# Classical statistical inference
## *Regression and Model fitting*

Associated notebook:

[05-MLE_and_regression/Regression_short.ipynb](05-MLE_and_regression/Regression_short.ipynb)
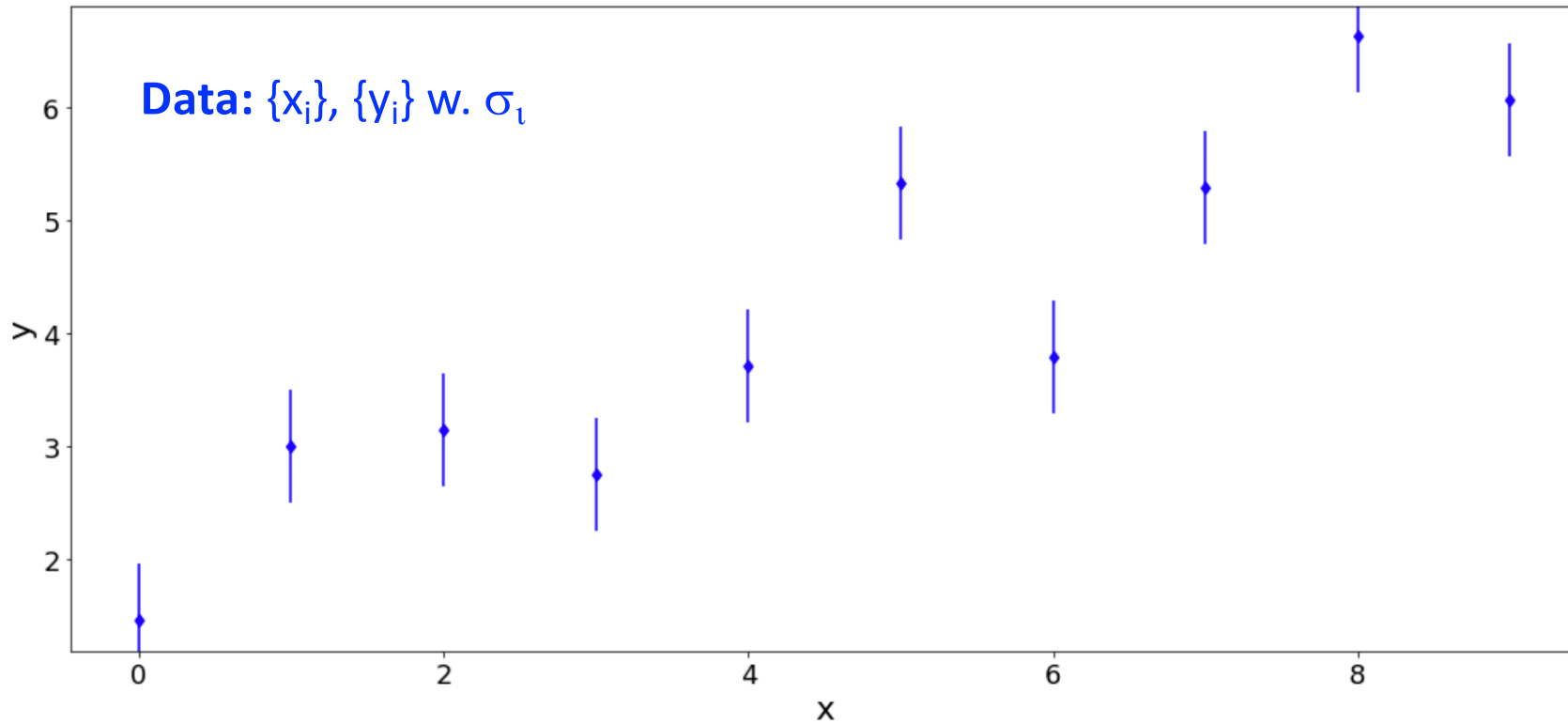
# Regression and model fitting

*Problem:* the quantities of interest are parameters of a model, not the RV that you measure

Examples

| Observation | Quantity of interest | Model |
|---|---|---|
| Position of a star: $\mathbf{x}(t)$ | Proper motion (velocity) of the star | $V = f(\mathbf{x}, t, …)$ |
| Photometry of an asteroid: mag(t) | P (period of rotation) | mag = f(t, P, …) |
| Transit of a planet: mag(t) | P (period), e (eccentricity), D (dist to star) | $\Delta m = f(t, P, e, D, …)$ |
| Spectrum of a QSO: F ($\lambda$) | $M_{BH}$ (Black hole mass of QSO) | FWHM = $f(M_{BH}, L, …)$ |

# Regression and model fitting

*Problem:* You measure $D \equiv (\{y_i\}, \{x_i\})$



**Data:** $\{x_i\}, \{y_i\}$ w. $\sigma_\iota$

**Model:** $M(\theta)$

# Regression and model fitting

*Problem:* You measure D ≡ ({$y_i$}, {$x_i$})



**Data:** {$x_i$}, {$y_i$} w. $\sigma_\iota$

M($\theta$): y = a * x + b

**Model:** M($\theta$)

y = a * x + b

$\theta$ = a, b

# Regression and model fitting

*How to find a good model?*
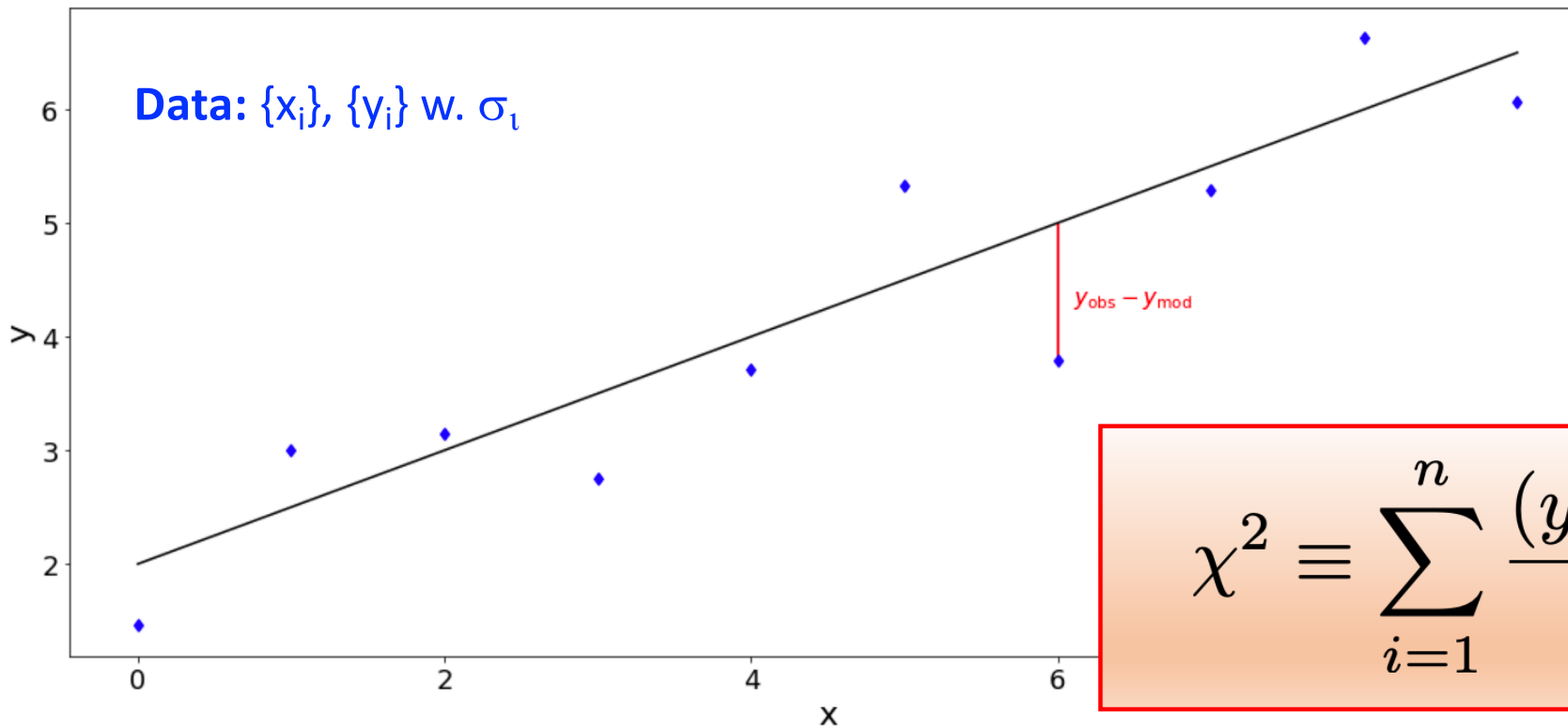


**Data:** $\{x_i\}$, $\{y_i\}$ w. $\sigma_\iota$

$y_{obs} - y_{mod}$

**Model:** $M(\theta)$

$y = a * x + b$

$= f(x \mid \theta)$

$\theta = a, b$

Minimize

$$\chi^2 \equiv \sum_{i=1}^{n} \frac{(y_i - y_{i,mod})^2}{\sigma_i^2}$$

# Regression and model fitting

$$\chi^2 \equiv \sum_{i=1}^{n} \frac{(y_i - y_{i,mod})^2}{\sigma_i^2}$$

If $\sigma_i = 1$ : Least square regression

If $\sigma_i \neq 1$ : chi-square regression

The $\chi^2$ is called a **merit** function

When uncertainties between variables are correlated, the $\chi^2$ is expressed:

$$\chi^2 = \sum_{i=1}^{n} \sum_{l=1}^{n} (y_i - y_{i,mod}) F_{i,l} (y_l - y_{l,mod})$$

Where *F* is the Fisher matrix

$$[F]^{-1} = [C] = \begin{bmatrix} \sigma_k^2 & \sigma_{kl} \\ \sigma_{kl} & \sigma_l^2 \end{bmatrix}$$

# Link between χ² and likelihood $\boxed{L \;=\; p(D \,|\, M(\boldsymbol{\theta}))}$

Case of a straight line:  $y_i = \theta_0 + \theta_1\, x_i + \epsilon_i$  with  $\epsilon_i \sim N(0, \sigma_i)$

For each $y_k$ we have:  $p(y_k \,|\, \mu, \sigma) = \dfrac{1}{\sigma \sqrt{2\,\pi}}\, \exp\left[-0.5 \left(\dfrac{y_k - \mu}{\sigma}\right)^2\right]$

Hence, we have for our data set D:

$$L \equiv p(\{y_i\} \,|\, \{x_i\}, \boldsymbol{\theta}) = \prod_{i=1}^{N} \frac{1}{\sigma_i \sqrt{2\,\pi}}\, \exp\left[\left(\frac{-(y_i - (\theta_0 + \theta_1\, x_i))^2}{2\,\sigma_i^2}\right)\right]$$

# Link between $\chi^2$ and likelihood

Hence, we have for our data set D:

$$L \equiv p(\{y_i\} \,|\, \{x_i\}, \boldsymbol{\theta}) = \prod_{i=1}^{N} \frac{1}{\sigma_i \sqrt{2\,\pi}} \, \exp\left[\left(\frac{-(y_i - (\theta_0 + \theta_1\, x_i))^2}{2\,\sigma_i^2}\right)\right]$$

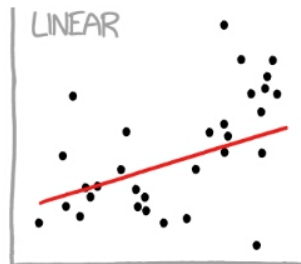$$\ln(L) \propto \sum_{i=1}^{N} \left(\frac{-(y_i - (\theta_0 + \theta_1\, x_i))^2}{2\,\sigma_i^2}\right)$$

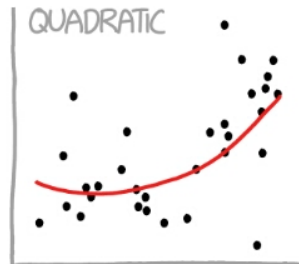$$\chi^2 \equiv \sum_{i=1}^{n} \frac{(y_i - y_{i,mod})^2}{\sigma_i^2}$$

$$\Rightarrow \quad L \propto \exp(-\chi^2/2)$$

Minimizing $\chi^2$ is equivalent to maximizing L

# Regression and model fitting



https://xkcd.com/2048/

# How to choose a suitable regression method ?

The way you will tackle a regression problem may depend on:

- **Linearity:** is the model linear *in its parameters?* $\quad f(x \,|\, \boldsymbol{\theta}) = \sum_{p=1}^{k} \theta_p g_p(x)$

- **Complexity:** large number of parameters increase complexity and covariance matrix on uncertainties

- **Error behaviour:** uncertainties on dependent and independent variable and their correlation.

# How to choose a suitable regression method ?

| Frequentist: (this lecture) | Bayes (future lecture): |
|---|---|
| *Optimization* with some merit function | *Sampling* of the likelihood |
| Search for *best* (fit) model *parameters* | PDF on parameters |
| Often when *simple* error behaviour | More *complex* error behaviour |

# Linear vs non linear regression

A model is **linear** if:

$$f(x \mid \boldsymbol{\theta}) = \sum_{p=1}^{k} \theta_p g_p(x)$$

$g_p(x)$ can be a non linear function of *x* BUT does not depend on any free parameter

In this case, the values of the parameters that yield $\dfrac{\partial \ln(L)}{\partial \theta_i} = 0$ (max. likelihood) can be found "analytically"

When the model is **not linear**, the minimization of the $\chi^2$ has to be performed *numerically*

# Linear vs non linear regression

**Linear** model fitting:   See Sect. IV.1

  Python implementation: `numpy.polyfit(x, y, deg=1, w=1/sigma)`

**NON Linear** model fitting:   See Sect. IV.3

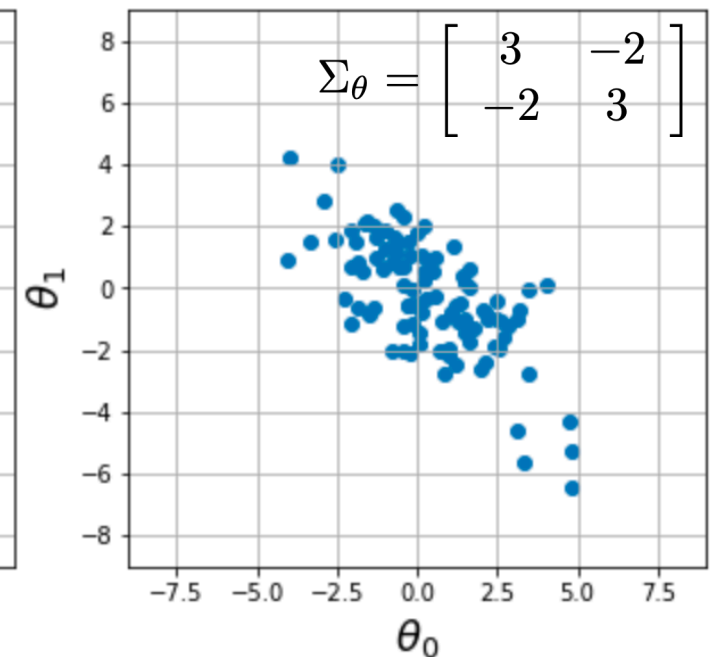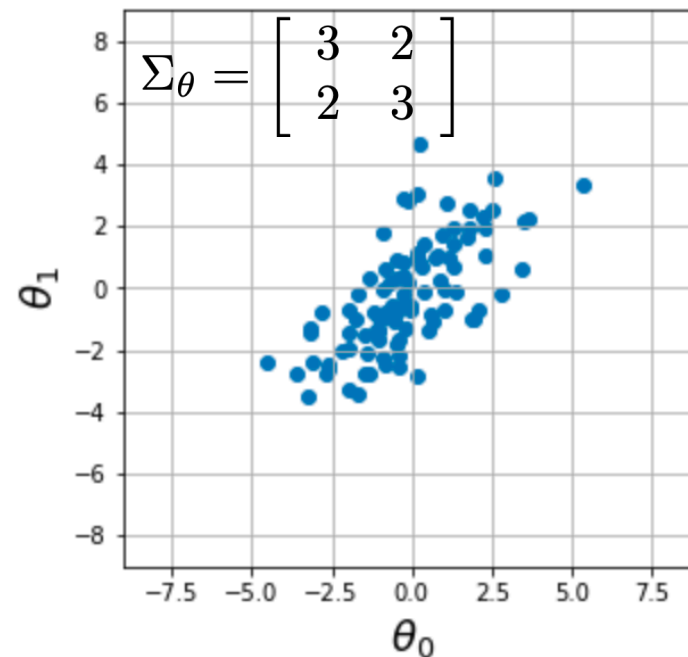  Python implementation: `scipy.optimize.curvefit()`

Go to Sect. IV.1.1 of the Notebook for practical example

# Uncertainty on the fitted parameters

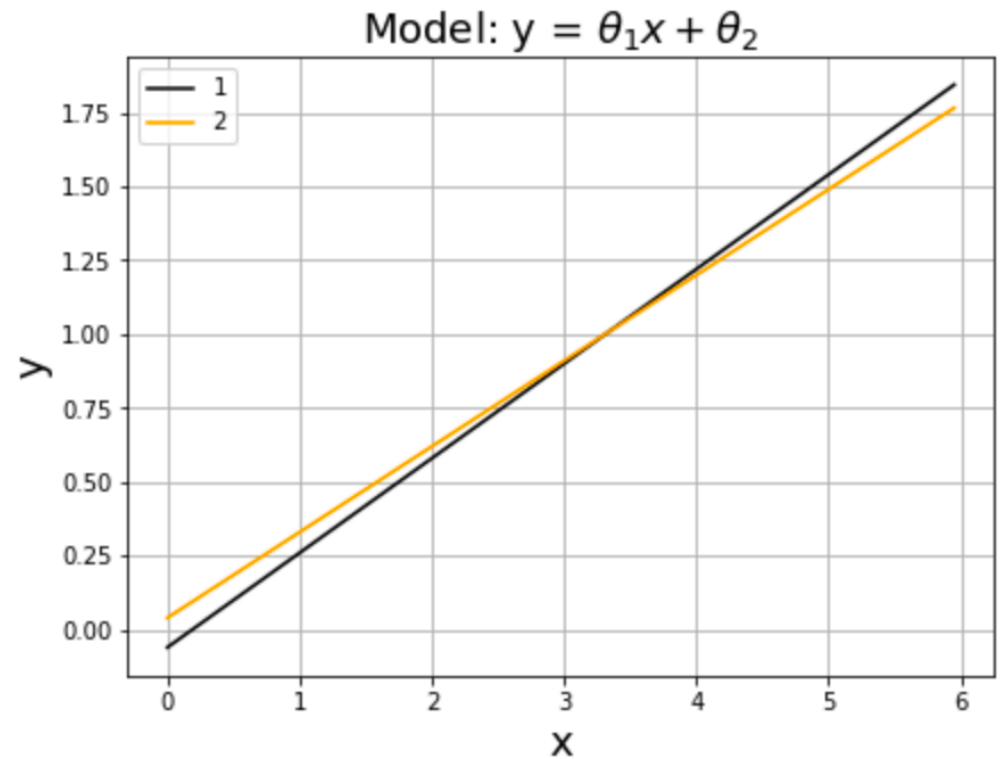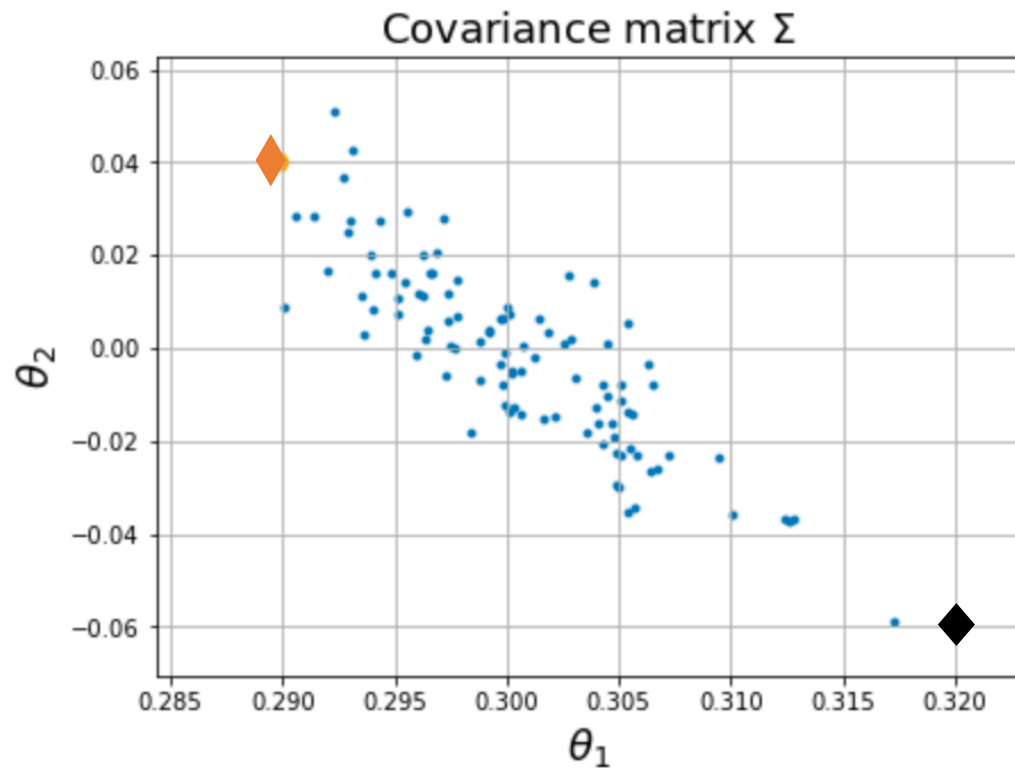The python functions return a covariance matrix (Warning : use arg. cov=True)

The diagonal elements of the matrix give the variance on the parameters (uncertainty$^2$)

$$\Sigma_\theta = \begin{bmatrix} \sigma^2_{\theta_0} & \sigma_{\theta_0\theta_1} \\ \sigma_{\theta_0\theta_1} & \sigma^2_{\theta_1} \end{bmatrix}$$

# Uncertainty on the fitted parameters

$$\Sigma_\theta = \begin{bmatrix} + & - \\ - & + \end{bmatrix}$$



Covariance matrix $\Sigma$

Model: $y = \theta_1 x + \theta_2$

# Uncertainty on the fitted parameters

$$\Sigma_\theta = \begin{bmatrix} \sigma^2_{\theta_0} & \sigma_{\theta_0\theta_1} \\ \sigma_{\theta_0\theta_1} & \sigma^2_{\theta_1} \end{bmatrix}$$
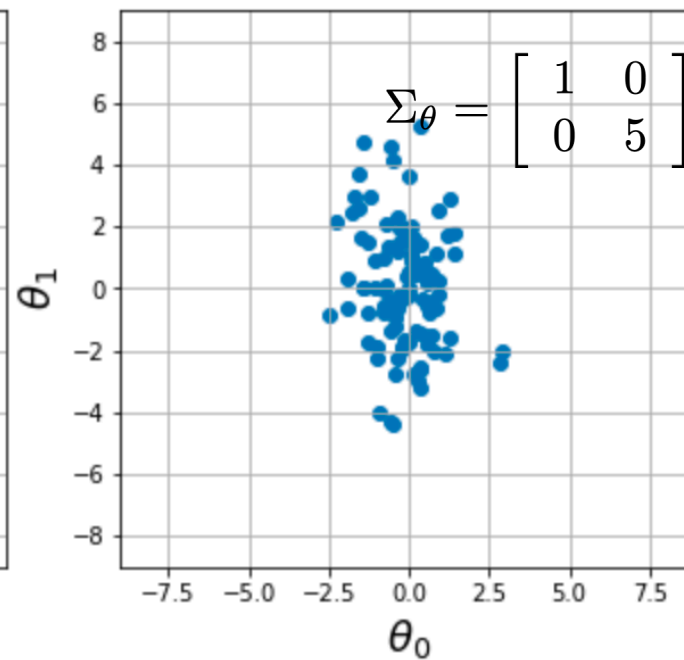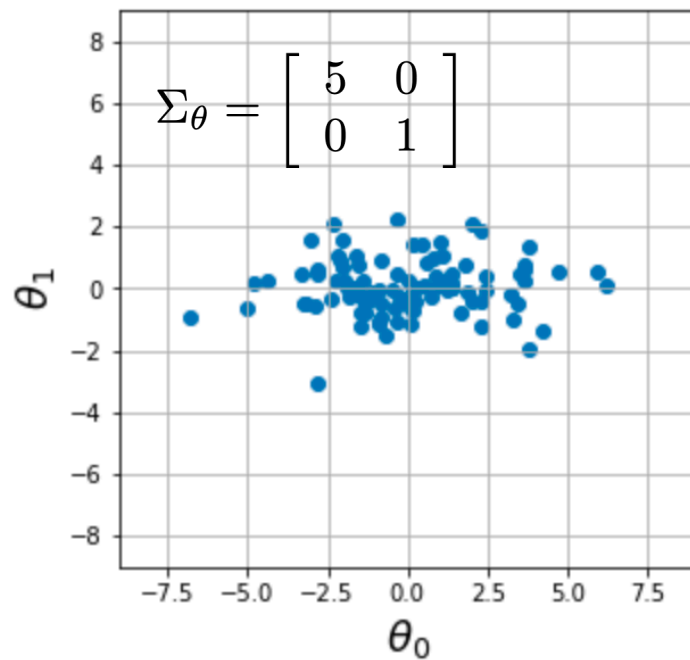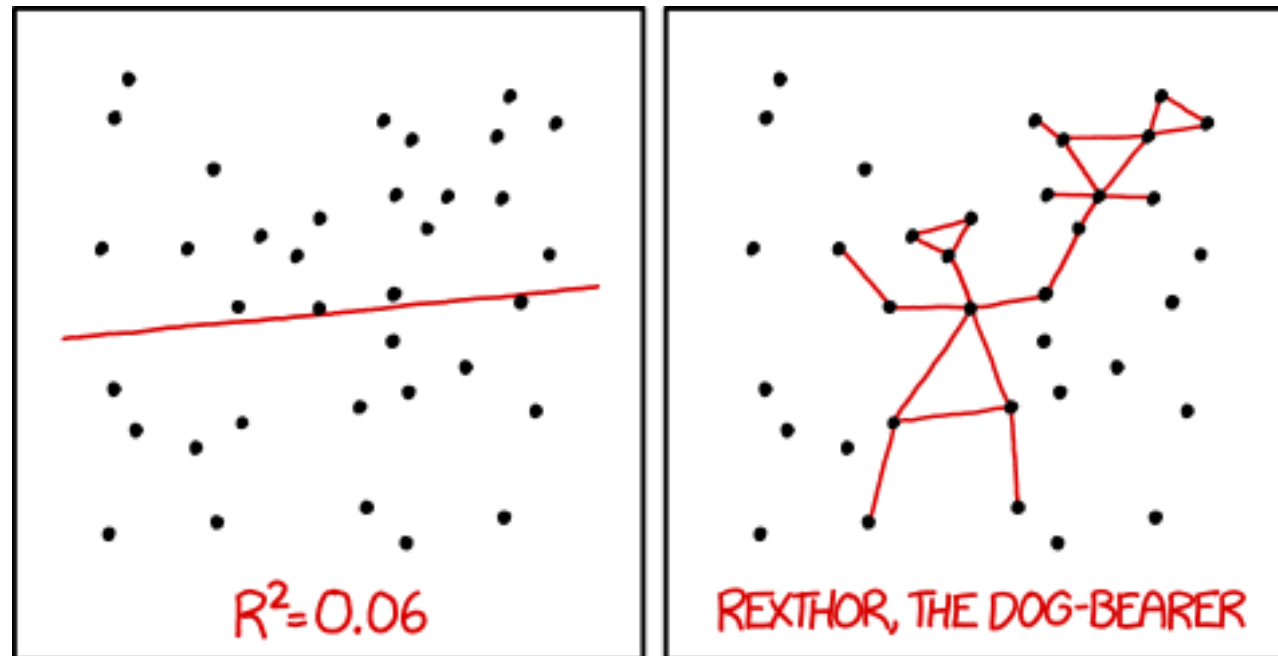
# Quality of the regression



R²=0.06

REXTHOR, THE DOG-BEARER

I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

xkcd.com/1725/

# Quality of the regression

Your $\chi^2$ is a random variable !

$$Q = \sum_{i=1}^{k} z_i^2 \rightarrow p(Q|k) = \frac{1}{(2\,\Gamma(k/2))}(Q/2)^{k/2-1}\exp(-Q/2)$$

k = **d**egree **o**f **f**reedom = N *points* − n *parameters*

If you fit a model with 2 parameters on a set of 100 points => 98 d.o.f.

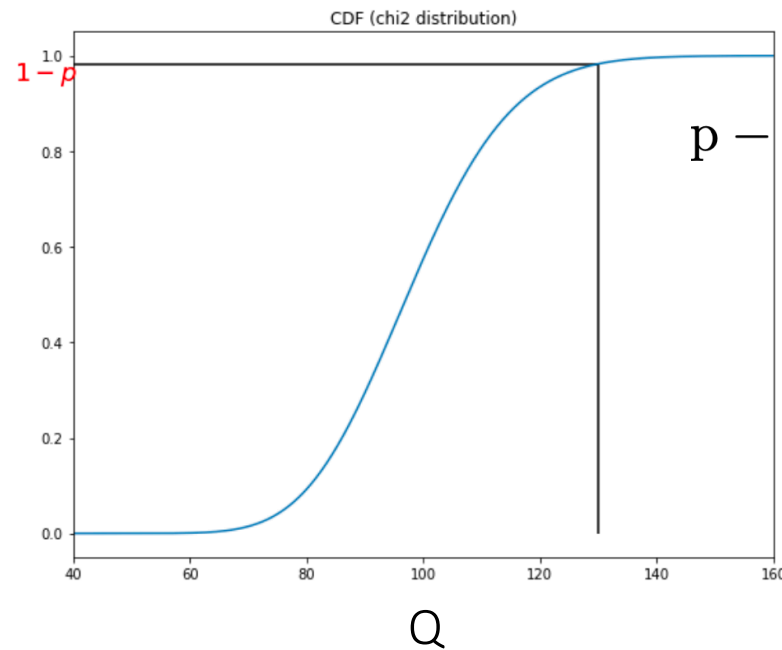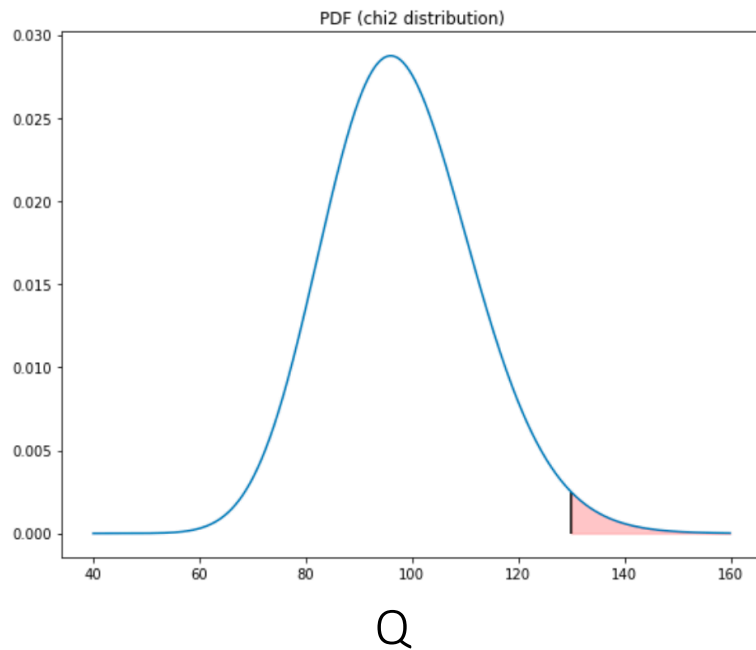Expectation $E(\chi^2) = 100 - 2 = 98$

Reduced $\chi^2$ : $\chi^2_{red} = \chi^2_{dof}$ / d.o.f.     ⇒     Reduced $\chi^2$ ≡ 1. if good fit

See also the Notebook 03-Basic_statistics_and_proba_concepts/Descriptive_statistics_02.ipynb

# Quality of the regression

$$p(Q|k) = \frac{1}{(2\,\Gamma(k/2))}(Q/2)^{k/2-1}\exp(-Q/2)$$



PDF (chi2 distribution)

Q



CDF (chi2 distribution)

Q

$$\mathrm{p-value} = \mathrm{p}(Q \geq \chi^2_{\mathrm{obs}})$$

$$= 1 - p(Q \leq \chi^2_{\mathrm{obs}})$$

CDF

Typically:
p-value < 0.05 : 🙁
0.05 < p-value < 1: 😀
p-value close to 1 : 🙁

```
1-scipy.stats.chi2.cdf(chi2_data, df= len(data)-nparam)
```

Go to Sect. IV.1.1 of the Notebook for practical example