# Probability of the data versus likelihood of the parameters

- Suppose you are counting how many cars pass in front of your window on Sundays between 9:00 and 9:02 am. Counting experiments are generally well described by the Poisson distribution. Therefore, if the mean counts are $\lambda$, the probability of counting n cars follows the distribution:
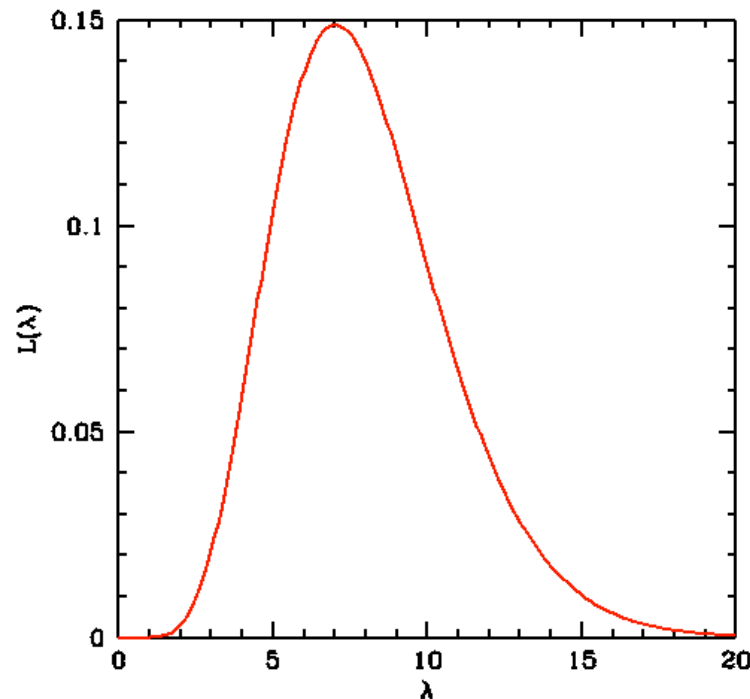
$$P(n \mid \lambda) = \frac{\lambda^n e^{-\lambda}}{n!}$$

- This means that if you repeat the experiment many times, you will measure different values of n following the frequency P(n). Note that the sum over all possible n is unity.

- Now suppose that you actually perform the experiment once and you count 7. Then, the likelihood for the model parameter $\lambda$ GIVEN the data is:

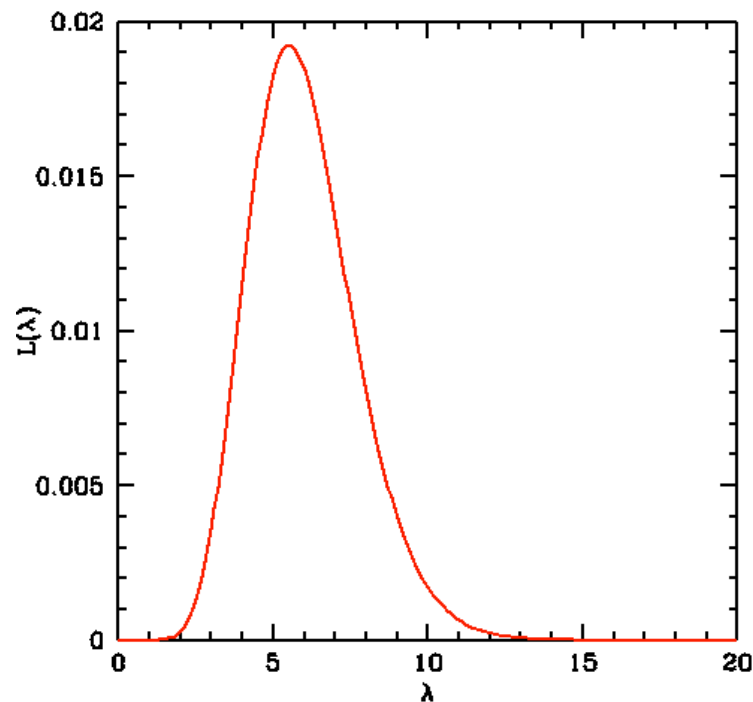$$L(\lambda) = P(7 \mid \lambda) = \frac{\lambda^7 e^{-\lambda}}{5040}$$

# The likelihood function

- This is a function of $\lambda$ only but it is NOT a probability distribution for $\lambda$! It simply says how likely it is that our measured value of n=7 is obtained by sampling a Poisson distribution of mean $\lambda$. It says something about the model parameter GIVEN the observed data.

# The likelihood function

- Let us suppose that after some time you repeat the experiment and count 4 cars. Since the two experiments are independent, you can multiply the likelihoods and obtain the curve below. Note that now the most likely value is $\lambda$=5.5 and the likelihood function is narrower than before, meaning that we know more about $\lambda$.

# Likelihood for Gaussian errors

- Often statistical measurement errors can be described by Gaussian distributions. If the errors $\sigma_i$ of different measurements $d_i$ are independent:

$$L(\theta) = P(d \mid \theta) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left[-\frac{(d_i - m_i(\theta))^2}{2\sigma_i^2}\right]$$

$$-\ln L(\theta) = \sum_{i=1}^{N} \frac{(d_i - m_i(\theta))^2}{2\sigma_i^2} + \text{const.} = \frac{\chi^2(\theta, d)}{2} + \text{const.}$$

- Maximizing the likelihood corresponds to finding the values of the parameters $\theta = \{\theta_1, ..., \theta_n\}$ which minimize the $\chi^2$ function (weighted least squares method).

# The general Gaussian case

- In general, errors are correlated and

$$-\ln L(\theta) = \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\left[d_i - m_i(\theta)\right] C_{ij}^{-1}\left[d_j - m_j(\theta)\right] + \text{const.} = \frac{\chi^2(\theta,d)}{2} + \text{const.}$$

where $C_{ij} = \langle \varepsilon_i \, \varepsilon_j \rangle$ is the covariance matrix of the errors.

- For uncorrelated errors the covariance matrix is diagonal and one reduces to the previous case.

- Note that the covariance matrix could also derive from a model and then depend on the model parameters. We will encounter some of these cases in the rest of the course.

# The Likelihood function: a summary

- In simple words, the likelihood of a model given a dataset is proportional to the probability of the data given the model

- The likelihood function supplies an order of preference or plausibility of the values of the free parameters $\theta_i$ by how probable they make the observed dataset

- The likelihood ratio between two models can then be used to prefer one to the other

- Another convenient feature of the likelihood function is that it is functionally invariant. This means that any quantitative statement about the $\theta_i$ implies a corresponding statements about any one to one function of the $\theta_i$ by direct algebraic substitution
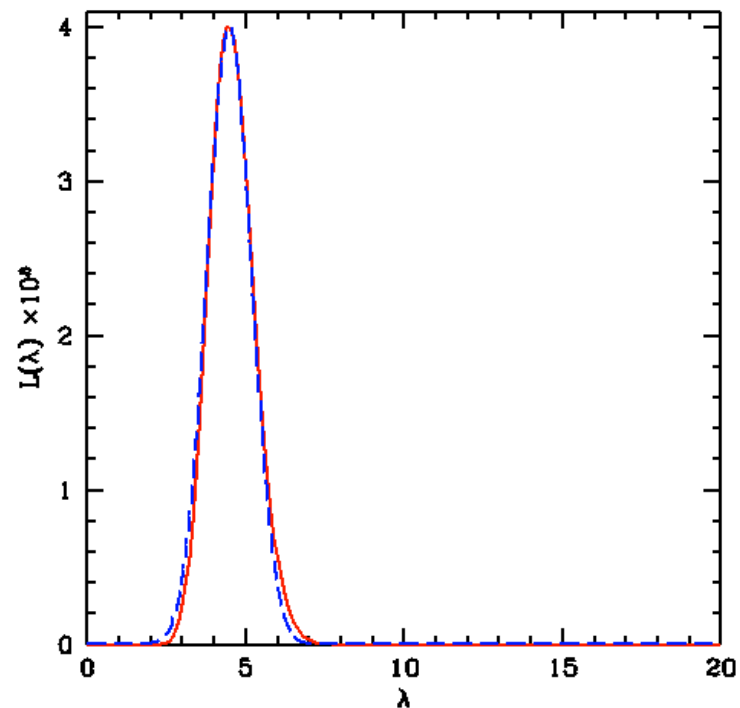
# Maximum Likelihood

- The likelihood function is a statistic (i.e. a function of the data) which gives the probability of obtaining that particular set of data, given the chosen parameters $\theta_1$, ... , $\theta_k$ of the model. It should be understood as a function of the unknown model parameters (but it is NOT a probability distribution for them)

- The values of these parameters that maximize the sample likelihood are known as the Maximum Likelihood Estimates or MLE's.

- Assuming that the likelihood function is differentiable, estimation is done by solving

$$\frac{\partial L(\theta_1,...,\theta_k)}{\partial \theta_i} = 0 \qquad \text{or} \qquad \frac{\partial \ln L(\theta_1,...,\theta_k)}{\partial \theta_i} = 0$$

- On the other hand, the maximum value may not exists at all.

# Back to counting cars

- After 9 experiments we collected the following data: 7, 4, 2, 6, 4, 5, 3, 4, 5. The new likelihood function is plotted below, together with a Gaussian function (dashed line) which matches the position and the curvature of the likelihood peak ($\lambda$ =4.44). Note that the 2 curves are very similar (especially close to the peak), and this is not by chance.

# Score and information matrix

- The first derivative of the log-likelihood function with respect to the different parameters is called the Fisher score function:

$$S_i = \frac{\partial \ln L(\theta)}{\partial \theta_i}$$

- The Fisher score vanishes at the MLE.

- The negative of the Hessian matrix of the log-likelihood function with respect to the different parameters is called the observed information matrix:

$$O_{ij} = -\frac{\partial^2 \ln L(\theta)}{\partial \theta_i \, \partial \theta_j}$$

- The observed information matrix is definite positive at the MLE. Its elements tell us how broad is the likelihood function close to its peak and thus with what accuracy we determined the model parameters.

# Example

**1 datapoint**
**Low information**
**Large uncertainty in $\lambda$**

**9 datapoints**
**High information**
**Small uncertainty in $\lambda$**

# Fisher information matrix

- If we took different data, then the likelihood function for the parameters would have been a bit different and so its score function and the observed information matrix.

- Fisher introduced the concept of information matrix by taking the ideal ensemble average (over all possible datasets of a given size) of the observed information matrix (evaluated at the true value of the parameters).

$$F_{ij} = -\left\langle \frac{\partial^2 \ln L(\theta)}{\partial \theta_i \, \partial \theta_j} \right\rangle$$

- Under mild regularity conditions, it can been shown that the Fisher information matrix also corresponds to

$$F_{ij} = \left\langle \frac{\partial \ln L(\theta)}{\partial \theta_i} \frac{\partial \ln L(\theta)}{\partial \theta_j} \right\rangle$$

i.e. to the covariance matrix of the scores at the MLE's.

# Cramér-Rao bound

- The Cramér-Rao bound states that, for ANY unbiased estimator of a model parameter $\theta_i$, the measurement error (keeping the other parameters constant) satisfies

$$\Delta\theta_i \geq \frac{1}{\sqrt{F_{ii}}}$$

- For marginal errors that also account for the variability of the other parameters (see slide 75 for a precise definition), instead, it is the inverse of the Fisher information matrix that matters and

$$\Delta\theta_i \geq \sqrt{F_{ii}^{-1}}$$

# Fisher matrix with Gaussian errors

- For data with Gaussian errors, the Fisher matrix assumes the form (the notation is the same as in slide 43)

$$F_{ij} = \frac{1}{2}\mathrm{Tr}\left[ C^{-1}C,_i\, C^{-1}C,_j + C^{-1}M_{ij} \right]$$

where

Information from the noise    Information from the signal

$$M_{ij} = m,_i\, m,_j^T + m,_j\, m,_i^T$$

(note that commas indicate derivatives with respect to the parameters while data indices are understood)

# Properties of MLE's

As the sample size increases to infinity (under weak regularity conditions):

- MLE's become asymptotically efficient and asymptotically unbiased
- MLE's asymptotically follow a normal distribution with covariance matrix (of the parameters) equal to the inverse of the Fisher's information matrix (that is determined by the covariance matrix of the data).

However, for small samples,

- MLE's can be heavily biased and the large-sample optimality does not apply

# Maximizing likelihood functions

- For models with a few parameters, it is possible to evaluate the likelihood function on a finely spaced grid and search for its minimum (or use a numerical minimisation algorithm).

- For a number of parameters >>2 it is NOT feasible to have a grid (e.g. 10 point in each parameter direction, 12 parameters = $10^{12}$ likelihood evaluations!!!)

- Special statistical and numerical methods needs to be used to perform model fitting.

- Note that typical cosmological problems consider models with a number of parameters ranging between 6 and 20.

# Forecasting

- Forecasting is the process of estimating the performance of future experiments for which data are not yet available

- It is a key step for the optimization of experimental design (e.g. how large must be my survey if I want to determine a particular parameter to 1% accuracy?)

- The basic formalism has been developed by Fisher in 1935

# Figure of merit



Figure of merit = 1 / (area of the ellipse)

# iCOSMO.org

# Open source Fisher matrices

# Fisher 4cast (Matlab toolbox)

# Counting cars, again

- In our study of the car counts we implicitly assumed that all the values of $\lambda$ are equally likely a priori (i.e. before we started taking the data). However, we didn't consider that an automatic gate regulates the traffic in our street and does not allow more than 8 cars to enter every 10 minutes. Therefore $\lambda$ cannot be larger than 8 and the likelihood derived from our counts should have been truncated at $\lambda$ =8.

- Also, we live close to a church and whenever there is a wedding the traffic is more intense than usual. This means that on wedding days a higher value of $\lambda$ is more likely than on non-wedding days.

- Moreover, a fellow that had been living in our flat before us did the same exercise and told us that he obtained $\lambda$ =4.2±0.5.

- Is there a way to account for all this information in our study?

# The Bayesian way



Bruno de Finetti (1906–1985)



Harold Jeffreys (1891–1989)

# What is probability?

- **<span style="color:red">Probability is a modern concept</span>** first discussed in a correspondence between Blaise Pascal and Pierre de Fermat in 1654

- There is no unique definition of probability, statisticians are divided into different schools with contrasting views

- <span style="color:red">Classic definition:</span> The probability of an event is the ratio of the number of cases favorable to it, to the number of all cases possible when nothing leads us to expect that any one of these cases should occur more than any other, which renders them, for us, equally possible. (Pierre-Simon de Laplace, *Théorie analytique des probabilités*, 1812)

- This is based on the "principle of insufficient reason" (or principle of indifference) which states that when cases are only distinguishable by their name they should be assigned the same probability

# What is probability?

- **Frequentist:** the long-run expected frequency of occurrence of a random event

- **Axiomatic:** given a sample space $\Omega$, a $\sigma$-algebra F of events E (a set of subsets of $\Omega$), we call probability measure a real function on F such that $P(E) \geq 0$, $P(\Omega) = 1$, and for any countable series of pairwise disjoint events $P(E_1 \cup E_2 \cup \ldots \cup E_N) = P(E_1) + P(E_2) + \ldots + P(E_N)$. These are known as Kolmogorov axioms.

- **Bayesian:** a measure of the degree of belief (the plausibility of an event given incomplete knowledge)

# Reasoning with beliefs

- There is 90% chance that today it will rain

- There is a 30% chance that my favourite football team will win the league this year

- There is a 10% chance that I will fail the observational cosmology examination

- There is a 0.1% chance that I will die before being 30

- There is 68.3% chance that $H_0$ lies between 67 and 73 km/s/Mpc

# De Finetti's game
# Can you measure degree of belief?

Suppose we are on a trip and you say that you are "pretty sure" you locked the door of your flat. I want to determine how sure you are.

- I offer you to play a game: I propose you to draw a marble from a bag containing 95 red and 5 blue marbles. If you pick at random a red marble, I give you one million euros. Alternatively, I offer you to go back home and check the door. If you choose this option and the door is locked indeed, I give you one million euros.

- If you choose to pick a marble, it means that your degree of belief is lower than 95%

- I can then propose many other rounds of the game by progressively reducing the fraction of red marbles until you choose to go back. This would measure your degree of belief.

# Posterior probability

$P(x|\theta)$: old name "direct probability"
It gives the probability of contingent events (i.e. observed data) for a given hypothesis (i.e. a model with known parameters $\theta$)

$L(\theta)=P(x|\theta)$: modern name "likelihood function" or simply "likelihood"
It quantifies the likelihood that the observed data would have been observed as a function of the unknown model parameters (it can be used to rank the plausibility of model parameters but it is not a probability density for $\theta$)

$P(\theta|x)$: old name "inverse probability"
        modern name "posterior probability"
Starting from observed events and a model, it gives the probability of the hypotheses that may explain the observed data (i.e. of the unknown model parameters)

# Bayes theorem


Rev. Thomas Bayes (1702–1761)


Pierre Simon (Marquis de) Laplace (1749–1827)

$$p(\theta \,|\, x) = \frac{p(x \,|\, \theta)\, p(\theta)}{p(x)}$$

**Prior probability** for the parameters (what we know before performing the experiment)

**Posterior probability** for the parameters given the data

**Likelihood function**

$$p(x \,|\, \theta) = L(x \,|\, \theta)$$

**Evidence** (normalization constant useful for Bayesian model selection)

$$p(x) = \int p(x \,|\, \theta)\, p(\theta)\, d\theta$$

# Bayes theorem

Let us re-write Bayes theorem emphasizing that all the probabilities that we generically called "p" are actually different functions:

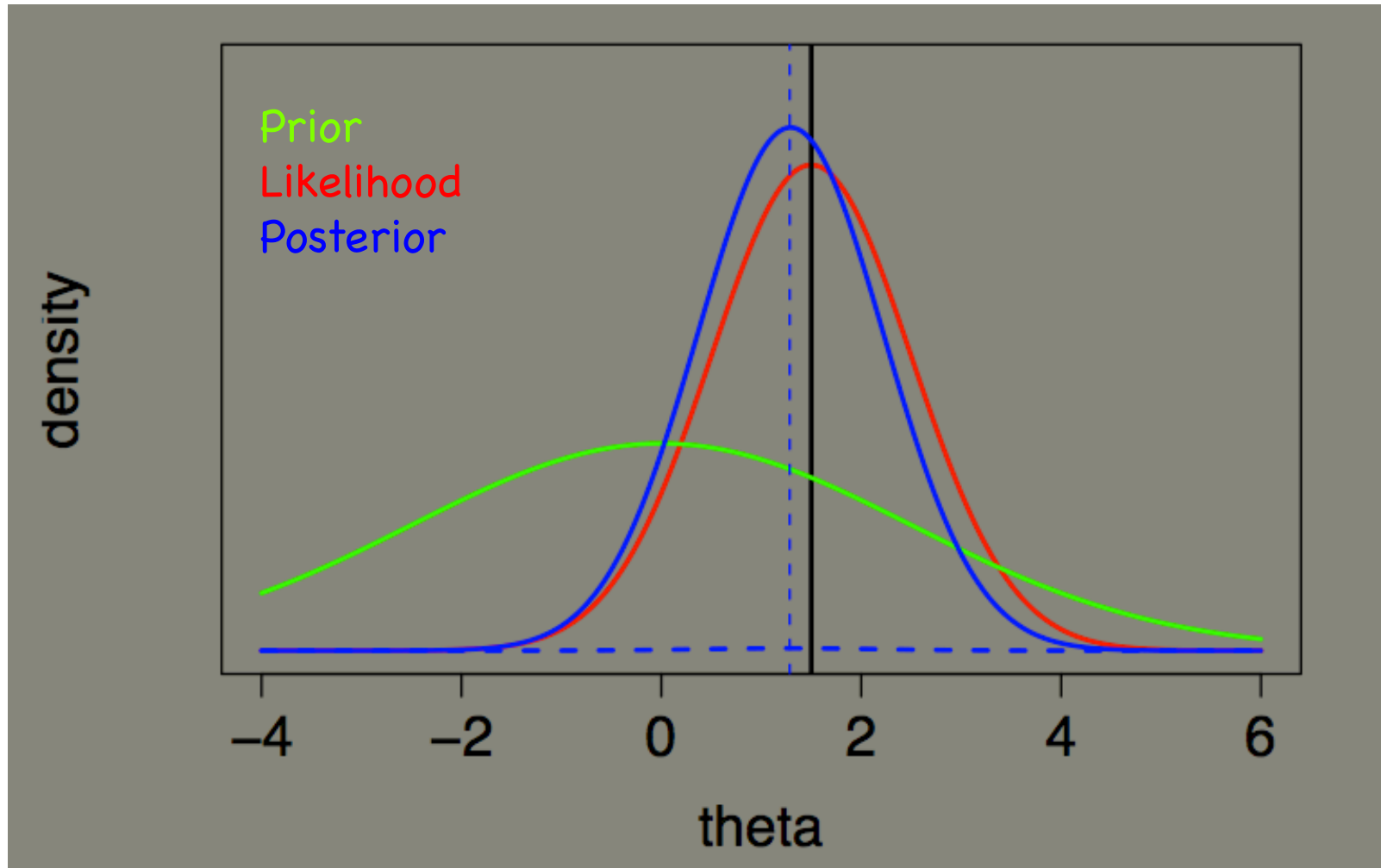$$P(\theta \mid x) = \frac{L(x \mid \theta)\,\pi(\theta)}{E(x)}$$
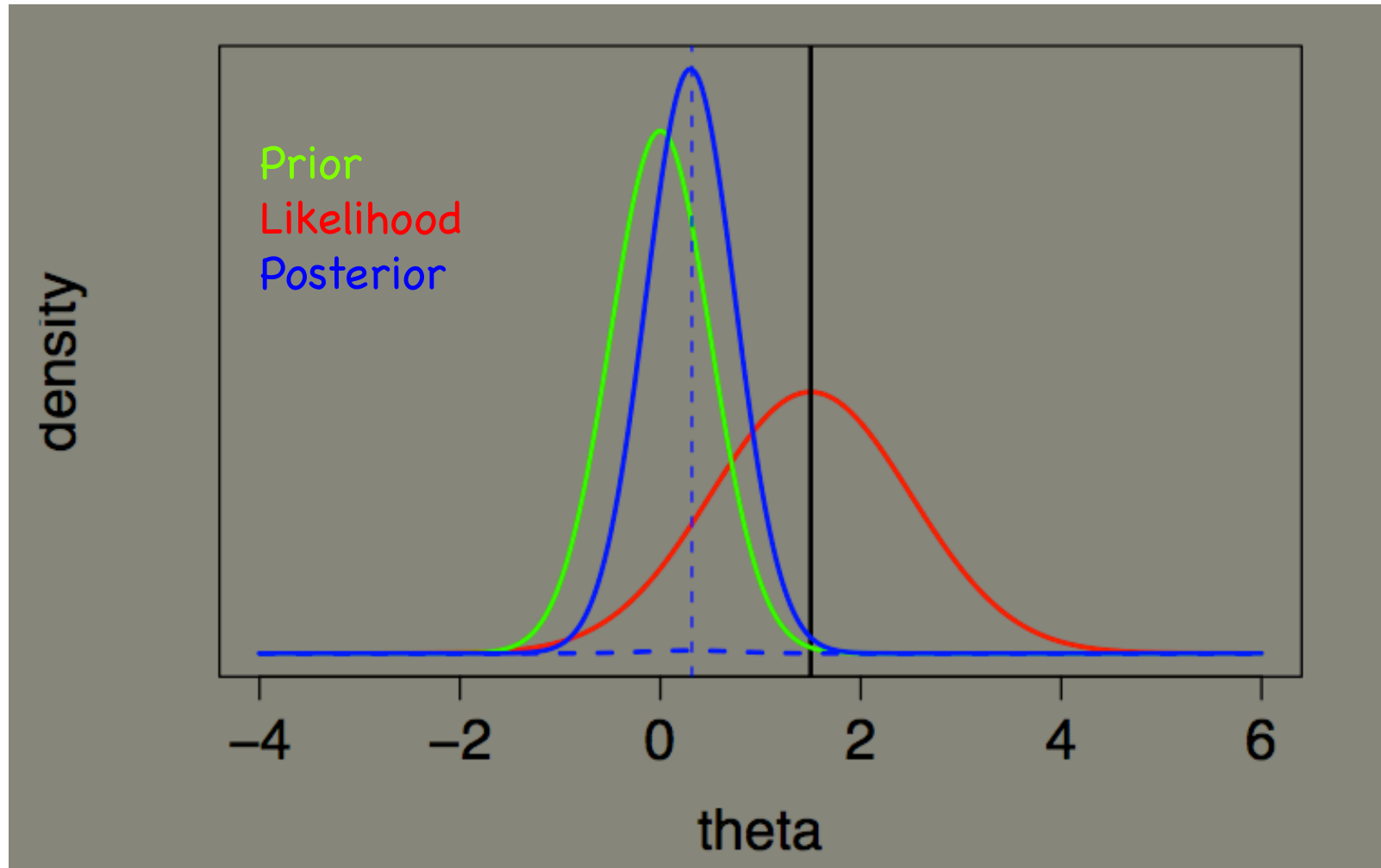
# Bayes theorem visually

# Priors

- Non-informative: if it has a minimal impact on the posterior distribution (i.e. it is "flat" with respect to the likelihood function). Non-informative priors are also called "vague", "diffuse" or "flat".

- Improper: if $\int \pi(\vartheta) d\theta = \infty$ e.g. uniform prior on the real line. Generally leads to a proper posterior but, if also the posterior is improper, inference is invalid.

- Informative: a prior which is not dominated by the likelihood. Must be handled with care in actual practice. On the other hand, illustrates the power of the Bayesian method: information gathered from previous study, past experience, or expert opinion can be combined with current information in a natural way

# Not very informative prior



density

Prior
Likelihood
Posterior

theta

# Very informative prior

# Bayesian estimation

- In the Bayesian approach to statistics, population parameters are associated with a posterior probability which quantifies our DEGREE OF BELIEF in the different values

- Sometimes it is convenient to introduce estimators obtained by minimizing the posterior expected value of a loss function

- For instance one might want to minimize the mean square error, which leads to using the mean value of the posterior distribution as an estimator

- If, instead one prefers to keep functional invariance, the median of the posterior distribution has to be chosen

- Remember, however, that whatever choice you make is somewhat arbitrary as the relevant information is the entire posterior probability density.

# Estimation: frequentist vs Bayesian

- Frequentist: there are TRUE population parameters that are unknown and can only be estimated by the data

- Bayesian: only data are real. The population parameters are an abstraction, and as such some values are more believable than others based on the data and on prior beliefs.

# Confidence vs. credibility intervals

- **Confidence intervals** (Frequentist): measure the variability due to sampling from a fixed distribution with the TRUE parameter values. If I repeat the experiment many times, what is the range within which 95% of the results will contain the true values?

- **Credibility interval** (Bayesian): For a given significance level, what is the range I believe the parameters of a model can assume given the data we have measured?

- They are profoundly **DIFFERENT** things even though they are often confused. Sometimes practitioners tend use the term "confidence intervals" in all cases and this is ok because they understand what they mean but this might be confusing for the less experienced readers of their papers. PAY ATTENTION!