

Lecture 6

- Fitting a line to data

$$\begin{aligned} \text{Data} &= y_i \\ \text{Signal} &= s_i = ax_i + b \\ \text{Noise} &= r_i \quad (\text{or } "n_i" \dots \text{doesn't matter}) \end{aligned}$$

$$\boxed{\text{Data} = \text{Signal} + \text{Noise}}$$

⇒ what makes this a probabilistic problem?
↳ noise!

$$p(r_i) = N(0, \sigma)$$

$$\therefore p(y_i - s_i) = N(0, \sigma)$$

$$p(y_i) = N(s_i, \sigma)$$

$$\Rightarrow p(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(y_i - s_i)^2}{\sigma^2}}$$

BUT s_i is now a LINE RELATIONSHIP not just the single " μ " as in the last lecture.

$$\ln L = \ln [\prod_i p(y_i)]$$

$$= \text{constant} - \underbrace{\frac{1}{2} \sum_i \frac{(y_i - ax_i - b)^2}{\sigma^2}}_{\chi^2}$$

(NOTE) Also possible to modify the likelihood to reduce the impact of outliers (see the notebook).

◦ Goodness of Fit

- MLE gives best-fit parameters of model.
- How do we know if the model is any good?

⇒ Max-likelihood should not be an unlikely occurrence ... need to know L distribution.

$$z_i = \frac{(y_i - s_i)}{\sigma_i}$$

$$\ln L = \text{constant} - \frac{1}{2} \sum_i z_i^2$$

$$= \text{constant} - \frac{1}{2} \chi^2$$

— chi-squared

Mean of χ^2 distribution = $N - k$

Std of χ^2 distribution = $\sqrt{2(N - k)}$

data-points # params

χ^2 per degree of freedom

$$\left. \begin{aligned} \chi^2_{\text{dof}} &= \frac{1}{N-k} \sum_i z_i^2 \end{aligned} \right\} \begin{aligned} &\text{GOOD FIT} \\ &\parallel \\ &\chi^2_{\text{dof}} \approx 1 \end{aligned}$$

- Model Comparison — much more on this later in the course.

$$AIC_M = \underbrace{-2 \ln[L^{\text{Max}}(M)]}_{\propto \chi^2} + \underbrace{2k + \frac{2k(k+1)}{N-k-1}}_{\text{penalty on model complexity.}}$$

⇒ Model with lowest AIC wins!

- Confidence Estimation — Fisher matrix approach is an approximation

⇒ simulate NEW data by resampling from the data you have — NOT MAGIC!!

BOOTSTRAP ⇒ resample data (with replacement) and compute statistics on each new data realization ... EASY.

JACKKNIFE ⇒ remove one data point at a time, and re-compute statistics
 ↳ gives (N-1) new data realizations ... EASY.