

Statistical analysis

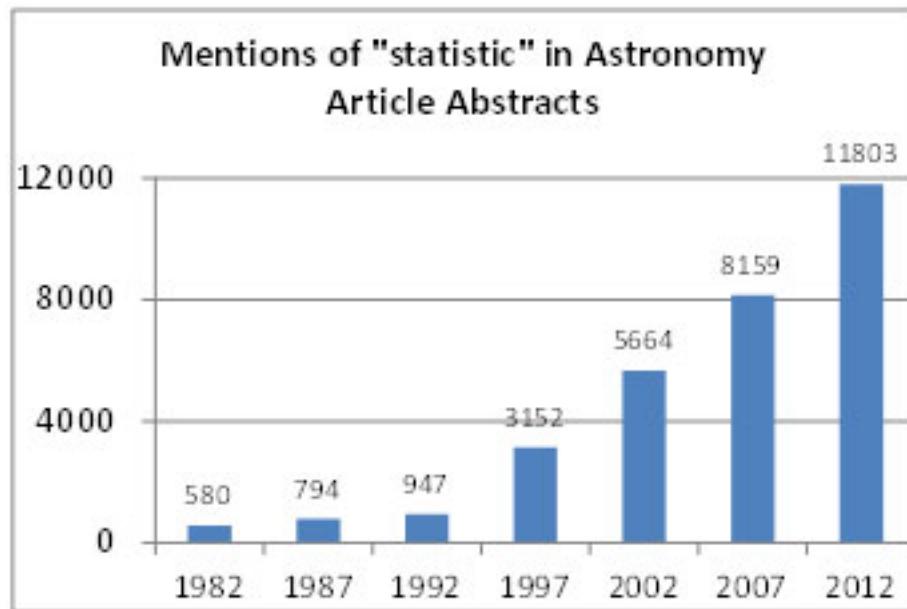
Amy Lien
Goddard Space Flight Center

Acknowledgement

Class material from:

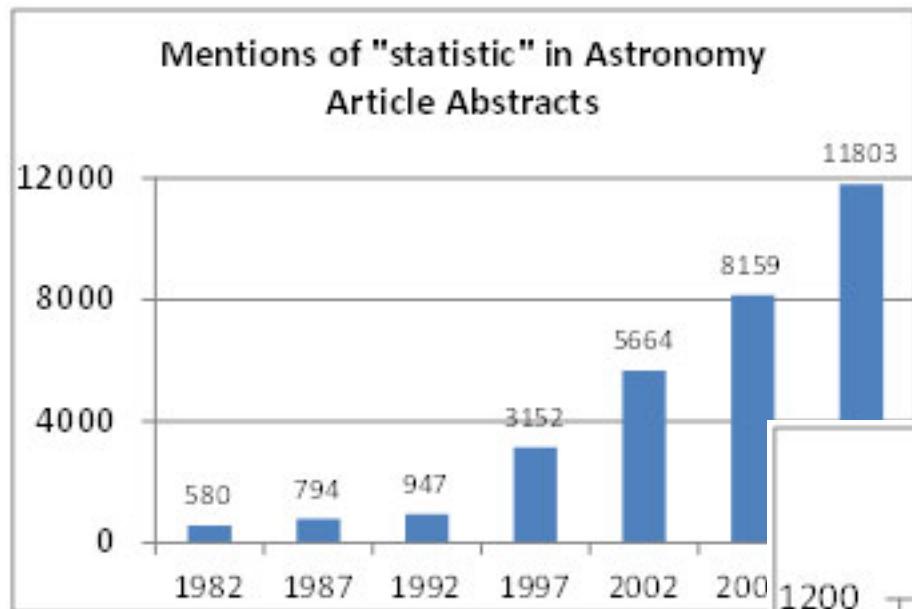
- **Slides from Prof. M.A. Thomson (Univ. of Cambridge)**
<http://www.hep.phy.cam.ac.uk/~thomson/lectures/lectures.html>
- **Prof. Sylvain Guiriec (George Washington Univ.)**

Why statistics?

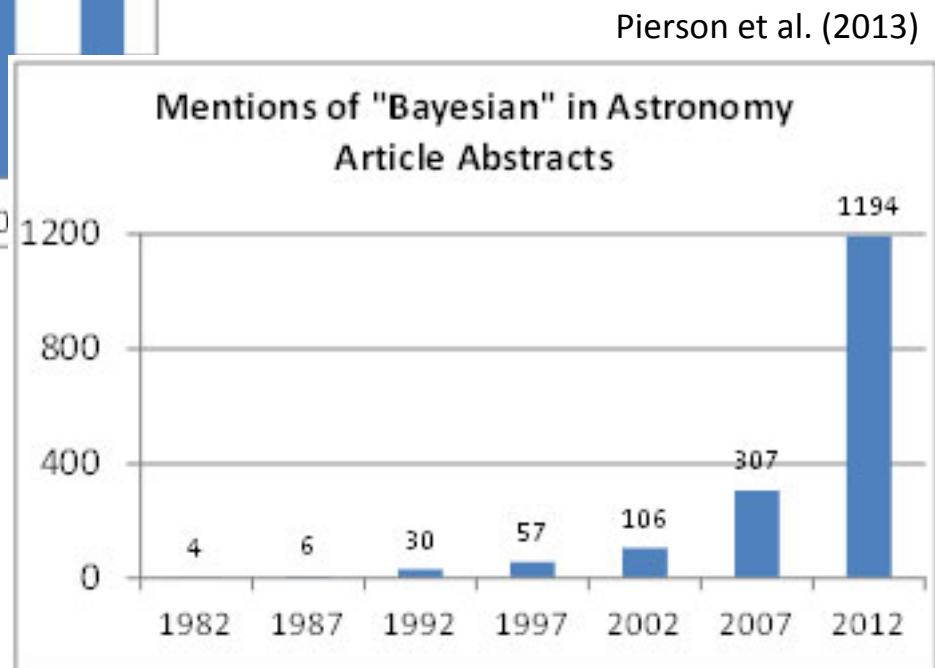


Pierson et al. (2013)

Why statistics?



Pierson et al. (2013)



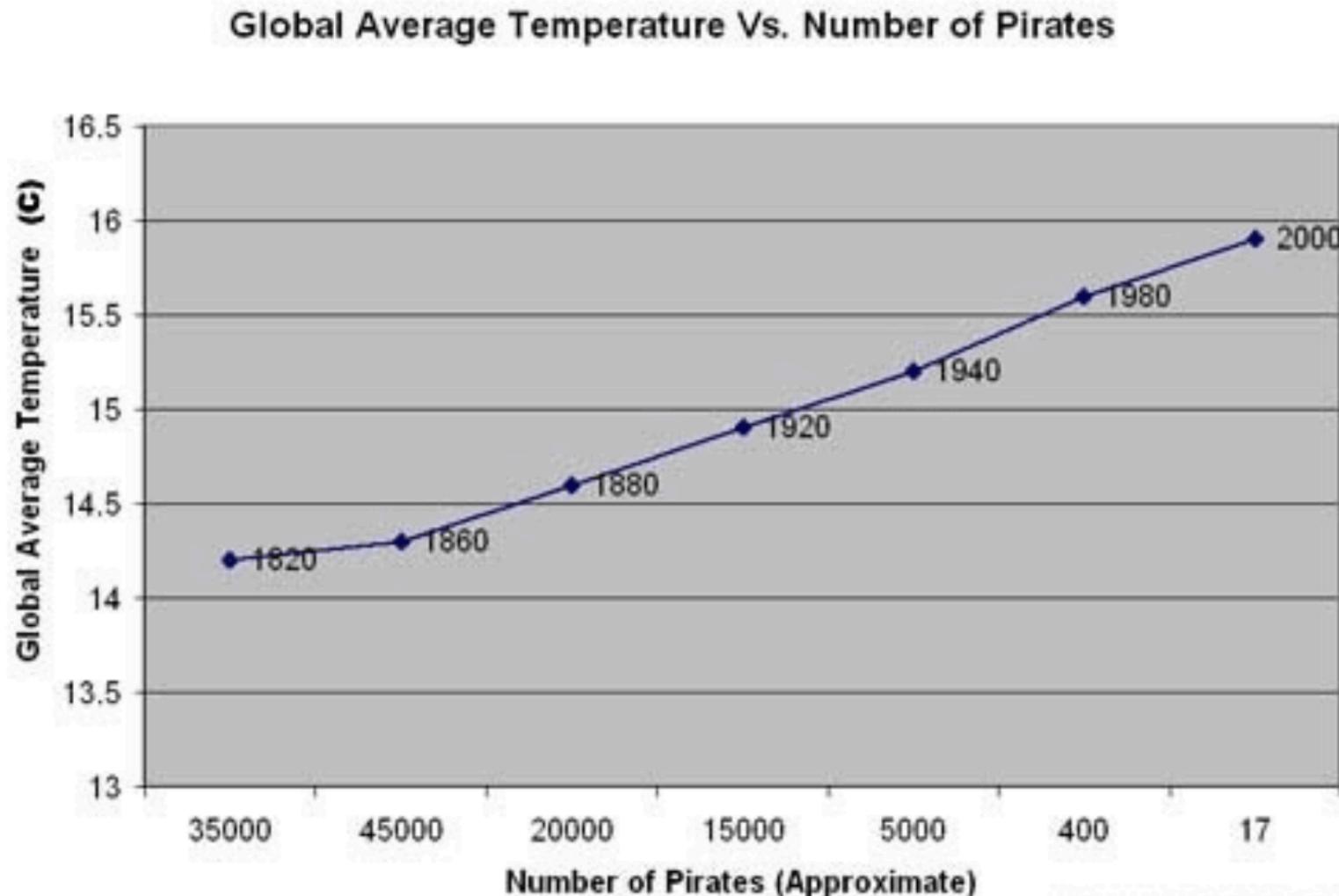
Pierson et al. (2013)

Most astronomers and physicists don't understand statistics

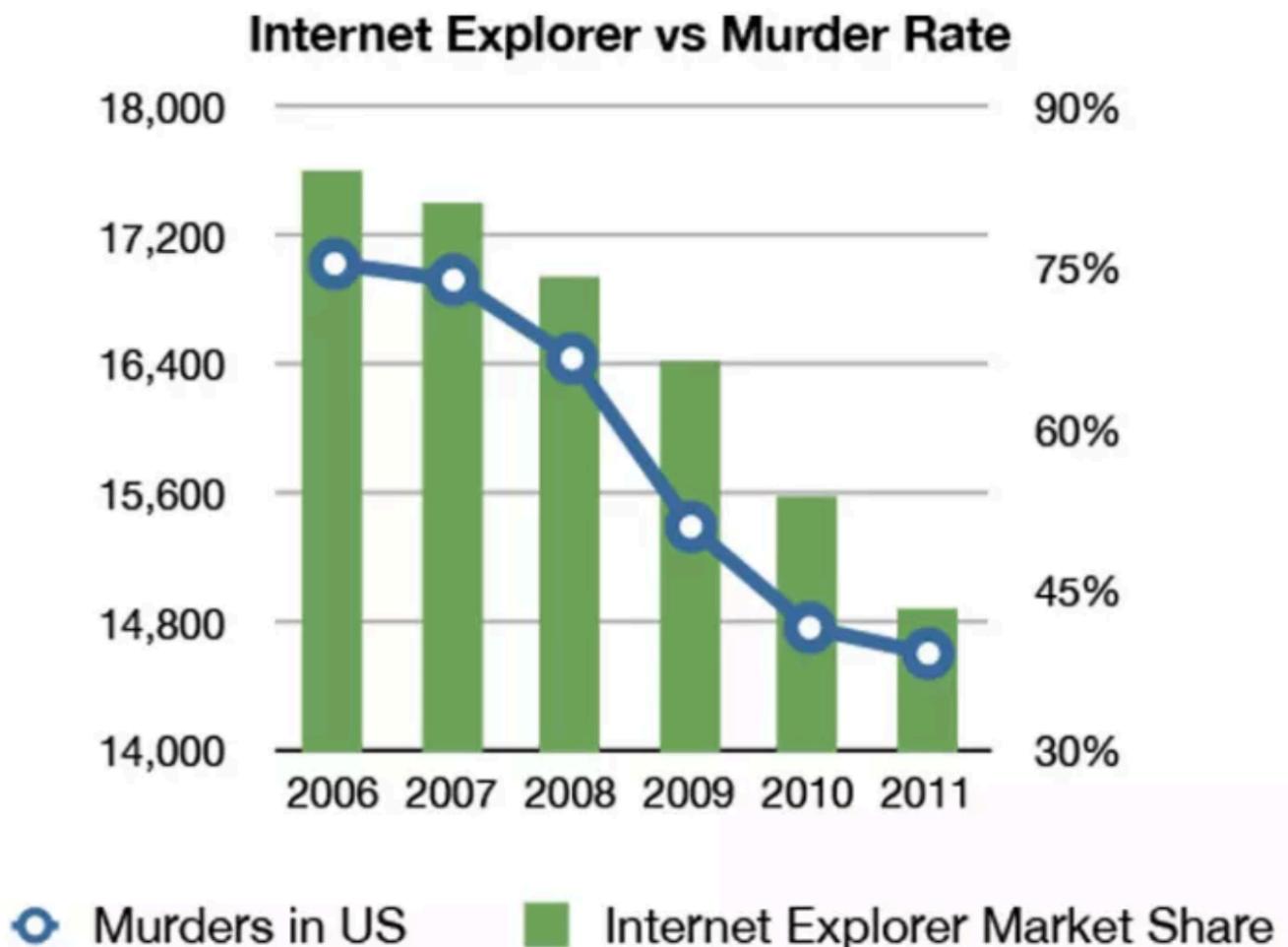
Suppose you derive a 90% confidence region on some parameter of interest. What does this mean ?

Does it tell you that there is a probability of 0.9 that the true value of the parameter lies in the range calculated ?

Be careful with statistics analysis



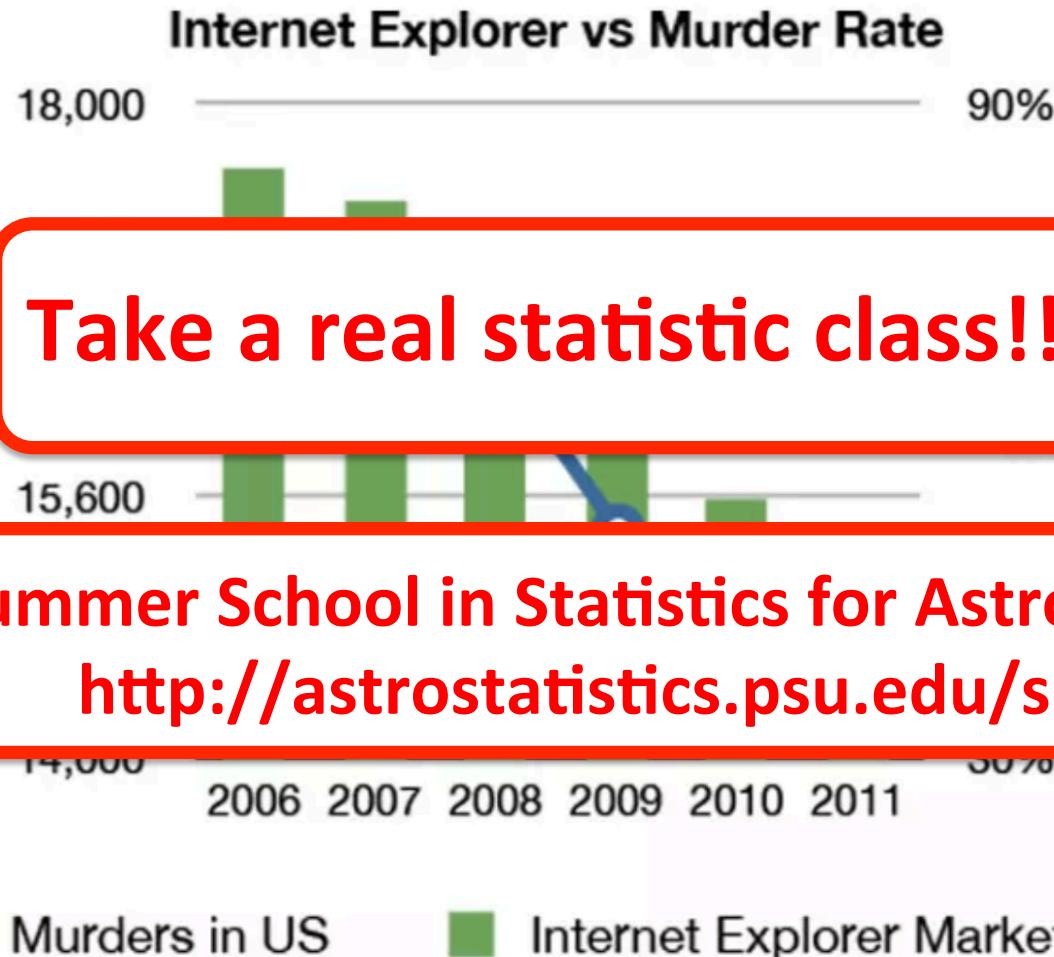
Be careful with statistics analysis



Be careful with statistics analysis



Be careful with statistics analysis

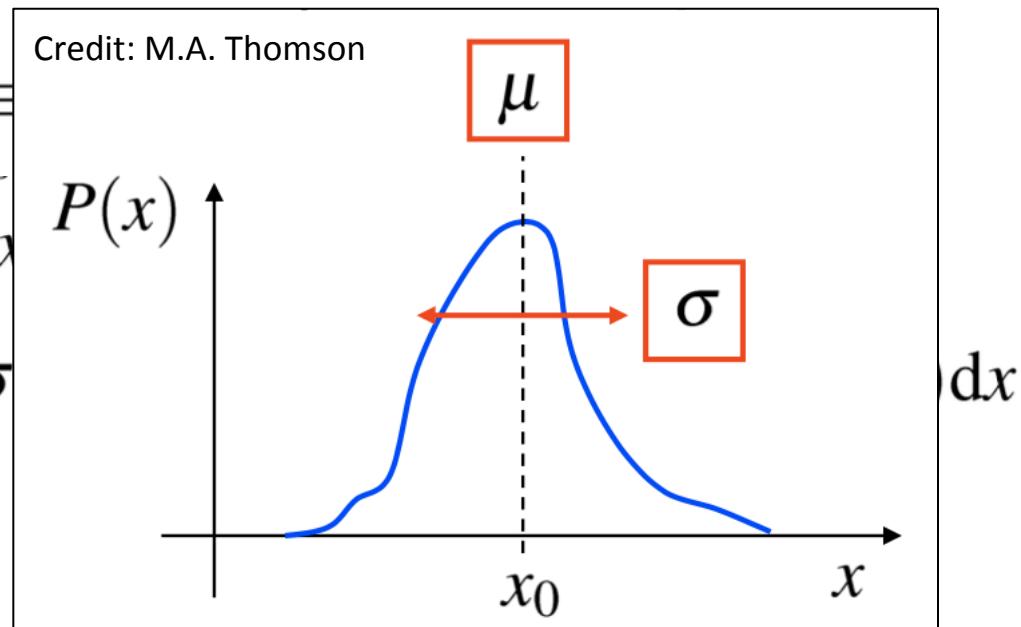


Basic statistical terms

- Mean (average): $\mu \equiv \langle x \rangle = \int x P(x) dx$
- Mean of squares: $\langle x^2 \rangle = \int x^2 P(x) dx$
- Variance: $Var(x) \equiv \sigma^2 \equiv \langle (x - \mu)^2 \rangle = \int (x - \mu)^2 P(x) dx$
 $= \langle x^2 \rangle - \mu^2$
- Standard deviation $\sigma = \text{sqrt}(\text{variance})$
- Degrees of freedom
 $= \text{number of data} - \text{number of variables}$

Basic statistical term

- Mean (average): $\mu \equiv \langle x \rangle$
- Mean of squares: $\langle x^2 \rangle$
- Variance: $Var(x) \equiv \sigma^2$



- Standard deviation $\sigma = \sqrt{\text{variance}}$
- Degrees of freedom
= number of data – number of variables

Poisson distribution

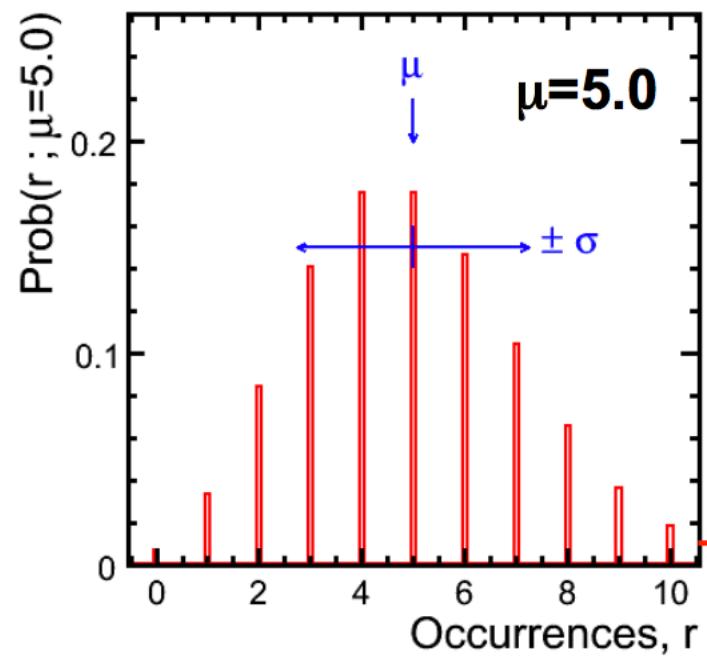
- Most of the high-energy instruments simply count photons → Poisson distribution
- Poisson distribution:

$$P(x) = \frac{\mu^x e^{-\mu}}{x!}$$

$$\langle x \rangle = \mu$$

$$\sigma^2 = \mu$$

μ : Average of the distribution
 σ : Standard deviation



Poisson distribution

- Most of the high-energy instruments simply count photons → Poisson distribution
- Poisson distribution:

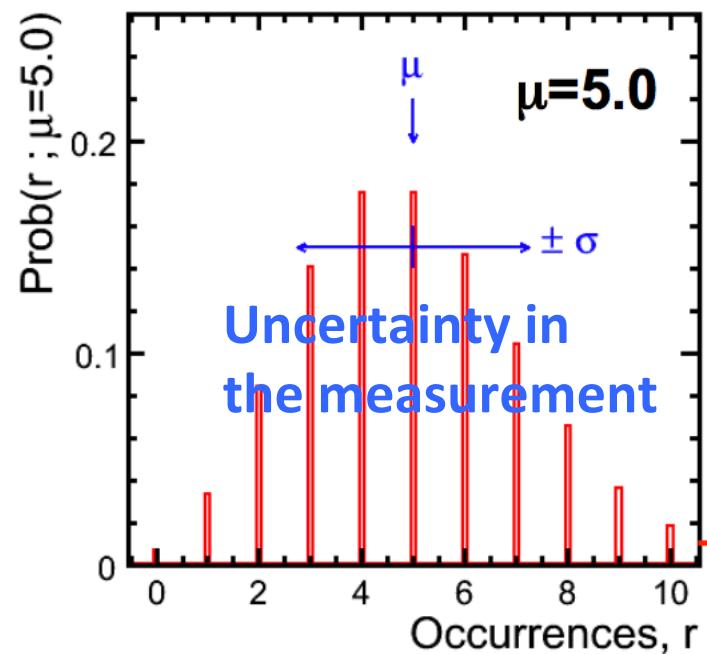
$$P(x) = \frac{\mu^x e^{-\mu}}{x!}$$

$$\langle x \rangle = \mu$$

$$\sigma^2 = \mu$$

μ : Average of the distribution

σ : Standard deviation



Gaussian distribution

- When the number of data (μ) is large,

Poisson distribution

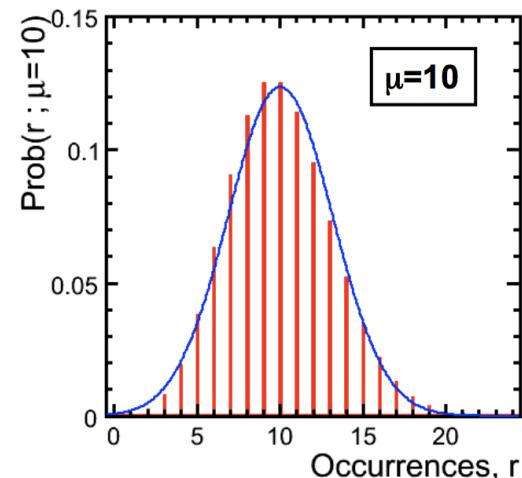
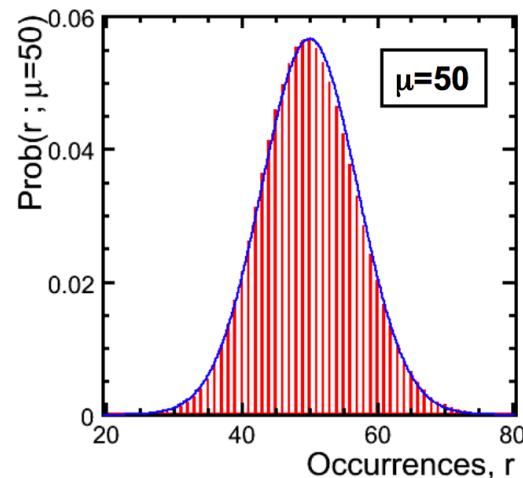
Gaussian distribution

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$P(x; \mu) \approx ke^{-\frac{(x-\mu)^2}{2\mu}}$$

$$\langle x \rangle = \mu$$

$$\langle (x - \mu)^2 \rangle = \sigma^2$$



Gaussian distribution

- When the number of data (μ) is large,

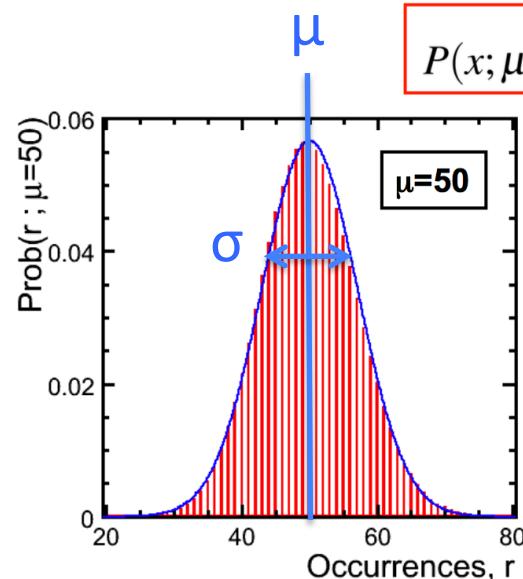
Poisson distribution

Gaussian distribution

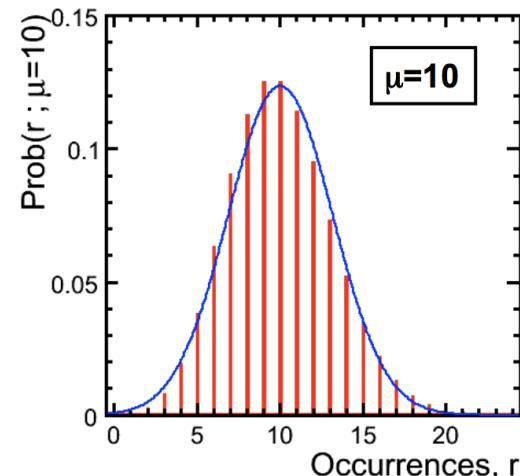
$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\langle x \rangle = \mu$$

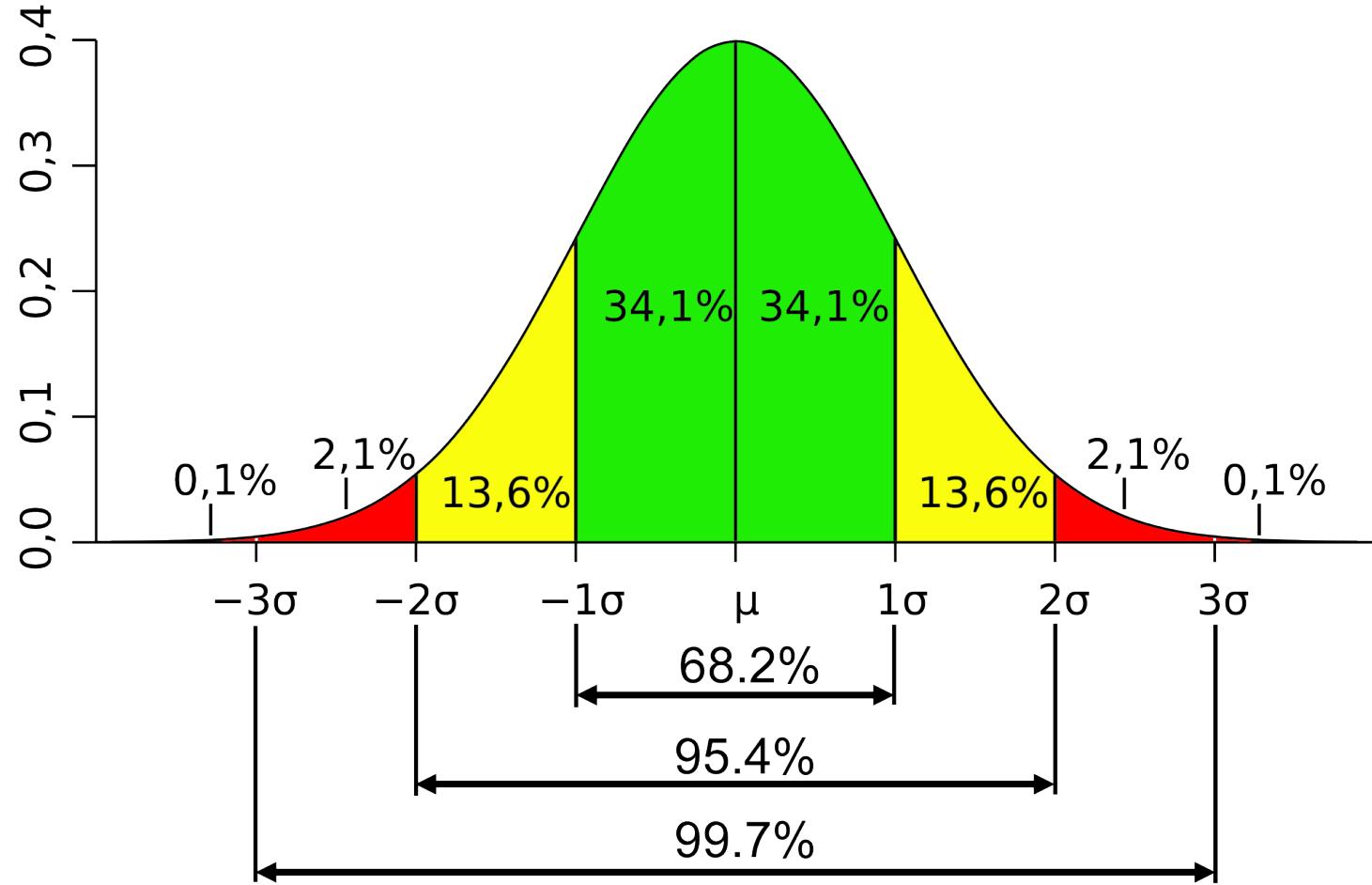
$$\langle (x - \mu)^2 \rangle = \sigma^2$$



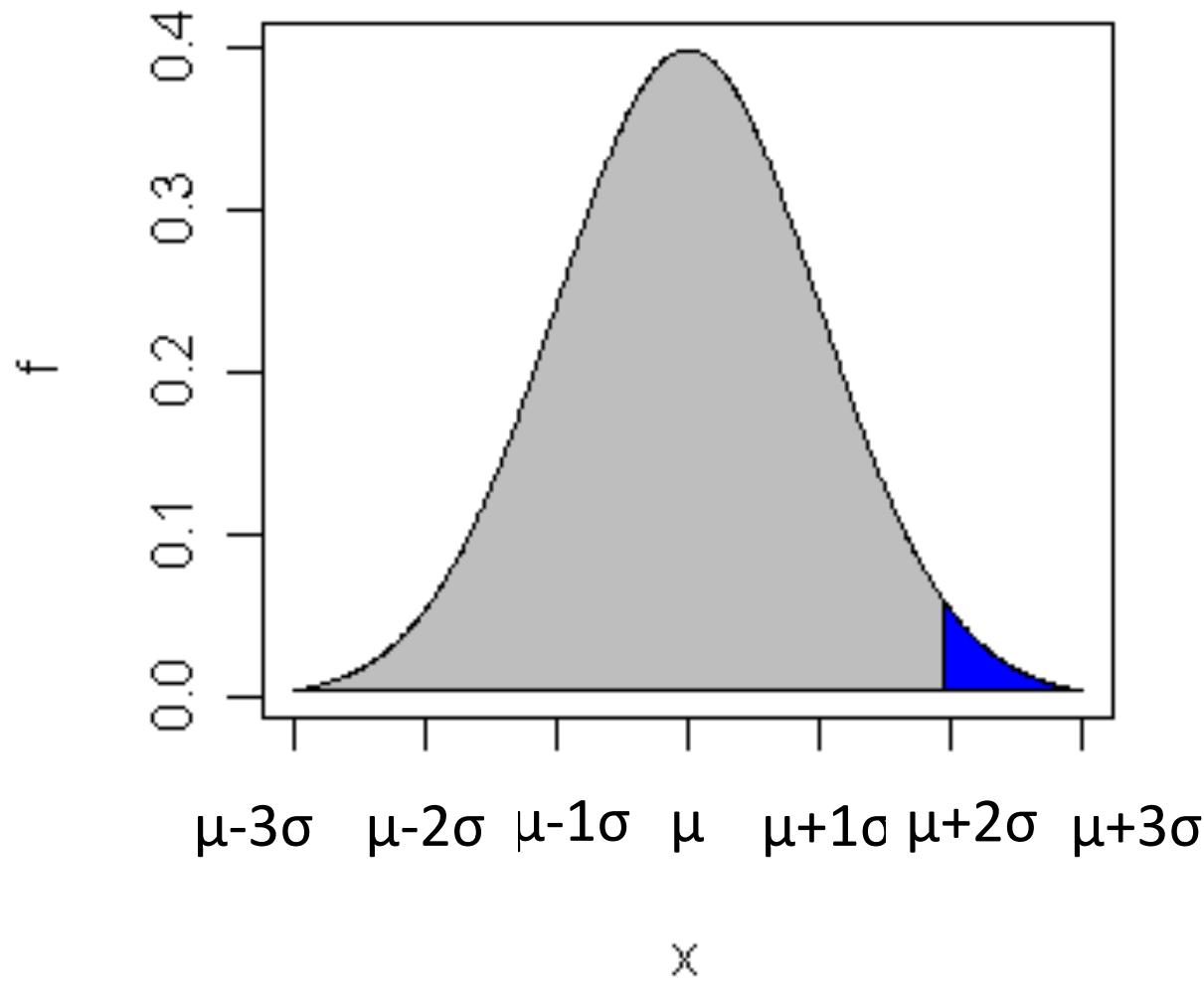
$$P(x; \mu) \approx ke^{-\frac{(x-\mu)^2}{2\mu}}$$



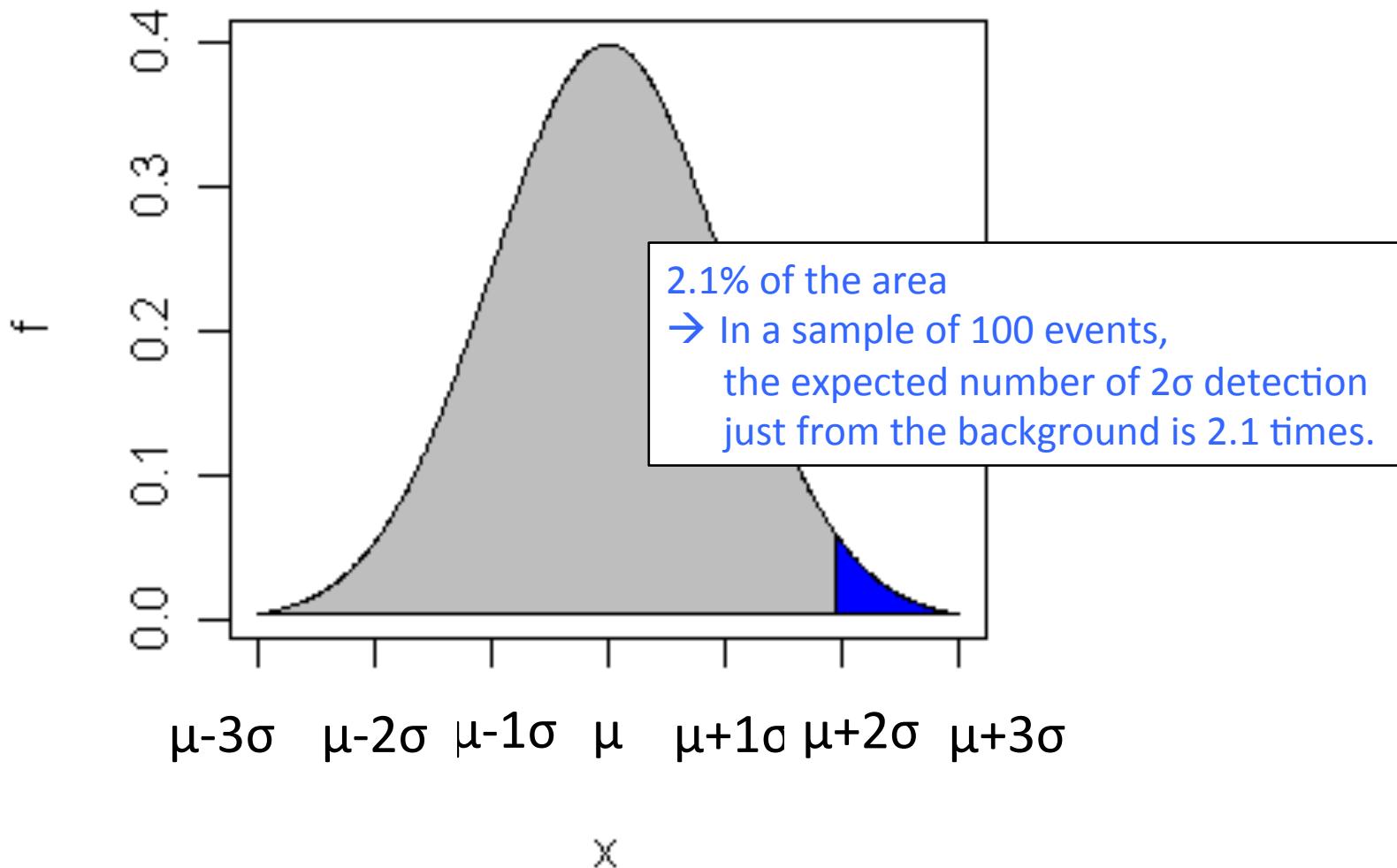
Detection significance vs false-detection rate



Detection significance vs false-detection rate



Detection significance vs false-detection rate



Data fitting

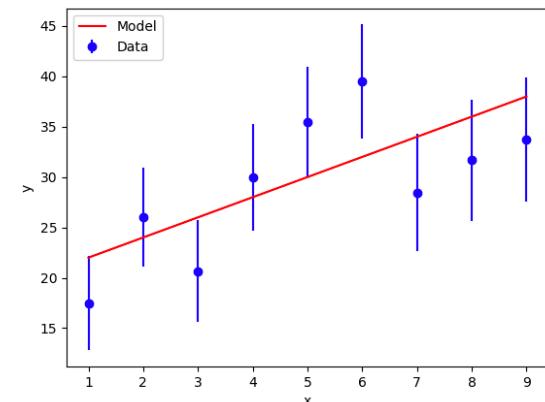
Maximum Likelihood as estimation for data with Gaussian distribution

The model parameters that have the best-chance to explain the observed data

- The likelihood function:

$$L(\theta) = \prod_{i=1}^n p(x_i; \theta)$$

The probability of measuring x_i , given a model θ



Data fitting

Maximum Likelihood as estimation for data with Gaussian distribution

The model parameters that have the best-chance to explain the observed data

- The likelihood function:

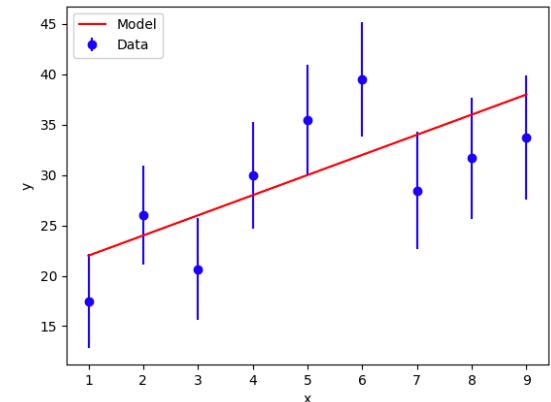
$$L(\theta) = \prod_{i=1}^n p(x_i; \theta)$$

The probability of measuring x_i , given a model θ

- If all your data have Gaussian distribution

$$L = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i}} e^{-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}}$$

$$\ln(L) = -\frac{N}{2} \ln(2\pi\sigma_i^2) - \sum_{i=1}^N \frac{(x_i - \mu_i)^2}{2\sigma_i^2}$$



Data fitting

Maximum Likelihood as estimation for data with Gaussian distribution

The model parameters that have the best-chance to explain the observed data

- The likelihood function:

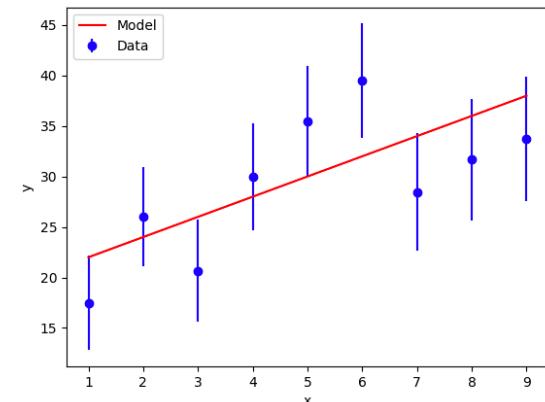
$$L(\theta) = \prod_{i=1}^n p(x_i; \theta)$$

The probability of measuring x_i , given a model θ

- If all your data have Gaussian distribution

Only related to data, not model
→ don't need for maximizing L when finding the right model

$$\ln(L) = -\frac{N}{2} \ln(2\pi\sigma_i^2) - \sum_{i=1}^N \frac{(x_i - \mu_i)^2}{2\sigma_i^2}$$



Data fitting

Maximum Likelihood as estimation for data with Gaussian distribution

The model parameters that have the best-chance to explain the observed data

- The likelihood function:

$$L(\theta) = \prod_{i=1}^n p(x_i; \theta)$$

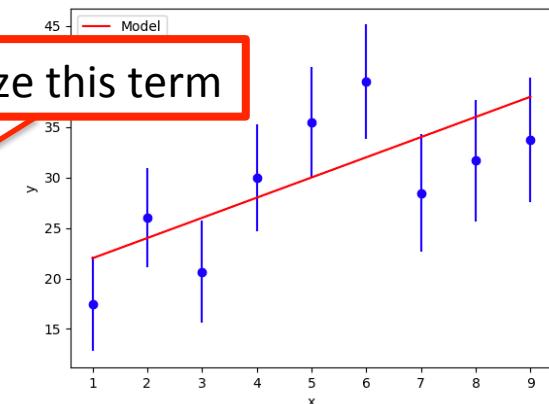
The probability of measuring x_i , given a model θ

- If all your data have Gaussian distribution

$$L = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i}} e^{-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}}$$

Maximize this term

$$\ln(L) = -\frac{N}{2} \ln(2\pi\sigma_i^2) - \sum_{i=1}^N \frac{(x_i - \mu_i)^2}{2\sigma_i^2}$$



Data fitting

Maximum Likelihood as estimation for data with Gaussian distribution

The model parameters that have the best-chance to explain the observed data

Maximize this term

$$\sum_{i=1}^N \frac{(x_i - \mu_i)^2}{2\sigma_i^2}$$



Minimize this term

$$\sum_{i=1}^N \frac{(x_i - \mu_i)^2}{\sigma_i^2}$$

Data fitting

Maximum Likelihood as estimation for data with Gaussian distribution

The model parameters that have the best-chance to explain the observed data

Maximize this term

$$\sum_{i=1}^N \frac{(x_i - \mu_i)^2}{2\sigma_i^2}$$



Minimize this term

$$\sum_{i=1}^N \frac{(x_i - \mu_i)^2}{\sigma_i^2}$$

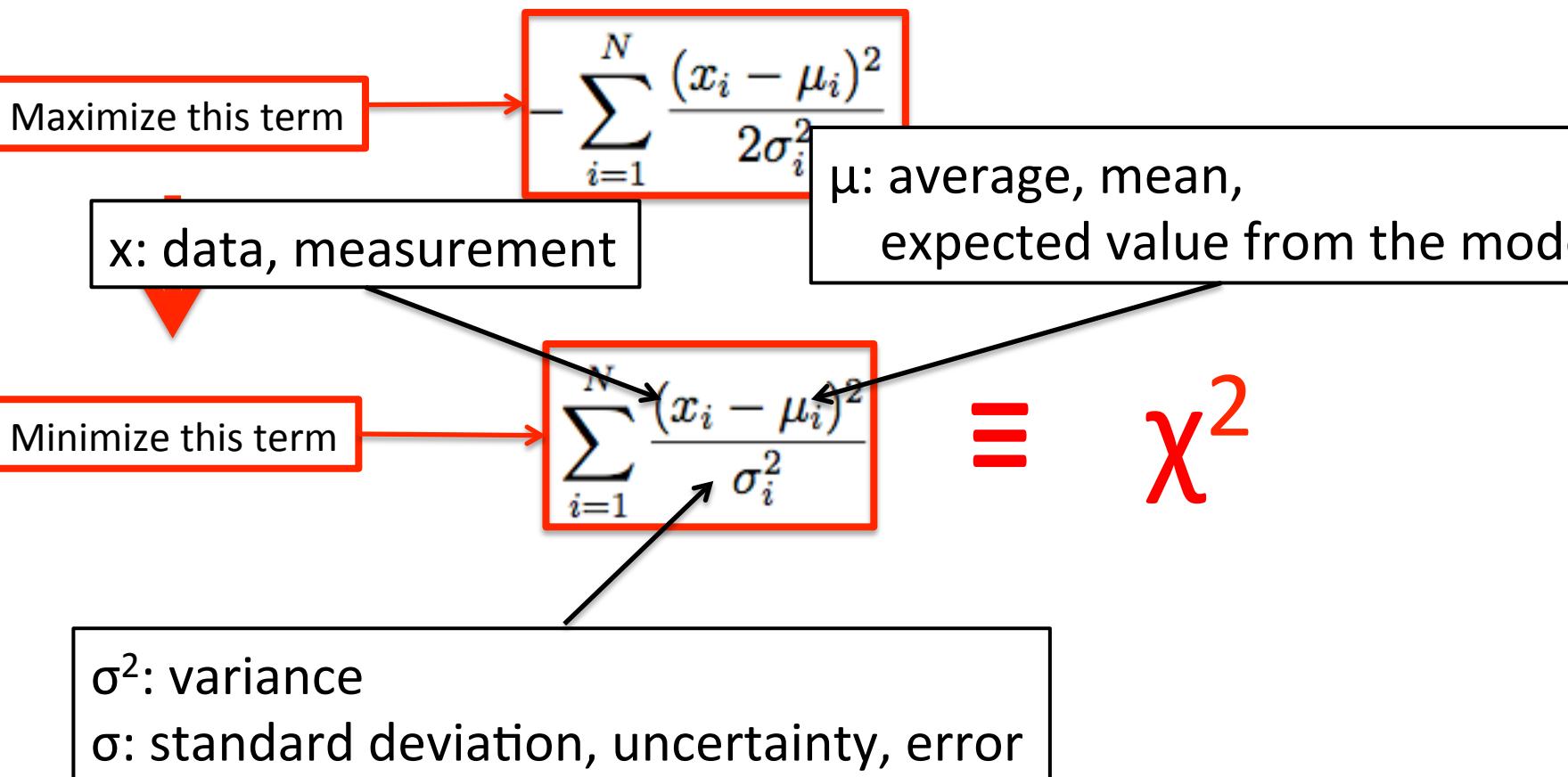
\equiv

χ^2

Data fitting

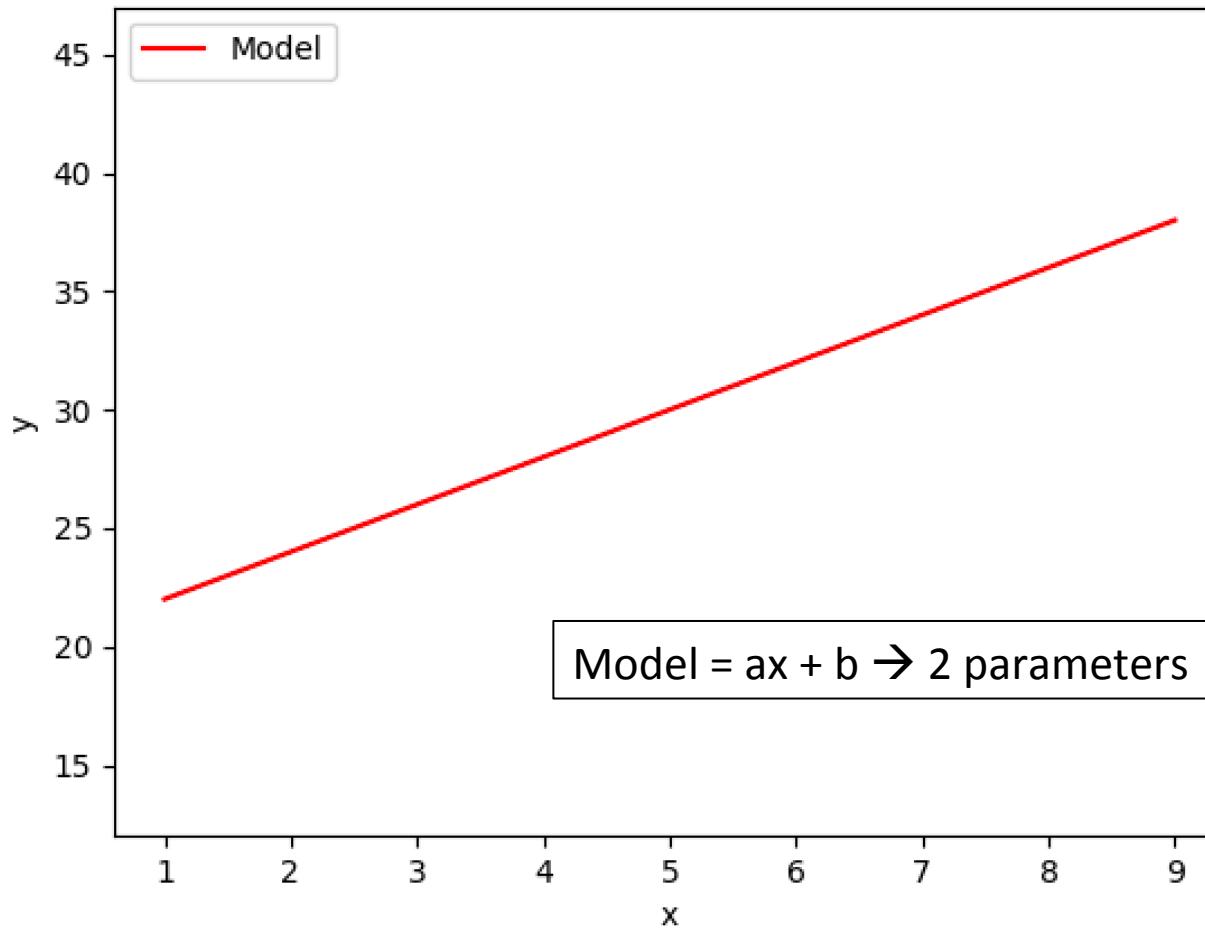
Maximum Likelihood as estimation for data with Gaussian distribution

The model parameters that have the best-chance to explain the observed data



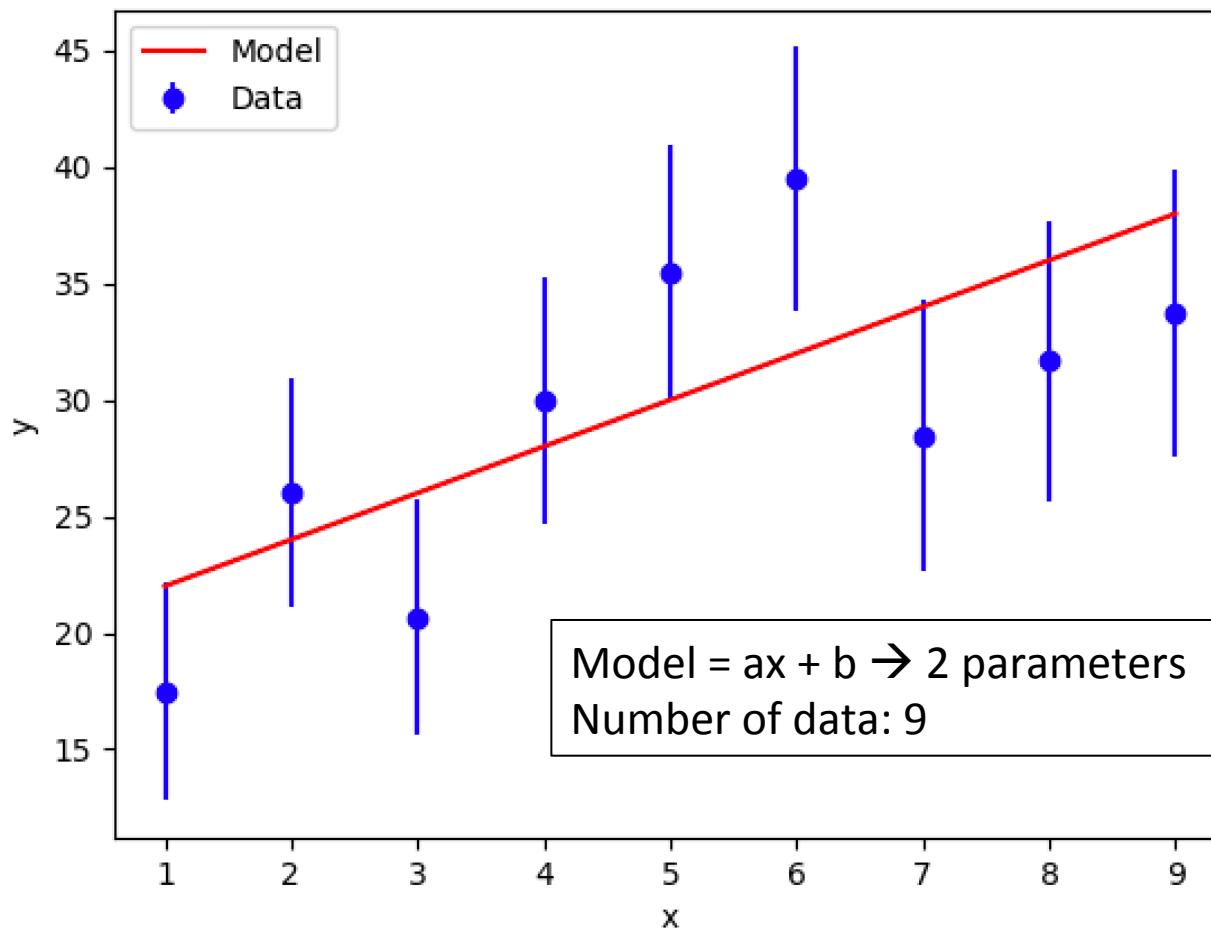
Data fitting

Maximum Likelihood as estimation for data
with Gaussian distribution



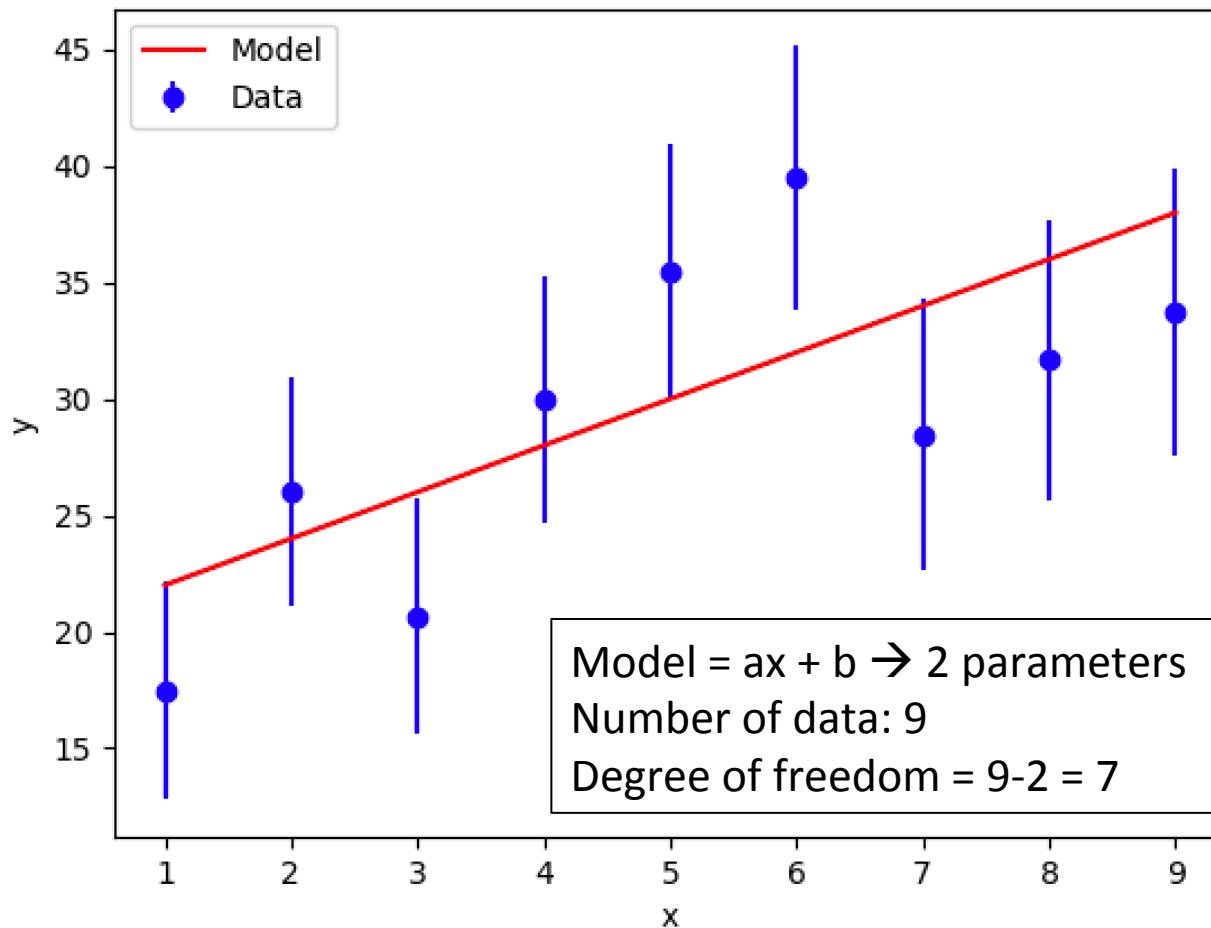
Data fitting

Maximum Likelihood as estimation for data
with Gaussian distribution



Data fitting

Maximum Likelihood as estimation for data
with Gaussian distribution



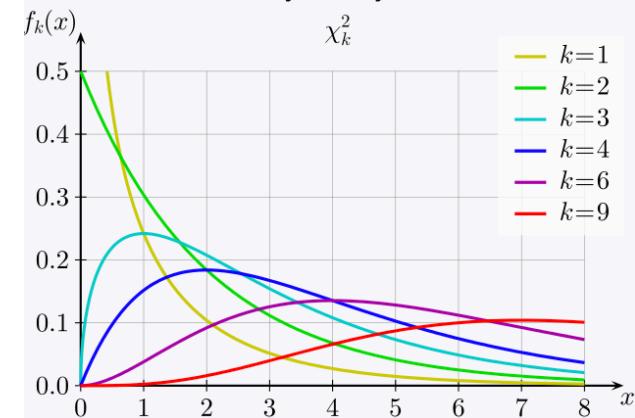
Chi-square distribution

- If χ^2 is a sum of the squares of independent normal/Gaussian random variables.

$$f(x; k) = \begin{cases} \frac{x^{\frac{k}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{k}{2}} \Gamma\left(\frac{k}{2}\right)}, & x > 0; \\ 0, & \text{otherwise.} \end{cases}$$

x: value of χ^2

k: degree of freedom



$$\langle x \rangle = k$$

→ IF you have enough data points, the reduced $\chi^2 = \chi^2/k \sim 1$