

2 | Contrast

From now on the true parameter is denoted θ_0

it is unknown
from the
statistician

Remember last time,

$$\begin{aligned} \text{MLE} \quad \hat{\theta} &= \underset{\theta \in \mathcal{U}}{\operatorname{argmax}} \quad l_{\theta}(x) \\ &= \underset{\theta \in \mathcal{U}}{\operatorname{argmin}} \quad -l_{\theta}(x) \end{aligned}$$

In the gaussian linear models

$$\hat{m} = \underset{m \in V}{\operatorname{argmin}} \quad \|Y - m\|^2$$

$$\mathcal{N}(m, \sigma^2)$$

In general a contrast is a function $C(\theta, X)$ where $\theta \in \Theta$ and X is the observation and such that

$$E_{\theta_0}(C(\theta, X)) \text{ is minimal in } \theta = \theta_0$$

Then the estimator defined by the minimum of the contrast is

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} C(\theta, X)$$

a) -log-likelihood

Let us prove that $C(\theta, X) = -l_\theta(X)$ is a contrast

$$E_{\theta_0}(-l_\theta(X)) = E_{\theta_0}(-\log f_\theta(X)) \quad \text{where } f_\theta \text{ is the density of your model with parameter } \theta.$$

$$= \int -\log[f_\theta(x)] f_{\theta_0}(x) \underbrace{dx}$$

↳ is a canonical notation for a $x \in \mathbb{R}^d$
 $dx = dx_1 \dots dx_d$

The Kullback Leibler divergence is defined by

f, g are two densities (or p.d.f.)
In fact

$$K(f, g) = \int \log \left(\frac{f(x)}{g(x)} \right) f(x) dx$$

$$= E_{X \sim f} \left[\log \left(\frac{f(X)}{g(X)} \right) \right]$$

$$= \int \underbrace{\left[\frac{g(x)}{f(x)} - \log \left(\frac{g(x)}{f(x)} \right) - 1 \right]}_{=0} f(x) dx$$

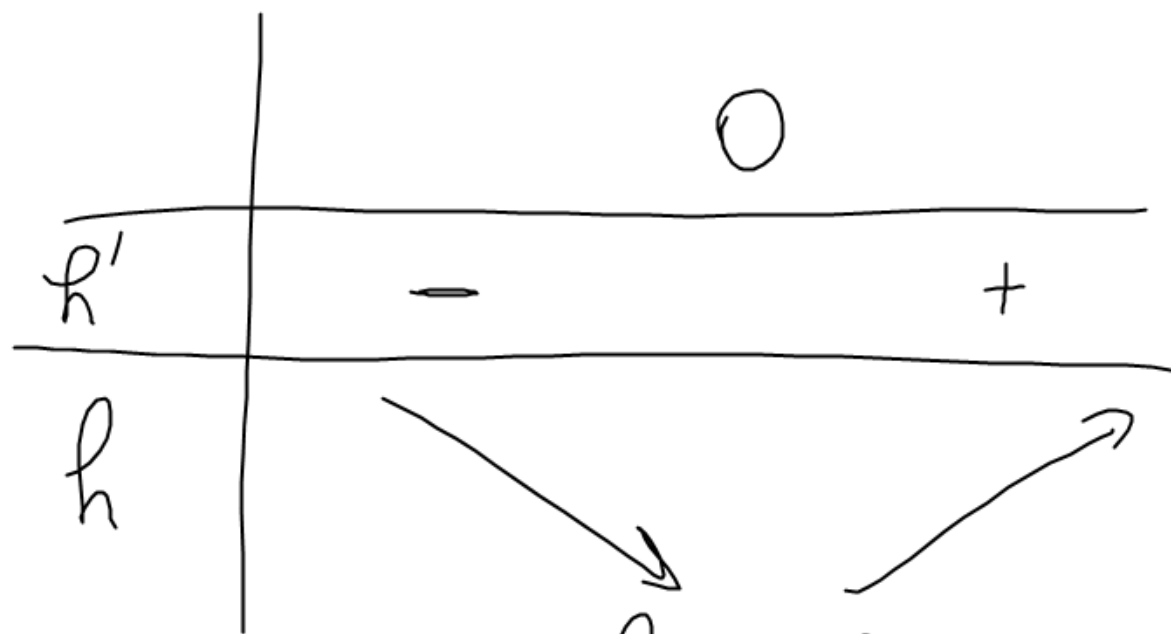
$$= \underbrace{\int \frac{g}{f} f}_{\int g = 1} + \underbrace{\int \left(-\log \frac{g}{f} \right) f}_{-\int f = -1} + \underbrace{\int (-1) f}_{-\int f = -1} = \int \log \frac{f}{g} f = K(f, g)$$

$\int f = 1$
 $\int g = \int \frac{g}{f} f = 1$

$$\frac{g(x)}{f(x)} - \log \frac{g(x)}{f(x)} - 1 = e^u - u - 1 \quad \text{with } u = \log \frac{g(x)}{f(x)}$$

But $u \mapsto e^u - u - 1 = h(u)$

$$h'(u) = e^u - 1$$



$$K(f, g) = \int h\left(\log \left[\frac{g(x)}{f(x)} \right]\right) f(x) dx$$

$$= E_{x \sim f} \left(h\left(\log \frac{g(x)}{f(x)}\right) \right)$$

So $\int K(f, g)$ is always ≥ 0
 $\int K(f, g) = 0$ iff $\forall x, \log \frac{g(x)}{f(x)} = 0$

which means $\frac{g(x)}{f(x)} = 1$ or $g(x) = f(x)$

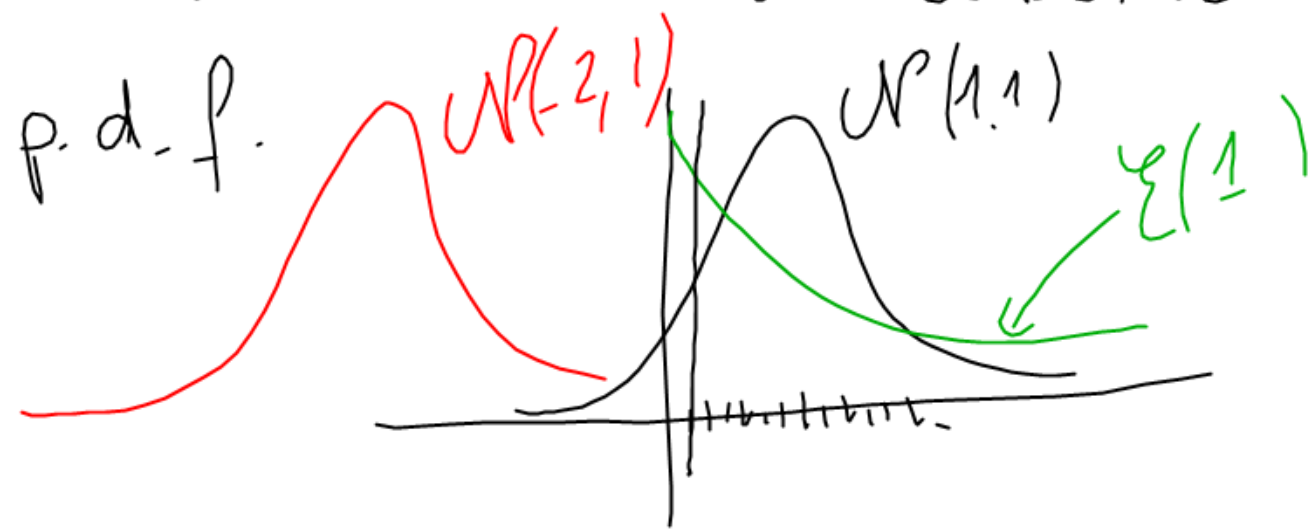
So

$$K(f, g) = E_{x \sim f} \left(\log \frac{f(x)}{g(x)} \right)$$

is always ≥ 0

$K(f, g) = 0$ iff $f(x) = g(x) \forall x$

So the Kullback Leibler divergence measures a distance between densities. The same applies to p.d.f.



Let's go back to $C(\theta, x) = -\log(x)$

$$\begin{aligned}\theta \mapsto \mathbb{E}_{\theta_0}(C(\theta, x)) &= \mathbb{E}_{\theta_0}(-\log f_{\theta}(x)) \\ &= \underbrace{\mathbb{E}_{\theta_0}(\log(f_{\theta_0}(x))) + \mathbb{E}_{\theta_0}(-\log f_{\theta}(x))}_{\mathbb{E}_{\theta_0}(\log \frac{f_{\theta_0}(x)}{f_{\theta}(x)}) = K(f_{\theta_0}, f_{\theta})} - \underbrace{\mathbb{E}_{\theta_0}(\log f_{\theta_0}(x))}_{\text{this doesn't depend on } \theta}\end{aligned}$$

$\int_0 \mathbb{E}_{\theta_0}(-\log(x))$ is minimal when $K(f_{\theta_0}, f_{\theta})$ is null, hence when $f_{\theta} = f_{\theta_0}$ $\forall x$
 \neg there is no problem of identifiability (meaning that 2 different θ 's encode 2 different densities)
 then this implies that $\theta = \theta_0$.

$$X \sim P \begin{cases} \rightarrow \text{disc} & f(x) \text{ } x \in X \text{ discrete} \\ \rightarrow c^0 & f(x) \text{ density.} \end{cases}$$

$$E(h(X)) = \begin{cases} \rightarrow \text{disc} & \sum_{x \in X} h(x) \underbrace{f(x)}_{P(X=x)} \\ \rightarrow c^0 & \int h(x) f(x) dx \end{cases}$$

b) least-square contrast

$$\mathbb{R}^d \ni X = \underbrace{\theta_0}_{\in \mathbb{R}^d} + \varepsilon \quad \varepsilon \text{ noise} \quad \frac{E(\varepsilon) = 0}{E_{\theta_0}(X) = \theta_0} \quad \left(\begin{array}{l} \text{for each} \\ \text{Coordinate} \end{array} \right)$$

$\varepsilon \in V \not\subset \mathbb{R}^d$ (for instance regression)
 $V \subset \mathbb{R}^d$ and $V \neq \mathbb{R}^d$

If we do not specify the dist^o of ε we cannot compute \mathbb{E}
 But you can use the constraint:

$$C(\theta, X) = \|X - \theta\|^2$$

Let us verify that this is a constraint.

$$\begin{aligned} \mathbb{E}_{\theta_0}(C(\theta, X)) &= \mathbb{E}_{\theta_0}(\|X - \theta\|^2) = \mathbb{E}_{\theta_0}(\|X\|^2 - 2\langle \theta, X \rangle + \|\theta\|^2) \\ &= \underbrace{\mathbb{E}_{\theta_0}(\|X\|^2)}_{\text{does not depend on } \theta} - 2\langle \theta, \underbrace{\mathbb{E}_{\theta_0}(X)}_{\theta_0} \rangle + \|\theta\|^2 \end{aligned}$$

2, $\sum \theta, X:$
"

$$E_{\theta_0}(\|X - \theta\|^2) = \underbrace{\cancel{\|X\|^2}}_{E_{\theta_0}(\|X\|^2)} - 2\langle \theta, \theta_0 \rangle + \|\theta\|^2$$

$$= \|\theta - \theta_0\|^2 - \underbrace{\|\theta_0\|^2}_{\text{no } \theta \text{ here}} + (\text{no } \theta \text{ here})$$

$$\underbrace{\|\theta\|^2 - 2\langle \theta, \theta_0 \rangle + \|\theta_0\|^2}_{\text{no } \theta \text{ here either}}$$

So this is minimal when $\theta = \theta_0$ and $\theta \mapsto \|X - \theta\|^2$ is a contrast

(it's known as the
least square contrast
for vectors)

So I can define

$$\hat{\theta} = \underset{\theta \in V}{\operatorname{argmin}} \|X - \theta\|^2$$

$$= \Pi_V X$$

c) least-square contrast for density

$X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$ X_1, \dots, X_n iid with density f_0 .

Let f be a candidate density and let us look at

$$C(f, X) = -\frac{2}{n} \sum_{i=1}^n f(X_i) + \int [f(x)]^2 dx$$

$$\mathbb{E}_{\text{all } X_i \sim f_0} (C(f, X)) = -\frac{2}{n} \sum_{i=1}^n \int f(x) f_0(x) dx + \int [f(x)]^2 dx$$

$$= -2 \int f(x) f_0(x) dx + \int [f(x)]^2 dx$$

$$= \int (f(x) - f_0(x))^2 dx - \int [f_0(x)]^2 dx$$

So $C(f, X)$ is minimal when $\int (f(x) - f_0(x))^2 dx$ is minimal

But $\int (f(x) - f_0(x))^2 dx$ is ≥ 0

and $= 0$ iff $\forall x, f(x) = f_0(x)$.

So $C(f, X)$ is a contrast, called the least-square contrast for density

III Choice of models

Ex (from Neuroscienza)

$$\rightarrow \underbrace{Y_i}_{\text{firing rate}} = f_0(\underbrace{W_i}_{\text{weight}}) + \varepsilon_i$$

$$\varepsilon_i \text{ iid } \mathcal{N}(0, \sigma^2)$$

model 1

$$f(w) = a + b w, \quad a, b \text{ unknown}$$

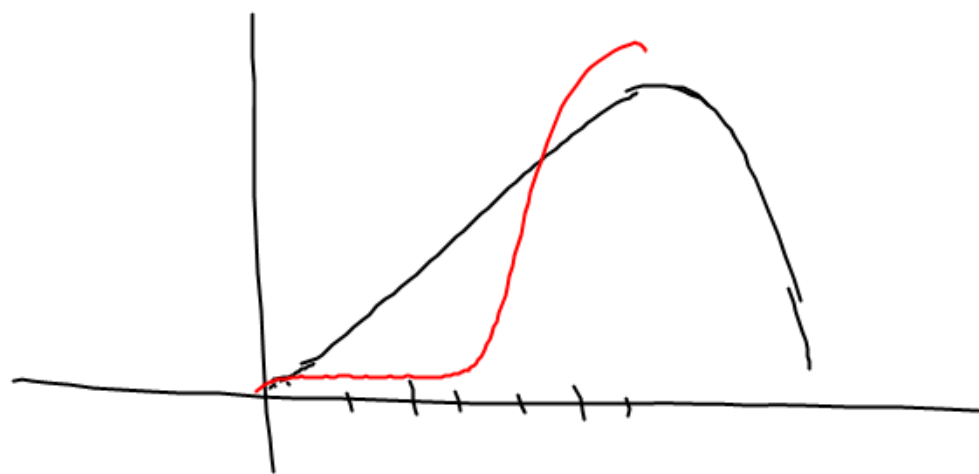
$$f(w) = a + b w + c w^2$$

$a, b, c \text{ unknown}$

\vdots

$$f(w) = a_0 + a_1 w + a_2 w^2 + \dots + a_d w^d$$

(model of dim $d+1$)



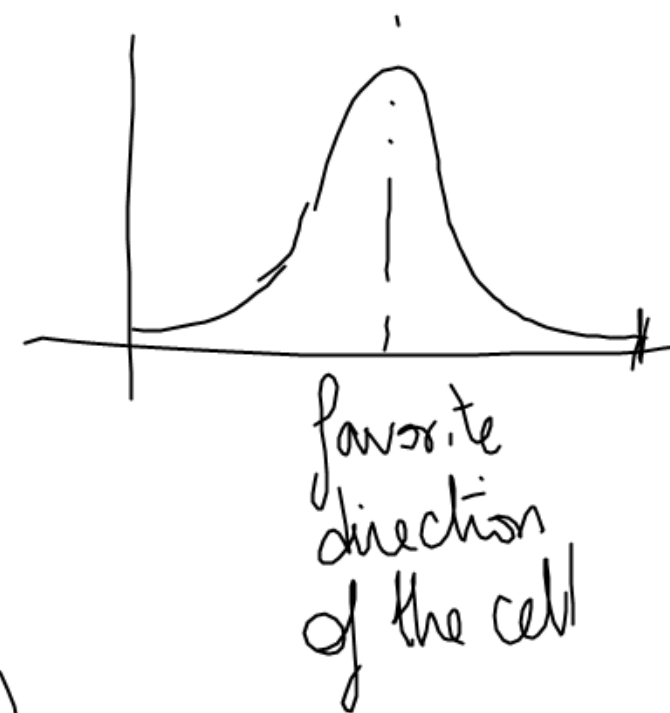
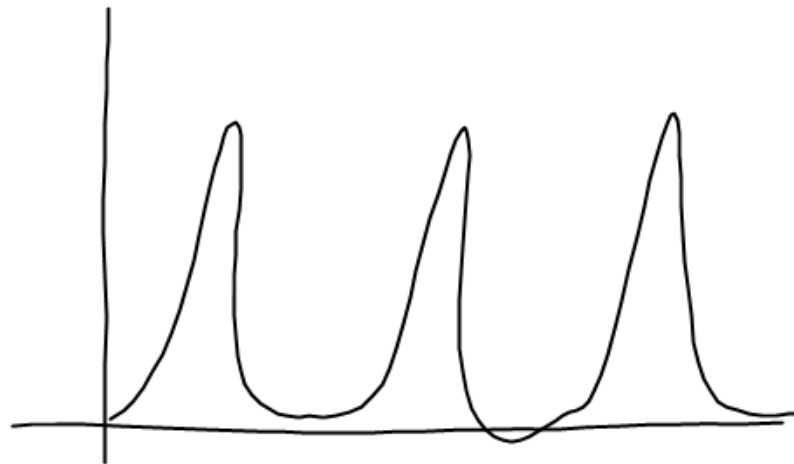
$$\rightarrow Y_i = f_0(U_i) + \varepsilon_i \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

U_i is the angle of the movement

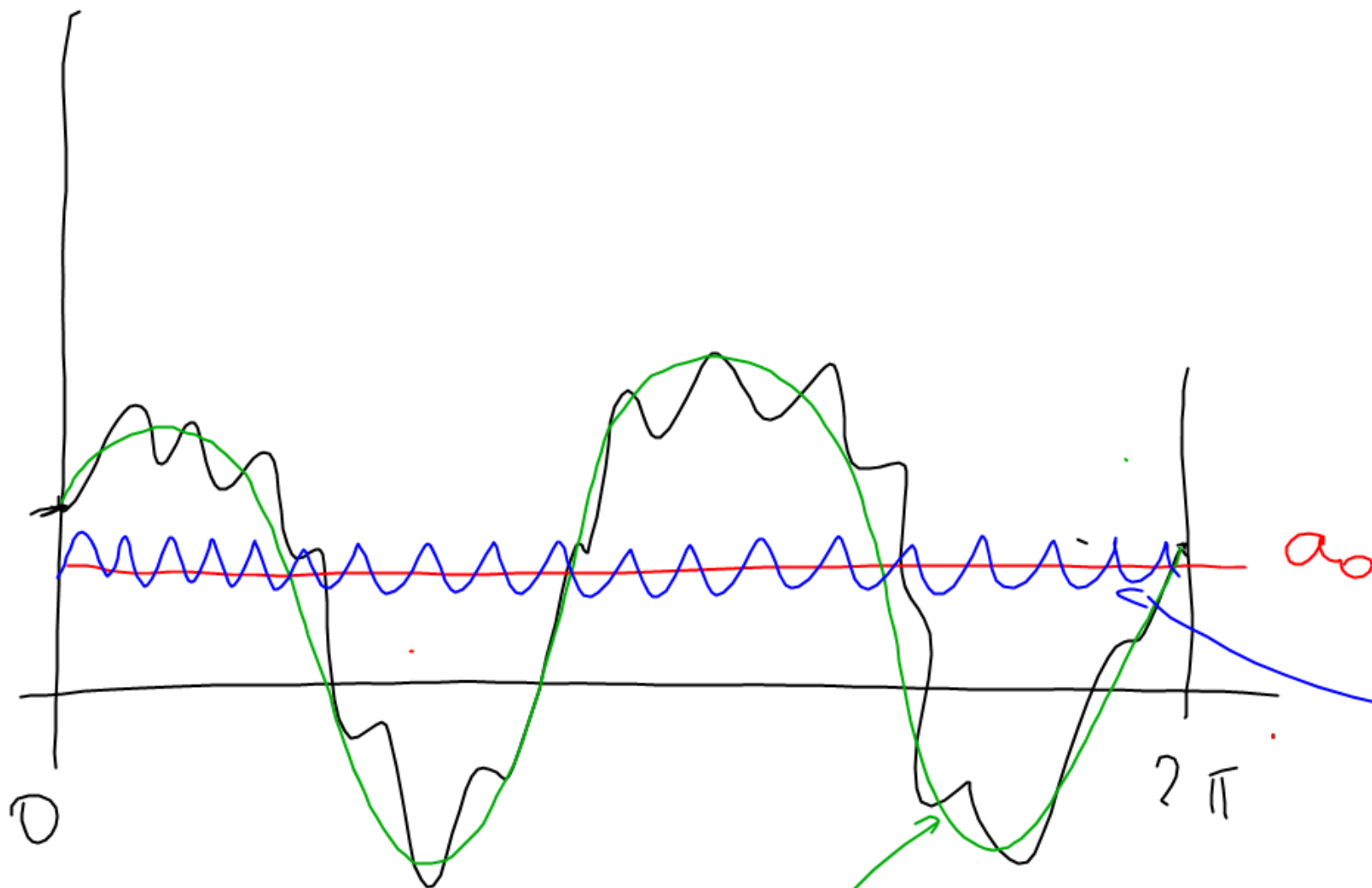
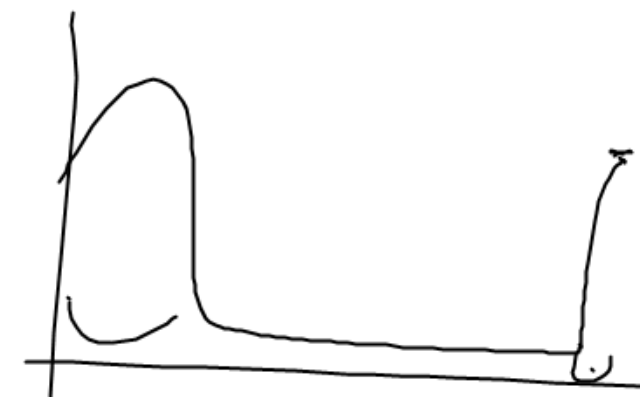
model 1

$$f(U_i) = a + b \cos(U_i)$$

$$U_i \in [0, 2\pi]$$



$$f(U_i) = a_0 + a_1 \cos(U_i) + a_{-1} \sin(U_i) \\ + \dots + a_d \cos(dU_i) + a_{-d} \sin(dU_i)$$



f (smooth periodic)
 $d \rightarrow +\infty$

$$f(u) = a_0 + \sum_{k=1}^d a_k \cos(ku) + a_{-k} \sin(ku)$$

$$a_i = \int \frac{f(u) \cos(iu) du}{2\pi}$$

$$a_{-i} = \int \frac{f(u) \sin(iu) du}{2\pi} \quad a_k \cos(ku) + a_{-k} \sin(ku)$$

higher frequency
 $a_k \cos(ku) + a_{-k} \sin(ku)$

hard = red + green + blue

In general we have a bunch of linearly independent functions

$$\varphi_1(x) \dots \varphi_d(x) \quad i=1 \dots n$$

and you look at the problem $y_i = f_0(x_i) + \varepsilon_i$
 $\varepsilon_i \sim N(0, \sigma^2)$

model of dim d is $f \in \text{Vect}(\varphi_1(x) \dots \varphi_d(x)) = V \subseteq \mathbb{R}^n$

$$f(x_i) = a_1 \varphi_1(x_i) + \dots + a_d \varphi_d(x_i) \quad \forall i.$$

Of course you will stop before n but when?

for this model, you know that $\begin{pmatrix} \hat{f}(x_1) \\ \vdots \\ \hat{f}(x_n) \end{pmatrix} = \Pi_V(y)$

$$V = \text{Vect} \left(\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} / \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \right)$$

$$\begin{aligned} \varphi_1(x) &= 1 \\ \varphi_2(x) &= x \end{aligned}$$

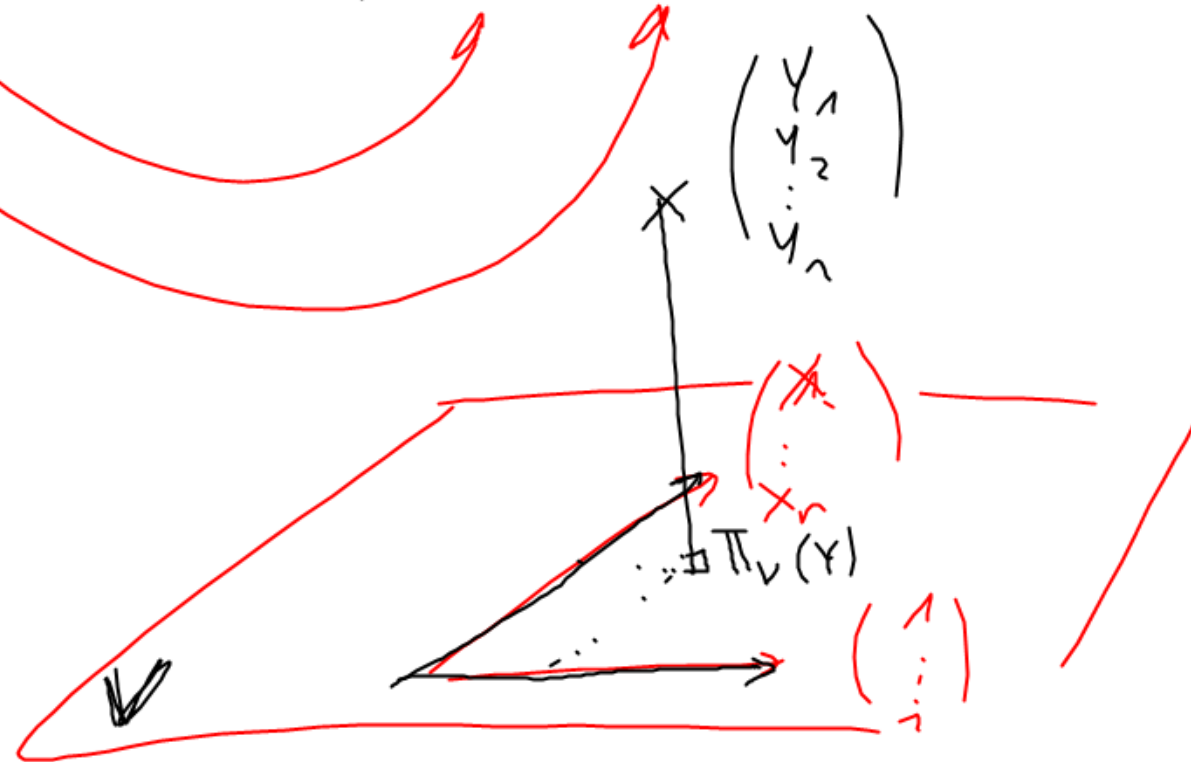
$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$f(x) = \underline{a} + \underline{b}x$$

$$\pi_V Y = \hat{a} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + \hat{b} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

orthogonal projection
of Y on V

$$\pi_V Y \in \mathbb{R}^n$$



1/ Bias - Variance decomposition

$$\text{So } Y_i = f_0(X_i) + \varepsilon_i$$

→ for model of dim d I have $\begin{pmatrix} \hat{f}_d(X_1) \\ \vdots \\ \hat{f}_d(X_n) \end{pmatrix} = \hat{f}_d(X)$

What is the d for which

$$E_{Y \sim f_0} \left(\| f_0(X) - \hat{f}_d(X) \|^2 \right) \text{ is the smallest?}$$

→ this should give me the best d .

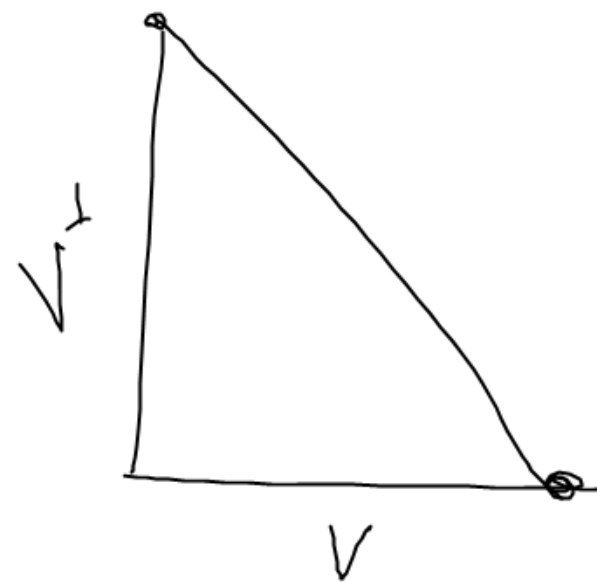
$$E_0 \left(\| f_0(X) - \hat{f}_d(X) \|^2 \right)$$

$$= E_0 \left(\| f_0(X) - \pi_V Y \|^2 \right)$$

$$= E_0 \left(\left\| \underbrace{f_0(X) - \pi_V f_0(X)}_{\perp V} + \underbrace{\pi_V f_0(X)}_{\in V} \right\|^2 \right)$$

$$= E_0 \left(\underbrace{\| f_0(X) - \pi_V f_0(X) \|^2}_{\text{BIAS TERM}} \right) + E_0 \left(\| \pi_V \varepsilon \|^2 \right)$$

BIAS TERM
decreases when $d \uparrow$



$$\underbrace{\sigma^2 d}_{\text{VARIANCE TERM}} \left(\text{since } \| \pi_V \varepsilon \|^2 \sim \sigma^2 d \right)$$

VARIANCE TERM increases with d

To choose a good d you need a trade off between

- complexity of the model (to have a small bias) \longrightarrow depends on f_0
- variability of each of the coefficients
(\rightarrow want a small variance)

\Rightarrow We define an oracle \tilde{d} which is a benchmark

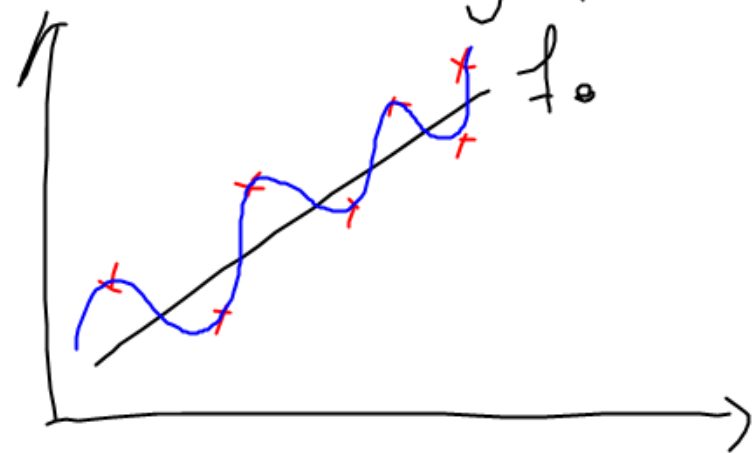
$$\tilde{d} = \underset{\substack{d \in [1..n-1]}}{\operatorname{argmin}} \left(\|f_0(X) - \Pi_{V^d} f_0(X)\|^2 + \sigma^2 d \right)$$

2/ Mallon's C_p

This consists in choosing

$$\hat{d} = \underset{d}{\operatorname{argmin}} \left(\underbrace{\|Y - \Pi_{V_d} Y\|^2}_{\text{least-square}} + \underbrace{2\sigma^2 d}_{\text{penalty}} \right) \quad \left(\text{we assume that } \sigma \text{ is known} \right)$$

If there wasn't any penalty $\rightarrow \hat{\hat{d}} = \underset{d}{\operatorname{argmin}} (\|Y - \Pi_{V_d} Y\|^2)$
 $= \text{largest } d$



\rightarrow you are overfitting the data

A penalty is always there to avoid overfitting
 Mallon's C_p is a particular case which satisfies
 an oracle inequality

$$\mathbb{E} \left(\left\| f_0(x) - \hat{f}_{\hat{d}}(x) \right\|^2 \right) \leq C \underbrace{\min_d \left[\mathbb{E} \left(\left\| f_0(x) - \Pi_{V_d} f_0(x) \right\|^2 \right) + \sigma^2 d \right]}_{\text{Oracle risk}}$$

\uparrow
 the choice
 with Mallon's C_p

$$\hat{d} = \arg \min_d \left(\left\| Y - \Pi_{V_d} Y \right\|^2 + 2d\sigma^2 \right)$$

C is very close to 1.

3) Akaike's criterion (AIC)

Akaike information criterion

model M parametrized by $\Theta_M \rightarrow \text{MLE } \hat{\Theta}_M \rightarrow f_{\hat{\Theta}_M}(x)$
 (eg $\mathcal{NP}(m, \sigma^2)$)
 $(\hat{m}, \hat{\sigma}^2) \rightarrow f_{\hat{m}, \hat{\sigma}^2}(x)$

density of X
when param = MLE

the likelihood
on the model M
at the param $\hat{\Theta}_M$

AIC criterion is

$$\hat{M} = \underset{M \in \mathcal{M}}{\text{argmin}} \left\{ \left(-\log f_{\hat{\Theta}_M}(x) \right) + \underbrace{\dim(M)}_{\text{nbof parameters inside } M} \right\}$$

oracle inequalities exist for that too, as long as there are few models with the same number of parameters

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \text{ iid } \sim g \quad \text{model } \Pi \text{ model } g^M \text{ of } g$$

$$f(X) = g(x_1) \cdots g(x_n)$$

$$\hat{\theta}_M = \underset{\theta \in \Theta_M}{\operatorname{argmax}} \sum_{i=1}^n g_{\theta}^M(x_i)$$

$$\hat{M} = \underset{M}{\operatorname{argmin}} \left(- \sum_{i=1}^n g_{\hat{\theta}_M}^M(x_i) + \underbrace{\dim(M)}_{\substack{\rightarrow \text{CP}(m, \sigma^2) \rightarrow 2 \\ \rightarrow \mathcal{E}(1) \rightarrow 1}} \right)$$

$$M \leftrightarrow \text{CP}(m, \sigma^2) \quad \theta_M = \begin{pmatrix} m \\ \sigma^2 \end{pmatrix}$$

$$\hat{\theta}_M = \begin{pmatrix} \bar{X} \\ \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2 \end{pmatrix}$$

$$M \leftrightarrow \mathcal{E}(1) \quad \theta_M = 1$$

$$\hat{\theta}_M = \frac{1}{\bar{X}}$$

4/ Bic criterion

$$Y_i = a_0 + a_1 X_i^1 + \dots + a_p X_i^p + \varepsilon_i \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

To do variable selection I could put in competition

$$V_0 = \text{Vect} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$V_1 = \text{Vect}(X^1) \quad \dots \quad V_p = \text{Vect}(X^p)$$

$$V_{01} = \text{Vect} \left(\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, X^1 \right)$$

$$V_{ij} = \text{Vect}(X^i, X^j) \leftarrow$$

$$V_{0\dots d} = \text{Vect} \left(\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, X^1, X^2, \dots, X^d \right)$$

number of models
with dim 2
 $\frac{p(p-1)}{2}$

write them all

→ AIC and Mallon's cannot work: too many models with the same dimension

$$BIC$$

$$\hat{M} = \underset{M \in \mathcal{M}}{\operatorname{argmin}} \left[\frac{\ln(n) \dim(M)}{2} - \log(\hat{f}_{\hat{\theta}_M}(x)) \right]$$

→ we need to penalize more...

Be careful

the more models in competition, the larger the penalty, and you may end with a model of small dimension.
 whereas if the nb of models ^{was} not too big, you could have use AIC and select it.

Also there are other penalties

→ Lasso

→ Ridge

→ slope heuristic

→ for people who want more details in math:

Bergé and Massart

"Gaussian model selection"

