

# Basics of statistics

## Getting used to noise

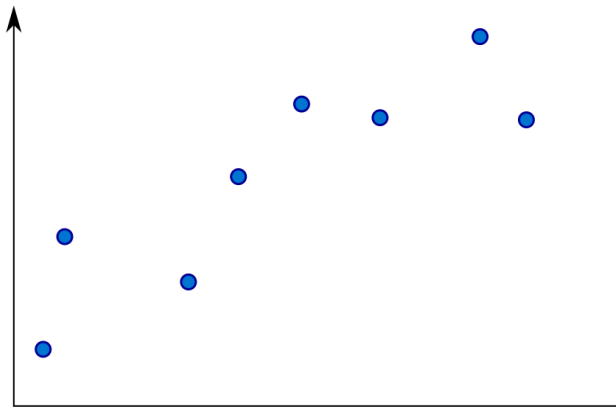
Elena Sellentin

Université de Genève  
Département de Physique Théorique

MACSS Summer school  
8th-12th May 2017

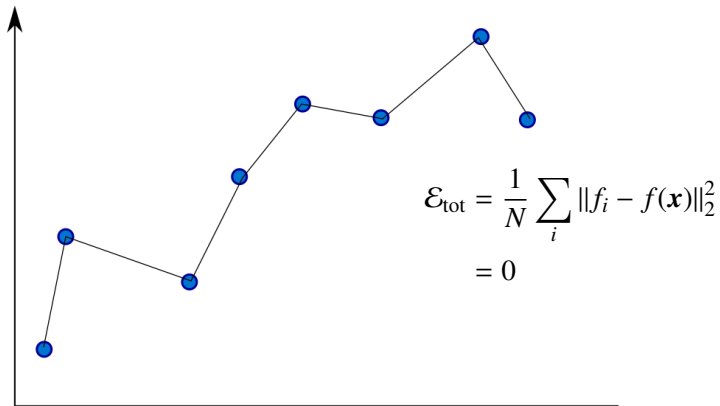
# Cross validation

Given the following data points, how do we prove which curve describes the data best?

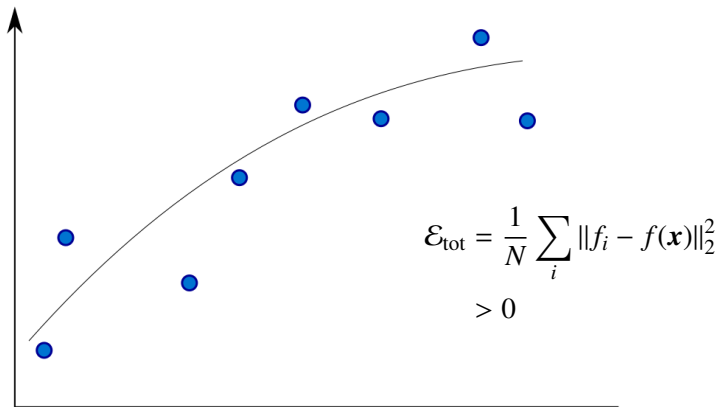


# Cross validation

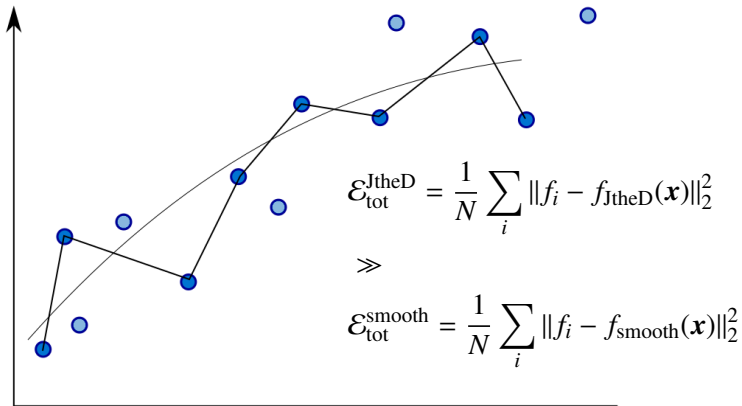
We all know this is wrong, but how do we prove it?



# Cross validation



# Cross validation



# Cross validation

- We saw that join-the-dots ‘perfectly’ explains one data set, but then failed catastrophically on the repeated measurement.
- The fitted curve explained the first dataset somewhat ‘worse’, but then correctly predicted the outcome of the second measurement.

The best model minimizes...

$$\mathcal{E}_{\text{tot}} = \frac{1}{N} \sum_i \|f_i - f(\mathbf{x})\|_2^2 \quad (1)$$

... for current and future measurements.

⇒ It minimizes the distance to taken data, **and** not taken but statistically iid data.

# Measuring in noisy data space

We know how to calculate the distance with a Euclidean metric.

$$\mathbf{x}_1, \mathbf{x}_2 \Rightarrow \sqrt{D(\mathbf{x}_1, \mathbf{x}_2)} = \sqrt{(x_1^1 - x_2^1)^2 + (x_1^2 - x_2^2)^2 + \dots + (x_1^n - x_2^n)^2} \quad (2)$$

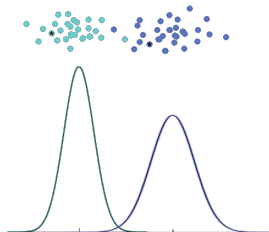
$$D(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 - \mathbf{x}_2)^T \mathbb{I} (\mathbf{x}_1 - \mathbf{x}_2). \quad (3)$$

Now generalize this concept to a noisy space....

# Measuring in noisy data space

Noise  $\Rightarrow$  Distances must become uncertain.

Dimensionful parameters  $\Rightarrow$  metric must measure in those units.



The covariance matrix:

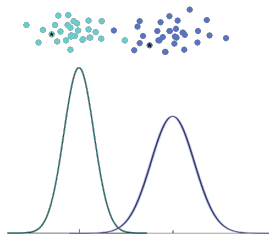
$$C = \langle (x - \langle x \rangle)(x - \langle x \rangle)^T \rangle \quad (4)$$

$$D(x, \mu) = (x - \mu)^T C^{-1} (x - \mu). \quad (5)$$

Mahalanobis distance: measure in units of the expected scatter.



# Measuring in noisy data space



$$D(x, \mu) = (x - \mu)^T C^{-1} (x - \mu). \quad (6)$$

⇒ ‘Distance’ now measures compatibility, expectations and surprises.

Imagine you have taken one measurement  $x_1$  and worked out your measurement errors. Now you repeat and get  $x_2$ .

- Huge  $D(x_1, x_2)$ : Eeew, that came unexpectedly!
- Tiny  $D(x_1, x_2)$ : Someone must be fooling you...

# Consequences of Mahalanobis distances

Introduce parameters: trying to explain the data with a parametric model.

$$D(\mathbf{x}, \boldsymbol{\mu}(\mathbf{p})) = (\mathbf{x} - \boldsymbol{\mu}(\mathbf{p}))^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}(\mathbf{p})). \quad (7)$$

- Which values do the parameters need to take, such that my model is most compatible with the data?  
 $\Rightarrow$  minimize  $D(\mathbf{x}, \boldsymbol{\mu}(\mathbf{p})) \Rightarrow$  'Fitting' procedure.
- How much wiggle room is then still left?  $\Rightarrow \hat{\mathbf{p}} \pm \Delta \mathbf{p}$ .
- Which  $D_{\text{crit}} = D_{\min}(\mathbf{x}, \boldsymbol{\mu}(\hat{\mathbf{p}})) + \Delta D(\mathbf{x}, \boldsymbol{\mu}(\mathbf{p}))$  is impossibly far away?

$\rightarrow$  where  $\Delta \chi^2$  comes from (just **one** example).

$\rightarrow$  depends on the chosen metric:  $\mathbf{C}, \mathbb{I}, \mathbf{F}^{-1}$

# Likelihoods

If  $\mathbf{x}$  is random, then any  $f(\mathbf{x})$  is random, hence  $D(\mathbf{x}, \mu(\mathbf{p}))$  is also a random variable.

⇒ If possible, work with a full **likelihood**.

Most famously:

$$\mathcal{G}(\mathbf{x}, \mu(\mathbf{p}), \mathbf{C}) = \frac{1}{\sqrt{|2\pi\mathbf{C}|}} \exp\left(-\frac{1}{2}[\mathbf{x} - \mu(\mathbf{p})]^T \mathbf{C}^{-1}[\mathbf{x} - \mu(\mathbf{p})]\right). \quad (8)$$

Minimizing  $\chi^2 \leftrightarrow$  maximizing a Gaussian likelihood.

# Rules for the Game of Noise

- Random variables are drawn from probability distribution functions. We write  $x \sim \mathcal{P}(x)$  in the univariate case,  $\mathbf{x} \sim \mathcal{P}(\mathbf{x})$  in the multivariate case and  $\mathbf{S} \sim \mathcal{P}(\mathbf{S})$  in the matrix-variate case.
- For any function  $f$ , if  $x$  is a random variable, so is  $f(x)$ .
- $P \geq 0$
- $\sum_i P_i = 1$ , especially  $P(A) + P(\bar{A}) = 1$ , or  $\int \mathcal{P}(x) dx = 1$
- $P(A, B|C) = P(A|C)P(B|A, C) = P(B|C)P(A|B, C)$ ,
- Exchange of variables:

$$g(\mathbf{y}) = f(\mathbf{x} = \tilde{\mathbf{y}}(\mathbf{x}))|J| \quad (9)$$

where the Jacobian is

$$|J| = \det \left( \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right). \quad (10)$$

- Bayes' Theorem:  $P(A|B, C) = \frac{P(A|C)P(B|A, C)}{P(B|C)}$

# Consequences

- Estimators, e.g. best-fitting parameters, are random variables (they are a function of the data).
- Sometimes you must get, and will get, ‘untypical’ measurements. They come from the tail of the distributions  $\mathcal{P}(x)$ .
- Non-linear functions of Gaussian random variables follow non-Gaussian distributions, and the CLT agrees with this.
  - the ratio of two Gaussian rv follows a Cauchy distribution
  - the exponential of a Gaussian rv is log-normally distributed
  - the absolute value of Fourier modes from a Gaussian random field follow a Rayleigh distribution
  - many more  $\Rightarrow$  Exercises.

## Bayesian Inference

# Bayesian Inference

**Inverse Problem:** Given data  $x$ , how do we explain them?

- Given: a dataset  $x$  and a model  $\mathcal{M}(p)$
- To which posterior do the data  $x$  constrain the parameters  $p$ ?
- Given two models  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , which one explains the data better?

$$\mathcal{P}(\theta, \mathcal{M}|x) = \frac{L(x|\theta, \mathcal{M})\pi(\theta)}{\epsilon(x|\mathcal{M})} \quad (11)$$

# Bayesian Inference

$$\mathcal{P}(\theta, \mathcal{M}|x) = \frac{L(x|\theta, \mathcal{M})\pi(\theta)}{\epsilon(x|\mathcal{M})} \quad (12)$$

- $\mathcal{P}(\theta, \mathcal{M}|x)$ : the posterior.
- $L(x|\theta, \mathcal{M})$ : the likelihood.
- $\pi(\theta)$ : the priors.
- $\epsilon(x|\mathcal{M})$ : the evidence ('marginal likelihood').

→ Even for a known likelihood, the posterior can be difficult to obtain. Need sampling techniques (MCMC, Gibbs...) & DALI.