

Fisher information

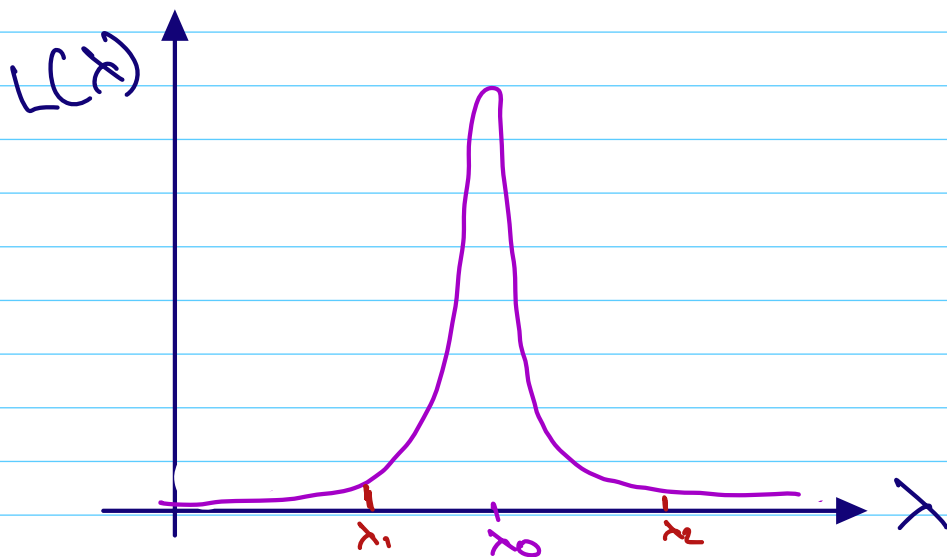
likelihood function is the basis of a lot of modern statistics and science and it's defined as follows.

$$L(\theta | \text{Data}) \equiv P(\text{Data} | \theta) = P(\text{Data} | \text{theory})$$

So we are asking having made a measurement how likely is it that we got measurement, given that the universe is described by certain theory (θ)

For example:

$$L(\lambda | x_1, x_2, x_3, \dots, x_n) = P(x_1, x_2, \dots, x_n | \lambda)$$

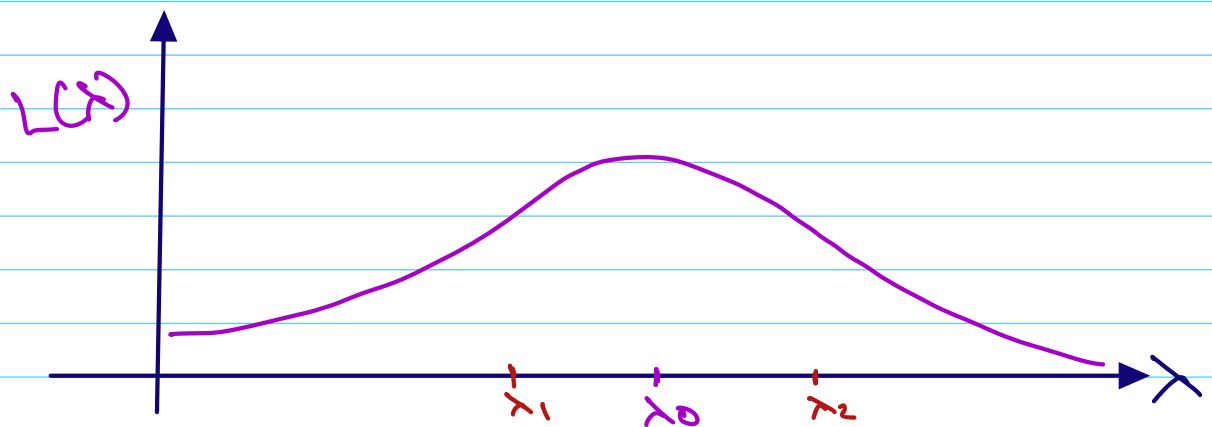


* The above Likelihood function have very sharp Peak in it. meaning at some value λ_0 , a true underlying value of our Parameter λ , we were very likely to get our data.

* Our data x_1, x_2, \dots, x_n is inconsistent with λ being $\lambda < \lambda_1$ or $\lambda > \lambda_2$. But it extremely likely we would have our data if our $\lambda = \lambda_0$

\Rightarrow The true Parameter λ^* will be very close to λ_0 .

Now Imagine another Likelihood function, where the distribution is very broad



in the above likelihood function, λ_0 is where the Peak Likelihood is i.e

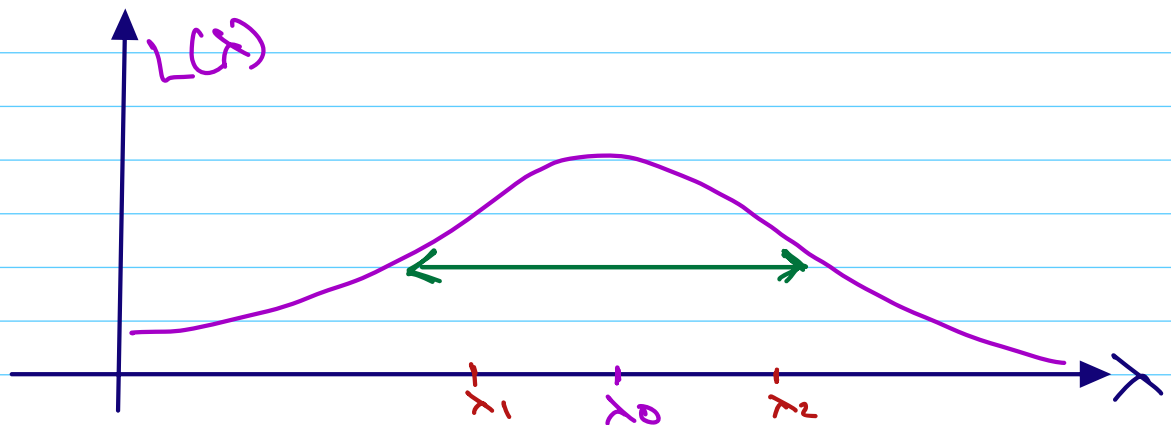
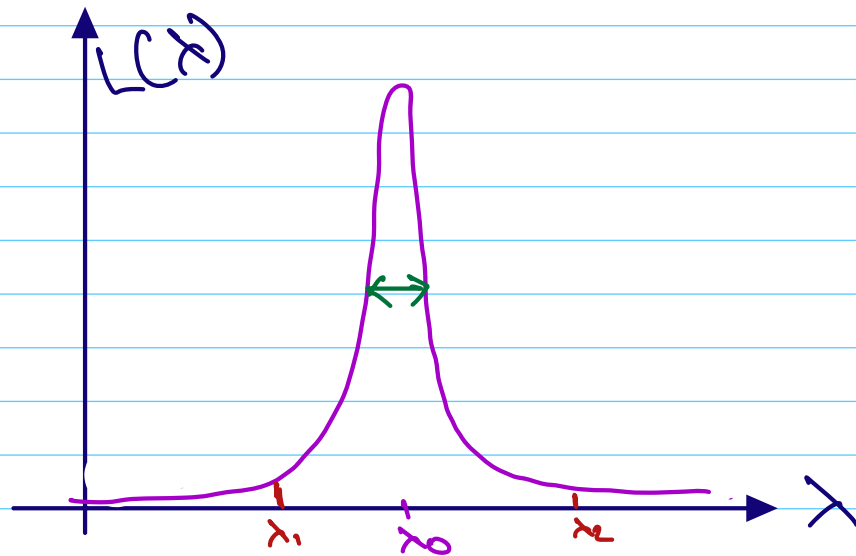
$$\lambda_0 = \underset{\lambda}{\operatorname{argmax}} L(\lambda | \text{Data})$$

\Rightarrow There is good chance that the data is generated using λ_0 , But there is also a good deal of likelihood at some distance away from λ_0 , saying that while we are most likely to recover data with the value λ_0 , there is a good deal of likelihood we will acquire data with another value, even if the λ_0 is the true value.

* So, we are looking here some way to quantify the difference b/w the above likelihood function's. and ask given a set of measurements, (Data) how constraining

are we on the underlying theory?

so, the clear sign we are looking for is the width of this likelihood function, or some measure of how fast it falls off from the maximum.



if it falls off very quickly, then our data is relatively constraining, however if it falls off very slowly, then our data doesn't do much to constrain the underlying theory.

* So, to come up with a more quantitative framework for it. we will Taylor expand the Likelihood function, around its peak

$$L(\lambda) \approx L(\lambda_0) + \frac{\partial L(\lambda)}{\partial \lambda} \bigg|_{\lambda=\lambda_0} (\lambda - \lambda_0) + \frac{1}{2} \frac{\partial^2 L(\lambda)}{\partial \lambda^2} \bigg|_{\lambda=\lambda_0} (\lambda - \lambda_0)^2$$

we know that $\frac{\partial L}{\partial \lambda} \bigg|_{\lambda=\lambda_0} = 0$ because it's Peak.

$$\Rightarrow L(\lambda) \approx L(\lambda_0) + \frac{1}{2} \frac{\partial^2 L}{\partial \lambda^2} \bigg|_{\lambda=\lambda_0} (\lambda - \lambda_0)^2$$

we are approximating our likelihood function as parabola. This is generally not a good approximation.

The Likelihood can have any shape, so better approximation is to look at 2nd derivative of log-Likelihood.

De-tour

will come back after some analysis.

why Parabola Approximation fails,
and Gaussian Approximation works.

①

The Likelihood is Always positive,
But a Parabola is not.

if we expand the Likelihood function $L(\theta)$ directly around its maximum
② :

$$L(\theta) \approx L(\hat{\theta}) + \frac{1}{2} L''(\hat{\theta}) (\theta - \hat{\theta})^2$$

① Since $L''(\hat{\theta}) < 0$, this is downward opening Parabola.

② But this approximation will eventually become negative for sufficiently large $|\theta - \hat{\theta}|$, impossible for Likelihood, because $L(\theta) \geq 0$

So the Parabola approximation is only valid infinitesimally close to the maximum, and has no probabilistic interpretation.

② The Log-Likelihood is the right Object to expand.

The log-Likelihood $Q(\theta) = \log L(\theta)$

is a smooth, concave function, and
 Taylor - expanding it ensures that
 the exponentiated form remains
 positive and normalizable

$$l(\theta) = \log L(\theta)$$

Peak. $\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} L(\theta | \text{Data}) = \hat{\theta}_n^{\text{MLE}}$

$$l(\theta) \approx l(\hat{\theta}) + \cancel{l'(\hat{\theta})} (\theta - \hat{\theta}) + \frac{1}{2} \ddot{l}(\hat{\theta}) (\theta - \hat{\theta})^2$$

$\hat{\theta} \leftarrow \text{MLE}$

Exponentiate

$$\Rightarrow L(\theta) = e^{l(\theta)} \approx e^{l(\hat{\theta}) + \frac{1}{2} \ddot{l}(\hat{\theta}) (\theta - \hat{\theta})^2}$$

$$\Rightarrow L(\theta) \approx e^{l(\hat{\theta})} \exp \left\{ -\frac{1}{2} I(\hat{\theta}) (\theta - \hat{\theta})^2 \right\}$$

where $I(\hat{\theta}) = -\ddot{Q}(\hat{\theta})$ is the

Observed Fisher information.

This gives a Gaussian kernel,
always positive, has a clear Probabilistic interpretation.

Expanding Log-likelihood is equivalent to assuming: "Around the peak, the information content changes quadratically"

This means:

- * The likelihood drops off exponentially with squared distance from the maximum.
- * The exponential decay reflects how evidence weakens as we move

away from the best parameter

- * The Curvature (Fisher information) determines the rate of decay, which directly measures how strongly data constrain θ .

$$L(\theta) \approx e^{l(\hat{\theta})} \cdot e^{\frac{1}{2}(\theta - \hat{\theta})^2 l''(\hat{\theta})}$$

Define $I(\hat{\theta}) = -l''(\hat{\theta})$

(That's the observed Fisher information, the curvature at the peak).

$$\Rightarrow L(\theta) \approx e^{l(\hat{\theta})} \cdot e^{-\frac{1}{2} I(\hat{\theta}) (\theta - \hat{\theta})^2}$$

We recognize the Gaussian form

The term $e^{-\frac{1}{2} I(\hat{\theta}) (\theta - \hat{\theta})^2}$

is exactly the kernel of a Normal distribution with

$$\theta \sim \mathcal{N}(\hat{\theta}, \frac{1}{I(\hat{\theta})})$$

So around its maximum, the likelihood behaves like a Gaussian centered the MLE with variance equal to the inverse of the curvature.

Back to our
content

in multi parameter space, the

2nd derivative is Hessian

$$\begin{bmatrix} \frac{\partial^2 L}{\partial \lambda^2} & \frac{\partial L}{\partial \lambda} \frac{\partial L}{\partial \beta} \\ \frac{\partial L}{\partial \beta} \frac{\partial L}{\partial \lambda} & \frac{\partial^2 L}{\partial \beta^2} \end{bmatrix} = -F$$

This we define $-F$, where F is the Fisher or Curvature matrix

→ The Fisher matrix is sometimes called Curvature matrix, because it is 2nd derivative of log-likelihood function

→ It tells how curved the likelihood function is around its maximum.

So the bigger the values in the

fisher matrix, the more curved it is, meaning the more peaked it is, meaning the more constraining our data is for that particular parameter.

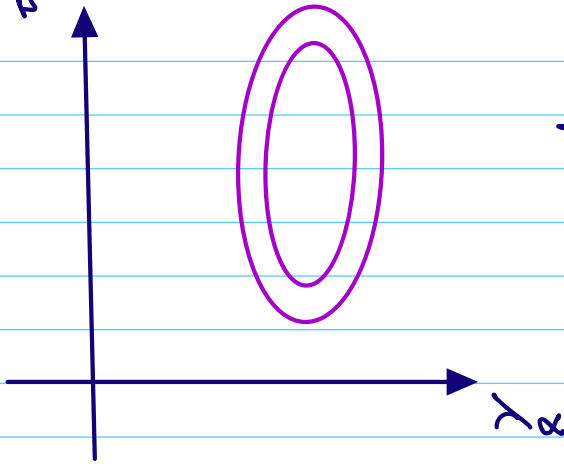
Another expression in statistics, the co-variance matrix, which is related to the fisher matrix as follows.

$$\Rightarrow F_{\alpha\beta}^{-1} = C_{\alpha\beta}$$

Co-variance is inverse of Fisher matrix.

$$C_{\alpha\beta} = \begin{bmatrix} \sigma_{\alpha}^2 & \sigma_{\alpha\beta} \\ \sigma_{\alpha\beta} & \sigma_{\beta}^2 \end{bmatrix}$$

(i) γ_B



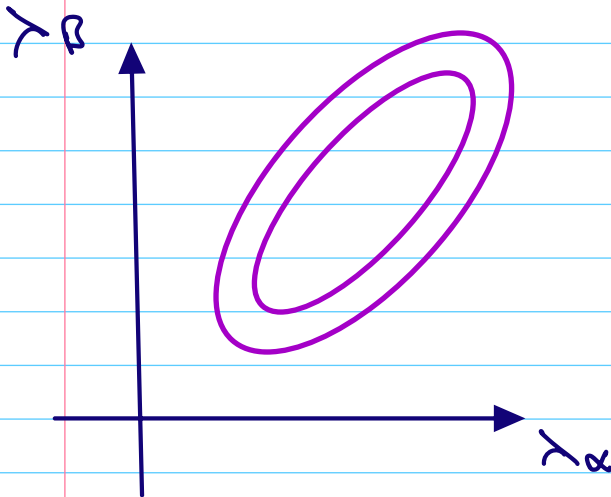
\Rightarrow

$$\begin{bmatrix} \sigma_\alpha^2 & 0 \\ 0 & \sigma_B^2 \end{bmatrix}$$

$$\sigma_B^2 > \sigma_\alpha^2$$

$$\sigma_{\alpha B} = \sigma_{B\alpha} = 0 \text{ (No dependence)}$$

(ii)



are correlated.

another interpretation.

$$x \sim f_x(x|\theta) \quad \theta \in \Theta$$

we don't know the true parameter θ^*

\Rightarrow let's assume we observed one realization $x \sim f_x(x|\theta^*)$ and want to know how much this observation tells us about the true parameter.