

Information geometry

J.-F. Cardoso, C.N.R.S.
Institut d'Astrophysique de Paris

Peyresq 2018

Note

These slides are supporting material for the lecture.

They do not provide a self-contained exposition.

Intro: Geometry in spaces of probability distributions

From Shun-ichi Amari, *Differential-geometrical methods in statistics*.

One may ask why geometry, in particular differential geometry, is useful for statistics. The reason seems very simple and strong. A statistical model is a set of probability distributions. In particular, a parametric model usually forms a finite-dimensional manifold embedded in the set of all possible probability distributions. [...] One often uses a statistical model to carry out statistical inference, assuming that the true distribution is included in the model. However, a model is merely a hypothesis. The true distribution may not be in the model but be only close to it. Therefore, in order to evaluate statistical inference procedures, it is important to know what part the statistical model occupies in the entire set of probability distributions and what shape the statistical model has in the entire set. This is the problem of geometry of statistical models. It is therefore expected that a fundamental role is played in statistics by geometrical quantities such as the distance or divergence between two probability distributions, the flatness or curvature of a statistical model, etc...

Bad news from Wikipedia

Not logged in [Talk](#) [Contributions](#) [Create account](#)

Article

[Talk](#)

Read

[Edit](#)

[View history](#)

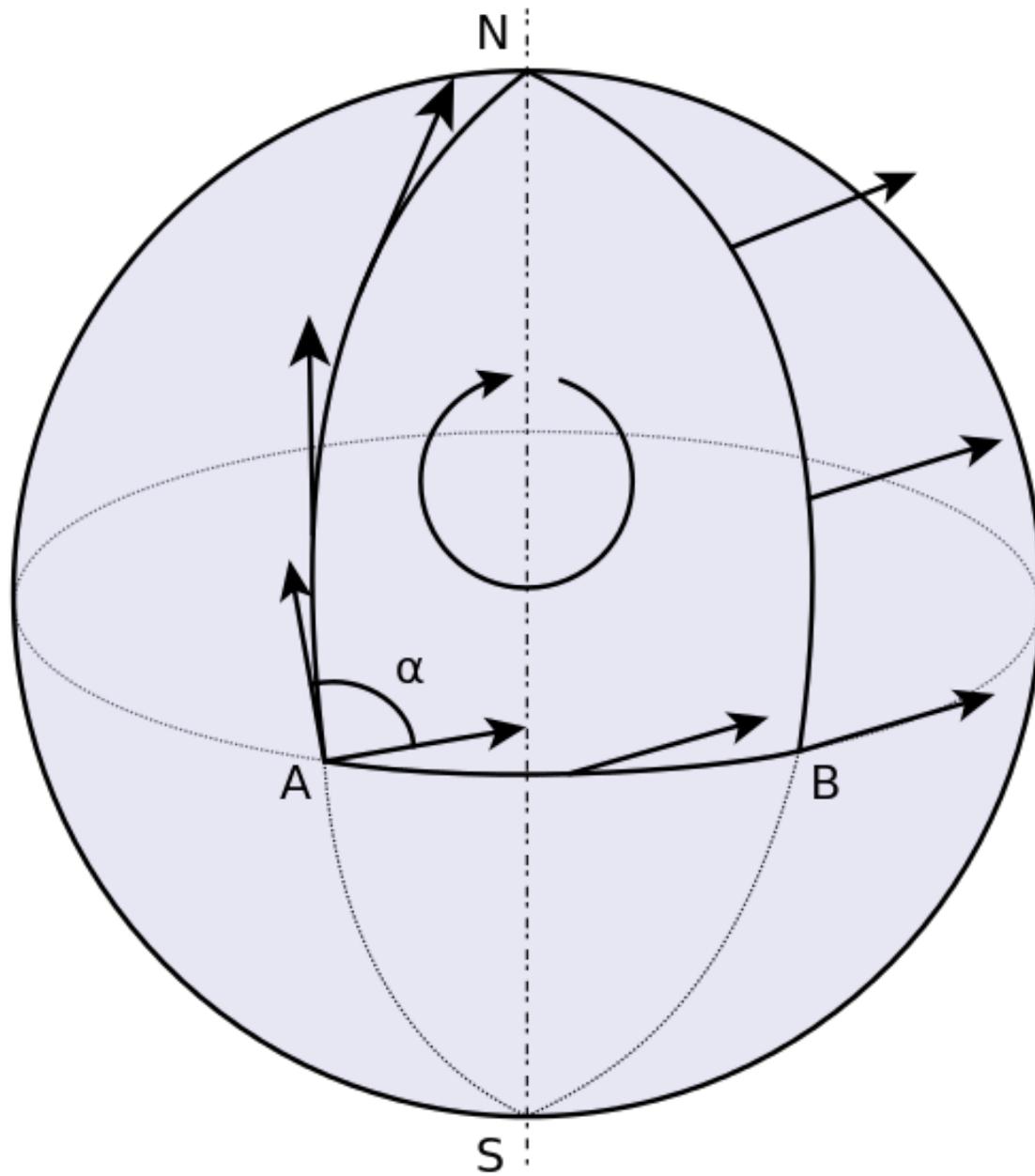
Information geometry

From Wikipedia, the free encyclopedia



This article may need to be **rewritten entirely** to comply with Wikipedia's [quality standards](#), as What an unreadable mess of math equations and name lists, which needs to be put into shape for a general audience encyclopedia.. [You can help](#). The [discussion page](#) may contain suggestions. *(May 2013)*

Information geometry is a branch of [mathematics](#) that applies the techniques of [differential geometry](#) to the field of [probability theory](#). This is done by taking [probability distributions](#) for a [statistical model](#) as the points of a [Riemannian manifold](#), forming a [statistical manifold](#). The F

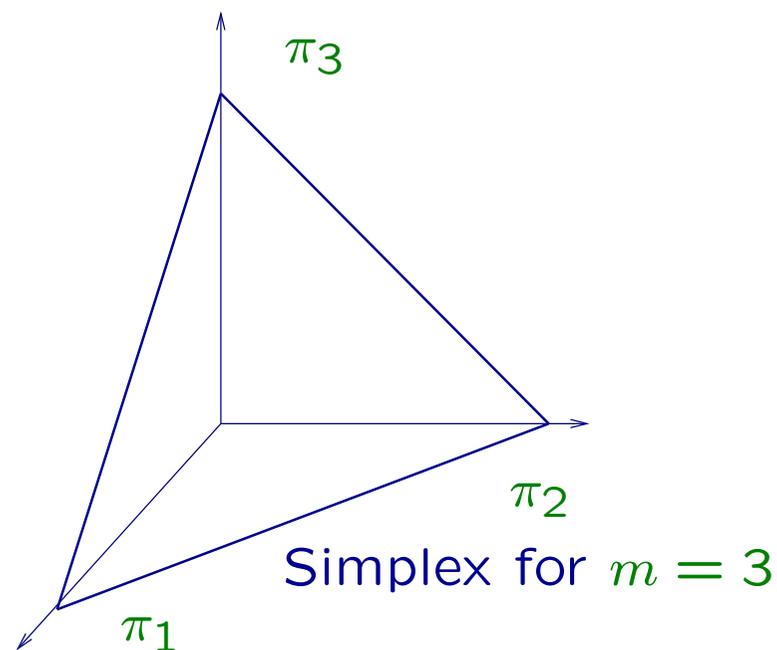


Points, manifolds and curves

Points I: Probability distributions for categorical variables

A random variable taking only X a discrete variable taking m distinct values $X \in \mathcal{X} = \{S_1, \dots, S_m\}$ with probability $\pi = (\pi_1, \dots, \pi_m)$.

Categorical distributions span the probability simplex: the subset of \mathbb{R}^m such that $\sum_i \pi_i = 1$ and $0 \leq \pi_i \leq 1$.



Multinomial: n independent realizations of a categorical variable.

An n -sample projects right away on the simplex by the sample frequencies:

$$\hat{\pi} = (\hat{\pi}_1, \dots, \hat{\pi}_m) = \left(\frac{n_1}{n}, \dots, \frac{n_m}{n} \right)$$

Points II: Probability distributions for continuous variables

Continuous: a random variable taking values in a measurable space.

We shall only consider families of probability distributions having a density with respect to a common (dominating) measure. And we further require common support (or trouble will happen).

Parametric families $\{p(x; \theta) \mid \theta \in \Theta \subset \mathbb{R}^p\}$ as p -dimensional manifolds.

The empirical distribution \hat{p} of n samples

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i)$$

does not live on the manifold.

Curves in distribution space: one-dimensional families

Such as simple transformation models:

- Location $p(x; \theta) = q(x - \theta)$
- Scale $p(x; \theta) = \frac{1}{\theta} q\left(\frac{x}{\theta}\right)$

But such curves do not exist for a general random variable, for which addition or multiplication are not defined.

We are looking for a **generic**, way of drawing a 'line' between two distributions.

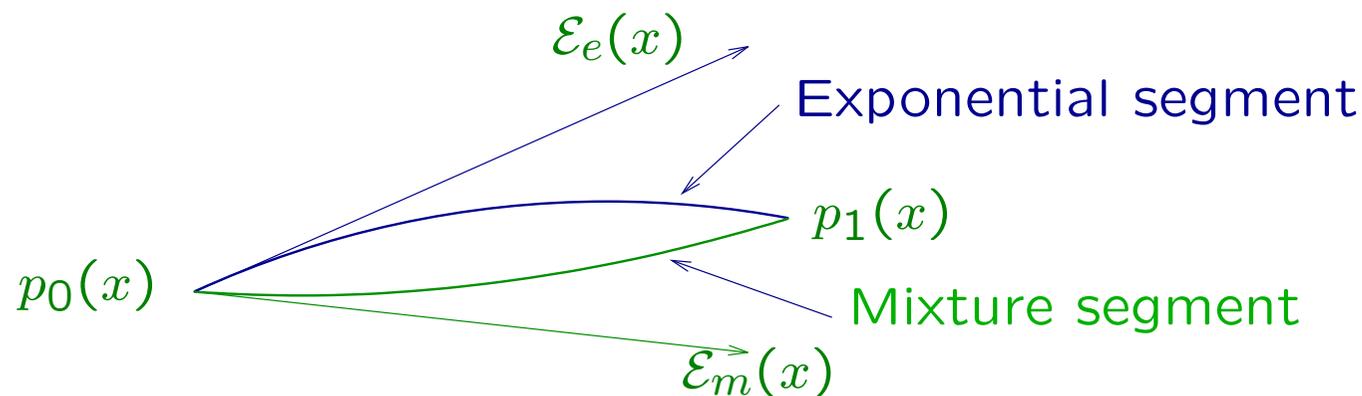
Statistics suggest at least two segments:

- Mixture segment
- Exponential segment

See next slide

Straight segment from p_0 to p_1 ? Two statistical ideas

- Defining 'segments' between two points $p_0(x)$ et $p_1(x)$.



- A 'mixture segment' from $p_0(x)$ to $p_1(x)$:

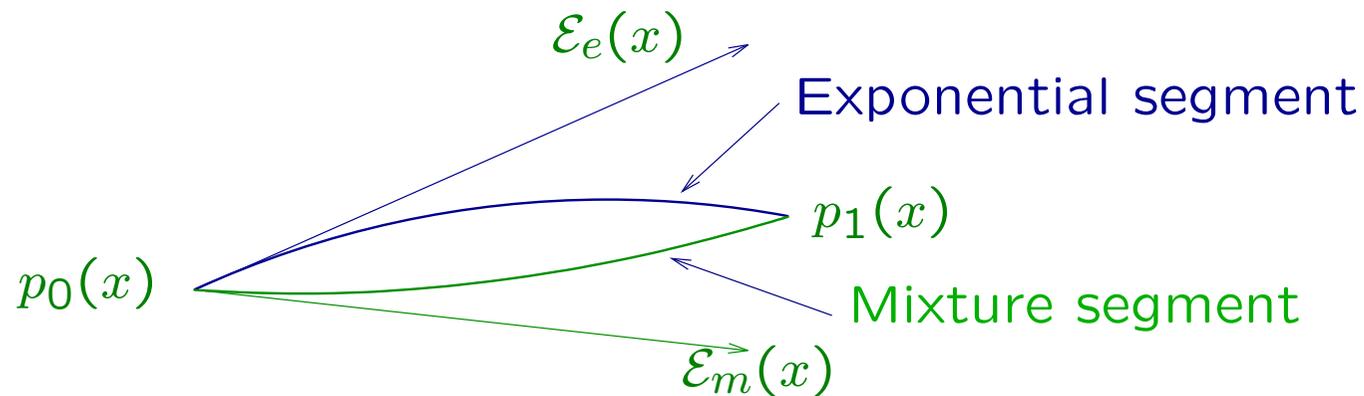
$$p_m(x; t) = (1 - t)p_0(x) + tp_1(x), \quad t \in [0, 1]$$

- An 'exponential segment' from $p_0(x)$ to $p_1(x)$:

$$\log p_e(x; t) = (1 - t) \log p_0(x) + t \log p_1(x) - \log \phi(t), \quad t \in [0, 1]$$

Do they look good in a simple Gaussian case ?

Tangent vectors (?)



The mixture and exponential segments from $p_0(x)$ to $p_1(x)$ also read

$$p_m(x; t) = p_0(x) (1 + t \mathcal{E}_m(x)) \quad \mathcal{E}_m(x) = \frac{p_1(x) - p_0(x)}{p_0(x)}$$

$$p_e(x; t) = p_0(x) \exp(t \mathcal{E}_e(X) - \psi(t)) \quad \mathcal{E}_e(x) = \log \frac{p_1(x)}{p_0(x)} - \mathbf{E}_0 \log \frac{p_1(x)}{p_0(x)}$$

We may interpret $\mathcal{E}_e(X)$ and $\mathcal{E}_m(X)$ as tangent vectors at point p_0 .

Note that they are zero-mean there:

$$\mathbf{E}_0 \mathcal{E}_e(X) = \mathbf{E}_0 \mathcal{E}_m(X) = 0$$

Affine space

Affine space: a vector space which has forgotten its origin (and does not care).

In Euclidean space, 3 points A, B, C define a plane \mathcal{P} .

Pick up an origin, say point C , and represent any point M in the plane \mathcal{P} as

$$M = M(\alpha, \beta) : \quad \vec{CM} = \alpha\vec{CA} + \beta\vec{CB}$$

From segments to affine manifolds

Def. An exponential (resp. mixture) family contains all the exponential (resp. mixture) segments between any two of its points.

So we start with $n + 1$ distributions p_0, \dots, p_n and build the n -dimensional model 'spanned' by them.

In the exponential case, that would be:

$$\log p(x; \theta) = \log p_0(x; \theta) + \sum_{i=1}^n \theta_i \log \frac{p_i(x)}{p_0(x)} - \psi(\theta) \quad \theta \in \Theta \subset \mathbb{R}^n$$

This is an exponential family. The usual, more general (?), definition:

$$p(x; \theta) = g(x) e^{\theta^\dagger S(x) - \psi(\theta)} \quad \theta \in \Theta \subset \mathbb{R}^n$$

where $g(x)$ is some positive measure and $S(x)$ an $n \times 1$ function of x , called the *sufficient statistic*. Then

$$\psi(\theta) = \log \int g(x) e^{\theta^\dagger S(x)} dx$$

This means bliss. Note that this is an 'affine' model. More about it soon.

Statistical manifolds

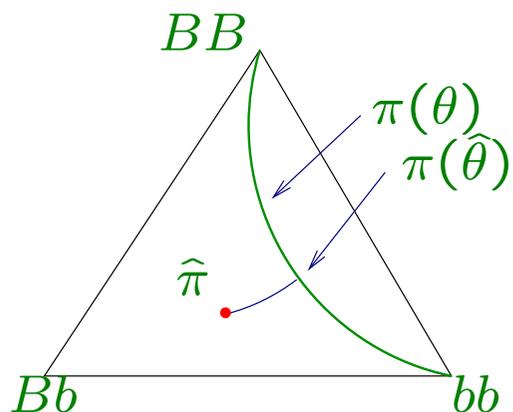
Hardy-Weinberg equilibrium: a manifold for a 3-categorical variable

In a population, a gene exists in two forms B and b in proportions θ and $1 - \theta$ respectively.

Individuals can carry pairs (BB) , (Bb) , (bb) .

At equilibrium, these are found with probabilities:

(BB)	(Bb)	(bb)
θ^2	$2\theta(1 - \theta)$	$(1 - \theta)^2$



With just a little bit of work, the Hardy-Weinberg model is seen to be an exponential family. Hint:

$$\log P_{BB} + \log P_{bb} - 2 \log p_{Bb} = -\log 4$$

Statistical manifolds: smooth parametric families of distributions

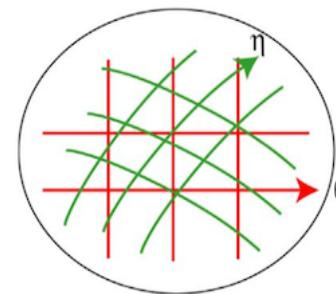
Parametric model: $p(x; \theta)$ for the distribution of X . $X \in \mathcal{X}$, $\theta \in \Theta \subset \mathbb{R}^p$.

Assume smoothness and a dominating measure.

The parameter θ is used to label the distributions in the parametric model but there are (infinitely) many ways to do it.

For instance, the normal family $\mathcal{N}(\mu, \sigma^2)$ may be parameterized by

$$\theta = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix} \quad \theta = \begin{bmatrix} \mu \\ \log \sigma \end{bmatrix} \quad \theta = \begin{bmatrix} \mu \\ \mu^2 + \sigma^2 \end{bmatrix} \quad \theta = \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix}$$



Statistical inference should be invariant with respect to parameterization

→ we mean *geometry*.

Still, geometry may help us pick ‘preferred’ parameterizations. . .

In the above example, which one do you think is ‘best’ ?

And for the Hardy-Weinberg model ?

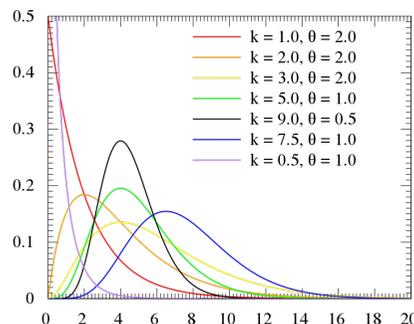
Some parametric models. I: Generative and transformation models

- Location $p(x; \theta) = q(x - \theta)$
- Scale $p(x; \theta) = \frac{1}{\theta} q\left(\frac{x}{\theta}\right)$, $\theta \in \mathbb{R}^+$
- Location-scale: $p(x; \theta) = \frac{1}{\theta_1} q\left(\frac{x - \theta_2}{\theta_1}\right)$
- Multi-dimensional scale: $p(x; \theta) = \frac{1}{|\det \theta|} q(\theta^{-1}x)$, with $x \in \mathbb{R}^n$, $\theta \in \text{GL}(n)$.
- Contamination: the distribution of $\frac{X + \lambda Y}{\sqrt{\sigma_Y^2 + \lambda^2 \sigma_Y^2}}$ when λ is varied over $[0, +\infty]$
- Mixture of Gaussians
- Graphical models
- From statistical physics
- Any stochastic physical model with adjustable parameters. . .

Some parametric models. II: From the big book

Such as, say, the Gamma distribution, with support on R^+ . Scale parameter θ and shape parameter k :

$$p(x; \theta, k) = \frac{1}{\Gamma(k)\theta^k} x^{k-1} e^{-x/\theta}$$



For integer k , it is the distribution of the sum of k independent exponentially distributed random variables, each with mean θ (Erlang distribution).

And many, many more...

Some parametric models. III: Exponential families. Oh so nice!

- A p -dimensional family of distributions for a random X (of any kind):

$$p(x; \theta) = g(x) e^{\theta^\dagger S(x) - \psi(\theta)} \quad \theta \in \Theta \subset \mathbb{R}^p$$

where $g(x)$ is some positive measure, $S(x)$ is a $p \times 1$ function of x .

$S(x)$: sufficient statistic, $\psi(\theta) = \log \int g(x) e^{\theta^\dagger S(x)} dx$: potential function.

- Differentiate $\mathbf{E}_\theta \mathbf{1} = \mathbf{1}$ twice to get:

$$\frac{\partial \psi}{\partial \theta} = \mathbf{E}_\theta S(x), \quad \frac{\partial^2 \psi}{\partial \theta^2} = \mathbf{Cov}_\theta S(x) \quad \text{Convexity!}$$

Differentiate further: ψ is the cumulant (of $S(X)$) generating function.

- Optional normalizations:

– Take $\int g(x) = 1$ to have $p(x; 0) = g(x)$ then ...

$$\psi(0) = 0.$$

– Replace $S(x)$ by $S(x) - \mathbf{E}_0 S(x)$ then ...

$$\frac{\partial \psi}{\partial \theta}(0) = 0.$$

– Replace $S(x)$ by $[\mathbf{Cov}_0 S(x)]^{-1/2} S(x)$ then ...

$$\frac{\partial^2 \psi}{\partial \theta^2}(0) = I_p.$$

– Change the origin and the basis vectors as you please.

- See the affine structure for the **canonical parameter** θ .

Who is exponential ?

The key thing in the exponential structure is the sufficient statistic

$$p(x; \theta) = g(x) e^{\theta^\dagger S(x) - \psi(\theta)} \quad \theta \in \Theta \subset \mathbb{R}^p$$

and that data and parameters interact only in the log through $\theta^\dagger S(x)$.

Is that so common ? Maybe after some massaging...

- The normal family $\mathcal{N}(\mu, R)$ for $X \in \mathbb{R}^d$ is a p -dimensional exponential family $p = d + d(d + 1)/2$ since $-(x - \mu)^\dagger R^{-1} (x - \mu)/2 = \theta^\dagger S(x) - \mu^\dagger R^{-1} \mu/2$ with

$$S(x) = [x_i, x_i x_j]_{1 \leq i \leq j \leq d} \quad \theta = [(R^{-1} \mu)_i, -(R^{-1})_{ij}/2]_{1 \leq i \leq j \leq d}$$

- The set of all multinomial distributions.
- Look in the big book: Normal, Gamma, Wishart, Poisson, ...
- *Maxent* distributions: distributions of maximal entropy with given moments.
- Curved exponential families as universal approximants to smooth models.

Kullback-Leibler divergence and the likelihood

Likelihood for multinomials

Data $X^n = [X_1, \dots, X_n]$ where $X_i \in \mathcal{X} = [S_1, \dots, S_p]$ modelled as n independent samples, each drawn with probability $\pi = [\pi_1, \dots, \pi_p]$ i.e. $\text{Prob}(X_i = S_k) = \pi_k$.

Probability of an n -sequence under a distribution categorical distribution π .

$$\begin{aligned}\log p(X^n; \pi) &= \sum_{k=1}^n \log \text{Prob}(X_k = S_k) \\ &= \sum_{i=1}^p \#(X_k = S_i) \log \pi_i \\ &= n \sum_{i=1}^p \hat{\pi}_i \log \pi_i\end{aligned}$$

where $\hat{\pi}_i = \#(X_k = S_i)/n$ is the frequency of S_i in the sequence X^n .

Key point: the probability of the sample depends only on the empirical distribution $\hat{\pi} = [\hat{\pi}_1, \dots, \hat{\pi}_p]$ which thus is a *sufficient statistic*.

Likelihood and Kulback matching

An enlightning decomposition:

$$-\frac{1}{n} \log p(X^n; \pi) = - \sum_{i=1}^p \hat{\pi}_i \log \pi_i = \sum_{i=1}^p \left[\hat{\pi}_i \log \frac{\hat{\pi}_i}{\pi_i} - \hat{\pi}_i \log \hat{\pi}_i \right]$$

For two categorical distributions p and q , define the KL divergence:

$$K [p | q] \stackrel{\text{def}}{=} \sum_i p_i \log \frac{p_i}{q_i} \quad \text{and} \quad H(p) \stackrel{\text{def}}{=} - \sum_i p_i \log p_i \quad \text{Shannon entropy.}$$

The Kullback-Leibler divergence from p to q also reads

$$K [p | q] = \sum_i p_i k\left(\frac{q_i}{p_i}\right) \quad k(u) \stackrel{\text{def}}{=} u - 1 - \log u \geq 0$$

This (non symmetric) measure of discrepancy between discrete distributions allows us to write the likelihood as

$$p(X^n; \pi) = e^{-n(K[\hat{\pi} | \pi] + H(\hat{\pi}))}$$

The likelihood measures goodness of fit in the Kullback sense.

Maximum likelihood estimation is a Kullback-projection onto the model.

Distribution of the empirical distribution

If n samples are i.i.d. distributed according to π and if X^n has type $\hat{\pi}$.

$$p(X^n; \pi) = e^{-n\{K[\hat{\pi}|\pi]+H(\hat{\pi})\}}$$

The set $\mathcal{T}(\hat{\pi})$ of n -samples with an empirical distribution $\hat{\pi}$ has cardinality:

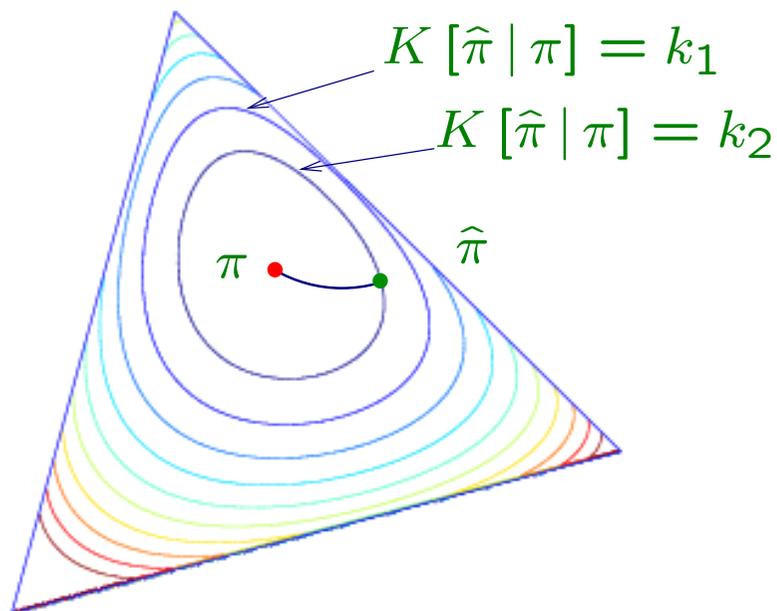
$$|\mathcal{T}(\hat{\pi})| = \frac{n!}{\prod_i n_i!} \approx e^{nH(\hat{\pi})} \quad (\text{via Stirling approx.})$$

Then, the probability of observing a particular type $\hat{\pi}$ is

$$p(\hat{\pi}; \pi) = \sum_{X^n \in \mathcal{T}(\hat{\pi})} p(X^n; \pi) = |\mathcal{T}(\hat{\pi})| e^{-n\{K[\hat{\pi}|\pi]+H(\hat{\pi})\}} \approx e^{-nK[\hat{\pi}|\pi]}$$

Thus, the sample distribution $\hat{\pi}$ is concentrated in a ‘Kullback ball’ centered at the true distribution π with a ‘radius’ of order $1/n$.

Distribution of sample distributions (ctd)



- $p(\hat{\pi}) \approx e^{-n K[\hat{\pi} | \pi]}$
- Most of the probability is concentrated in a Kullback-ball of center π and of radius $1/n$
- or is it $1/\sqrt{n}$?

For $\hat{\pi}$ close to π , the KLD reduces to the (symmetric) quadratic χ^2 distance:

$$K[\hat{\pi} | \pi] \approx \frac{1}{2} \sum_i \frac{(\hat{\pi}_i - \pi_i)^2}{\pi_i} \approx K[\pi | \hat{\pi}]$$

and small Kullback balls become ellipsoids with covariance matrices:

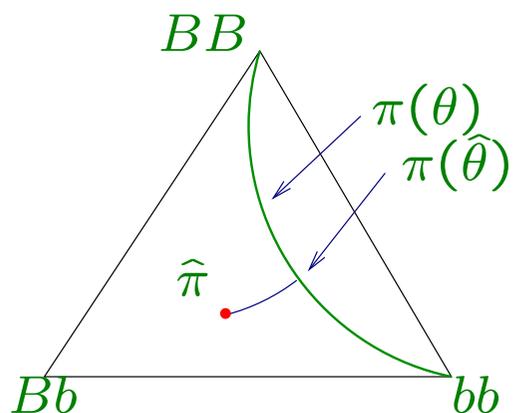
$$\text{Cov}(\hat{\pi}_i, \hat{\pi}_j) = \frac{1}{n} (\delta_{ij} \pi_i - \pi_i \pi_j)$$

Can you see why?

Back to Hardy-Weinberg

If $\pi = \pi(\theta)$, the maximum likelihood estimate of θ is the minimizer of $K[\hat{\pi} | \pi(\theta)]$.
Kullback matching.

Remember Hardy-Weinberg



(BB)	(Bb)	(bb)
θ^2	$2\theta(1 - \theta)$	$(1 - \theta)^2$

Finding the maximum likelihood solution is Kullback-projecting the sample distribution onto the model.

We'll see later why this picture is a bit wrong.

Kullback matching for continuous variables. I

- Kullback-Leibler divergence *from* a distribution P to a distribution Q :

$$K [P | Q] = \mathbf{E}_P \log \frac{P(x)}{Q(x)} = \int P(x) \log \frac{P(x)}{Q(x)} dx$$

- It is *not* symmetric: in general, $K [P | Q] \neq K [Q | P]$
- It is strictly positive unless P and Q agree P -almost everywhere because $K [P | Q] = \mathbf{E}_P [k (Q(x)/P(x))]$ with $k(u) = u - 1 - \log(u) \geq 0$.
- It is invariant: under any invertible transform $f(\cdot)$,

$$K [P_X | P_Y] = K [P_{f(X)} | P_{f(Y)}]$$

Kullback matching for continuous variables. II

Question: With a continuous variable X and a parametric model $p(x; \theta)$, can we play the same trick : maximizing the likelihood is finding the model distribution which is the closest in the Kullback sense to the sample distribution ?

No, in general: the sample distribution for n independent realizations x_1, \dots, x_n :

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i)$$

is at infinite distance from the model $K[\hat{p}(x) | p(x; \theta)] = +\infty$.

But considering the likelihood loss:

$$\hat{\mathcal{L}}(\theta) = -\frac{1}{n} \sum_i \log p(x_i; \theta)$$

as the estimate of the expected loss, with $p_*(x)$ the true distribution of X :

$$\mathcal{L}(\theta) = -\mathbf{E}_* \hat{\mathcal{L}}(\theta) = -\int p_*(x) \log p(x; \theta) = K[p_*(x) | p(x; \theta)] + H(p_*)$$

gives us the **expected shape of the likelihood**.

Kullback matching in exponential families

Recall the generic exponential family:

$$p(x; \theta) = g(x) e^{\theta^\dagger S(x) - \psi(\theta)} \quad \psi(\theta) = \log \int g(x) e^{\theta^\dagger S(x)} dx$$
$$\frac{\partial \psi(\theta)}{\partial \theta} = \mathbf{E}_\theta S(X) \stackrel{\text{def}}{=} \eta(\theta) \quad \frac{\partial^2 \psi(\theta)}{\partial \theta^2} = \mathbf{Cov}_\theta S(X)$$

Since $\psi(\theta)$ is convex, there is a 1-to-1 mapping $\theta \longleftrightarrow \eta = \frac{\partial \psi(\theta)}{\partial \theta}$.
Any distribution can be uniquely labeled with $\eta = \eta(\theta) = \mathbf{E}_\theta S(x)$.

$$\frac{\partial \log p(x; \theta)}{\partial \theta} = S(x) - \frac{\partial \psi(\theta)}{\partial \theta} = S(x) - \eta(\theta)$$

So the maximum likelihood estimate of η based on x is $\hat{\eta}_{\text{ML}} = S(x)$.

Based on an n -sample, it is $\hat{\eta}_{\text{ML}} = \hat{\mathbf{E}}S(x)$.

With $\hat{\theta}_{\text{ML}} = \theta(\hat{\eta}_{\text{ML}})$, one has exactly

$$\hat{\mathcal{L}}(\theta) = -\frac{1}{n} \sum_i \log p(x_i; \theta) = K [\hat{\theta}_{\text{ML}} | \theta] + \text{data only term}$$

yielding the exact shape of the likelihood in geometric terms.

Summary

- There is a natural divergence in spaces of probability distributions.
- Maximum likelihood is Kullback matching.
- So we found a (non)-distance for information geometry.
- But differential geometry needs
a **local (infinitesimal) quadratic distance**.

The matrix and the metric

Estimation theory gives us a metric

Parametric model $p(x; \theta)$ for the distribution of $X \in \mathcal{X}$ with $\theta \in \Theta \subset \mathbb{R}^p$.

Estimator: a function $\mathcal{T} : \mathcal{X} \mapsto \mathbb{R}^p$ of the data, returns an estimate $\hat{\theta} = \mathcal{T}(X)$.

Cramér-Rao bound: The variance of an unbiased estimator is lower bounded:

$$\text{Cov}_{\theta}(\mathcal{T}(X)) \geq F(\theta)^{-1} \quad F(\theta): \text{ the } p \times p \text{ Fisher information matrix}$$

Hence, in essence, $p(x; \theta)$ and $p(x; \theta + \Delta\theta)$ are undistinguishable if, say,

$$\Delta\theta^{\dagger} F(\theta) \Delta\theta < 1 \quad \text{A resolution cell.}$$

Infinitesimal distance between $p(x; \theta)$ and $p(x; \theta + d\theta)$?

The CR bound makes it statistically unavoidable (Rao 1945) to define it as

$$\sqrt{d\theta^{\dagger} F(\theta) d\theta}$$

This is parameterization invariant, of course.

Note: for n independent observations, the Fisher matrix is $nF(\theta)$.

The resolution cell shrinks accordingly.

Score function and Fisher matrix

The score function is the vector-valued function

$$\varphi(x; \theta) = \frac{\partial \log p(x; \theta)}{\partial \theta}$$

and the score $\varphi(X; \theta)$ is a random vector which depends on θ .

A trivial but striking property: for any smooth function $h(x)$,

$$\frac{\partial}{\partial \theta} \mathbf{E}_\theta(h(X)) = \mathbf{E}_\theta(h(X)\varphi(X; \theta)).$$

In particular, the score $\varphi(X; \theta)$ has zero mean when $X \sim p(x; \theta)$:

$$\mathbf{E}_\theta \varphi(X; \theta) = 0$$

The ‘Fisher information matrix’ is the covariance matrix of the score:

$$F(\theta) = \text{Cov}_\theta(\varphi(X; \theta)) = \mathbf{E}_\theta \varphi(X; \theta)\varphi(X; \theta)^\dagger$$

Reminder: some (statistical) Euclidean geometry

- In a vector space \mathbb{R}^n : consider the approximation of a vector y as a linear combination of p vectors $\sum_{i=1}^p \Phi_i \alpha_i = \Phi \alpha$. Best ($\min_{\alpha} \|y - \Phi \alpha\|^2$) solution:

$$y^* = \Phi \alpha^* = \Phi (\Phi^\dagger \Phi)^{-1} \Phi^\dagger y = \Pi_{\Phi} y \quad \Pi_{\Phi} \stackrel{\text{def}}{=} \Phi (\Phi^\dagger \Phi)^{-1} \Phi^\dagger = \text{Proj}(\text{Span}(\Phi))$$

Basis-free characterization: $\Pi_{\Phi} y^* = y^*$ and $\Pi_{\Phi} (y - y^*) = 0$ (orthogonality).

Reminder: some (statistical) Euclidean geometry

- In a vector space \mathbb{R}^n : consider the approximation of a vector y as a linear combination of p vectors $\sum_{i=1}^p \Phi_i \alpha_i = \Phi \alpha$. Best ($\min_{\alpha} \|y - \Phi \alpha\|^2$) solution:

$$y^* = \Phi \alpha^* = \Phi (\Phi^\dagger \Phi)^{-1} \Phi^\dagger y = \Pi_{\Phi} y \quad \Pi_{\Phi} \stackrel{\text{def}}{=} \Phi (\Phi^\dagger \Phi)^{-1} \Phi^\dagger = \text{Proj}(\text{Span}(\Phi))$$

Basis-free characterization: $\Pi_{\Phi} y^* = y^*$ and $\Pi_{\Phi} (y - y^*) = 0$ (orthogonality).

- Prediction \hat{Y} of a random **scalar** Y by a linear combination of p zero-mean random variables: $\hat{Y} = \sum_{i=1}^p \alpha_i \xi_i$. The best ($\min \mathbf{E}(\hat{Y} - Y)^2$) prediction is

$$\hat{Y} = X^\dagger \mathbf{E}(X X^\dagger)^{-1} \mathbf{E}(X Y) \quad X = [\xi_1, \dots, \xi_p]^\dagger$$

Actually, we are working in a vector space of scalar random variables with scalar product $\langle \xi_i, \xi_j \rangle = \mathbf{E}(\xi_i \xi_j)$ and (squared) norm $\|\xi\|^2 = \langle \xi, \xi \rangle = \mathbf{E} \xi^2$.

Orthogonal decomposition as $Y = \hat{Y} + (Y - \hat{Y})$ since $\mathbf{E}((\hat{Y} - Y)\hat{Y}) = 0$. It entails $\mathbf{E} Y^2 = \mathbf{E} \hat{Y}^2 + \mathbf{E} (Y - \hat{Y})^2$.

Spoiler: tangent planes can be constructed like that.

The Frechet-Darmois-Cramér-Rao bound a.k.a. the CRB...

Decompose the error $e = \hat{\theta} - \theta$ made by an estimator $\hat{\theta} = \mathcal{T}(X)$ as

$$\hat{\theta} - \theta = e = e_{//} + e_{\perp}$$

where $e_{//}$ is the projection of e onto the linear span of $\varphi(X; \theta)$ and e_{\perp} is the orthogonal (uncorrelated) part. The projected error $e_{//}$ is:

$$e_{//} = \mathbf{E}_{\theta} (e \varphi^{\dagger}) \mathbf{E}_{\theta} (\varphi \varphi^{\dagger})^{-1} \varphi \quad \text{with } \varphi = \varphi(X; \theta)$$

If $\mathcal{T}(x)$ is unbiased, $\mathbf{E}_{\theta} (e \varphi^{\dagger})$ must be the identity matrix:

$$\mathbf{E}_{\theta} (e \varphi^{\dagger}) = \mathbf{E}_{\theta} ((\hat{\theta} - \theta) \varphi^{\dagger}) = \frac{\partial}{\partial \theta} \mathbf{E}_{\theta} \hat{\theta} - \theta \mathbf{E}_{\theta} \varphi^{\dagger} = I - 0.$$

Hence, its parallel error is $e_{//} = F_{\theta}^{-1} \varphi(X; \theta)$. It does *not* depend on \mathcal{T} !!

$$\hat{\theta} = \theta + F_{\theta}^{-1} \varphi(X; \theta) + e_{\perp}$$

Hence the Cramér-Rao bound: covariance of the unavoidable error:

$$\mathbf{Cov}_{\theta}(\hat{\theta}) = \mathbf{Cov}(e_{//}) + \mathbf{Cov}(e_{\perp}) \geq \mathbf{Cov}(e_{//}) = \mathbf{Cov}_{\theta} (F_{\theta}^{-1} \varphi) = F_{\theta}^{-1}$$

An *efficient estimator* is unique: it must behave as $\mathcal{T}(X) = \theta + F_{\theta}^{-1} \varphi(X; \theta)$.

Summary

We have a statistically well founded Riemmanian metric.

Further, it is **Kullback-compatible**: For infinitesimal $d\theta$,

$$K[\theta | \theta + d\theta] = K[\theta + d\theta | \theta] = \frac{1}{2} d\theta^\dagger F(\theta) d\theta$$

So, we are happy.

Tangent planes

Representing the tangent plane with random variables

- Tangent plane at point θ_* = span of the scores there.

$$\mathcal{T}(\theta_*) = \left\{ \alpha^\dagger S(x; \theta_*) = \sum_i \alpha_i S_i(x; \theta_*) = \sum_i \alpha_i \frac{\partial \log p(X; \theta_*)}{\partial \theta_i} \right\}$$

So the tangent plane is a vector space of zero-mean scalar random variables with the norm controlled by the Fisher matrix

$$\|\alpha\|^2 = \mathbf{E}(\alpha^\dagger S(x; \theta_*))^2 = \alpha^\dagger F(\theta_*) \alpha$$

- More generally, for $P(x)$ a prob. dist. and $d(X)$ a zero-mean variable $\mathbf{E}_P d(X) = 0$, one can define a ‘translation’ $P \rightarrow Q : P \vec{Q} = \vec{d}$ by $\vec{d} = d(X)$ by:

$$Q(x) = P(x) e^{d(x) - \psi(\vec{d})} \quad \psi(\vec{d}) = \log \mathbf{E}_P e^{d(X)} \approx \frac{1}{2} \mathbf{E}_P d^2(X) = \frac{1}{2} \|\vec{d}\|_P^2.$$

Note that $K [P | P + \vec{d}] = K [P | Q] = \mathbf{E}_P (\log P/Q) = \psi(\vec{d})$.

If $P(x) = p(x; \theta)$ and $d(x) = \alpha^\dagger S(x; \theta)$, the translated distribution stays in the model approximately. It stays in it **exactly** for an exponential family.

The exponential family tangent to a smooth model

- 2nd-order expansion around point θ_* of the log-density in a smooth model:

$$\log p(x; \theta) \approx \log p(x; \theta_*) + (\theta - \theta_*)^\dagger \frac{\partial \log p(x; \theta_*)}{\partial \theta} + \frac{1}{2} (\theta - \theta_*)^\dagger \frac{\partial^2 \log p(x; \theta_*)}{\partial \theta^2} (\theta - \theta_*)$$

- Exponential family based at θ_* with sufficient statistic $\varphi(x; \theta_*) = \frac{\partial \log p(x; \theta_*)}{\partial \theta}$:

$$p(x; \theta) = p(x; \theta_*) e^{(\theta - \theta_*)^\dagger \varphi(x; \theta_*) - \psi(\theta - \theta_*)}.$$

The potential ψ satisfies $\psi(0) = 0$ and $\frac{\partial \psi(0)}{\partial \theta} = \mathbf{E}_* \varphi(X; \theta_*) = 0$ and

$$\frac{\partial^2 \psi(0)}{\partial \theta^2} = \mathbf{Cov}_*(\varphi(X, \theta_*)) = F(\theta_*)$$

so, at second order,

$$\psi(\theta - \theta_*) \approx \frac{1}{2} (\theta - \theta_*)^\dagger F(\theta_*) (\theta - \theta_*).$$

Hence the two models are almost identical at second-order since

$$-\frac{\partial^2 \log p(x; \theta_*)}{\partial \theta^2} \approx -\mathbf{E}_* \frac{\partial^2 \log p(x; \theta_*)}{\partial \theta^2} = \mathbf{E}_* \frac{\partial \log p(x; \theta_*)}{\partial \theta} \frac{\partial \log p(x; \theta_*)}{\partial \theta}^\dagger = F(\theta_*)$$

and they become equivalent if many independent samples are available.

Smooth models as curved exponential families

For a smooth model, define $\ell(x; \theta) = \log p(x; \theta)$, $\theta \in \mathbb{R}^n$, pick a θ_* and define the higher order scores:

$$\ell_i(x) = \frac{\partial \ell(x; \theta_*)}{\partial \theta_i}, \quad \ell_{ij}(x) = \frac{\partial^2 \ell(x; \theta_*)}{\partial \theta_i \partial \theta_j}, \quad \ell_{ijk}(x) = \frac{\partial^3 \ell(x; \theta_*)}{\partial \theta_i \partial \theta_j \partial \theta_k},$$

In general $d = n + n(n + 1)/2 + n(n + 1)(n + 2)/6$ such independent scores.

Then, build with them a d -dimensional exponential family $p_3(x; T)$, $T \in \mathbb{R}^d$

$$p_3(x; T) = p(x; \theta_*) e^{T^\dagger S(x) - \psi(T)} \quad \text{with} \quad S(x) \stackrel{\text{def}}{=} [\ell_i(x) \ \ell_{ij}(x) \ \ell_{ijk}(x)]_{1 \leq i \leq j \leq k \leq n}^\dagger.$$

The third order expansion of the log-density $\log p(x; \theta_* + t)$ reads

$$\log p(x; \theta_* + t) \approx \log p(x; \theta_*) + \sum_i t_i \ell_i(x) + \sum_{ij} t_i t_j \ell_{ij}(x) + \sum_{ijk} t_i t_j t_k \ell_{ijk}(x) + \dots$$

so the log-density $\log p(x; \theta)$ is approximated at third-order around θ_* by

$$\log p(x; \theta_* + t) \approx \log p_3(x; T(t)) \quad T(t) \stackrel{\text{def}}{=} [t_i \ t_{ij} \ t_{ijk}]_{1 \leq i \leq j \leq k \leq n}^\dagger.$$

We have approximated a smooth model by a **curved exponential family**.

It is curved because the $\mathbb{R}^n \mapsto \mathbb{R}^d$ mapping $t \rightarrow T(t)$ is not linear.

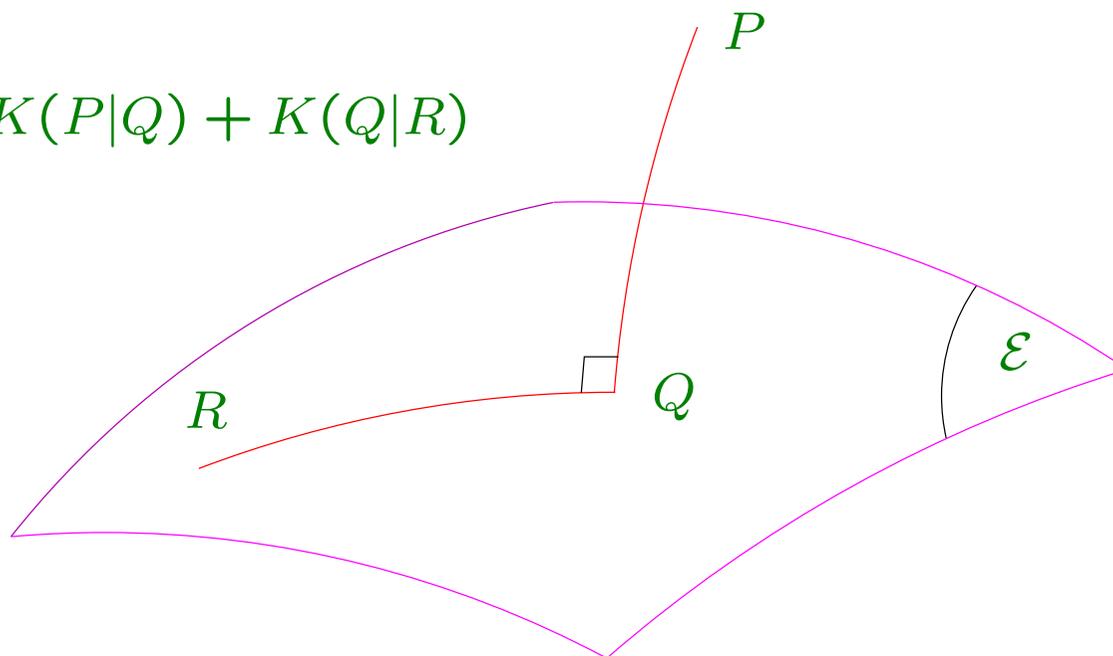
Summary

- Tangent plane to $p(x; \theta)$ at θ_* as a vector space . . .
- of zero-mean random variables . . .
- spanned by the scores at θ_* .
- Exponential model tangent to a smooth family
- ‘Osculating’ higher-dimensional exponential models to embed (approximately) a smooth manifold as a curved submodel.

The Pythagorean theorem of information geometry.

The Pythagorean theorem of information geometry

$$K(P|R) = K(P|Q) + K(Q|R)$$



- \mathcal{E} : an exponential family of distributions.
- Q : the minimizer in \mathcal{E} of $K(P\|\cdot)$. It is the m -projection of P onto \mathcal{E} .
- Then, for any other distribution R of \mathcal{E}

$$K(P\|R) = K(P\|Q) + K(Q\|R).$$

MLE and foliationfoliation here, maybe with maxent

Kullback divergence in an exponential family

Recall the generic exponential family:

$$p(x; \theta) = g(x) e^{\theta^\dagger S(x) - \psi(\theta)} \quad \psi(\theta) = \log \int g(x) e^{\theta^\dagger S(x)} dx$$

$$\frac{\partial \psi(\theta)}{\partial \theta} = \mathbf{E}_\theta S(x) \stackrel{\text{def}}{=} \eta(\theta) \quad \frac{\partial^2 \psi(\theta)}{\partial \theta^2} = \mathbf{Cov}_\theta S(x)$$

Since $\psi(\theta)$ is convex, there is a 1-to-1 mapping $\theta \longleftrightarrow \eta = \frac{\partial \psi(\theta)}{\partial \theta}$. Any point is uniquely labeled with $\eta = \eta(\theta) = \mathbf{E}_\theta S(x)$, the conjugate variable.

The convex conjugate of $\psi(\theta)$ is $\varphi(\eta)$ which is minus the intercept of the tangent plane with slope η , so $\psi(\theta) = -\varphi(\eta) + \theta^\dagger \eta$ or

$$\varphi(\eta) = \theta^\dagger \eta - \psi(\theta) \quad \text{for } \theta \text{ s.t. } \frac{\partial \psi(\theta)}{\partial \theta} = \eta$$

Then $\psi(\theta) \geq -\varphi(\eta) + \eta^\dagger \theta$ for all θ, η with equality if $\frac{\partial \psi(\theta)}{\partial \theta} = \eta$. Then

$$K [1 | 2] = \varphi(\eta_1) + \psi(\theta_2) - \eta_1^\dagger \theta_2 \quad \text{Kullback as Bregman}$$

The dual potential functions can also be expressed as

$$\varphi(\eta) = K [p(x; \eta) | g(x)] \quad \psi(\theta) = K [g(x) | p(x; \theta)]$$

Euclidean vs Information geometry

- Do we have, *non locally in distribution space* something like this?

$$\|\vec{AC}\|^2 = \|\vec{AB} + \vec{BC}\|^2 = \|\vec{AB}\|^2 + \|\vec{BC}\|^2 - 2\vec{BA} \cdot \vec{BC}$$

- Take any three distributions P, Q, R , and check:

$$K[P|R] - K[P|Q] - K[Q|R] = \int P \log \frac{P}{R} - \int P \log \frac{P}{Q} - \int Q \log \frac{Q}{R} = \int (P - Q) \log \frac{Q}{R}$$

- Making sense of the difference:

$$\int (P - Q) \log \frac{Q}{R} = - \int Q \left(\frac{P}{Q} - 1 \right) \left(\log \frac{R}{Q} + K[Q|R] \right)$$

- Everything is *beautiful* with mixture and exponential vectors:

$$K[P|R] = K[P|Q] + K[Q|R] - \langle \vec{QP}^m, \vec{QR}^e \rangle_Q$$

$$\vec{QP}^m = \frac{P(X)}{Q(X)} - 1 \quad \vec{QR}^e = \log \frac{R(X)}{Q(X)} + K[Q|R]$$

and $\langle \vec{QP}^m, \vec{QR}^e \rangle_Q = 0$ for the Pythagorean theorem.

A dual mixture/exponential structure is hidden in the KLD.

Back to Hardy-Weinberg

Parallel transport and connections

Connecting tangent planes, the regular way

Let $V_i(\theta)$ be a basis for the tangent plane at $\theta = [\theta^1, \dots, \theta^n]$.

The scalar product between two vectors at θ :

$$\langle X|Y \rangle = \langle \sum_i x^i V_i | \sum_j y^j V_j \rangle = \sum_{ij} x^i y^j \langle V_i | V_j \rangle \stackrel{\text{def}}{=} \sum_{ij} x^i y^j g_{ij}(\theta)$$

Variation of a vector field $X(\theta) = \sum_i x^i(\theta) V_i(\theta)$:

$$\delta X(\theta) = \sum_i \delta [x^i V_i] = \sum_i \sum_j \left[\frac{\partial x^i}{\partial \theta^j} V_i + x^i \frac{\partial V_i}{\partial \theta^j} \right] \delta \theta^j$$

First order variation of the basis controlled by n^3 Christoffel symbols Γ_{ij}^k :

$$\frac{\partial V_i}{\partial \theta^j} = \sum_k \Gamma_{ij}^k V_k$$

The metric g_{ij} and the Christoffel symbols Γ_{ij}^k fully determine the geometry.

You/the math/the physics/the stats decide what to pick for $g_{ij}(\theta)$ and $\Gamma_{ij}^k(\theta)$...

... but a unique connection can be derived from the metric.

Statistical connections

Connections can be determined from statistical principles. Define:

$$\ell_i = \ell_i(x; \theta) = \frac{\partial \log p(x; \theta)}{\partial \theta_i} \quad \ell_{ij} = \ell_{ij}(x; \theta) = \frac{\partial^2 \log p(x; \theta)}{\partial \theta_i \partial \theta_j}$$

The metric (Levi-Civita) connection for the Fisher metric $g_{ij}(\theta) = \mathbf{E}_\theta \ell_i \ell_j$ is

$$\Gamma_{ij,k}^0(\theta) = \mathbf{E}_\theta \left((\ell_{ij} + \frac{1}{2} \ell_i \ell_j) \ell_k \right)$$

but we do not want it. Amari introduced a pair of dual connections:*

$$\Gamma_{ij,k}^{(e)}(\theta) = \mathbf{E}_\theta (\ell_{ij} \ell_k) \quad \Gamma_{ij,k}^{(m)}(\theta) = \mathbf{E}_\theta ((\ell_{ij} + \ell_i \ell_j) \ell_k)$$

$\Gamma_{ij,k}^{(e)} = 0$ in an exponential model. $\Gamma_{ij,k}^{(m)} = 0$ in a mixture model.

The scalar product of two vectors is preserved during parallel transport along any curve if one is e -transported while the other is m -transported.

But we are not going to do that.

*Actually, he introduced a 1D family of such pairs of dual connections.

Some parallel transport

Dual parallel transport preserves the scalar product.

Let a and b be two vectors of the tangent plane at P , represented by two random variables with zero mean under P .

$$\mathbf{E}_P a(X) = \mathbf{E}_P b(X) = 0 \text{ with scalar product } \langle a|b \rangle_P = \mathbf{E}_P (a(X)b(X))$$

The vectors are parallelly transported to Q .

Vector a undergoes an e -parallel transport and becomes a^e .

Vector b undergoes an m -parallel transport and becomes b^m .

How does parallel transport work here ?

$$a^e(X) = a(X) - \mathbf{E}_Q a(X) \quad b^m(X) = b(X) \frac{p(X)}{q(X)}$$

Check that $a^e(X)$ and $b^m(X)$ have zero-mean under Q and that

$$\langle a|b \rangle_P = \langle a^e|b^m \rangle_Q$$

Information geometry of ICA

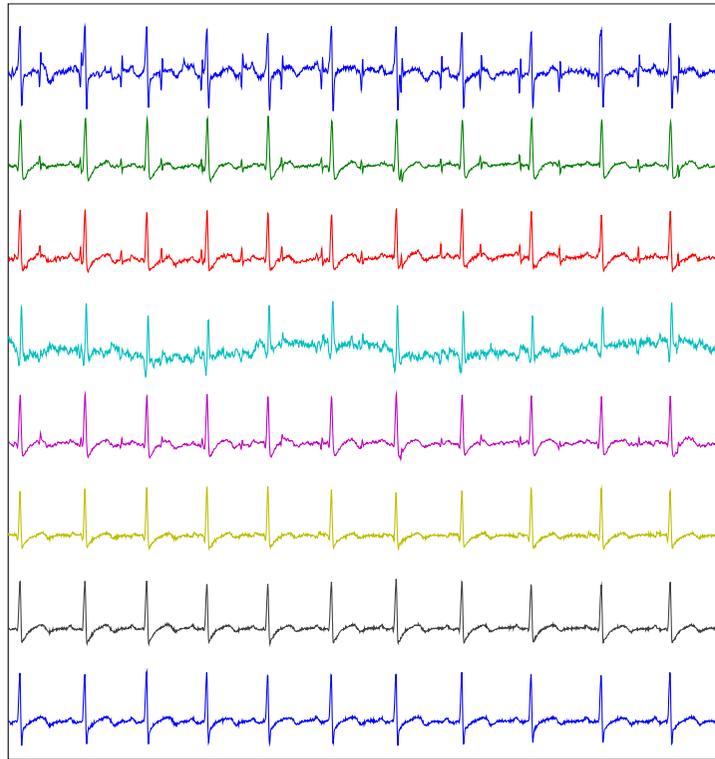
Multi-channel time series: a biomedical example



8 ECG electrodes located on the thorax and the abdomen of a pregnant woman.

Looking for linear decompositions: $\text{Data} = \text{Mixing matrix} \times \text{Sources}$.

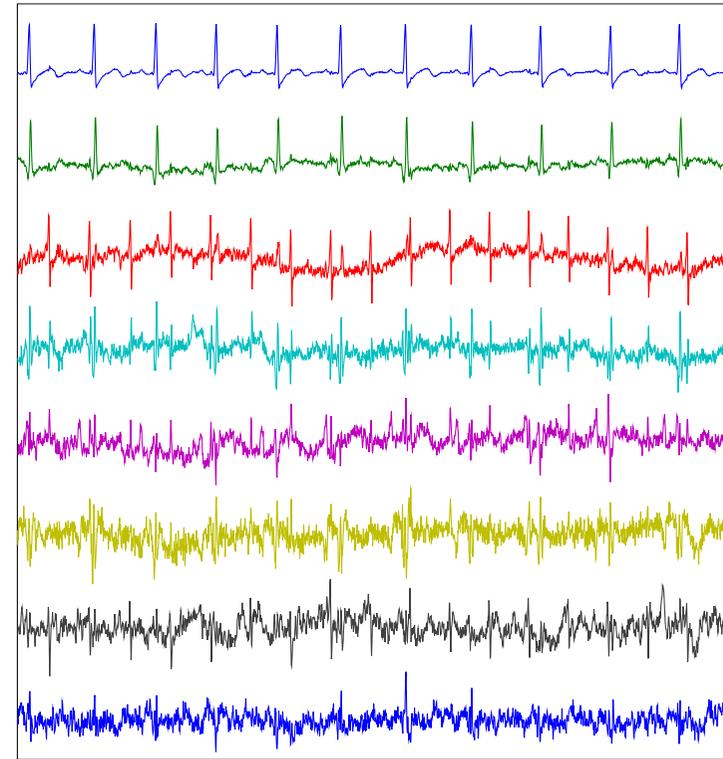
Principal component analysis



=

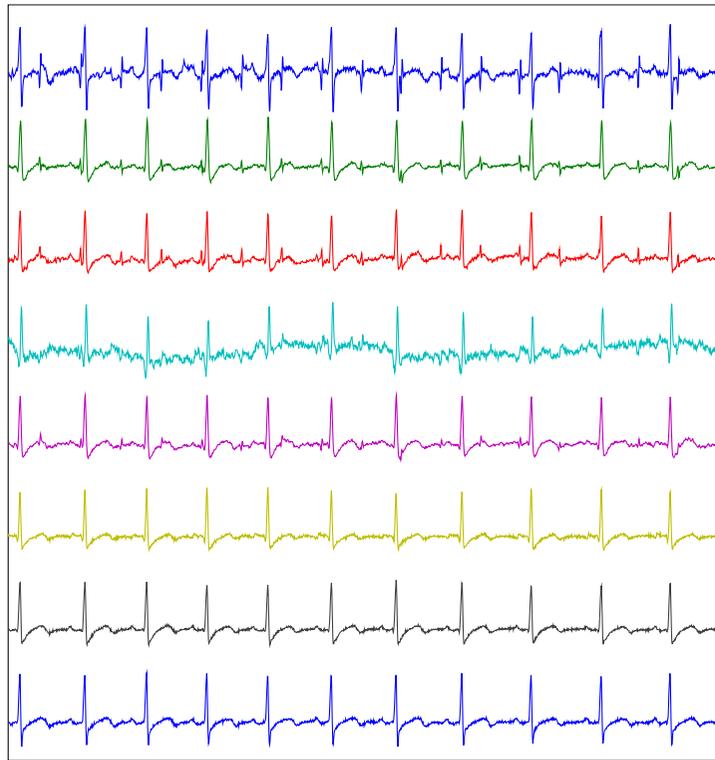
95	-51	-33	8	2	0	2	2
87	53	16	6	-2	-5	4	4
94	17	14	0	11	6	1	1
-9	105	-28	-5	1	0	0	0
97	26	0	5	0	-1	-9	-9
89	7	1	0	-11	8	0	0
97	-6	1	-4	-1	-4	1	1
92	-32	-3	-15	0	-1	0	0

×



- Orthogonal mixture, uncorrelated components $\frac{1}{T} \sum_t y_i(t)y_j(t) = 0$ for $i \neq j$
- Decorrelation is weak (always possible), orthogonality is implausible.

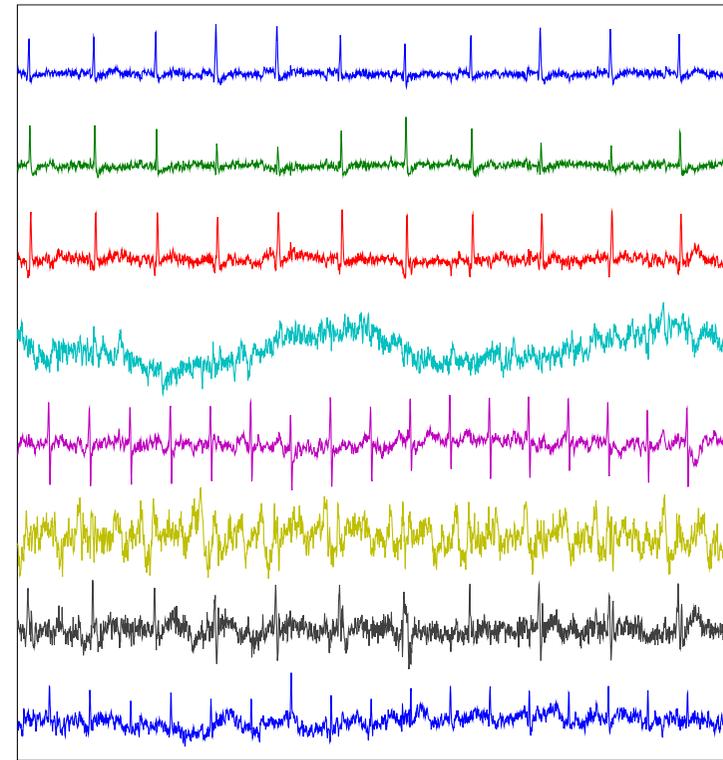
Independent component analysis



=

98	60	-71	5	-34	17	3	3
93	114	37	0	10	5	6	6
101	99	-4	0	13	6	6	6
-17	47	79	33	-19	-13	-1	-1
106	102	7	0	-10	4	10	10
110	71	-8	0	-1	1	0	0
107	82	-22	0	0	7	17	17
112	42	-43	0	-1	11	20	20

×



- Linear decomposition into “the most independent sources”
- Blind: only independence is at work but it must go beyond decorrelation.
- Independence is statistically very strong but often physically plausible.
- Weak assumptions → wide applicability

The basic ICA model

- An $n \times T$ data set $X = \{x_i(t) \mid 1 \leq i \leq n, 1 \leq t \leq T\} \dots$

... modelled as $X = AS$ with an $n \times T$ source matrix of *independent* rows.

$$\boxed{X} = \boxed{A} \times \begin{array}{|c|} \hline \dots S_1 \dots \\ \hline \vdots \\ \hline \dots S_n \dots \\ \hline \end{array}$$

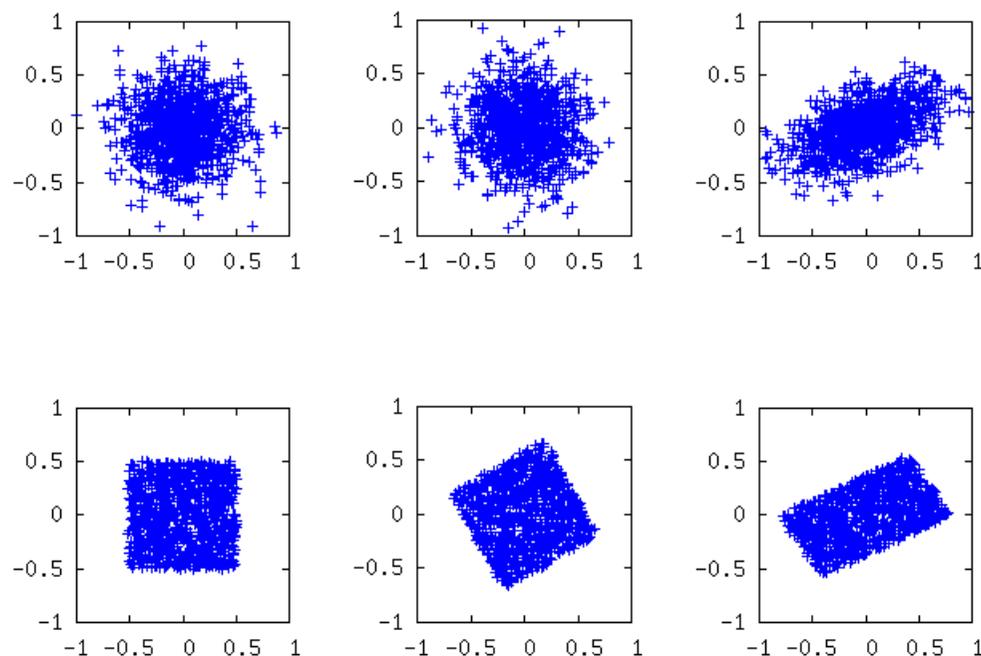
Assuming $X = AS$, is the assumption of statistical independence sufficient to recover the sources *blindly* and *efficiently* ?

Blindness: why not? If mixing always destroys independence, then recovering independence should unmix, *shouldn't it?*

The crux of the matter

Blind methods good statistical descriptors.

Mixing Gaussian (top row) or uniform (bottom) random variables with the identity (left), a rotation (center), a generic transform (right):



Rotation is invisible unless the data **and** the model are non Gaussian.

Mere decorrelation, $\frac{1}{T} \sum_t y_i(t) y_j(t) = 0$, does not cut it.

How to do it ?

- Ignore time/space structure but use non Gaussianity ...
 - Minimize dependence between recovered sources ...
 - ... as measured by mutual information, *why ? hard !*
 - ... or approximated using high order cumulants, *clumsy ?*
 - or find the most non Gaussian sources *why ? how ?*
 - or find sources uncorrelated through non linear functions:
e.g $\frac{1}{T} \sum_t \phi(y_i(t)) \psi(y_j(t)) = 0$ for $i \neq j$. Which functions ϕ, ψ ?
- ... or assume Gaussianity but use temporal/spatial structure, such as
 - ... correlations, or spectral diversity,
 - ... or non stationarity / inhomogeneities

But how to make sense of all that ? How to do it properly ? *Statistics!*

Likelihood and Kullback matching

The distribution of $X = AS$ is specified by $\theta = (A, P_{S_1}, \dots, P_{S_n})$.

$$-\log P_\theta(X) \quad \text{- log-likelihood}$$

$$= -E_X \log P_\theta(X) + \text{stoch.} \quad \text{The average likelihood landscape}$$

$$\stackrel{c}{=} K [P_X | P_\theta] + H(X) \quad H(X) = \text{entropy} = \text{cst}$$

$$\stackrel{c}{=} K [P_X | P_{AS}] \quad \text{The model is } X = AS.$$

$$= K [P_{A^{-1}X} | P_S] \quad \text{Invariance under invertible transforms.}$$

$$= K [P_Y | P_S] \quad \text{No } X, \text{ no } A, \text{ only } Y \stackrel{\text{def}}{=} A^{-1}X.$$

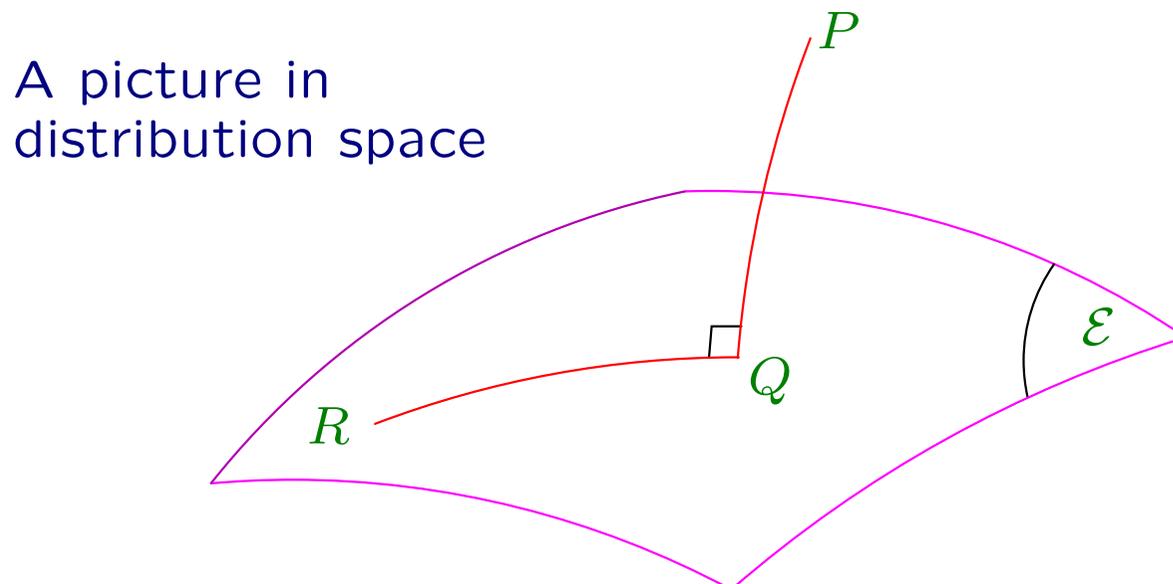
The likeliest A is the one making P_Y the closest to $P_S = \prod_i P_{S_i}$. Nice !

$$P_X \longrightarrow \boxed{A^{-1}} \longrightarrow P_Y \stackrel{\mathcal{D}}{\approx} P_S = \prod_i P_{S_i}$$

Data X and mixing A enter only through $Y = A^{-1}X$: equivariance.

The Pythagorean theorem of information geometry

The Kullback divergence may not be a distance, it still has its own private Pythagoras theorem.



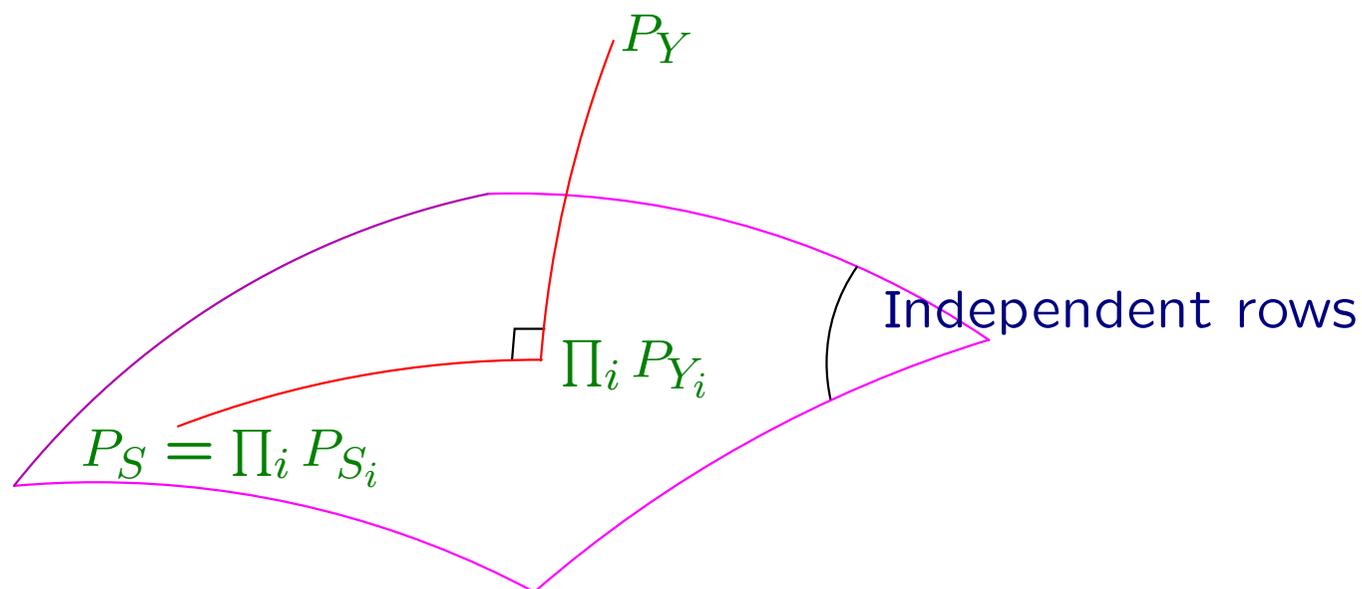
- \mathcal{E} : an *exponential* family of distributions.
- Q : the minimizer in \mathcal{E} of $K[P|\cdot]$. It is a projection of P onto \mathcal{E} .
- Then, for any other distribution R of \mathcal{E}

$$K[P|R] = K[P|Q] + K[Q|R]$$

The Kullback divergence behaves here as a squared Euclidean distance.

Likelihood and mutual information

- Maximizing the likelihood $p(X|A)$ for $X = AS$ and $S \sim P_S$ is equivalent to minimizing the (empirical) Kullback mismatch $K[P_Y | P_S]$ with $Y = A^{-1}X$.
- The target source distribution P_S is usually unknown except for $P_S = \prod_i P_{S_i}$.



$$\begin{aligned} K[P_Y | P_S] &= K[P_Y | \prod_i P_{Y_i}] + K[\prod_i P_{Y_i} | \prod_i P_{S_i}] \\ &= I(Y) + \sum_i K[P_{Y_i} | P_{S_i}] \\ &= \text{dependence} + \text{sum of marginal mismatches} \end{aligned}$$

The likelihood suggests to measure dependence by the mutual information $I(Y) = K[P_Y | \prod_i P_{Y_i}]$.

Form Kullback matching to mutual information

- Maximizing the ICA likelihood is equivalent to minimizing $K [P_Y | P_S]$.
- The target source distribution P_S is usually unknown except for $P_S = \prod_i P_{S_i}$.

$$K [P_Y | P_S] = I(Y) + \sum_i K [P_{Y_i} | P_{S_i}]$$

Optimizing over nuisance parameters P_{S_i} kills each $K [P_{Y_i} | P_{S_i}]$ and leaves us with *(in)dependence*:

$$I(Y) \stackrel{\text{def}}{=} K [P_Y | \prod_i P_{Y_i}] \quad \text{a.k.a. } \textit{mutual information}$$

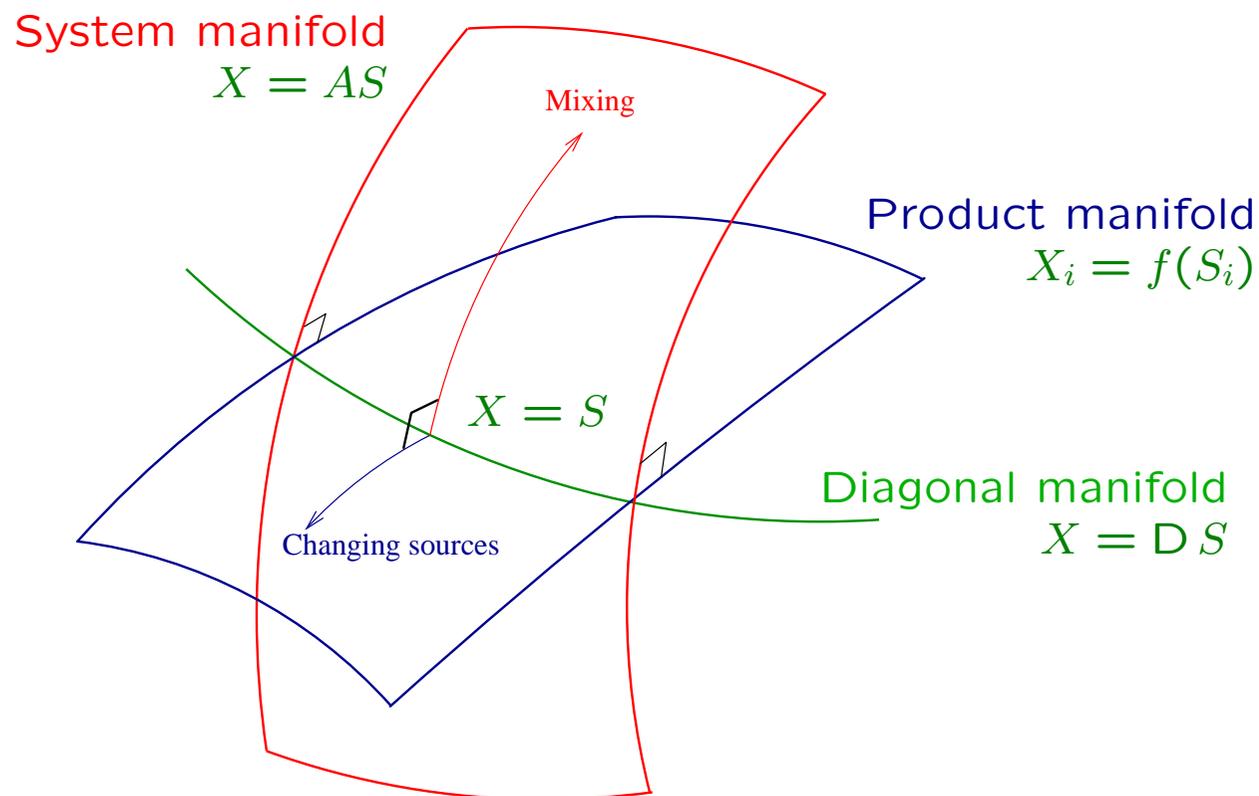
- Maximum likelihood leads to maximum independence. . .
- . . . and provides a definition for it.

A key property:

mixing matrix and source distributions are orthogonal parameters.

- Black board time.

Some statistical manifolds of ICA



- $X = S$: the distribution of a vector S with independent entries. Move from there.
- The *system manifold*: the family of distributions of $X = AS$ for all invertible matrices A .
- The *product manifold*: distort the marginals of S , retaining independence.
- The *diagonal manifold*: change only the scales of the entries of S .

→ ICA is doable: the system manifold intersects the product manifold along the diagonal manifold *only*. And even easy by nuisance orthogonality. But breaks for Gaussian S ...

Classic ICA: mining non Gaussianity

Classic ICA ignores any 'time' structure and thus **must** rely on non Gaussianity. It assumes (explicitly or not) i.i.d. sequences:

$$\begin{aligned} P(Y) &= P(Y(1), Y(2), \dots, Y(T)) \\ &= \prod_t P_t(Y(t)) \quad \text{Independently and ...} \\ &= \prod_t p(Y(t)) \quad \dots \text{ identically distributed.} \end{aligned}$$

where $p()$ is the n -dimensional pdf common to all $Y(t)$, $1 \leq t \leq T$.

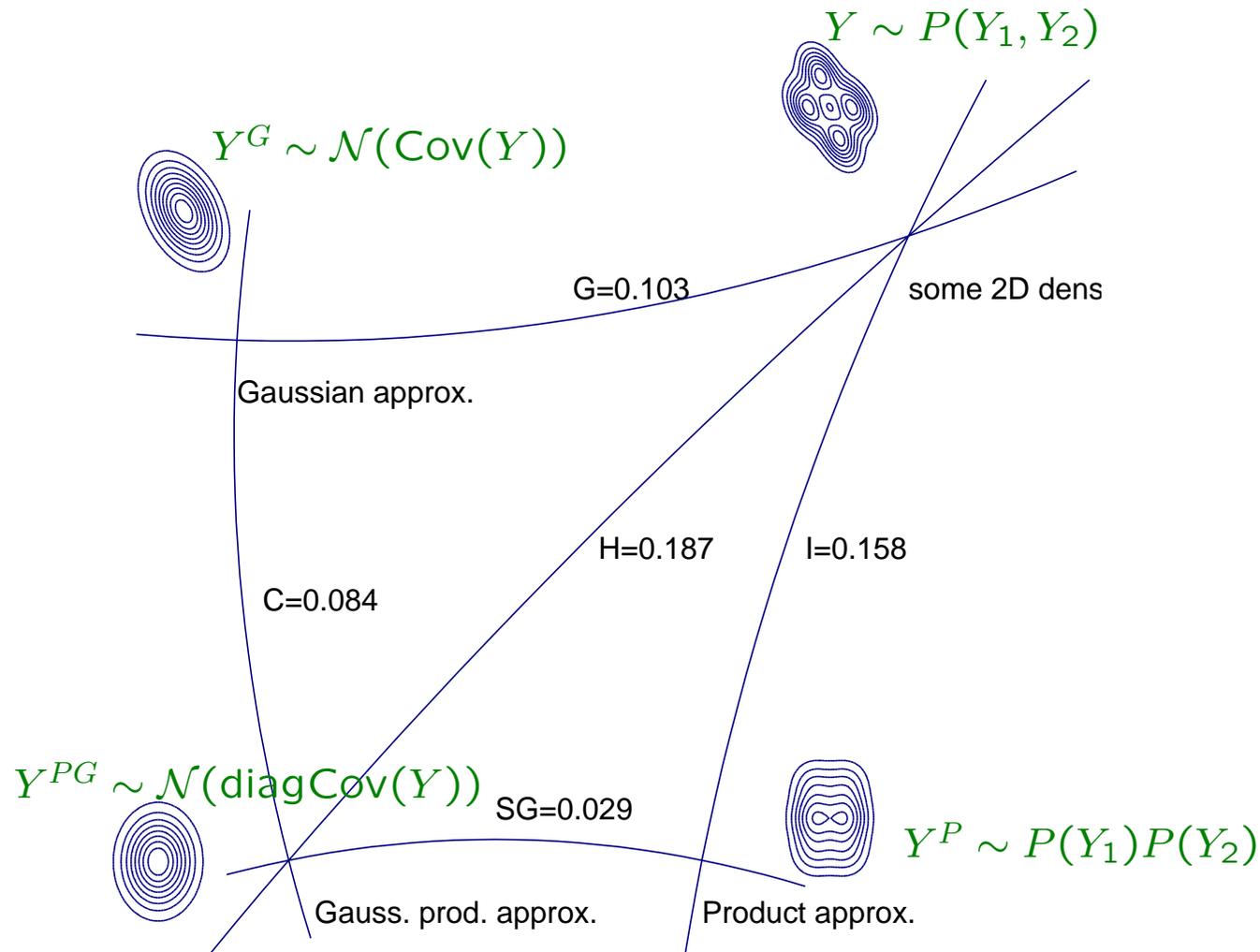
Then, divergence between T -long time series distributions reduces to

$$K [P | Q] = T \times K [p | q]$$

depending only on the n -dimensional distributions p and q .

Note: for readability, we drop the T factor in the following.

Geometry of non Gaussian dependence

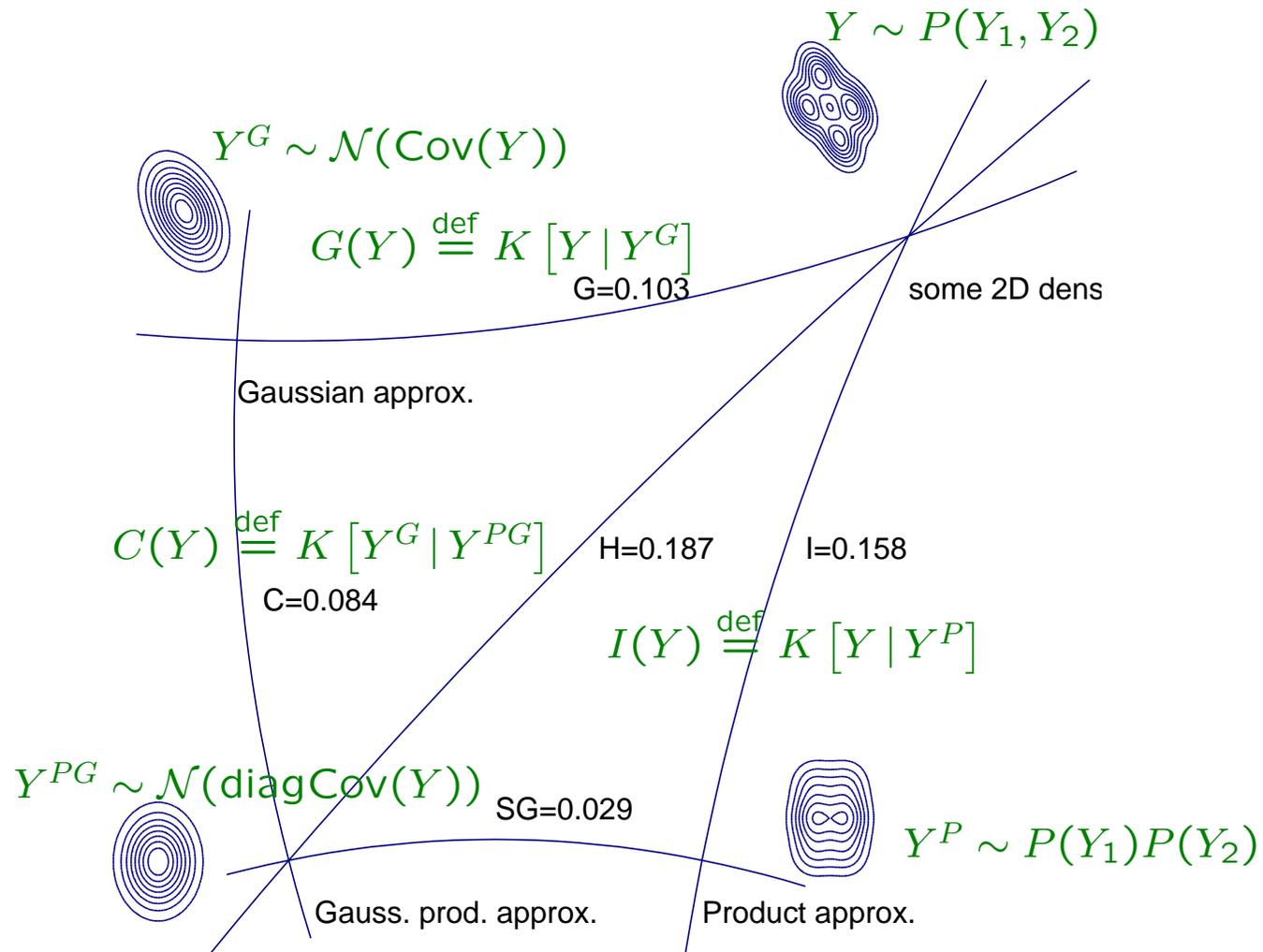


Two familiar statistical approximations can be seen as projections:

Oh, let's assume it's Gaussian...

Ho, let's assume they're independent...

Geometry of non Gaussian dependence



Orthogonal projections \rightarrow two right triangles \rightarrow two Pythagorean theorems

Dependence and non Gaussianity

- *Non Gaussianity*. Define the non Gaussianity $G(Y)$ of Y as

$$G(Y) = K [P_Y | \mathcal{N}(R_Y)]$$

i.e. how much the best Gaussian approx. fails to mimic the distrib. of Y .

- The *correlation* $C(Y)$ of Y

$$C(Y) = K [\mathcal{N}(R_Y) | \mathcal{N}(\text{diag}R_Y)]$$

i.e. how much the covariance matrix R_Y of Y fails to be diagonal.

- All these are (geometrically) connected by

$$I(Y) + \sum_i G(Y_i) = C(Y) + G(Y)$$

- Under *linear* transforms, $G(Y)$ is constant. The mutual information then is

$$I(Y) = C(Y) - \sum_i G(Y_i) + \text{cst}$$

→ Under linear transforms, making the entries of Y as independent as possible ($\min I(Y)$) is *identical* to making (as much as possible) Y uncorrelated and *each of its entries* non Gaussian.

Non Gaussianity (cont.)

Repeat: Under linear transforms, making the entries of Y as independent as possible ($\min I(Y)$) is *identical* to making (as much as possible) Y uncorrelated and *each of its entries* non Gaussian.

- Note 1: The relation $I(Y) = C(Y) - \sum_i G(Y_i) + G(Y)$ also reads

Complicated = Simple - Simple + Complicated constant

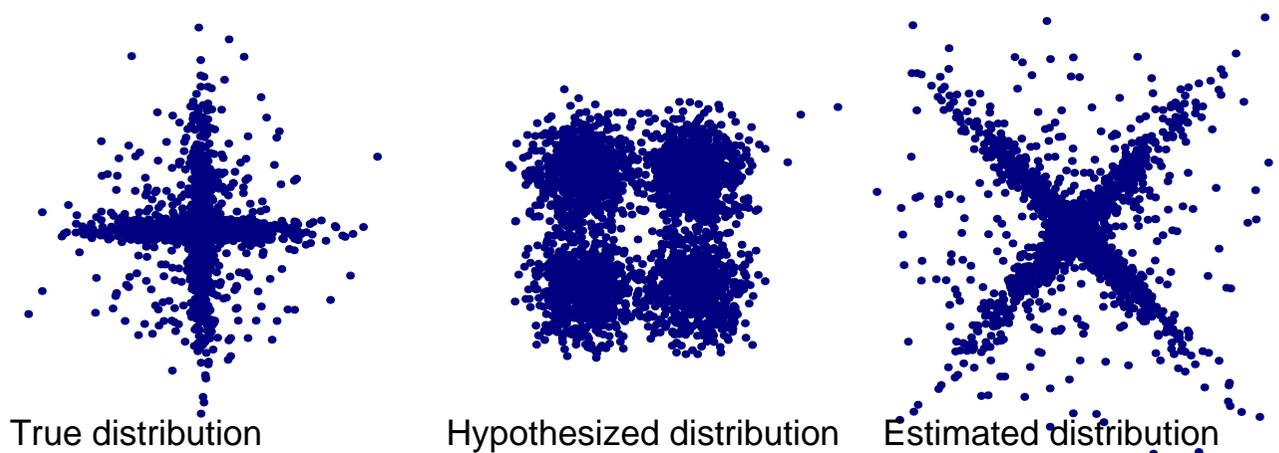
- Note 2: Some people/algos *enforce* $C(Y) = 0$ (pre-whitening) but even that leaves us with the not-so-easy-after-all measure of marginal non-Gaussianity $G(Y_i)$. So, how do we do that?
- Note 3: Connection between non Gaussianity and sparsity.

Non Gaussianity, sparsity and bad luck

To avoid estimating (explicitly or implicitly) the source distributions, one may use fixed targets $P_{S_i} = Q_i$ for each component and just minimize $K [P_Y | \prod_i Q_i]$.

If the pdf's of the sources are known, one should do just that, indeed. If the pdf's are *not* known but are expected to be “sparse”, then one could use *some* sparse guess Q_i in place of the true but unknown distribution,

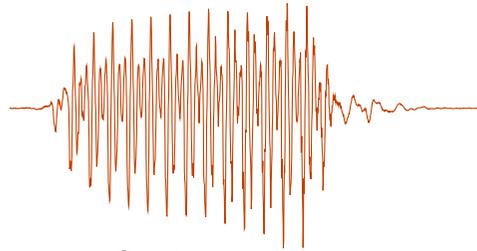
This works (provably) provided the mismatch between P_{S_i} and Q_i is not too large. It does *not* work if a sparse model is used and the true sources are “anti-sparse” or vice-versa.



It's not only wrong; it's maximally wrong!

Other geometries: three points of view on a time series.

A random (!) sequence



Marginal probability density

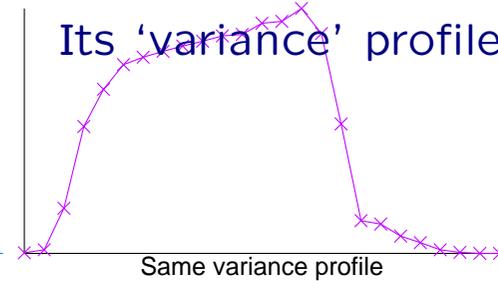
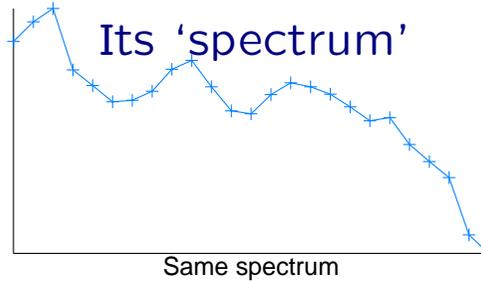
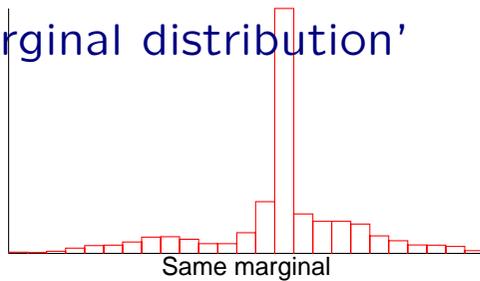
Spectral energy density

Temporal energy density

Its 'marginal distribution'

Its 'spectrum'

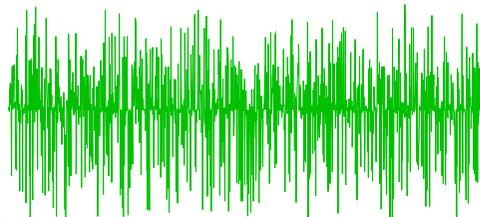
Its 'variance' profile



Same marginal

Same spectrum

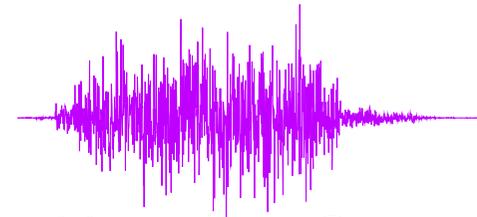
Same variance profile



Non Gaussian i.i.d



Gaussian stationary



Modulated Gaussian i.i.d.

- *All models are wrong, but some are useful* — George Box

Other instances of mutual information by projecting the data onto simple Gaussian models for time series:
stationary but colored, white but non stationary.

Conclusion

Some nice conclusion