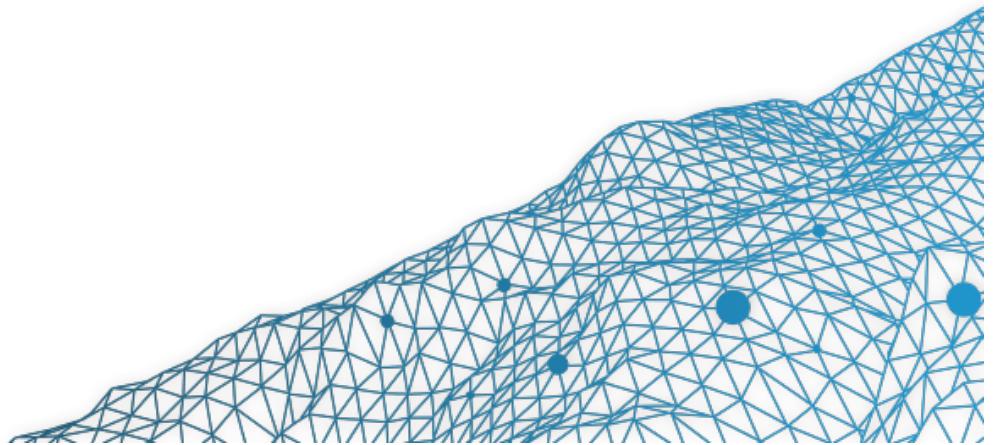


1 Introduction & Geometry of Least Squares

Arthur Charpentier (Université du Québec à Montréal)

Machine Learning & Econometrics

SIDE Summer School - July 2019



Preamble : complex does not mean better...

Econometrics

Machine Learning

Econometrics Courses in a Nutshell

Haavelmo (1944, **The Probabilistic Approach in 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

CHAPTER III

STOCHASTICAL SCHEMES AS A BASIS FOR ECONOMETRICS

As far as is known, the scheme of probability and random variables is, at least for the time being, the only scheme suitable for formulating such theories. We may have objections to using this scheme, but among these objections there is at least one that can be safely dismissed, viz., the objection that the scheme of probability and random variables is not general enough for application to economic data. Since, however, this is apparently not commonly accepted by economists we find ourselves justified in starting our discussion in this chapter with a brief outline of the modern theory of stochastical variables, with particular emphasis on certain points that seem relevant to economics.

The more recent developments in statistical theory are based upon the so-called modernized classical theory of probability. Here “probability” is defined as an absolutely additive and nonnegative *set-function*,¹ satisfying certain formal properties.²

Let us first take an example to illustrate this probability concept.

Data (y_i, x_i) are seen as realizations of (iid) random variables (Y, \mathbf{X}) on some probabilistic space $(\Omega, \mathcal{F}, \mathbb{P})$.

The Linear Regression

Consider some sample $\mathcal{S}_n = \{(y_i, \mathbf{x}_i)\}$ where $y_i \in \mathbb{R}$ and $\mathbf{x}_i \in \mathbb{R}^p$.

\mathcal{S}_n is the realization of n i.i.d. random vectors $(Y_i, X_{1,i}, \dots, X_{p,i})$ with unknown distribution \mathbb{P}

Assume that $Y_i = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_p x_{p,i} + \varepsilon_i$, or $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

where ε_i 's satisfy $\mathbb{E}[\varepsilon_i] = 0$, $\text{Var}[\varepsilon_i] = \sigma^2$, $\text{Var}[\varepsilon_i, \varepsilon_j] = 0$

Least square estimator of (unknown) $\boldsymbol{\beta}$ is

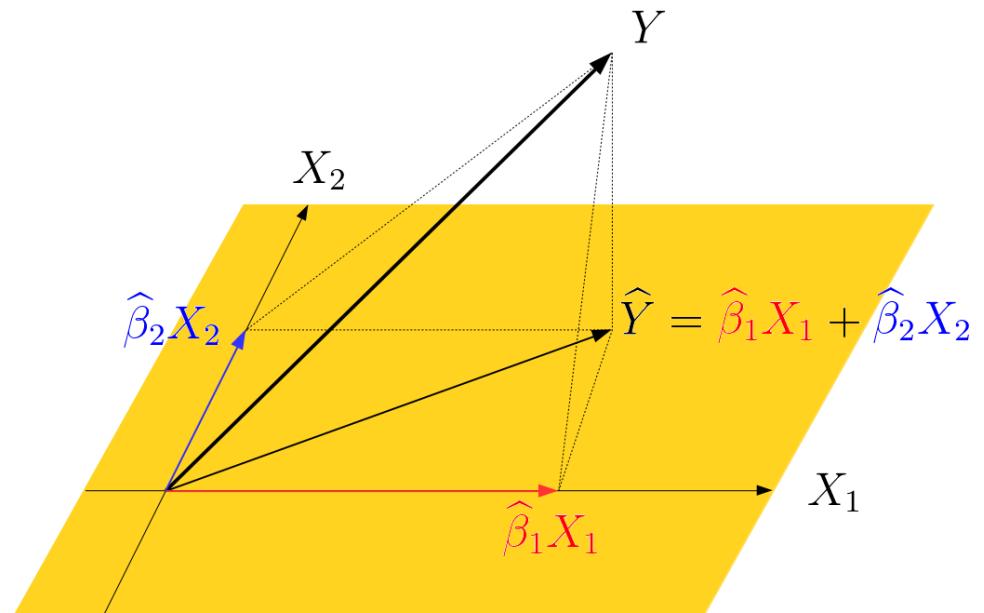
$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \right\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

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 θ between \mathbf{y} and $\Pi_{\mathbf{X}}\mathbf{y}$,

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

see Davidson & MacKinnon (2003)

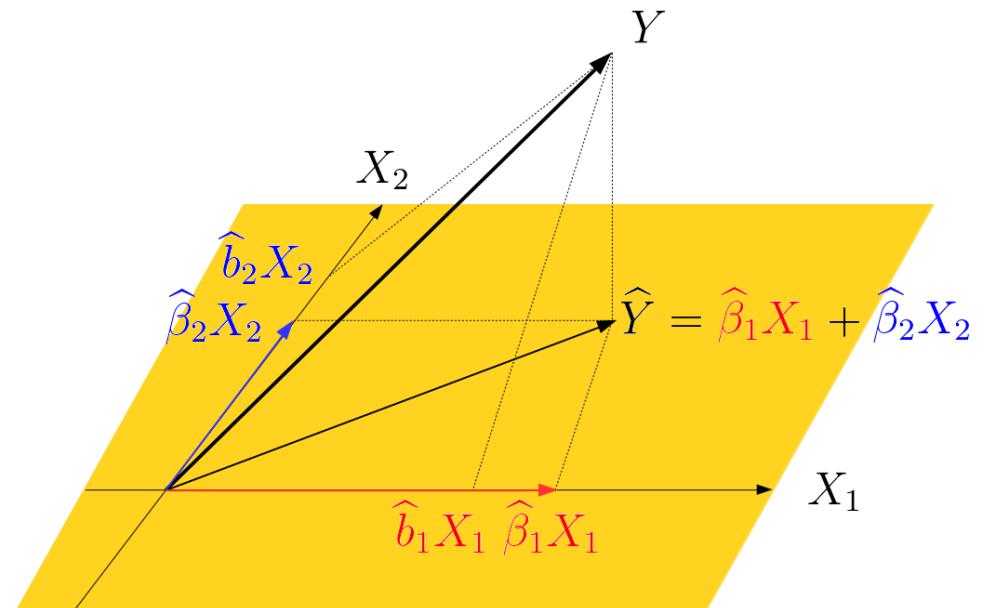
$$\mathbf{y} = \beta_0 + \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \varepsilon$$

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

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

$\mathbf{X}_2^* = \mathbf{X}_2$ if $\mathbf{X}_1 \perp \mathbf{X}_2$,

Frisch-Waugh theorem.



Least Squares ? The probabilistic interpretation

Recall that we assumed that (Y_i, \mathbf{X}_i) were i.i.d. with unknown distribution \mathbb{P} .

In the (standard) linear model, we assume that $(Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{N}(\mathbf{x}^\top \boldsymbol{\beta}, \sigma^2)$.

The maximum-likelihood estimator of $\boldsymbol{\beta}$ is then

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \left\{ - \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\}$$

For the logistic regression, $(Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{B}(\text{logit}^{-1}(\mathbf{x}^\top \boldsymbol{\beta}))$.

The maximum-likelihood estimator of $\boldsymbol{\beta}$ is then

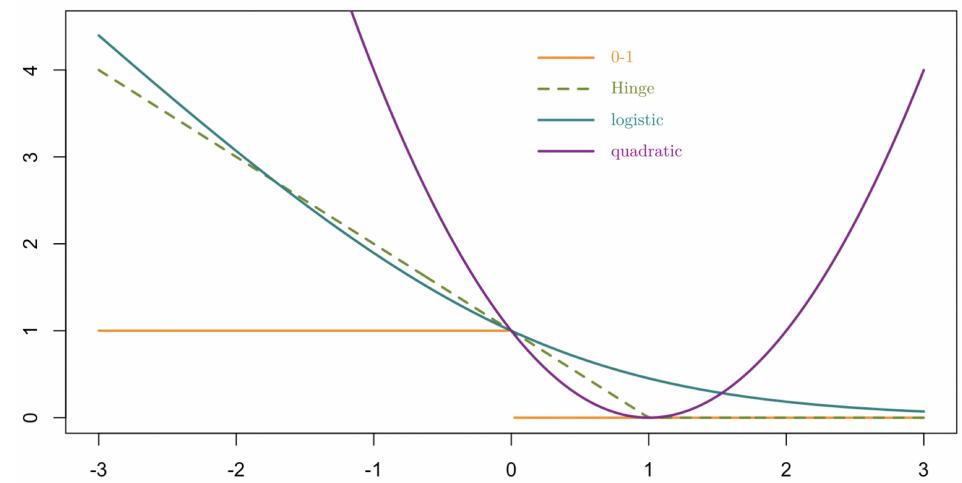
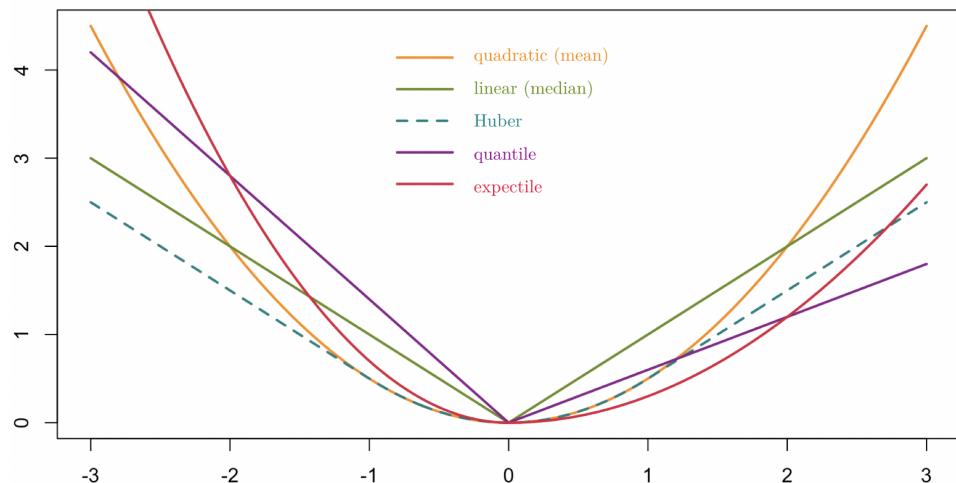
$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i \cdot \log(\text{logit}^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta})) + (1 - y_i) \cdot \log(1 - \text{logit}^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}))) \right\}$$

Least Squares ?

Consider a linear model, $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$.

In a general setting, consider some **loss function** $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$, e.g.

$$\ell(y, m) = (y - m)^2 \text{ (quadratic } \ell_2 \text{ norm)}$$



For **inference**, we try to minimize the empirical risk associated with ℓ ,

$$\hat{\boldsymbol{\beta}} \in \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^{p+1}} \left\{ \sum_{i=1}^n \ell(y_i, \mathbf{x}_i^\top \boldsymbol{\beta}) \right\}$$

where $\mathbf{x}_i = (1, x_{1,i}, x_{2,i}, \dots, x_{p,i})$. But cannot be used for **validation**.

Least Squares and Quadratic Risk

The risk associated with model m is $\mathcal{R}_{\mathbb{P}}(m) = \mathbb{E}_{\mathbb{P}}[\ell(Y, m(\mathbf{X}))]$. Hence

$$\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = \operatorname{argmin} \left\{ \mathcal{R}_{\mathbb{P}}(m) = \mathbb{E}_{\mathbb{P}}[\ell_2(Y, m(\mathbf{X}))] \right\}$$

where $\ell_2(y, m) = (y - m)^2$, so m^* depends on (unknown) \mathbb{P} .

Average Risks

The average risk associated with \hat{m}_n is $\mathcal{R}_{\mathbb{P}}(\hat{m}_n) = \mathbb{E}_{\mathbb{P}}[\ell(Y, \hat{m}_n(\mathbf{X}))]$

Empirical Risks

The empirical risk associated with \hat{m}_n is $\widehat{\mathcal{R}}_n(\hat{m}_n) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}_n(\mathbf{x}_i))$.

If we minimize the average risk, we overfit...

Cross Validation

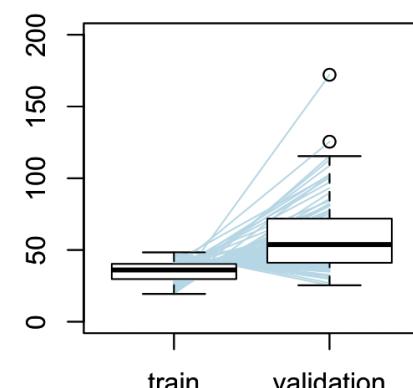
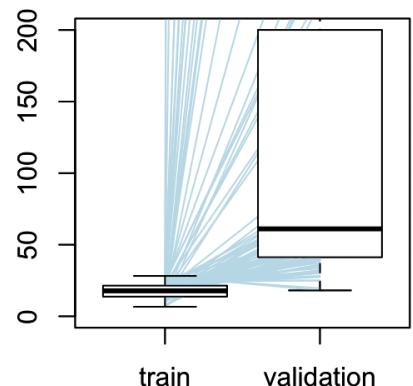
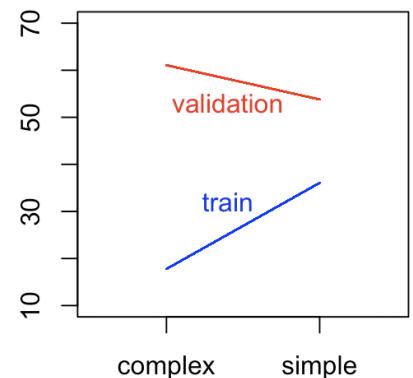
Classical Approach : split the sample \mathcal{S}_n in two parts

Hold-Out Cross Validation

1. Split $\{1, 2, \dots, n\}$ in T (training) and V (validation)
- 2 . Estimate \hat{m} on sample (y_i, \mathbf{x}_i) , $i \in T$: \hat{m}_T
3. Compute $\frac{1}{|V|} \sum_{i \in V} \ell(y_i, \hat{m}_T(x_{1,i}, \dots, X_{p,i}))$

```

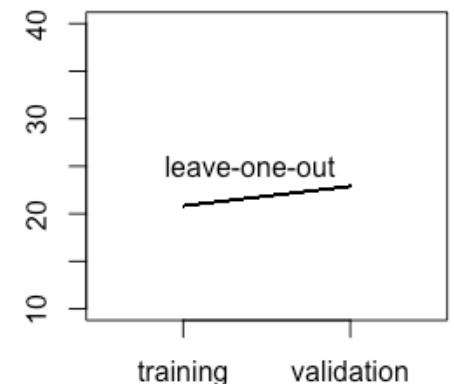
1 chicago <- read.table("http://freakonometrics.free.fr/
  chicago.txt", header=TRUE, sep=";")
2 idx <- sample(1:nrow(chicago), nrow(chicago)*.7)
3 train <- chicago[idx,]
4 valid <- chicago[-idx,]
```



Cross Validation

Leave-one-Out Cross Validation

1. Estimate n models : \hat{m}_{-j} on sample (y_i, \mathbf{x}_i) , $i \in \{1, \dots, n\} \setminus \{j\}$
2. Compute $\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}_{-i}(\mathbf{x}_i))$



Can be computationally intensive...

In the case of a linear regression, there is a simple formula to compute $\hat{\beta}_{-j}$ when observation j is removed. Let $\mathbf{H} = \Pi_x = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

$$(\mathbf{X}_{(j)}^\top \mathbf{X}_{(j)})^{-1} \mathbf{X}_{(j)}^\top = \frac{1}{1 - H_{j,j}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_{(j)}^\top$$

Cross Validation

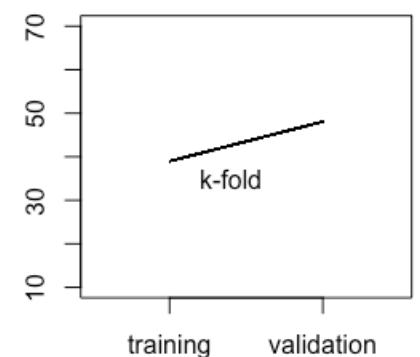
Instead of fitting n models, fit only K

K -Fold Cross Validation

1. Split $\{1, 2, \dots, n\}$ in K groups V_1, \dots, V_K
2. Estimate K models : \hat{m}_k on sample (y_i, \mathbf{x}_i) , $i \in \{1, \dots, n\} \setminus V_k$
3. Compute $\frac{1}{K} \sum_{k=1}^K \frac{1}{|V_k|} \sum_{i \in V_k} \ell(y_i, \hat{m}_k(x_{1,i}, \dots, X_{p,i}))$

If $K = 10$ we fit on 90% of the observation, and validate on the remaining 10%

Here the K groups should be created randomly...



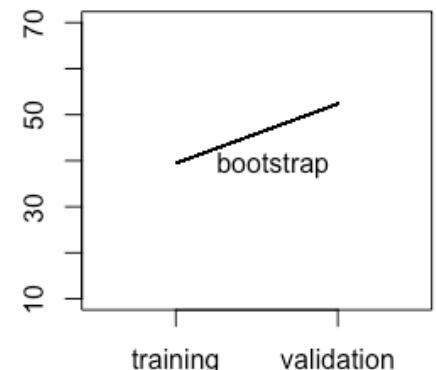
Cross Validation

Bootstrap Cross Validation

1. Generate B bootstrap samples from $\{1, 2, \dots, n\}$, I_1, \dots, I_B
2. Estimate B models : \hat{m}_b on sample (y_i, \mathbf{x}_i) , $i \in I_b$
3. Compute $\frac{1}{B} \sum_{b=1}^B \frac{1}{n - |I_b|} \sum_{i \notin I_b} \ell(y_i, \hat{m}_b(x_{1,i}, \dots, X_{p,i}))$

The probability that $i \notin I_b$ is $\left(1 - \frac{1}{n}\right)^n \sim e^{-1} (= 36.78\%)$

At stage b , we validate on $\sim 36.78\%$ of the dataset.



Norms, Inner Products and Kernels

Most properties are related to the geometry the Euclidean space of \mathbb{R}^n .

Inner Product

An inner product on a vector space \mathcal{H} is the application $(f, g) \mapsto \langle f, g \rangle_{\mathcal{H}}$ (taking value in \mathbb{R}) bilinear, symmetric, definite positive:

- $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$
- $\langle \alpha f + \beta g, h \rangle_{\mathcal{H}} = \alpha \langle f, h \rangle_{\mathcal{H}} + \beta \langle g, h \rangle_{\mathcal{H}}$
- $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ if and only if $f = 0$.

Example : $\mathcal{H} = \mathbb{R}^n$, $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$

Example : $\mathcal{H} = \ell_2 = \left\{ u : \sum_{i=1}^{\infty} u_i^2 < \infty \right\}$, $\langle u, v \rangle = \sum_{i=1}^{\infty} u_i v_i$

Example : $\mathcal{H} = L_2(\mu) = \left\{ f : \int f(x)^2 d\mu(x) < \infty \right\}$, $\langle f, g \rangle = \int f(x)g(x)d\mu(x)$

Norms, Inner Products and Kernels

If \mathcal{H} is finite, $\mathcal{H} = \{h_1, \dots, h_d\}$, $\langle x, y \rangle$ take value $K_{i,j}$ if $x = h_i$ and $y = h_j$. Let $\mathbf{K} = [K_{i,j}]$

\mathbf{K} is a symmetric $d \times d$ matrix, $\mathbf{K} = \mathbf{V}\Lambda\mathbf{V}^\top$ for some orthogonal matrix \mathbf{V} where columns are eigenvectors, and $\Lambda = \text{diag}[\lambda_i]$ (positive values). Let

$$\Phi(x) = (\sqrt{\lambda_1}V_{i,1}, \sqrt{\lambda_2}V_{2,i}, \dots, \sqrt{\lambda_d}V_{d,i}) \text{ if } x = h_i$$

Note that

$$K_{i,j} = [\mathbf{K} = \mathbf{V}\Lambda\mathbf{V}^\top]_{i,j} = \sum_{l=1}^d \lambda_l V_{i,l} V_{l,j} = \langle \Phi(h_i), \Phi(h_j) \rangle$$

Matrix K defines an inner product, it is called a **kernel**. It is symmetric, associated with a positive semi-definite matrix.

Then $K(u, u) \geq 0$ and $K(u, v) \leq \sqrt{K(u, u) \cdot K(v, v)}$.

Norms, Inner Products and Kernels

Let $\varphi : u \mapsto K(\cdot, u)$, then $K(x, y) = \langle \varphi(x), \varphi(y) \rangle$

In a general setting, let $\|f\|_{\mathcal{H}} = \langle f, f \rangle_{\mathcal{H}}^{1/2}$, and define the distance from f to $\mathcal{G} \subset \mathcal{H}$

$$d(f, \mathcal{G}) = \inf_{g \in \mathcal{G}} \{\|f - g\|_{\mathcal{H}}\} = d(f, g^*) \text{ where } g^* \in \mathcal{G}$$

Note that $\langle g, f - g^* \rangle_{\mathcal{H}} = 0, \forall g \in \mathcal{G}$. And $\mathcal{H} = \mathcal{G} \oplus \mathcal{G}^\perp$.

Riesz representation theorem

For any continuous linear functionals L from \mathcal{H} into the field \mathbb{R} , there exists a unique $g_L \in \mathcal{H}$ such that $\forall f \in \mathcal{H}, \langle g_L, f \rangle_{\mathcal{H}} = Lf$.

Consider the case where $\mathcal{H} = \mathbb{R}^n$. Let Σ denote some symmetric $n \times n$ positive definite matrix. Then

$$\langle \mathbf{f}, \mathbf{g} \rangle_{\Sigma} = \mathbf{f}^T \Sigma^{-1} \mathbf{g}$$
 is an inner product on \mathbb{R}^n .

Norms, Inner Products and Kernels

Note that if $\boldsymbol{\sigma}_i$ denote columns of Σ $\langle \boldsymbol{\sigma}_i, \boldsymbol{\sigma}_j \rangle_{\Sigma} = \boldsymbol{\sigma}_i^{\top} \Sigma^{-1} \boldsymbol{\sigma}_j = \Sigma_{i,j}$, and more generally, $\langle \boldsymbol{\sigma}_i, \mathbf{f} \rangle_{\Sigma} = f_i$

The space \mathcal{H} of functions $\mathbb{R}^p \rightarrow \mathbb{R}$ is a Reproducing Kernel Hilbert Space (**RKHS**) if there is an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ such that \mathcal{H} with $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is an Hilbert space, and for all $\mathbf{x} \in \mathbb{R}^p$, linear functional $\delta_{\mathbf{x}} : \mathcal{H} \rightarrow \mathbb{R}$ defined as $\delta_{\mathbf{x}}(f) = f(\mathbf{x})$ is bounded.

Thus, \mathcal{H} is a RKHS if and only if $\forall f \in \mathcal{H}$ and $\mathbf{x} \in \mathbb{R}^p$, there exists $M_{\mathbf{x}}$ such that $|f(\mathbf{x})| \leq M_{\mathbf{x}} \cdot \|f\|_{\mathcal{H}}$.

From Riesz theorem, there exists a unique $\zeta_{\mathbf{x}} \in \mathcal{H}$ associated with $\delta_{\mathbf{x}}$, i.e.
 $\langle \zeta_{\mathbf{x}}, f \rangle_{\mathcal{H}} = f(\mathbf{x})$

Reproducing Kernel of \mathcal{H}

Function $\mathbf{x} \mapsto \zeta_{\mathbf{x}}$ is called reproducing function in \mathbf{x} and $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ defined as $K(\mathbf{x}, \mathbf{y}) = \zeta_{\mathbf{x}}(\mathbf{y})$ is the reproducing kernel of \mathcal{H} .

Norms, Inner Products and Kernels

Observe that $\langle K(\mathbf{x}, \cdot), K(\mathbf{y}, \cdot) \rangle_{\mathcal{H}} = K(\mathbf{x}, \mathbf{y})$.

The kernel is unique, and is (semi-)definite positive.

If \mathcal{H} is a closed subspace of Hilbert space \mathcal{X} . For any function $f \in \mathcal{X}$,
 $\mathbf{x} \mapsto \langle f, K(\cdot, \mathbf{x}) \rangle_{\mathcal{X}}$ is the projection of f on \mathcal{H} .

Note that conversely, **Moore-Aronszajn's theorem** allows to create a RKHS from a definite positive kernel K .

Norms, Inner Products and Kernels

Mercer's kernel

Let μ denote some measure on \mathbb{R}^p and $\mathcal{H} = L^2(\mathbb{R}^p, \mu)$, define

$$(L_K f)(\mathbf{x}) = \int K(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mu(\mathbf{y})$$

which is a compact bounded linear operator, self-adjoint and positive. Let $\lambda_1 \geq \lambda_2 \geq \dots$ denote eigenvalues of L_k , with (orthonormal) eigenvectors ψ_1, ψ_2, \dots , then

$$K(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^p \lambda_k \psi_k(\mathbf{x}) \psi_k(\mathbf{y}) = \Psi(\mathbf{x})^\top \Psi(\mathbf{y}) = \langle \Psi(\mathbf{x}), \Psi(\mathbf{y}) \rangle_{L^2}$$

where $\Psi(\mathbf{x}) = (\sqrt{\lambda_k} \psi_k(\mathbf{x}))$.

Norms, Inner Products and Kernels

Example : Consider the space \mathcal{H} defined as

$$\mathcal{H}_1 = \{f : [0, 1] \rightarrow \mathbb{R} \text{ continuously differentiable, with } f' \in L^2([0, 1]) \text{ and } f(0) = 0\}$$

\mathcal{H}_1 is an Hilbert space on $[0, 1]$ with inner product

$$\langle f, g \rangle_{\mathcal{H}_1} = \int_0^1 f'(t)g'(t)dt$$

with (definite positive) kernel $K_1(x, y) = \min\{x, y\}$:

$$\begin{aligned} \langle f, K(x, \cdot) \rangle_{\mathcal{H}_1} &= \int_0^1 f'(t) \underbrace{\frac{\partial K_1(t, x)}{\partial x}}_{= \mathbf{1}_{[0,x]}(t)} dt = \int_0^x f'(t)dt = f(x) \end{aligned}$$

Norms, Inner Products and Kernels

Example : Consider the Sobolev space $W^1([0, 1])$ defined as

$$W^1([0, 1]) = \{f : [0, 1] \rightarrow \mathbb{R} \text{ continuously differentiable, with } f' \in L^2([0, 1])\}$$

Observe that $W^1([0, 1]) = \mathcal{H}_0 \oplus \mathcal{H}_1$ where

$$\mathcal{H}_0 = \{f : [0, 1] \rightarrow \mathbb{R} \text{ continuously differentiable, with } f' = 0\}$$

The later is an Hilbert space with kernel $K_0(x, y) = 1$.

One can consider kernel $K(x, y) = K_0(x, y) + K_1(x, y)$ (related to linear splines).

More generally, consider

$$\mathcal{H}_2 = \{f : [0, 1] \rightarrow \mathbb{R} \text{ twice cont. diff., with } f'' \in L^2([0, 1] \text{ with } f'(0) = 0)\}$$

Then $\langle f, g \rangle_{\mathcal{H}_2} = \int_0^1 f''(t)g''(t)dt$ is an inner product, with kernel

$$K_2(x, y) = \int_0^1 (x - t)_+(y - t)_+ dt$$

Norms, Inner Products and Kernels

Consider some Hilbert space \mathcal{H} with kernel K and some functional $\Psi : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$, increasing in its last argument.

Representation theorem

Given $\mathbf{x}_1, \dots, \mathbf{x}_n$, $\min_{f \in \mathcal{H}} \{\Psi(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \|f\|_{\mathcal{H}})\}$ admits solution

$$\forall \mathbf{x}, \quad f(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x})$$

A classical expression for Ψ is, for some convex function ψ ,

$$\Psi(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \|f\|_{\mathcal{H}}) = \psi(\mathbf{y}, f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) + \lambda \|f\|_{\mathcal{H}}$$

$$\Psi(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \|f\|_{\mathcal{H}}) = \sum_{i=1}^n \ell(\mathbf{y}, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}}$$

Norms, Inner Products and Kernels

Assume that $y_i = m(x_i) + \varepsilon_i$, where $m \in W_2([0, 1])$, then polynomial splines of degree 2 is the solution of

$$\min_{m \in W_2} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - m(x_i))^2 + \nu \int_0^1 [m''(t)]^2 dt \right\}$$

then $m^*(x) = \beta_0 + \beta_1 x + \sum_{i=1}^n \gamma_i K_2(x_i, x)$ Note that one can use a matrix representation

$$\min \left\{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Q}\boldsymbol{\gamma})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Q}\boldsymbol{\gamma}) + n\nu\boldsymbol{\gamma}^\top \mathbf{Q}\boldsymbol{\gamma} \right\}$$

where $\mathbf{Q} = [K_1(x_i, x_j)]$. If $\mathbf{M} = \mathbf{Q} + n\nu\mathbb{I}$,

$$\boldsymbol{\beta}^* = (\mathbf{X}^\top \mathbf{M}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}^{-1} \mathbf{y} \text{ and } \boldsymbol{\gamma}^* = \mathbf{M}^{-1} (\mathbb{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{M}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}^{-1}) \mathbf{y}$$

Norms, Inner Products and Kernels

Kimeldorf & Wahba's representation theorem

Consider a kernel K and \mathcal{H}_K the associated RKHS. For any (convex) loss function $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$, the solution

$$m^* \in \operatorname{argmin}_{m \in \mathcal{H}_K} \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) + \|m\|_{\mathcal{H}_K}^2$$

can be expressed

$$m^*(\cdot) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \cdot)$$

See Kimeldorf & Wahba (1971, [Some results on Tchebycheffian spline functions](#)), some [slides](#) and later on, on SVM.

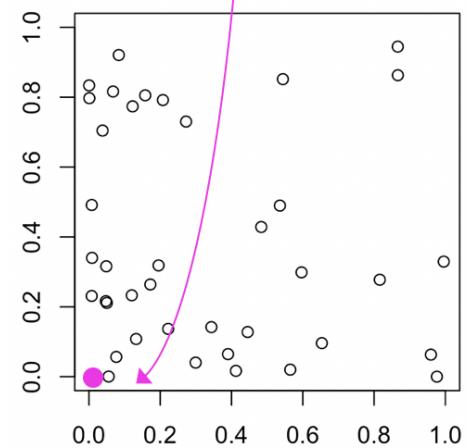
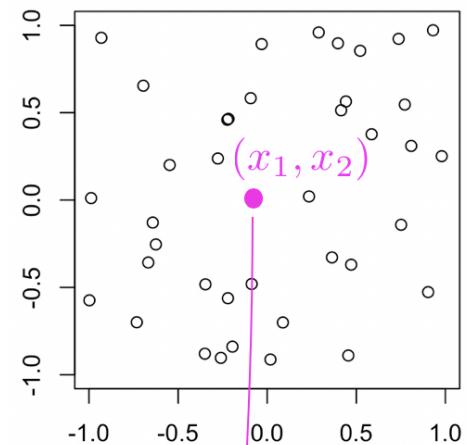
What's going on, here ?

For more technical (mathematical) results, see
Wahba (1990, [Spline Models for Observational Data](#))

We've seen that in many cases, $K(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^\top \varphi(\mathbf{y})$
for some $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^q$ (here $q = p$)

Consider for example $\varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}$

so that $K(\mathbf{x}, \mathbf{y}) = x_1^2 y_1^2 + x_2^2 y_2^2$

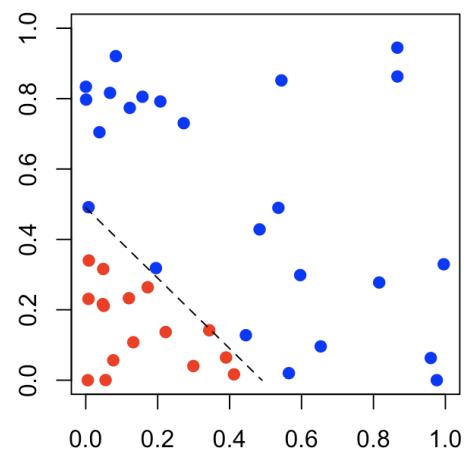
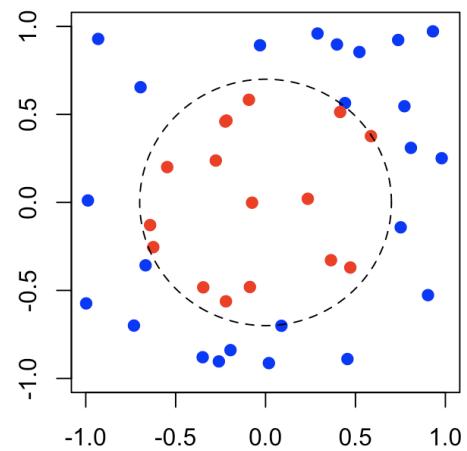


What's going on, here ?

Consider $\varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}$

$$K(\mathbf{x}, \mathbf{y}) = x_1^2 y_1^2 + x_2^2 y_2^2$$

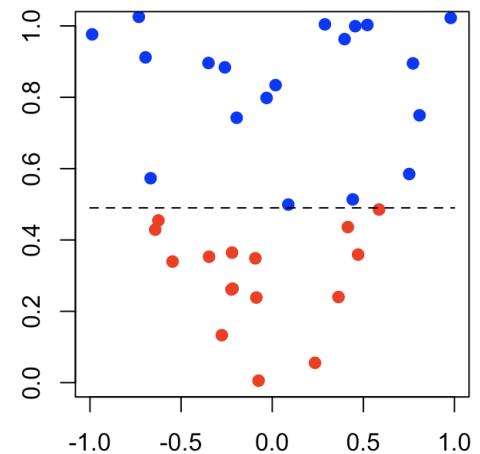
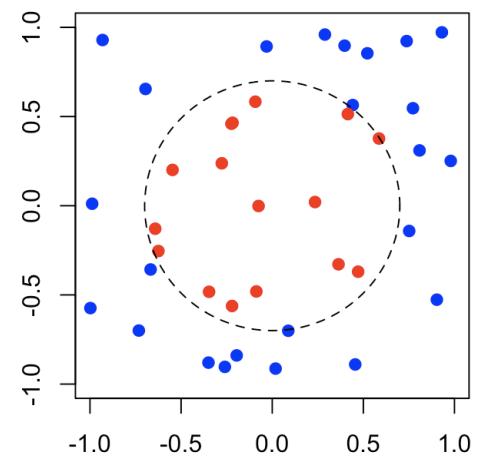
From data (y_i, \mathbf{x}_i) , transform the covariates into $(y_i, \phi(\mathbf{x}_i))$, and use a (classical) linear model



What's going on, here ?

$$\text{Consider } \varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1 \\ x_1^2 + x_2^2 \end{pmatrix}$$

$$K(\mathbf{x}, \mathbf{y}) = x_1 y_1 + (x_1^2 + x_2^2)(y_1^2 + y_2^2)$$

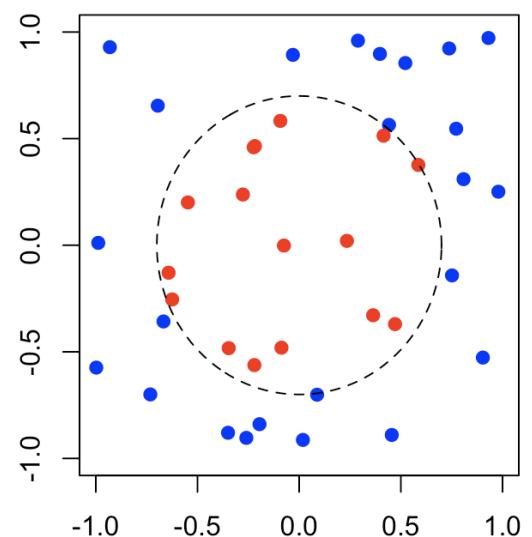


What's going on, here ?

But we can have $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^q$ (with $q \neq p$)

$$\text{Consider } \varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1 \\ x_2 \\ x_1 \cdot x_2 \end{pmatrix}$$

A classical idea with SVM will be to consider $q > p$
 to be able to find a linear separator of points
 red and blue



What's going on, here ?

$$\text{Consider } \varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}$$

$K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^\top \mathbf{y})^2$ is called **polynomial kernel** (of order 2). More generally

$$K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^\top \mathbf{y})^d, \text{ with } d \in \mathbb{N}.$$

For any degree $d \geq 2$, $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^q$ with $q = \binom{p+d}{d}$ (default in R, $d = 3$)

How to get those kernels ?

If K_1 and K_2 are two kernels, so are

$$K(\mathbf{x}, \mathbf{y}) = a_1 K_1(\mathbf{x}, \mathbf{y}) + a_2 K_2(\mathbf{x}, \mathbf{y}) \text{ and } K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) \cdot K_2(\mathbf{x}, \mathbf{y})$$

If h is some $\mathbb{R}^n \rightarrow \mathbb{R}^n$ function, $K(\mathbf{x}, \mathbf{y}) = K_1(h(\mathbf{x}), h(\mathbf{y}))$

If h is some $\mathbb{R}^n \rightarrow \mathbb{R}$ function, $K(\mathbf{x}, \mathbf{y}) = h(\mathbf{x}) \cdot h(\mathbf{y})$

If P is a polynomial with positive coefficients, $k(\mathbf{x}, \mathbf{y}) = P(K_1(\mathbf{x}, \mathbf{y}))$ as well as
 $K(\mathbf{x}, \mathbf{y}) = \exp[k_1(\mathbf{x}, \mathbf{y})]$