

## ADA 07 - 9am Tue 27 Sep 2022

Maximum Likelihood Estimation  
Error bars are Model parameters  
Fitting Poisson Data  
Noise Model Parameters

Conditional Probabilities  
Bayes Theorem  
Bayesian Inference

135

## Example: Correct the Bias in $(S^2)^{1/2}$

Define  $y(x) = x^b$ ,

Derivatives:  $y'(x) = b x^{b-1}$ ,  $y''(x) = b(b-1)x^{b-2}$

Evaluate the bias:

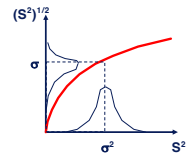
$$\langle (S^2)^b \rangle = y(\langle S^2 \rangle) + \frac{y''(\langle S^2 \rangle)}{2} \text{Var}(S^2) + \dots$$

$$= y(\sigma^2) + \frac{y''(\sigma^2)}{2} \frac{2\sigma^4}{N-1} + \dots$$

$$= \sigma^{2b} + \frac{b(b-1)\sigma^{2(b-2)}}{2} \frac{2\sigma^4}{N-1} + \dots = \sigma^{2b} \left( 1 + \frac{b(b-1)}{N-1} + \dots \right)$$

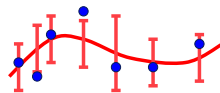
$$\langle (S^2)^{p/2} \rangle = \sigma^p \left( 1 + \frac{p(p-2)}{4(N-1)} + \dots \right)$$

$$\text{Bias-corrected: } \bar{S} = \frac{\sqrt{S^2}}{\left( 1 + \frac{p(p-2)}{4(N-1)} \right)^{1/p}} \quad \langle \bar{S}^p \rangle = \sigma^p$$

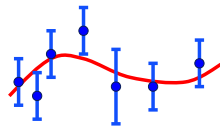


136

## Error Bars live with the Model



## Not with the Data



Usually the distinction is unimportant.  
But sometimes *it is important.*

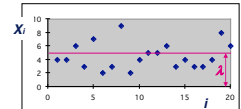
137

## Error bars live with the model, not the data!

Example: **Poisson data:**

$$\text{Prob}(x = n | \lambda) = \frac{\lambda^n e^{-\lambda}}{n!} \quad n = 0, 1, 2, \dots$$

$$\langle X_i \rangle = \lambda, \quad \sigma^2(X_i) = \lambda$$



How to attach error bars to the data points?

The **wrong way:** If  $\sigma(X_i) = \sqrt{X_i}$ , then  $1/\sigma^2 = \infty$  when  $X_i = 0$

$$\text{and } \hat{X} = \frac{\sum X_i / \sigma_i^2}{\sum 1/\sigma_i^2} = \frac{0 \cdot \infty}{\infty} = 0, \text{ clearly wrong!}$$

Assigning  $\sigma(X_i) = \sqrt{X_i}$  gives a **downward bias**. Points lower than average by chance are given smaller error bars, and hence more weight than they deserve.

The **right way:**

Assign  $\sigma = \sqrt{\lambda}$ , where  $\lambda =$  mean count rate *predicted by the model*.

138

## Maximum Likelihood (ML) Estimation

**Likelihood** of parameters  $\alpha$  for a given dataset:

$$L(\alpha) = P(X | \alpha) = P(X_1 | \alpha) \times P(X_2 | \alpha) \times \dots \times P(X_N | \alpha)$$

$$= \prod_{i=1}^N P(X_i | \alpha)$$

**Maximum Likelihood Parameters**

$$\alpha_{\text{ML}} \text{ satisfies } 0 = \frac{\partial}{\partial \alpha} [-2 \ln L(\alpha)],$$

$$\text{Var}[\alpha_{\text{ML}}] \approx \frac{2}{\left( \frac{\partial^2}{\partial \alpha^2} [-2 \ln L(\alpha)] \right)_{\alpha = \alpha_{\text{ML}}}}$$

Example: **Gaussian errors:**

$$P(X_i | \alpha) = \frac{1}{\sqrt{2\pi} \sigma_i} \exp \left\{ -\frac{1}{2} \left( \frac{X_i - \mu_i(\alpha)}{\sigma_i} \right)^2 \right\}$$

$$L(\alpha) = \frac{\exp\{-\chi^2/2\}}{Z_D}, \quad Z_D = (2\pi)^{N/2} \prod_{i=1}^N \sigma_i$$

$$\text{BoF} = -2 \ln L = \chi^2 + \sum \ln \sigma_i^2 + N \ln(2\pi)$$

To maximise  $L(\alpha)$ , minimise  $\chi^2 + \sum \ln \sigma_i^2$

**Generalises  $\chi^2$  fitting.**

1. For parameters that affect  $\sigma$
2. For non-Gaussian errors

139

## Need ML when Parameters alter Error Bars

• Data points  $X_i$  with no error bars:

$$\chi^2 = \sum_{i=1}^N \left( \frac{X_i - \mu}{\sigma} \right)^2$$

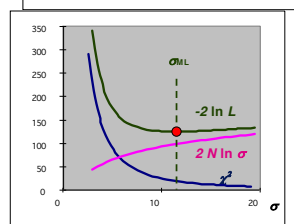
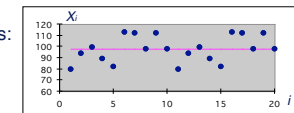
• To find  $\mu$ , minimise  $\chi^2$ .

• To find  $\sigma$ , minimising  $\chi^2$  fails!

$$\chi^2 \rightarrow 0 \text{ as } \sigma \rightarrow \infty$$

• ML method minimises

$$-2 \ln L = \chi^2 + N \ln \sigma^2$$



140

### Need ML to fit low-count Poisson Data

Example : **Poisson data** :

$$P(X = n | \lambda) = \frac{e^{-\lambda} \lambda^n}{n!} \quad n = 0, 1, \dots, \infty$$

Likelihood for  $N$  Poisson data points :

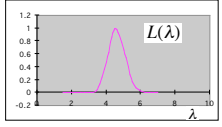
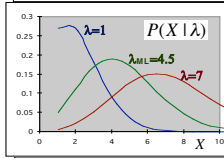
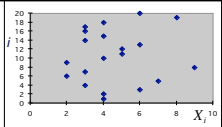
$$L(\lambda) = \prod_{i=1}^N P(X_i | \lambda) = \prod_{i=1}^N \frac{e^{-\lambda} \lambda^{X_i}}{X_i!}$$

$$\ln L = \sum_i (-\lambda + X_i \ln \lambda - \ln X_i!)$$

Maximum likelihood estimator of  $\lambda$  :

$$\frac{\partial \ln L}{\partial \lambda} = -N + \frac{1}{\lambda} \sum_i X_i = 0 \quad \text{at } \lambda = \lambda_{ML}$$

$$\therefore \lambda_{ML} = \frac{1}{N} \sum_i X_i$$



141

### Conditional Probabilities

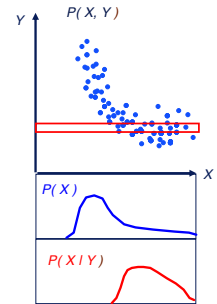
$P(X, Y)$  = joint probability density of  $X$  and  $Y$   
 $P(X)$  = projection of  $P(X, Y)$  onto  $X$  axis.

$$P(X) = \int P(X, Y) dY$$

**Conditional Probability:**

$P(X | Y)$  = "probability of  $X$  given  $Y$ "  
 = "normalised slice" of  $P(X, Y)$   
 at a fixed value of  $Y$ .

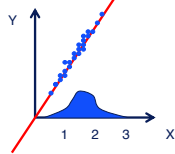
$$P(X | Y) = \frac{P(X, Y)}{P(Y)} = \frac{P(X, Y)}{\int P(X, Y) dX}$$



142

### Test Understanding

$Y = 3X$   
 $X = \text{Gaussian}$



$P(Y | X = 2) = ?$



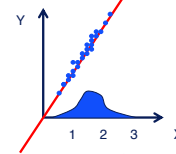
$P(Y | X > 2) = ?$



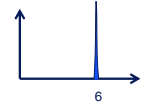
143

### Test Understanding

$Y = 3X$   
 $X = \text{Gaussian}$



$P(Y | X = 2) = ?$



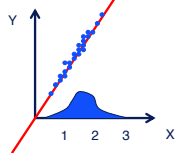
$P(Y | X > 2) = ?$



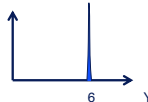
144

### Test Understanding

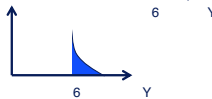
$Y = 3X$   
 $X = \text{Gaussian}$



$P(Y | X = 2) = ?$



$P(Y | X > 2) = ?$



145

### Conditional Probabilities

$P(X)$  = projection onto  $X$  axis.  
 $P(Y)$  = projection onto  $Y$  axis.

$$P(X) = \int P(X, Y) dY$$

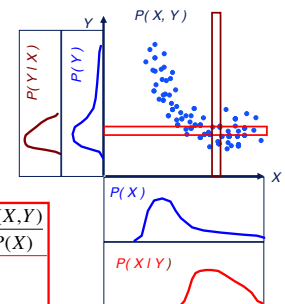
$$P(Y) = \int P(X, Y) dX$$

**Conditional Probability:**

$P(X | Y)$  = normalised slice at fixed  $Y$   
 $P(Y | X)$  = normalised slice at fixed  $X$

$$P(X | Y) = \frac{P(X, Y)}{P(Y)} \quad P(Y | X) = \frac{P(X, Y)}{P(X)}$$

$$P(X, Y) = P(X | Y) P(Y) = P(Y | X) P(X)$$



146

## Bayes' Theorem and Bayesian Inference

**Bayes' Theorem:**  $P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$

Since  $P(X,Y) = P(X|Y)P(Y) = P(Y|X)P(X)$   
 then  $P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)} = \frac{P(Y|X)P(X)}{\int P(Y|X)P(X)dX}$

**Bayesian Inference :**

$$P(\text{model} | \text{data}) = \frac{P(\text{data} | \text{model}) P(\text{model})}{P(\text{data})}$$

Shows us **how to change** our probability distribution  $P(\text{model}) \Rightarrow P(\text{model} | \text{data})$  over various models in light of new data.

147

## Inferences depend on Prior, not just Data

**Bayesian inference:** ( $M$  = model,  $D$  = data)

Posterior Probability = (Likelihood  $\times$  Prior Probability) / Evidence

$$P(M|D) = \frac{P(D|M)P(M)}{P(D)} = \frac{P(D|M)P(M)}{\int P(D|M)P(M)dM}$$

Relative probability of two models  $M_1$  and  $M_2$  :

$$\frac{P(M_1|D)}{P(M_2|D)} = \frac{P(D|M_1)}{P(D|M_2)} \times \frac{P(M_1)}{P(M_2)} \approx \exp\left(\frac{-\Delta\chi^2}{2}\right) \times \frac{P(M_1)}{P(M_2)}$$

- The **Likelihood**,  $P(\text{data} | \text{model})$ , is quantified by a "badness-of-fit" statistic. e.g.  $P(\text{data} | \text{model}) \sim \exp(-\chi^2/2)$
- The **Prior**,  $P(\text{model})$  expresses your **prejudice** (prior knowledge).
- The **Posterior**,  $P(\text{model} | \text{data})$ , gives your **inference**, the relative probabilities of different models (parameters), in light of the data.

**No absolute inferences !** New data **updates** your prior expectations, but your **conclusions depend also on your prior.**

148

## Choice of Prior

- A model for a set of data  $X$  depends on model parameters  $\alpha$ , and gives the Likelihood

$$L(\alpha) \equiv P(X|\alpha)$$

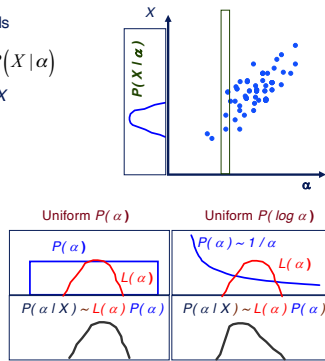
- Knowledge of  $\alpha$  before measuring  $X$  is quantified by the **prior**  $P(\alpha)$ .

- Choice of prior  $P(\alpha)$  is arbitrary, subject to common sense!

- After measuring  $X$ , Bayes theorem gives **posterior** :

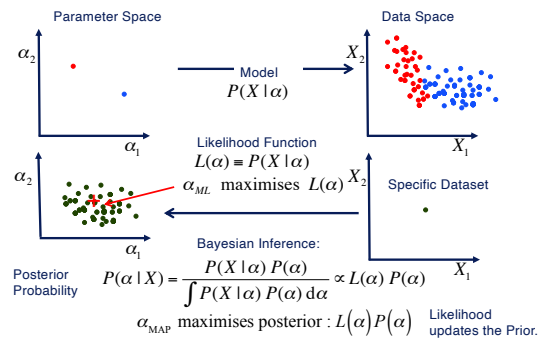
$$P(\alpha|X) \propto P(X|\alpha)P(\alpha) = L(\alpha)P(\alpha)$$

- Different priors  $P(\alpha)$  lead to different **inferences** :



149

## Max Likelihood and Bayesian Inference



150

## N=1 Gaussian Datum with Uniform Prior

Data :  $X \pm \sigma$  Model parameter :  $\mu$

Likelihood function :

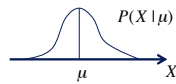
$$L(\mu) = P(X|\mu) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{X-\mu}{\sigma}\right)^2}$$

$\mu_{ML} = X$  maximises  $L(\mu)$ .

Posterior probability :

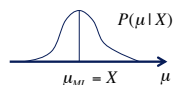
$$P(\mu|X) = \frac{P(X|\mu)P(\mu)}{P(X)}$$

$$P(X) = \int P(X|\mu)P(\mu)d\mu$$



Uniform prior:

$$P(\mu) = \text{constant}$$



Maximum Likelihood implicitly assumes a Uniform Prior

151

## N=1 Gaussian Datum with Gaussian Prior

Gaussian Data:  $X \pm \sigma$

$$\text{Likelihood: } L(\mu) = P(X|\mu) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{X-\mu}{\sigma}\right)^2}$$

$$\text{Prior: } P(\mu) = \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{1}{2}\left(\frac{\mu-\mu_0}{\sigma_0}\right)^2}$$

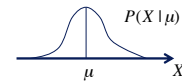
Posterior :  $P(\mu|X) \propto \text{Likelihood} \times \text{Prior}$

$$L(\mu)P(\mu) \propto e^{-\frac{1}{2}\left(\frac{X-\mu}{\sigma}\right)^2} e^{-\frac{1}{2}\left(\frac{\mu-\mu_0}{\sigma_0}\right)^2} \propto \exp\left\{-\frac{1}{2}\left(\frac{\mu-\mu_{MAP}}{\sigma(\mu_{MAP})}\right)^2\right\}$$

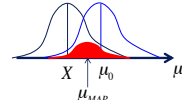
**Maximum Posterior (MAP) estimate:**

$$\mu_{MAP} = \frac{\frac{\mu_0}{\sigma_0^2} + \frac{X}{\sigma^2}}{\frac{1}{\sigma_0^2} + \frac{1}{\sigma^2}}, \quad \text{Var}(\mu_{MAP}) = \frac{1}{\frac{1}{\sigma_0^2} + \frac{1}{\sigma^2}}$$

Verify this result.



Likelihood  $\times$  Prior:



Same as Optimal Average !

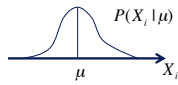
**Gaussian prior acts like 1 more data point.**

Data "pulls" the probability away from the prior, and vice-versa.

152

## N Gaussian Data with Gaussian Prior

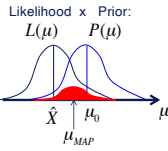
Likelihood:  $L(\mu) = P(X | \mu) = \prod_{i=1}^N P(X_i | \mu) = \frac{\exp\left\{-\frac{1}{2}\chi^2\right\}}{(2\pi)^{N/2} \prod_i \sigma_i}$



Prior:  $P(\mu) = \frac{1}{\sqrt{2\pi} \sigma_0} \exp\left\{-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right\}$

Posterior:  $P(\mu | X) \propto \text{Likelihood} \times \text{Prior}$

$L(\mu) P(\mu) \propto \exp\left\{-\frac{\chi^2}{2} - \frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right\} \propto \exp\left\{-\frac{1}{2}\left(\frac{\mu - \mu_{MAP}}{\sigma_{MAP}}\right)^2\right\}$



**Maximum Posterior (MAP) estimate:**

$$\mu_{MAP} = \frac{\frac{\mu_0}{\sigma_0^2} + \sum_{i=1}^N \frac{X_i}{\sigma_i^2}}{\frac{1}{\sigma_0^2} + \sum_{i=1}^N \frac{1}{\sigma_i^2}}, \quad \sigma^2(\mu_{MAP}) = \frac{1}{\frac{1}{\sigma_0^2} + \sum_{i=1}^N \frac{1}{\sigma_i^2}}$$

Same as Optimal Average !

Gaussian prior acts like 1 more data point.

153

## Summary:

1. Error bars live with the Model, not with the Data.
2. Bayes Theorem (**Bayesian Inference**)

$$P(\text{Model} | \text{Data}) = \frac{P(\text{Data} | \text{Model}) P(\text{Model})}{P(\text{Data})}$$

3. **Maximum Likelihood**,  $L(\text{Model}) \equiv P(\text{Data} | \text{Model})$

e.g. for Gaussian Data:

$$BoF = -2 \ln L = \chi^2 + \sum_{i=1}^N \ln \sigma_i^2 + const$$

4. Minimise  $\chi^2$  if Gaussian errors with known  $\sigma_i$ .
5. or Maximise likelihood ( e.g. minimise  $BoF = -2 \ln L$  ), if error bars unknown, or low-count Poisson data.
6. or full **Bayesian analysis**, including the prior:

e.g. for Gaussian Data:

$$BoF = -2 \ln P(\text{Model} | \text{Data}) = \chi^2 + \sum_{i=1}^N \ln \sigma_i^2 - 2 \ln P(\text{Model}) + const$$

154

Fini -- ADA 07

155