

# Cosmological Observations

**Cosmological parameter estimation**

## Statistical inference

Cosmological datasets are usually large, noisy and with systematics → the problem to solve is not a straightforward system of equations relating precise values of an observable with a cosmological function of the parameters.

**To solve the problem requires the use of random variables and the calculation of probabilities.**

The standard way to **estimate the values of the model free parameters** (the cosmological parameters and the nuisance parameters) in cosmological analyses is through [Bayesian inference](#).

### Additional Bibliography:

- David MacKay - “Information theory, Inference and Learning Algorithms” (chapter 29), CUP 2003
- Hobson, Jaffe, Liddle, Mukherjee and Parkinson - “Bayesian methods in Cosmology” includes chapters on MCMC (Lewis, Bridle 2002) and model selection (Mukherjee, Parkinson, Liddle 2005)
- Amendola and Tsujikawa - “Dark Energy” (chapter 13) includes sections on Fisher matrix and analytical marginalizations
- Dan Coe - “Fisher matrices - a quick-start guide and software” 2009 arXiv:0906.4123

## Forward Probability (the frequentist approach)

vs.

## Inverse Probability (the Bayesian approach)

The goal is to compute the probability distribution of the data given the fixed and true value of the parameters. **Data are random variables** with a **probability density function** (pdf). Their probability corresponds to the **frequency** with which its values occur in repetitions of the experiment.

A **statistic** is computed from the data and a pdf is derived for the statistic. From values obtained for a statistic with a known pdf, a rejection level may be assigned to a **hypothesis** (a parameter value).

**There is no probability distribution of the parameter values, they are absolute quantities.**

Here **the vector of parameters is a random variable**, and has a probability function. They are unobserved variables.

We want to compute that probability: a **conditional probability given the data**.

**Data are also random variables** and a **joint probability** may be defined:  $P(m,d)$

Note:

d = data (the estimated physical property)

m = model (the values of the parameters)

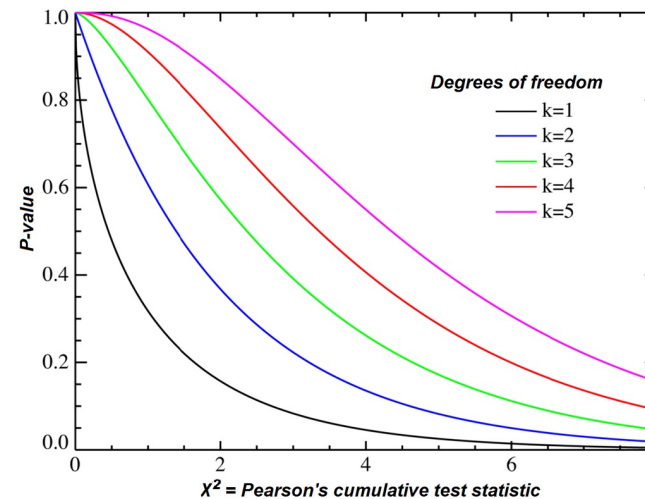
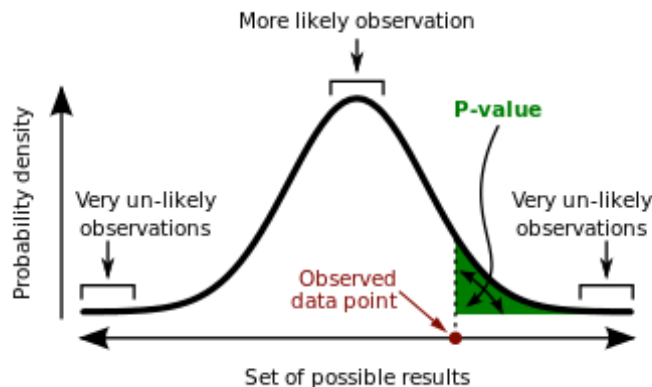
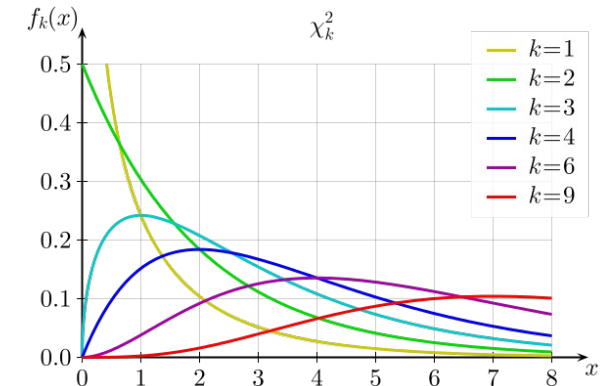
# Example : using the chi-square statistic in the frequentist approach

The chi-square is an example of a statistic

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

with known pdf (also called the chi-squared distribution)

Computing the chi-square value and knowing the pdf we can compute the **p-value** → if  $p\_value < threshold$  then the data rejects the hypothesis (the fact of the parameter value being the assumed one)



## Example : using the chi-square statistic in the **Bayesian approach**

The conditional probability of the *data given the model* is a Gaussian in the chi-square statistic.

Through Bayes theorem this implies that the conditional probability of the *model given the data* is well sampled by the values of that Gaussian.

Both methods use the chi-squared statistic, but in ***frequentist hypothesis testing*** the crucial information is the chi-squared distribution, while in ***Bayesian parameter inference*** the crucial information is the theoretical  $d(m)$  expression.

The joint probability may be written in terms of the Probability of m conditional to d (the probability of m to equal  $m_i$ , given that the data equals  $d_i$ ), and the intrinsic probability of the data to be equal to  $d_i$  :

$$P(m,d) = P(m|d) P(d)$$

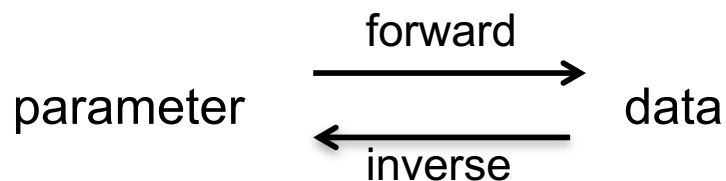
or also the other way around :

$$P(m,d) = P(d|m) P(m)$$

So, in Bayesian Inference we consider 2 spaces:

the **data space** (where random variables d live)  $P(d|m)$

the **parameter space** (where random variables m live)  $P(m|d)$



The two spaces are related through **Bayes theorem**, which we can obtain by equating the two expressions:

$$P(m|d) = \frac{P(d|m) P(m)}{P(d)}$$

$P(m|d)$  is the probability of the parameter values given the data. It is a distribution in the parameter space.

It is known as the **posterior** distribution → this is what we want to get.

$P(m)$  is the probability of the parameter values independently of these data → it can be something we know beforehand from another experiment, or from some intrinsic property of the model, or it may just be flat (no special restriction on that parameter).

It is known as the **prior**.

$P(d|m)$  is the probability of getting the measured data given the parameter values. It is a distribution in the data space.

Remember, in inverse probability we do not want to study the properties of the data space but rather we want to use  $P(d|m)$  to **compute/infer**  $P(m|d)$ .

Now, it is reasonable to assume that if the probability of the observed data, given a parameter value, is low (high), then that value is unlikely (likely) to occur.

**For these two reasons  $P(d|m)$  is named the **likelihood** of the *parameters*,  $L(m)$ , even though it is a quantity in the *data* space.**



$P(d)$  is the probability of the data independently of the parameter values. It may be obtained from the joint probability by integrating over the full range of parameter values, i.e., **marginalizing** over *all* parameters.

$$P(d) = \int_m P(d|m) P(m)$$

Being independent of  $m$ , it is a **normalization** constant for  $P(m|d)$ , in Bayes theorem.

Note that it is independent of the parameter values, but not on the modeling, and its value may be used as a criteria for **model comparison**.

For this reason, it is known as the **evidence**.

Note that for any **parametrization/theory/model** (e.g.:  $\Omega_m$ ,  $\Omega_\wedge$ ), (e.g.  $\Omega_m$ ,  $\Omega_\wedge$ ,  $\Omega_v$ ) the whole universe of possible **models/parameter values** has a total probability of 1.

Thus, when working within one case,  $P(m|d)$  may be renormalized to 1 and the evidence is not needed. But when comparing two cases, the absolute value has valuable information: the highest absolute value is the preferred case, hence the name  $\rightarrow$  there is highest evidence for that case.

## The data space

We know many things about the data space:

- we have a sample of the distribution there (the measured data)
- we know **moments of the distribution** - the mean, the variance-covariances (either from computing the average and dispersion of the measured sample, or computing from theory e.g.  $d(m)$  )

For most practical applications, data is large (even one measurement of SN magnitude involves a large number of independent photons) and the central limit theorem tells us that the full distribution (for which we just have a sample) must be a **Gaussian**.

$$P(d|m) = L(m) = \frac{1}{(2\pi)^n |C|^{1/2}} \exp \left[ -\frac{1}{2} (\bar{d} - d(m)) C^{-1} (\bar{d} - d(m))^T \right]$$

ex: 2D

$$P(d|m) = A \exp \left[ -\frac{1}{2} \frac{1}{(1-\rho^2)} \left( \frac{(\bar{d}_1 - d_1(m))^2}{\sigma_1^2} + \frac{(\bar{d}_2 - d_2(m))^2}{\sigma_2^2} - \frac{2\rho (\bar{d}_1 - d_1(m))(\bar{d}_2 - d_2(m))}{\sigma_1 \sigma_2} \right) \right]$$

data is 2D  $[C] = 2 \times 2 = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$  Covariance matrix

Correlation matrix is the covariance normalized by each sigma, such as to leave only the correlation information but not the absolute values

$$= \begin{bmatrix} \frac{\sigma_{11}}{\sigma_{11}} & \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}} \\ \frac{\sigma_{21}}{\sqrt{\sigma_{11}\sigma_{22}}} & \frac{\sigma_{22}}{\sigma_{22}} \end{bmatrix} = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

Notation: 2 indices mean a square quantity

$$\sigma_{ii} \text{ is } \sigma^2$$

$$\sigma_i \text{ is } \sqrt{\sigma^2} = \sigma$$

Correlation  $\rho = \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}} = \frac{\sigma_{12}}{\sigma_1\sigma_2}$

Note that each data point  $d_i$  is usually obtained from many measurements.

It is a random variable with mean  $\bar{d}_i$  and variance  $\sigma_{ii}$  (std error  $\sqrt{\sigma_{ii}} = \sigma_i$ )

Ex: SN  $d_i = \mu_{z_i}$

CMB  $d_i = C_{\ell_i}$

The full  $\underline{d}$  vector may have very large dimension.

Ex: SN in 100 redshift bins

$$\mu_{z=0.1}, \mu_{z=0.15}, \dots$$

$$[C] = 100 \times 100$$

CMB in 1000  $l$ -modes

$$[c] = 1000 \times 1000$$

$$C_1, C_2, \dots, C_{100}, \dots, C_{1000}$$

the number of correlations  $l$  is:

$$\frac{\overset{\text{total number}}{m^2} - m \rightarrow \text{diagonal}}{\underset{\text{symmetry}}{2}} = \frac{m(m-1)}{2}$$

So for SN we write:

$$P(d|m) = A \exp \left\{ -\frac{1}{2} \left[ \mu_{\text{obs}} - (5 \log_{10} D_L(m) + 25) \right]^T C^{-1} \left[ \mu_{\text{obs}} - (5 \log_{10} D_L(m) + 25) \right] \right\}$$

For  $N$  redshift bins  $C$  is a  $N \times N$  matrix, with  $N$  variances  
and  $\frac{N(N-1)}{2}$  covariances

and  $\mu$  and  $D_L$  are  $N$ -dimensional vectors

## The parameter space

From Bayes theorem we can compute the posterior from the likelihood, if we know the prior and if we renormalize the evidence.

Notice that a Gaussian likelihood does not necessarily lead to a Gaussian posterior (even in the case of a flat prior), because changing from one space to the other involves an inversion  $d(m) \rightarrow m(d)$

Only in the case that the response of the observable is linear in the parameter values, will the posterior also be a Gaussian.

→ this is the case for geometrical probes  $D(z;m)$ , i.e., for SN, BAO, but not for structure formation probes  $P(k;m)$ .

Assuming Gaussian,

$$P(m|d) = A \exp \left[ -\frac{1}{2} (\bar{m} - m(d))^T C^{-1} (\bar{m} - m(d)) \right]$$

Now, the dimension of  $C$  is the dimension of the parameter space

$[C] = p \times p \rightarrow$  number of parameters, not related with number of redshift bins or  $l$ -modes.

Usually it is much lower than the dimension of the data space.

In 2D it has also the general form

$$C = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

where 1, 2 are for example  $\Omega_m, \Omega_\Lambda$   
and not  $z_1, z_2$



Analogous to the data space, we could say that one value of the data vector will imply one value of the parameter vector  $m$ , which compared to the true (or fiducial) parameter vector  $\bar{m}$ , will give a value of  $P(m|d)$ , i.e., the likelihood of the data.

We usually do not say "the likelihood of the data" because the data is the quantity measured (known) and not the model.

To compute  $P(m|d)$  exactly and analytically we would need to write

$$m(d)$$

This is usually impossible to do.

A possible way forward is to notice that  $\frac{\partial^2 \ln P(m|d)}{\partial m^2} \Big|_{m_0} = C_{pp}^{-1}$

and use Bayes theorem to relate

$$\frac{\partial^2 \ln P(m|d)}{\partial m^2} \Big|_{m_0} \text{ with } \frac{\partial^2 \ln P(d|m)}{\partial d^2} \Big|_d = C_{dd}^{-1}$$

To estimate the parameter values and their uncertainties from data, we need to find their distribution in the parameters' space, i.e., the posterior distribution (or its moments, since in practice we do not need the full distribution).

There are two general ways of doing this:

- **Sampling the distribution** - we can get a sample of the posterior distribution by direct computation of the likelihood on a **grid**, or by using stochastic methods (**Monte Carlo**)
- **Fisher matrix** - we can compute a lower limit for the second-order moments of the posterior distribution (i.e., the variance of the parameters) in a deterministic way. However, we cannot compute the first-order moment in a similar way (i.e., the values of the parameters).

## Sampling the distribution: in a deterministic way

### Grid

Since the likelihood is proportional to the posterior distribution that we want to find, a direct way to sample it is to compute its values at some points in the parameter space:

compute the likelihood in a **grid** - a hypercube of likelihoods.

This does not give us a sample distributed as the posterior, but just give us some values of that function.

### Disadvantages:

- the resolution of the grid may be too low to make contours,
- will waste time computing in low likelihood places,
- the number of required points increase fast with dimension of the grid

**Maximization** means to reduce a dimension by fixing it in the grid.

**Marginalization** means to reduce a dimension by summing along it on the grid.

## Analytical marginalization (over nuisance parameters)

One way to decrease the dimension of the problem, making it possible to compute a grid of lower dimension is to marginalize the likelihood in advance,

i.e., to integrate the likelihood dependence on one or more parameters, obtaining a new likelihood with less dimensions. This should mainly be done to parameters we are not interested in.

Let us consider a data vector  $x_i$  (for example the distance modulus measurements at various redshifts, with associated error bars  $\sigma_i$ ), and the theoretical vector  $m_i$  (for example the distance modulus computed for the same redshifts, which is a function of the values of the cosmological parameters).

The Gaussian likelihood of a theoretical model given the data vector is:

$$L(p) = N \exp \left[ -\frac{1}{2} \sum_z \frac{(\bar{d}_z - d_z(p))^2}{\sigma_z^2} \right]$$

## Marginalization with an additive bias

Now, consider that there is a systematic effect contributing to the distance modulus in an additive way, parameterized by a parameter  $\alpha$ .

Therefore, the theoretical prediction, now including that effect, is:

$$d_i \rightarrow d_i + \alpha$$

This means that the theoretical model that will be applied to fit the data gets an extra parameter:  $d(p_1, \dots, p_n) + \alpha$

We need to estimate the cosmological parameters ( $p_i$ ) in the presence of  $\alpha$ , i.e., allowing for all possible values of  $\alpha$ .

Instead of building a  $N+1$  dimension grid (and since we are not interesting in estimating  $\alpha$ , but only in including its impact on the estimation of  $p_i$ ), we can marginalize a priori over all possible values of  $\alpha$ .

Let us then marginalize over a generic additive parameter:

$$\text{notation } \left\{ \begin{array}{l} \text{observed} = \bar{d}_i \\ \text{theoretical} : d_i(m) = d_i + \alpha \end{array} \right.$$

To marginalize is to get the constraint on the parameter <sup>sub</sup>vector  $p$ , assuming all possible values of  $\alpha$ , i.e., to integrate the likelihood over  $\alpha$ , obtaining a new expression for the likelihood.

$$\text{So, } L(m) = A \exp \left[ -\frac{1}{2} \sum_i \frac{(\bar{d}_i - d_i - \alpha)^2}{\sigma_i^2} \right]$$

Marginalize over  $\alpha$

$$\rightarrow L_{\text{New}}(m) = A \int d\alpha L(m) =$$

$$= A \int_{-\infty}^{+\infty} d\alpha \exp \left( -\frac{1}{2} \sum_i \frac{(\bar{d}_i - d_i)^2 + \alpha^2 - 2\alpha(\bar{d}_i - d_i)}{\sigma_i^2} \right) =$$

$$= A e^{-S_2/2} \int_{-\infty}^{+\infty} d\alpha e^{(\alpha S_1 - \alpha^2 S_0/2)}$$

part independent  
of  $\alpha$

where we defined

$$S_0 = \sum \frac{1}{\sigma_i^2}$$

$$S_1 = \sum \frac{(\bar{d}_i - d_i)}{\sigma_i^2}$$

$$S_2 = \sum \left( \frac{\bar{d}_i - d_i}{\sigma_i} \right)^2$$

Note that  $S_1$  and  $S_2$  depend on the model parameters (in  $d_i$ ) but not  $S_0$

$$= A e^{-\frac{S_2}{2}} \int_{-\infty}^{+\infty} d\alpha e^{-\frac{1}{2} S_0 \left( \alpha^2 - 2\alpha \frac{S_1}{S_0} \right)}$$

$$= A e^{-\frac{S_2}{2}} \int_{-\infty}^{+\infty} d\alpha e^{-\frac{1}{2} S_0 \left( \alpha - \frac{S_1}{S_0} \right)^2 + \frac{1}{2} \frac{S_1^2}{S_0}}$$

to complete the square

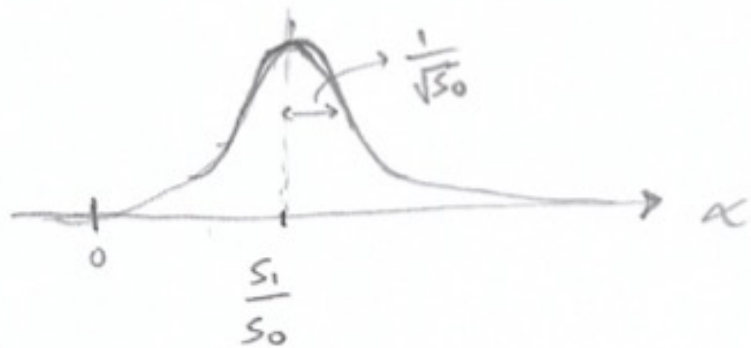
multiply this to cancel out the extra term that arises from the square.

$$= A e^{-\frac{1}{2} \left( S_2 - \frac{S_1^2}{S_0} \right)} \int_{-\infty}^{+\infty} d\alpha e^{-\frac{1}{2} S_0 \left( \alpha - \frac{S_1}{S_0} \right)^2}$$

This is a Gaussian in the variable  $\alpha$ , centered in  $\alpha = \frac{S_1}{S_0}$  with  $\sigma = \frac{1}{\sqrt{S_0}}$

Notice that the result of this integral only depends on the width of the Gaussian ( $S_0^{-1/2}$ ) and not on its central point  $S_1/S_0$ .

So it is just a constant, i.e., it is independent on the cosmological parameters contained in  $S_1$  and  $S_2$ .



The integral  $[-\infty, +\infty]$  of a Gaussian does not depend on the central point, only on the width:

$$\int_{-\infty}^{+\infty} e^{-\frac{1}{2\sigma^2}(x-x_0)^2} dx = \sqrt{(2\pi)^m \det C} = 2\pi\sigma = \frac{2\pi}{\sqrt{s_0}}$$

→ For 1 parameter ( $\alpha$ )  
 $m=1$ ,  $\det C = \sigma = \frac{1}{\sqrt{s_0}}$

⇒ The new likelihood (marginalized over  $\alpha$ ) is:

$$L_{\text{new}} = A_{\text{new}} e^{-\frac{1}{2} \left( s_2 - \frac{s_1^2}{s_0} \right)}$$

↓  
 new  
 constant  
 of normalization

$$A_{\text{new}} = \frac{1}{\sqrt{(2\pi)^m s_0 \det C}}$$



$$\Rightarrow \mathcal{L}_{\text{new}} = A_{\text{new}} \cdot e^{-\frac{1}{2} \sum \left( \frac{\bar{d}_i - d_i}{\sigma_i} \right)^2} \cdot \exp \frac{\left( \sum \frac{\bar{d}_i - d_i}{\sigma_i^2} \right)^2}{2 \left( \sum \frac{1}{\sigma_i^2} \right)}$$

The marginalized likelihood is like the normal likelihood with no  $\alpha$ , times a new likelihood factor.

The values of this likelihood on the points of a grid in the  $n$ -dimensional space of the cosmological parameters ( $p_1 \dots p_n$ ), are identical to the ones that would be obtained by first computing the original likelihood on the points of a grid in the  $(n+1)$ -dim space of the cosmological parameters ( $p_1 \dots p_n, \alpha$ ), and then summing up all the likelihood values along the  $\alpha$  dimension on each  $p$  (dim  $n$ ) point  $\rightarrow$  so this method only requires an  $n$ -dimension grid, instead of an  $(n+1)$ -dimension one.

## Marginalization with a multiplicative bias

If the systematic effect contributes to the distance modulus in a multiplicative way, parameterized by an  $\alpha$  parameter, the  $(n+1)$ -dim likelihood is:

$$L(p, \alpha) = N \exp \left[ -\frac{1}{2} \sum_z \frac{(\bar{d}_z - \alpha d_z(p))^2}{\sigma_z^2} \right]$$

Marginalizing over  $\alpha$ , (using the same approach as in the previous calculation) the  $n$ -dim likelihood becomes:

$$L_2(\theta) = N_2 \exp \left[ -\frac{1}{2} \left( \ln S_{02} - \frac{S_{11}^2}{S_{02}} \right) \right]$$

where

$$S_{ab} = \sum_z \frac{\bar{d}_z^a d_z(p)^b}{\sigma_z^2}$$

Note that the result depends on the way the effect is included in the modelling.

If the multiplicative bias parameter is applied to correct the data (instead of being included in the theoretical modelling), the  $(n+1)$ -dim likelihood is written as

$$L(p, \alpha) = N \exp \left[ -\frac{1}{2} \sum_z \frac{(\alpha \bar{d}_z - d_z(p))^2}{\sigma_i^2} \right]$$

In this case, after marginalizing over  $\alpha$  the resulting  $n$ -dim likelihood obtained is different.

It is given by

$$L_1(\theta) = N_1 \exp \left[ -\frac{1}{2} \left( S_{02} - \frac{S_{11}^2}{S_{20}} \right) \right]$$

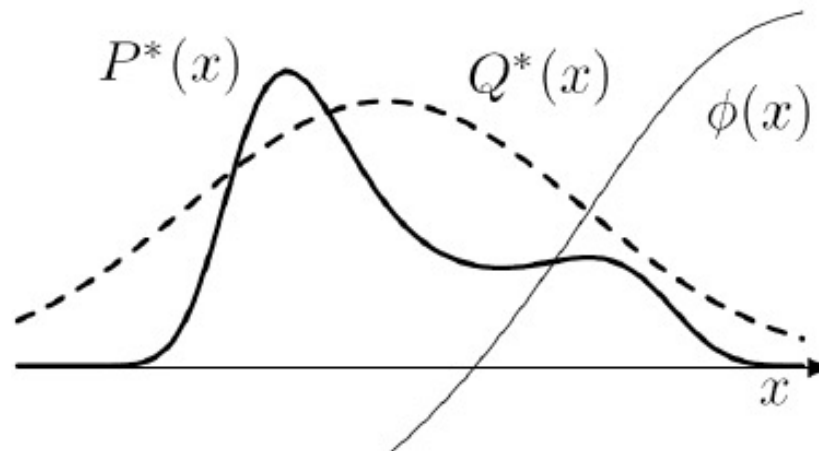
## Sampling the distribution: in a stochastic way

It a way is to get a **sample of the distribution**, i.e., a group of points in the parameter space in the same proportion as in the full distribution - and not just to know the values of the probability at certain points of the space (as with the grid).

### Importance Sampling

One possible way to do this is to sample from a known distribution (Q), that we assume will be similar to our **target distribution** (P).

(for example, Q may be a smooth approximation of P)



We need to define weights to get a true sample of P from a sample of Q.

Introduce weights

$$w_r = \frac{P^*(\mathbf{x}^{(r)})}{Q^*(\mathbf{x}^{(r)})}$$

This is impossible if we know nothing about P.

But if we suspect it may be similar to Q then this method is very useful, because we do not need to generate any points for P.

We just need to get the Q points and change their weights - for example computing the likelihood of those points (with our data).

The ratio between likelihood and their probability value under Q will be the new weight → the Q sample is changed into a P sample.

The advantage is that we did not need to compute likelihoods in points of a grid (which might be a bad coverage of the P sample) but on the sample points of Q (which are a better coverage of the P sample).

## Monte Carlo Markov chain (MCMC)

The Markov method is related with importance sampling in that we sample from an auxiliary distribution  $Q$ .

But now  $Q$  does not need to be similar to  $P$ .

We start sampling one point of  $P$  and center  $Q$  on that point.

Then we sample from  $Q$  - but  $Q$  depends on the current position in space.

$Q$  is not important - it may change from point to point.

This method builds correlated samples: each point depends on the previous one - this is the definition of a **Markovian process**.

This works if it fulfills the following properties:

It must be **irreducible** → there is a non-zero probability of reaching any model from any starting model.

For example, if the target distribution has several local maxima, it may happen that the chain cannot pass from one of those regions to another. In this case it may converge to different distributions, depending on the starting point of the chain.

It must be **aperiodic** → it must not oscillate between different sets of models in a periodic movement.

It must be **invariant** → once the chain follows the target distribution, all subsequent iterations will also have that same distribution.

The most used algorithm of Markov chain Monte Carlo (MCMC), has these three properties. It is called **Metropolis-Hastings**:

Given a point  $m$ , a candidate new point  $m'$  is generated from  $Q(m, m')$ .  
The point is accepted to be part of the sample with a certain probability:

$$\alpha(m, m') = \min \left( 1, \frac{P(m')Q(m'|m)}{P(m)Q(m|m')} \right)$$

$P$  are the likelihoods of the two points.

If  $Q$  are symmetric distributions  $\rightarrow Q(m'|m)=Q(m|m')$

In summary:

If  $P(m') > P(m) \rightarrow m'$  becomes a new point of the sample

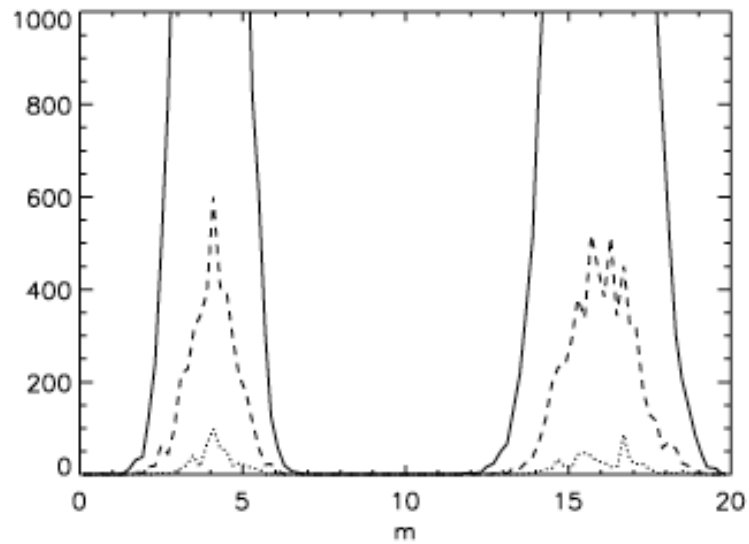
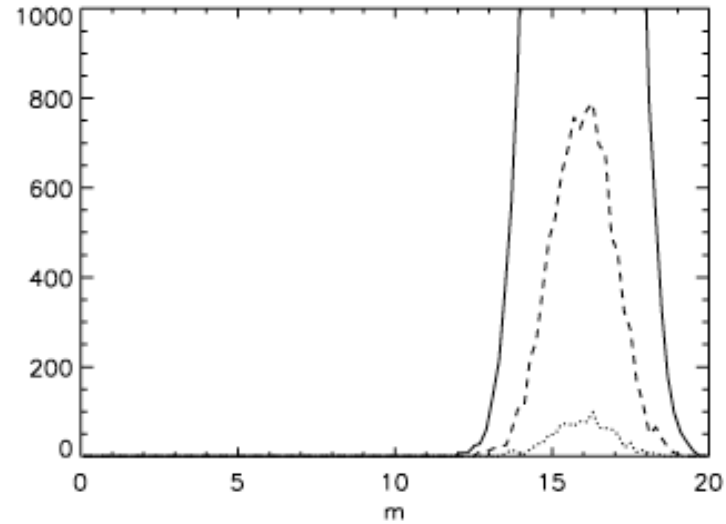
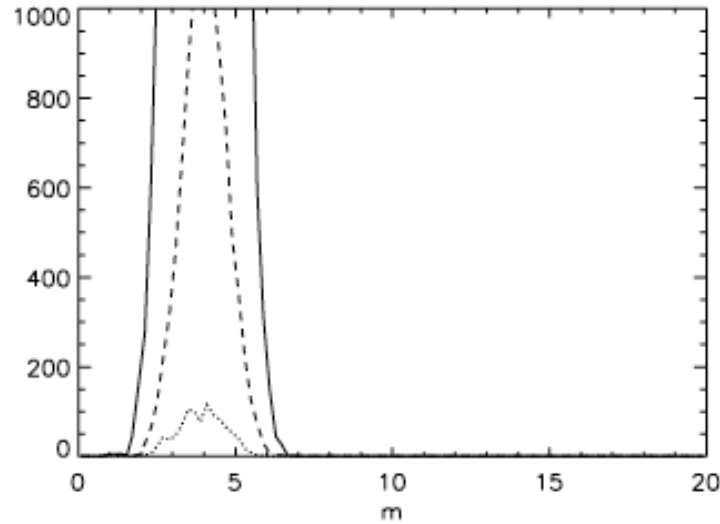
If  $P(m') < P(m) \rightarrow m'$  may or may not become a new point, with a probability  $P(m')/P(m)$ . The better it is, the better chance to be accepted.

When  $m'$  is not accepted, the chain stays at  $m \rightarrow$  the weight of  $m$  in the sample increases.



# Properties of MCMC

## Starting point



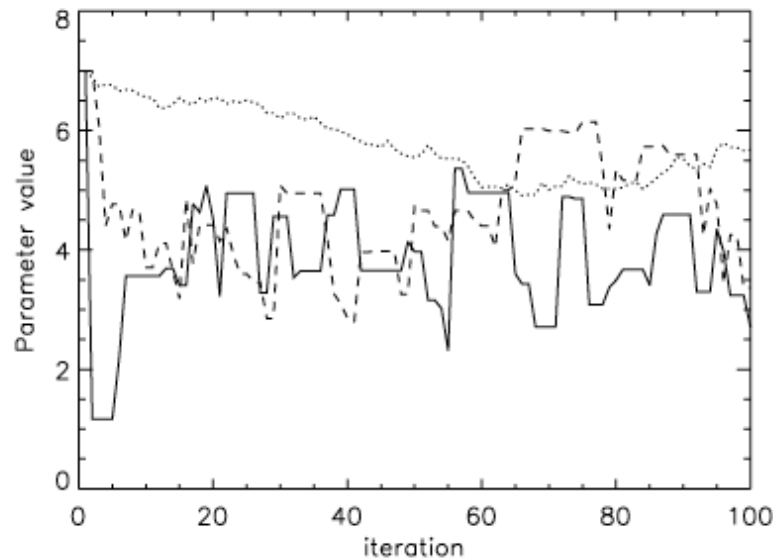
Sampling a binomial distribution with different choices of starting points.

## Step (scale of Q)

In comparison with the typical scale of P:

if too large → once the chain gets into a high posterior region, most of the subsequent proposed models will be in regions of lower posterior, and are likely to be rejected.

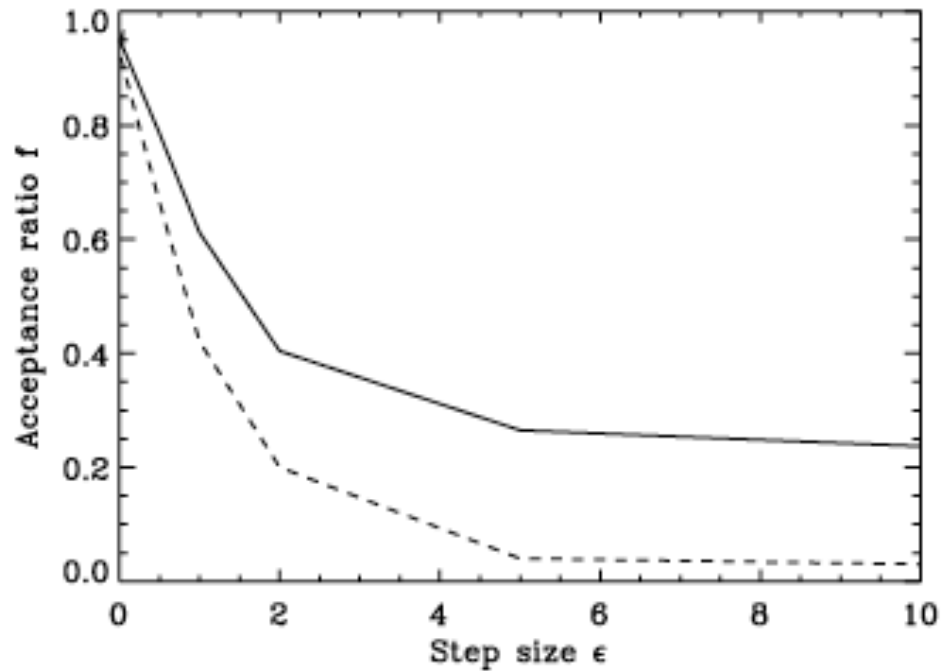
if too small → acceptance rate will be larger and the chain will move frequently. However, it will move in small steps, taking a long time to probe all space and being virtually non-irreducible.



increasing steps:

dot  
dash  
solid

## Acceptance rate



Each point in the line is made from a full chain.

**Optimal acceptance rate**  
0.3 - 0.5

**Optimal step size**  
 $2 \times \sigma_{\text{parameter}}$

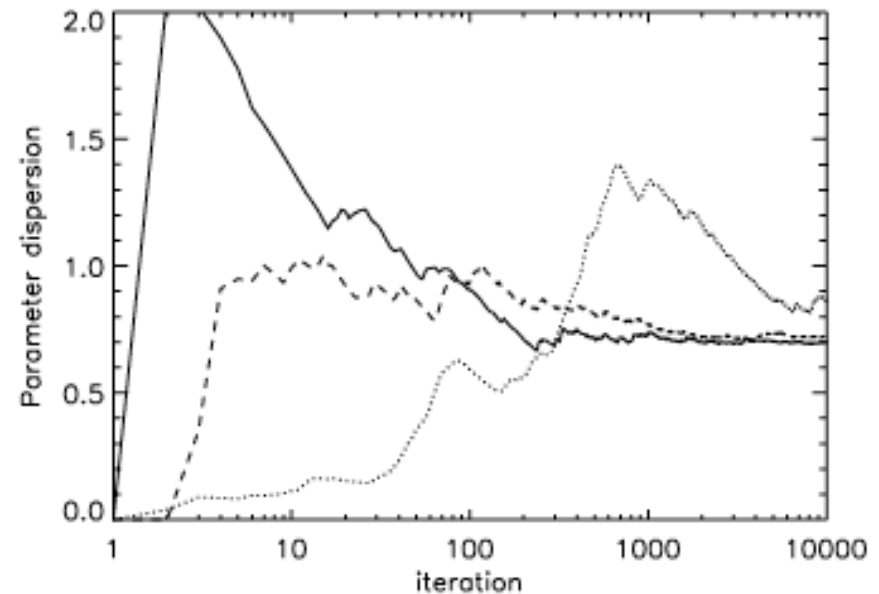
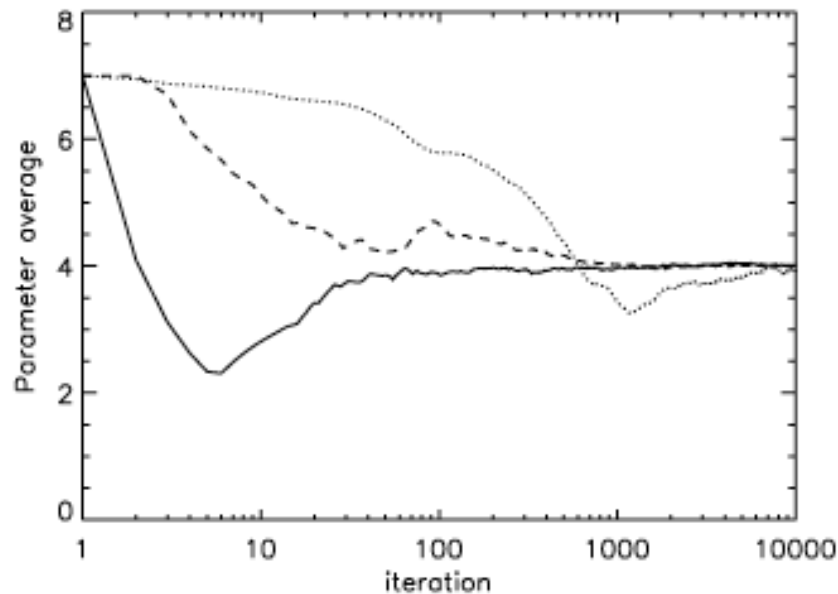
Q may have different scales on different directions.

Q may be chosen to be aligned with degeneracy directions for larger efficiency.

## Convergence

### How to assess convergence?

Convergence is related with the amount of time needed for the chain to start sampling from the target.



Comparing chains: the one with small step is the least efficient.

The part of the chain built before convergence need to be removed: the **burn-in** (it may be a large fraction of the chain).

## How to quantify convergence? the **Gelman-Rubin** convergence test

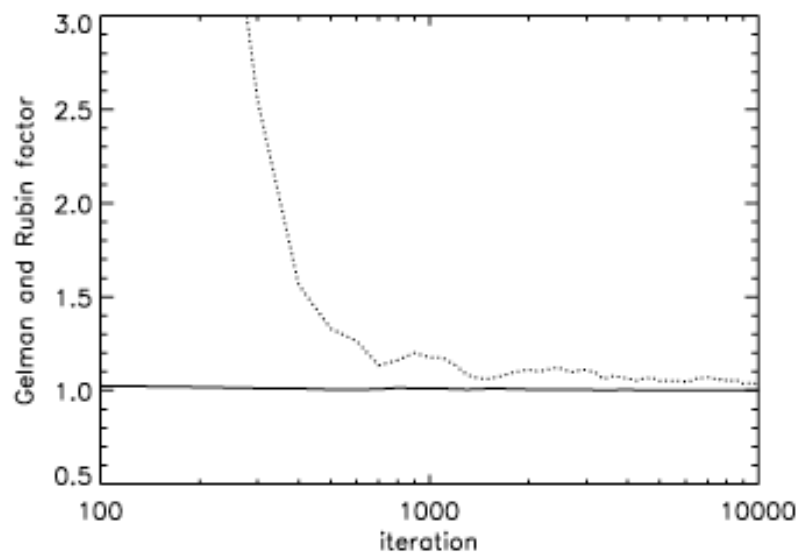
$$R = \frac{1}{W} \left( B + \frac{WI}{I-1} \right)$$

$$W(p) = \frac{1}{J} \sum_{j=1}^J \frac{1}{I-1} \sum_{i=I/2}^I [p_i(j) - \bar{p}(j)]^2 \quad \text{within-chain dispersion}$$

$$B(p) = \frac{1}{J-1} \sum_{i=1}^J [\bar{p}(j) - \bar{p}]^2 \quad \text{between-chain dispersion}$$

Convergence may require a long time  
(e.g. order  $10^6$  points)

For a large number of parameters ( $N > 4$ ),  
MCMC is usually faster than grid  
computation



## Correlation

The resulting chain is correlated  $\rightarrow$  samples are not independent

We can compute the **correlation between points as function of separation in the chain**: (values of a parameter  $p$  in positions  $i$  and  $i+j$ )

$$\text{variance} = \langle (p_i - p_0) (p_{i+j} - p_0) \rangle = \langle p_i p_{i+j} - p_i p_0 - p_0 p_{i+j} + p_0 p_0 \rangle$$

$$\text{since } \langle p_i \rangle = \langle p_{i+j} \rangle = p_0 \rightarrow \text{variance} = \langle p_i p_{i+j} \rangle - p_0^2$$

$$\text{covariance} = \langle (p_i - p_0) \rangle^2 = \langle p_i p_i \rangle - p_0^2$$

The correlation is the variance normalized by the covariance (and it has a value  $< 1$ ):

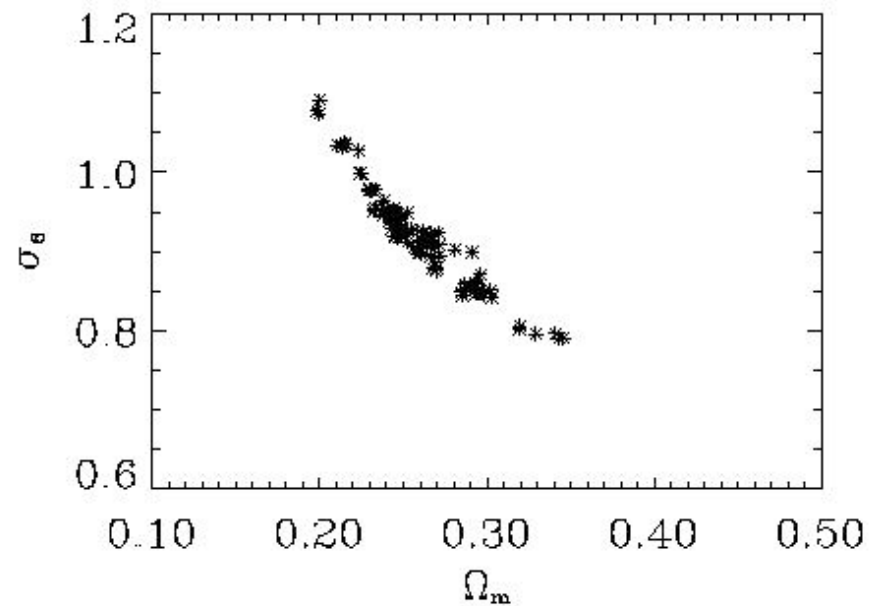
$$c_j = \frac{\langle \tau_i \tau_{i+j} \rangle - \langle \tau_i \rangle^2}{\langle \tau_i^2 \rangle - \langle \tau_i \rangle^2}$$

Reduce the **correlation length** of the chain : for each point of the chain remove the  $j$  subsequent points such that the  $c_j$  correlation is larger than a certain threshold (e.g.  $c_j > 0.5$ )  $\rightarrow$  **thin-out** the chain

## Output sample

The resulting chain - converged, with burn-in removed, and thinned-out - is a sample of the posterior in parameter space,  $P(m|d)$ .

A plot of the cloud of points directly shows the probability density of the sample  $P(m|d)$ .



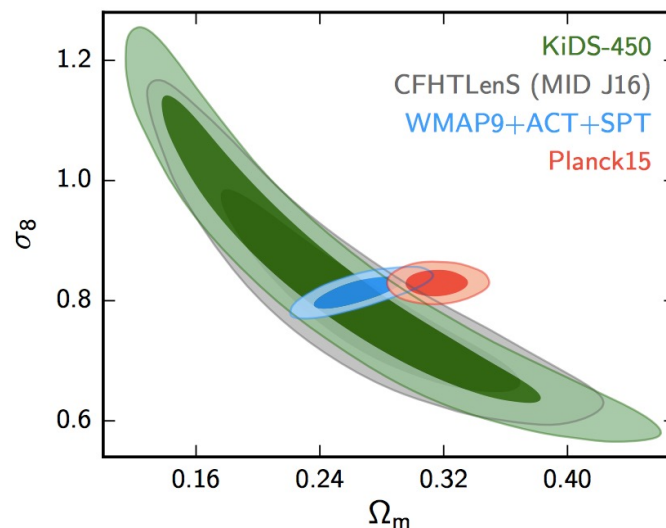
## Parameter constraints

Having obtained a converged representative sample of values of  $P(m|d)$ , the values of the likelihood are no longer needed to compute parameter constraints (they were only needed to build the chain).

We can compute **averages**, **dispersions** and **correlations** for all parameters directly from the chain  $\rightarrow$  they are moments of the  $P(m|d)$  distribution.

Notice that the results for each parameter are already marginalized over all the other parameters.

We can also draw **contour plots**: iso-probability contours, enclosing a given fraction of the total probability



The resulting  $P(m|d)$  is not necessarily Gaussian  $\rightarrow$  contours are not necessarily ellipses.



## Fisher Information Matrix

Assume the posterior is a Gaussian distribution in the parameters' space, centred in the best-fit  $m_0$

This method allows us to compute the covariance of the posterior distribution from the curvature matrix of the likelihood.

$$P(m|d) \propto \exp \left[ -\frac{1}{2} (m_0 - m) C^{-1} (m_0 - m)^T \right]$$

Since it is a Gaussian distribution, the covariance matrix in the parameters space (that represents the uncertainty of the parameters estimation) is computed from

$$C_{ij}^{-1}(m_0) = \left( \frac{\partial^2 (-\ln P)}{\partial p_i \partial p_j} \right) \Big|_{m_0}$$

where  $m = m_0$  is the parameter value with maximum probability, i.e., the peak of the posterior.

This is also the mean value, since it is a Gaussian (symmetric) distribution.

Inserting this in Bayes theorem we can write

$$C_{ij}^{-1}(m_0) = F_{ij} + \left( \frac{\partial^2(-\ln \text{Prior})}{\partial p_i \partial p_j} \right)_{|m_0}$$

i.e., the inverse covariance matrix of the posterior is the second-order derivative of the  $\ln(\text{prior})$  (which is zero for flat priors) + the Fisher matrix,

The **Fisher matrix** is defined as the second-order derivative of the chi-square function (the  $\ln(\text{likelihood})$ ) with respect to the parameters, averaged over all parameter values.

For a Gaussian or any symmetric distribution this is just the second-order derivative taken at the peak, known as the **curvature matrix** (also known as the Hessian), i.e.,

$$F_{ij} = \left\langle \frac{\partial^2 \mathcal{L}}{\partial p_i \partial p_j} \right\rangle \approx \left( \frac{\partial^2 \mathcal{L}}{\partial p_i \partial p_j} \right)_{|m_0} \quad \mathcal{L} = -\ln L$$

$F_{ij}$  is an  $N_p \times N_p$  matrix, where  $N_p$  is the number of parameters.

Naturally, the parameter values at the peak (the vector  $m_0$ ) are not known.

But the goal of this method is to find the uncertainty on the parameters  
(called the **credible intervals** in the Bayesian approach - also called the **confidence intervals** a name from the hypothesis testing in the frequentist approach)

and not the actual parameter values (the peak of the distribution - called the **best-fit** values).

So, any reasonable value may be chosen for  $m_0$  - this is called the **fiducial value**.  
The result we are looking for is  $F_{ij}$ , the inverse covariance matrix in the parameters' space → it will give us **the uncertainty of the parameters' values around the fiducial value**.

For this reason, this method is only used to make **forecasts** of the precision of achievable with future data, while with real data we want to find out not only the uncertainty of the estimates but the actual predictions for the parameters → likelihood sampling methods are used with real data (e.g. MCMC).

Now, even if the posterior distribution is not Gaussian in reality, it is still useful to consider the Fisher matrix method, because the **Rao-Cramer inequality** states that:

**the parameters confidence intervals obtained from a Fisher matrix analysis are a lower limit of the true ones.**

This result may be derived from the **Cauchy-Schwarz inequality**,

$$O_L(f^2) O_L(g^2) \geq [O_L(fg)]^2$$

( $O_L$  stands for linear operation and  $f$  and  $g$  are general functions).

If we choose the linear operator to be the **expectation value**, i.e. the ensemble average  $\langle \rangle$  :

$O_L(g^2) = E [(m-m_0)^2]$  , i.e., the variance (E denotes the **expectation value**  $\langle \rangle$ )

$O_L(f^2) = E [(d \ln L / dm)^2]$

the **Cauchy-Schwartz inequality** becomes: 
$$\text{Var}(m) \geq \frac{(E[(m - m_0) \partial \ln L / \partial m])^2}{E[(\partial \ln L / \partial m)^2]}$$

The numerator is

$$(E[(m)\partial \ln L/\partial m] - m_0 E[\partial \ln L/\partial m])^2 = \left( \int m \frac{\partial \ln L}{\partial m} L \right)^2 = \left( \frac{\partial E(m)}{\partial m} \right)^2 \geq 0.$$

while the denominator is:

$$\int \left( \frac{\partial \ln L}{\partial m} \right)^2 L = \int \frac{\partial \ln L}{\partial m} \frac{\partial L}{\partial m} = E \left[ \frac{\partial^2 \mathcal{L}}{\partial m^2} \right] = F.$$

which proves the Rao-Cramer inequality:  $\text{Var}(\mathbf{m}) \geq \mathbf{F}^{-1}$

In other words, the elements of the covariance matrix in the parameter space are larger than the elements of  $\mathbf{F}^{-1}$

i.e., ***the inverse Fisher matrix is a lower limit of the covariance matrix.***

Notice that **large values of Fisher matrix mean small uncertainties on the parameters' values**

(this is consistent with the Fisher matrix being the curvature matrix  $\rightarrow$  large curvature in the likelihood means a peaked distribution  $\rightarrow$  small sigma).

## Computing the Fisher matrix in practice

The Fisher matrix approximation is very useful because it is very fast to compute, being basically the derivative of the cosmological function with respect to the cosmological parameters.

Let us start from the log-likelihood:

$$-\ln L \propto \frac{1}{2} (\mu_{\text{obs}}(z) - \mu_{\text{th}}(z, \Omega_i))^T C_{zz}^{-1} (\mu_{\text{obs}}(z) - \mu_{\text{th}}(z, \Omega_i))$$

here written considering SN data  $\mu_{\text{obs}}(z)$  and the corresponding true value  $\mu_{\text{th}}$  that is function of the cosmological parameters:  $\mu_{\text{th}}(z; \Omega)$ ;  $C_{zz}$  is the covariance matrix of the data.

Now, remember that the Fisher matrix is computed from the second-order derivatives of the log-likelihood computed at the fiducial value, i.e., at the peak of the distribution, i.e., at  $\mu_{\text{obs}} = \mu_{\text{th}}$  (fiducial parameters).

This allows us to simplify the computation. Consider one element of the  $X^2$  sum and a diagonal covariance with elements  $\sigma^2$ . The derivation is:

$$\begin{aligned}
 \frac{\partial^2}{\partial \Omega_a \partial \Omega_b} \left( \frac{\mu_{\text{obs}} - \mu_{\text{th}}}{\sigma^2} \right)^2 &= \frac{\partial}{\partial \Omega_b} \left[ \frac{2(\mu_{\text{obs}} - \mu_{\text{th}})}{\sigma^2} \frac{\partial \mu_{\text{th}}}{\partial \Omega_a} \right] \\
 &= \frac{2}{\sigma^2} \frac{\partial}{\partial \Omega_b} \left[ (\mu_{\text{obs}} - \mu_{\text{th}}) \right] \frac{\partial \mu_{\text{th}}}{\partial \Omega_a} + \frac{2(\mu_{\text{obs}} - \mu_{\text{th}})}{\sigma^2} \frac{\partial^2 \mu_{\text{th}}}{\partial \Omega_a \partial \Omega_b} \\
 &= \frac{2}{\sigma^2} \frac{\partial \mu_{\text{th}}}{\partial \Omega_a} \frac{\partial \mu_{\text{th}}}{\partial \Omega_b} + 0
 \end{aligned}$$

**So, in practice we just need to compute the first derivative of the cosmological function  $\mu_{\text{th}}$  with respect to the cosmological parameters, due to the condition  $\mu_{\text{obs}} = \mu_{\text{th}}$ .**

This result is valid if the covariance matrix does not depend on the cosmological parameters. If this is not the case the covariance (i.e.  $1/\sigma^2$  in this example) also needs to be differentiated). Usually that dependence is weaker, and the above formula is a good approximation.

So, the general practical formula of the Fisher matrix is:

$$F_{ab} = \left( \frac{\partial \mu_{th}(z; \Omega)}{\partial \Omega_a} \right)^T C_{zz'}^{-1} \frac{\partial \mu_{th}(z'; \Omega)}{\partial \Omega_b}$$

(notice that the factor 2 cancels out with the  $\frac{1}{2}$  of the  $X^2$ )

We just need to **compute the derivatives of  $\mu$  with respect to all cosmological parameters**. For each parameter we will have a vector (a discretized function of  $z$ ), that contracted with the covariance matrix will produce a number for each pair of cosmological parameters.

In other words, the Fisher matrix has dimension of  $N_p \times N_p$  and is the sum of the products of two derivatives over the redshift range, normalized by the variances:

$$F_{ab} = \sum_z \left( \frac{\partial \mu_{th}(z)}{\partial \Omega_a} \frac{\partial \mu_{th}(z)}{\partial \Omega_b} \right) \frac{1}{\sigma_z^2}$$

(written here for the case of a diagonal covariance).



The diagonal terms of the Fisher matrix (for a diagonal covariance matrix) are just:

$$F_{\Omega_m \Omega_m} = \sum_z \left( \frac{\partial M_{\ell_h}(z)}{\partial \Omega_m} \right)^2 \frac{1}{\sigma_z^2}$$

If the **derivative of the cosmological function** with respect to a certain parameter is larger than with respect to another one, it means that the cosmological function is more sensitive to the first one  $\rightarrow$  the corresponding component of the Fisher matrix is larger  $\rightarrow$  the corresponding  $F^{-1}$  value is smaller  $\rightarrow$  the uncertainty on the first parameter is smaller than on the second one.

But what about the absolute value of the uncertainty? We saw it is smaller, but is it small? That depends on the data covariance matrix that equally affects the derivatives with respect to all parameters:

If **the data errors** are small  $\rightarrow$  the derivatives are divided by a small number  $\rightarrow$  the components of the Fisher matrix are larger  $\rightarrow$  the corresponding  $F^{-1}$  value are smaller  $\rightarrow$  the parameters are estimated with smaller uncertainties.

**In summary:**

**The inverse of the Fisher matrix is a covariance matrix in the parameters space.**

The square root of its diagonal gives the error bars on the estimated parameters. As in the data space, if the parameters are correlated, the full Fisher matrix is needed to quantify the errors.

**The error associated with the estimate of a cosmological parameter depends on two factors:**

- **the sensitivity of the cosmological function to the parameter** (the derivatives)
- **the precision and accuracy of the data** (the data covariance matrix)

## Finding the credible intervals in the parameters space (the contours)

The computation of the Fisher matrix we just did is exact, regardless of the posterior being a Gaussian or not.

Now, to plot the contours in the parameters' space we will consider the approximation that the posterior is a Gaussian and consider a Taylor expansion of the log-posterior in the parameters space:

$$\mathcal{L}(m) - \mathcal{L}(m_0) = 0 + (p_i - p_{i0})F_{ij}(p_j - p_{j0}) + \mathcal{O}(\Delta_p^3)$$

### Accuracy of the method:

- The expansion shows explicitly that the Fisher matrix method gives a lower limit for the parameters' variance. The result is only exact if higher-order derivatives are zero (which happens for a Gaussian, which is fully described by only two moments).
- Moreover, the result of this method is not accurate (even in the case of a Gaussian posterior) if the fiducial value chosen is not at the peak of the distribution.

The first term of the Taylor expansion,  $(\partial L/\partial m)|_{m_0}$  is zero, since the derivative is taken at the maximum of the likelihood (the peak).

This equation, to second order, is a **quadratic equation** in the variables  $\Delta p_i$ , with the center of the coordinates in  $m_0$ .

Note the Fisher matrix is **semi-definite positive** by construction from the derivatives of the likelihood (also from Cauchy-Schwartz).

$$\frac{\partial^2 \mathcal{L}}{\partial p_1^2} \frac{\partial^2 \mathcal{L}}{\partial p_2^2} \geq 2 \frac{\partial^2 \mathcal{L}}{\partial p_1 \partial p_2}$$

(i.e., the correlation coefficients are smaller than 1)

→ **the points of constant  $\Delta L$  define a (hyper)ellipse.**

A value of  $\Delta L = L_0 - L(p_i)$  gives a contour level, or (n-sigma) confidence interval, that connects all points  $p_i$  in the parameter space that have the same likelihood.

## 1D “contours”: (1 parameter)

In a 1D normalized Gaussian posterior distribution, consider the parameter values  $p_{\min}$  and  $p_{\max}$  (respectively to the left and the right of the peak at  $p_0$ ) such that

$$\ln L(p_0) - \ln L(p_{\min}) = \ln L(p_0) - \ln L(p_{\max}) = 1$$

i.e., the log-likelihood of those two points differs  $\Delta L=1$  from the log-likelihood of the peak. This is called the **1-sigma** level

For 1-sigma the quadratic equation is simply:

$$\Delta L = L_0 - L(p_i) = 1 \rightarrow F (p-p_0)^2 = 1 \rightarrow (p-p_0) = \text{sqrt}(1/F)$$

→ **The 1-sigma error is  $\text{sqrt}(1/F)$**

Incidentally, if we compute the integral of the normalized Gaussian from  $p_{\min}$  to  $p_{\max}$  the result is 0.683, meaning that the volume enclosed by the contour  $\Delta L=1$  contains 68.3% of the total probability.

Other probability levels are also usually defined:

$\Delta L=4$  is **2-sigma** → contains 95.4% of the total probability

$\Delta L=9$  is **3-sigma** → contains 99.7% of the total probability

## 2D contours: (2 parameters)

Integrating a 2D normalized Gaussian, we find that the 68.3%, 95.4% and 99.7% values correspond to different likelihood levels than in a 1D Gaussian.

The levels are now  $\Delta L = 2.3, 6.2, 11.8$ , respectively.  
Nevertheless, they are still called 1, 2 and 3-sigma levels.  
(For example, in 2D,  $\Delta L=1$  only encloses 40% of the probability).

The quadratic equation for a fixed  $\Delta L$  is

$$\Delta L = F_{xx} (x-x_0)^2 + 2F_{xy} (x-x_0)(y-y_0) + F_{yy} (y-y_0)^2 \rightarrow \text{this defines an ellipse.}$$

So iso-probability contours in the Fisher matrix method are ellipses. Larger ellipses correspond to larger probability volumes

In this way, the components of the inverse Fisher matrix give us directly:

$\text{sig}_{xx}^2$  - variance of parameter x  
 $\text{sig}_{xy}^2$  - covariance (correlation of x and y)  
 $\text{sig}_{yy}^2$  - variance of parameter y

If F is diagonal there is no correlation and the matrix axes are along the parameter axes x and y.

### 3D contours: (3 parameters)

The (hyper-)ellipse equation for a fixed  $\Delta L$  is

$$\Delta L = F_{xx} (x-x_0)^2 + 2F_{xy} (x-x_0)(y-y_0) + F_{yy} (y-y_0)^2 + 2F_{xz} (x-x_0)(z-z_0) + 2F_{yz} (y-y_0)(z-z_0) + F_{zz} (z-z_0)^2$$

The components of the inverse Fisher matrix give us directly:

$\text{sig}_{xx}^2$  - variance of parameter x

$\text{sig}_{xy}^2$  - covariance (correlation of x and y)

$\text{sig}_{yy}^2$  - variance of parameter y

$\text{sig}_{xz}^2$  - covariance between x and z

$\text{sig}_{yz}^2$  - covariance between y and z

$\text{sig}_{zz}^2$  - variance of parameter z

**But how can we plot ellipses (2D contours) in this case?**

There are two ways to plot the ellipses on each of the 3 planes (x,y), (x,z), (y,z).

Consider the contour in the (x,y) plane. The two ways are:

**Maximizing**: one of the parameters is kept fixed at the maximum (in the x,y case, we fix  $z = z_0$ ).

This corresponds to a 2D slice through the 3D hyper-ellipse.

**In practice → remove z line and column from F → use this reduced F to plot the contour (x,y) or invert it to read the uncertainties directly on the new  $F^{-1}$**



**Marginalizing**: integrating over the full range of the 3<sup>rd</sup> parameter.

This corresponds to projecting the hyper-ellipse on a 2D plane.

Integrating the likelihood will remove the dependence on the third parameter from the multivariate Gaussian, obtaining a Gaussian without that parameter that can be differentiated to get a (reduced) Fisher matrix.

**In practice → remove z-axis line and column from the covariance  $F^{-1}$ , obtaining a reduced  $F^{-1}$  → use it to read directly the uncertainties or invert it to insert in the ellipse equation and plot the contour.**

Notice that marginalizing results in a larger ellipse than maximizing.

The uncertainty volume can also be reduced by neglecting the axis that have small variance → **Principal Components Analysis**

Notice that it is also possible to **marginalize on a reduced interval** instead of integrating to infinity.

This is equivalent to introducing a **prior**, restricting the interval of a given parameter. In this case the prior contribution needs to be added to the Fisher matrix → this should result in larger error bars than the maximization but smaller than the full marginalization.

## Figure-of-Merit

The area of a 1-sigma ellipse is:

$$\pi a b = 2.3 \pi / \text{sqrt}(\det F)$$

**The square-root of the determinant of the 2D Fisher matrix is proportional to the inverse of the area of the ellipse.**

The **Figure-of-Merit (FoM)** is defined as

$$\text{FoM} = \text{sqrt}(\det F)$$

The FoM of the ellipse in the  $w_0, w_a$  plane is used to quantify the constraining power of cosmological surveys: it is called the **dark energy FoM**

## Cosmological parameters estimation: observations of SNe

The cosmological information of the measured distance modulus is contained in the luminosity distance:

$$\mu(z) = 5 \log_{10} (D_L(z; H_0, \Omega, w)) + 25$$

with  $D_L = (1+z) D_C$  (for a flat Universe)

Let us investigate what is the **cosmological information** that the distance-modulus contains:

First, the comoving distance from the observer at  $t_0$  to the source at  $t$  is computed from the metric as

$$D_C = \int_t^{t_0} \frac{1}{a} dt = \int_t^{t_0} (1+z) dt$$

$a(t)$  (or  $H(t)$ ) are usually computed using Einstein equations, and this expression will involve the density parameters (that define the cosmological model).

Alternatively, in order to access the cosmological information in a more fundamental **model-independent** way, let us consider first the **cosmographic approach**, where  $a(t)$  is expanded as:

$$a(t) = a_0 \left[ 1 + H_0 (t-t_0) - \frac{1}{2} q_0 H_0^2 (t-t_0)^2 + \frac{1}{3!} \dot{q}_0 H_0^3 (t-t_0)^3 + \frac{1}{4!} r_0 H_0^4 (t-t_0)^4 + \dots \right]$$

From here, we can also write an expansion for  $z(t)$ , since  $a^{-1} = -z/(z+1)$ :

$$z^{-1} = \frac{-z}{z+1} \Rightarrow z = \frac{1}{z^{-1} - 1} \Rightarrow z = (z+1) \left[ H_0 (t_0-t) + \frac{1}{2} q_0 H_0^2 (t_0-t)^2 + \dots \right]$$

$$\Rightarrow z = \underbrace{H_0 (t_0-t)}_{\substack{\text{the lowest-order term} \\ \text{is order } \mathcal{O}(t)}} + \frac{1}{2} q_0 H_0^2 (t_0-t)^2 + \underbrace{z H_0 (t_0-t)}_{\text{order } \mathcal{O}(t^2) + \dots} + \frac{1}{2} \underbrace{z q_0 H_0^2 (t_0-t)^2}_{\text{order } \mathcal{O}(t^3) + \dots} + \dots$$

Keeping only order  $\mathcal{O}(t^2)$  we may use  $z \sim H_0 (t_0-t)$  to insert here

$$\Rightarrow z = H_0 (t_0-t) + \left( \frac{1}{2} q_0 H_0^2 + H_0^2 \right) (t_0-t)^2 + \mathcal{O}(t^3)$$

Inserting the  $z(t)$  expansion in the comoving distance, we find, to second order:

$$D_c = \int_t^{t_0} [1 + H_0(t_0 - t)] dt$$

(to second order we just need to consider order  $t$  in the integrand, because the integral will be order  $t^2$ )

$$D_c = (t_0 - t) + \frac{H_0}{2} (t_0 - t)^2$$

At this point it would be useful to invert the expansion  $z(t)$ , to be able to find an expression for  $D_c(z)$  instead of  $D_c(t)$

Writing now for  $t$ :

$$\Rightarrow t_0 - t = \frac{z}{H_0} - \left( \frac{q_0 H_0^2 + H_0^2}{2} \right) \frac{1}{H_0} (t_0 - t)^2 + o(t^3)$$

again the same trick, using  $(t_0 - t) \sim \frac{z}{H_0}$

$$\Rightarrow t_0 - t = \frac{z}{H_0} - \left( q_0 \frac{H_0^2}{2} + H_0^2 \right) \frac{z^2}{H_0^3} + \dots$$

$$\text{or } t_0 - t = \frac{1}{H_0} z - \frac{1}{H_0} \left( 1 + \frac{q_0}{2} \right) z^2 + \sigma(z^3)$$

Now, Since  $D_c = (t_0 - t) + \frac{H_0}{2} (t_0 - t)^2 + \sigma(t^3),$

we get  $D_c^{(z)} = \frac{1}{H_0} z - \left( 1 + \frac{q_0}{2} \right) \frac{z^2}{H_0} + \frac{H_0}{2} \frac{z^2}{H_0^2} + \sigma(z^3)$

$$\text{or } D_c = \frac{1}{H_0} z - \left( 1 + \frac{q_0}{2} - \frac{1}{2} \right) \frac{z^2}{H_0} + \sigma(z^3)$$

$$\text{or } D_c^{(z)} = \frac{1}{H_0} \left[ z - \frac{1}{2} (1 + q_0) z^2 \right]$$

We arrive then at the expression for the luminosity distance that we were looking for:

$$D_L = \left[ \frac{z}{H_0} - \frac{z^2}{H_0} \left( \frac{1}{2} + \frac{1}{2} q_0 \right) \right] (1+z) + \mathcal{O}(z^3)$$

$$= \frac{z}{H_0} - \frac{z^2}{H_0} \left( \frac{1}{2} + \frac{1}{2} q_0 \right) + \frac{z^2}{H_0} + \mathcal{O}(z^3)$$

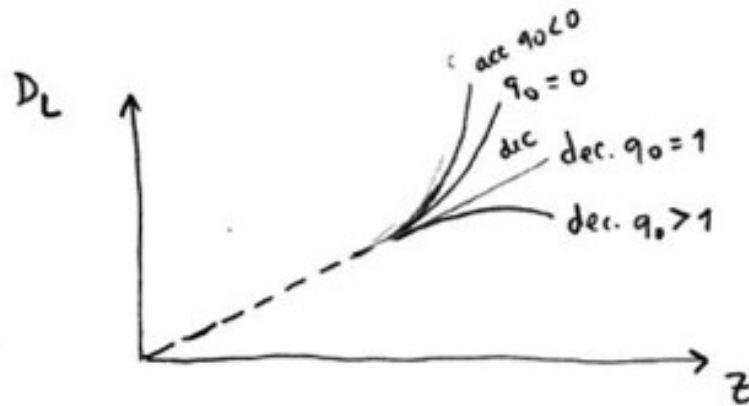
$$\Rightarrow \left[ D_L = \frac{1}{H_0} \left( z + \frac{1}{2} (1 - q_0) z^2 \right) \right] \quad (1+z)$$

**We see that, up to second-order, the luminosity distance:**

- **at low-z measures the (constant) velocity of the Universe** (with  $(c)z/H_0$  being the Doppler velocity)
- **at high-z measures the (constant) acceleration of the Universe** ( $q_0$ )

So, in order to detect an **acceleration** of the Universe we need to observe SNe at **high redshifts**.

SNe at **low redshifts** measure  $H_0$ , i.e., the slope of the  $D_L(z)$  straight line



To determine **the value of the acceleration**, both high and low redshift measurements are needed, to break the **degeneracy** between  $H_0$  (the Hubble law slope, important at lower  $z$ ) and the acceleration (important at higher  $z$ )

**Note that the evidence for acceleration is based only on the shape of the function.**

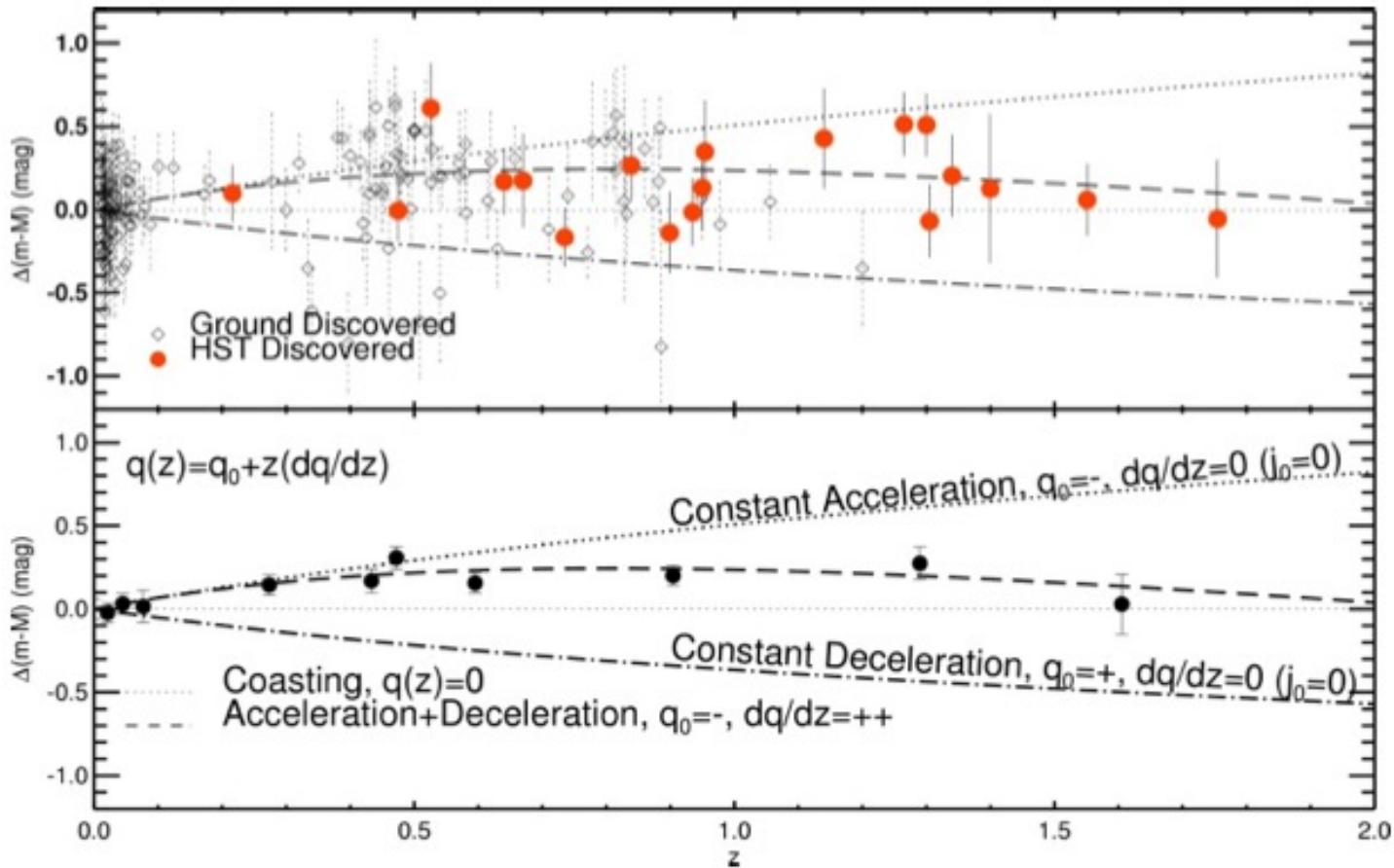
Even if the absolute values determined for the distances were not precise (i.e., if  $H_0$  was estimated with a large uncertainty) we could still find evidence for acceleration using only high- $z$  SNe.



This result is only an approximation. **If we consider higher-orders**, we find

$$D_c(z) = \frac{z}{H_0} \left[ 1 - \left(1 + \frac{q_0}{2}\right)z + \left(1 + q_0 + \frac{q_0^2}{2} - \frac{j_0}{6}\right)z^2 - \left(1 + \frac{3}{2}q_0(1+q_0) + \frac{5}{8}q_0^3 - \frac{1}{2}j_0 - \frac{5}{12}q_0j_0 - \frac{\Lambda_0}{24}\right)z^3 + O(z^4) \right]$$

The introduction of higher orders shows that the acceleration is not necessarily constant, i.e., the quadratic term depends also on  $j_0$ , i.e., there is a non-zero  $dq/dz$ .



The data seem to prefer a model with varying  $q_0$ , such that

$q_0 > 0$  at high- $z$

and

$q_0 < 0$  at low- $z$ .

This plot shows the dynamic behaviour of the Universe (independently of the values of the density parameters) → clear **evidence for a model with acceleration for  $z < 1$  and deceleration for  $z > 1$**  → proof of **late-time acceleration of the universe**

Now, the cosmographic analysis was very useful to get an insight of the dynamic behaviour of the Universe, but in order to estimate cosmological parameters this analysis is not needed.

**What we need is just to compute the observable cosmological function (i.e. the luminosity distance) from vectors of cosmological parameter values and compare the various theoretical  $D_L$  obtained with the observed one through the computation of likelihoods in the parameter space.**

The dependence of the luminosity distance on the cosmological parameters can be most easily seen by writing the distance as an integral over redshift:

In general,  $D_L(z) = (1+z) D_M$ ,

considering the case of flat Universe, we have

$$D_L(z) = (1+z) D_C(z) = (1+z) \int_t^{t_0} \frac{1}{a} dt = (1+z) \int_0^z \frac{cdz}{H(z)}$$

and the Hubble function is found in terms of the parameters of the cosmological fluid (densities of the various sources) through Friedmann's equation:

$$H^2(a) = H_0^2 \left( \Omega_r(1+z)^4 + \Omega_m(1+z)^3 + \Omega_K(1+z)^2 + \Omega_\Lambda \right) \quad (\text{here including a cosmological constant})$$

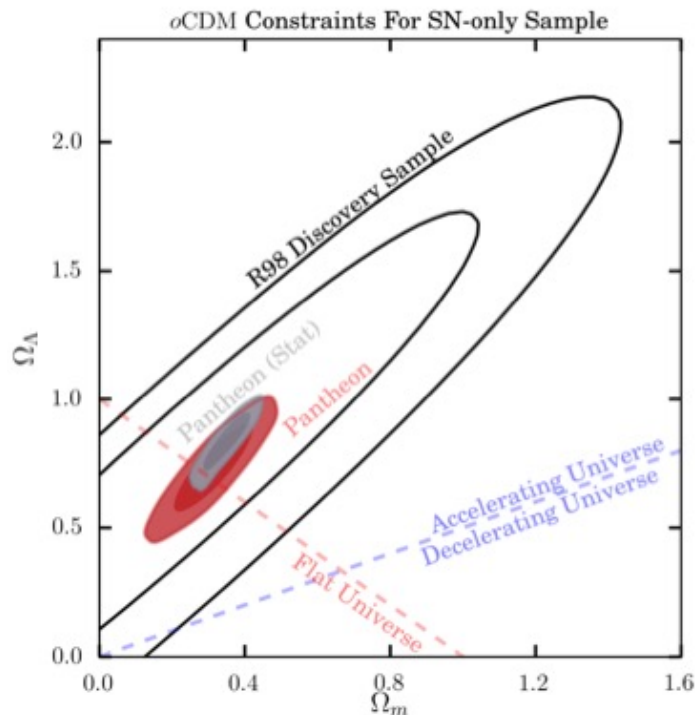
From this, the luminosity distance can be computed for any values of the vector of cosmological parameter, and for the redshifts of the various SN. Then all terms of the distance-modulus estimator are added (introducing a large number of nuisance parameters) and finally the theoretical distance-modulus  $\mu(z)$  is found. Its likelihood can then be computed by comparing with the distance-modulus data. In a sampling method, the procedure is then repeated for millions of points in the parameter space. The cosmological parameters estimates are finally found by marginalizing over the nuisance parameters.

It is important to realize that the results depend on the cosmology assumed (it is a working hypothesis).

Let us consider the  $\Lambda$ CDM scenario.

**flat  $\Lambda$ CDM**  $\rightarrow$  2 independent background cosmological parameters:  $H_0$ ,  $\Omega_m$ , since  $\Omega_r$  and  $\Omega_K$  are fixed and  $\Omega_\Lambda = 1 - \Omega_m$ . If curvature is not fixed *a priori* then there are 3 free parameters:  $H_0$ ,  $\Omega_m$ ,  $\Omega_\Lambda$ , and this is historically called **oCDM** (open CDM, even though the fit is free to have any curvature - flat, open or closed )

**Constraints in the  $(\Omega_m, \Omega_\Lambda)$  plane:** after marginalizing over the nuisance parameters and  $H_0$ , the constraints on the 2 density parameters are given as confidence contours in the 2D parameter space.



Some notes:

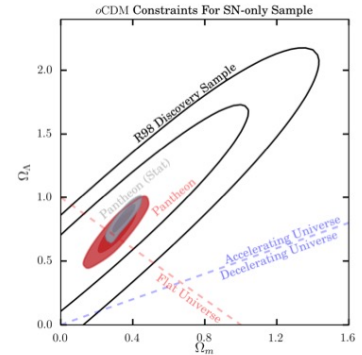
- The large contours are from the first SN results of 1998. They show quite large  $1\sigma$  and  $2\sigma$  probability contours, since the data has large error bars.

- The smaller contours (also showing  $1\sigma$  and  $2\sigma$  contours) are for the recent Pantheon results.

- Notice also the impact of considering or not the contribution of the systematic effects for the data error bars:

Grey contours - analysis done using Pantheon data with error bars including only the errors (statistical uncertainties)

Red contours - analysis done using Pantheon data with error bars including statistical + systematic uncertainties)



**So, with larger error bars, a larger region of the parameter space has a “good likelihood” and is included inside the confidence contours (red larger than grey).**

- If  $H_0$  was known (fixed in the analysis instead of marginalized), the  $(\Omega_m, \Omega_\Lambda)$  contours would be smaller (tighter constraints)

- Models such that  $\Omega_\Lambda = 1 - \Omega_m$  (i.e.,  $\Omega_K = 0$ ), lie on the straight line marked “flat”

- There is also a line dividing accelerating and decelerating models.

## The contours are all aligned on a preferred direction. Why is this?

To answer this question, let us remember that to first approximation, we are measuring the acceleration of the Universe, i.e., as we saw,  $D_L$  depends directly on  $q_0$

$$q_0 = - \frac{\ddot{a}}{a} \Big|_{t_0} \frac{1}{H_0^2}$$

From Raychadhuri's equation, we can write the acceleration parameter in terms of the source parameters:

$$\frac{\ddot{a}}{a} = - \frac{4\pi G}{3} (\rho + 3p)$$

$$= - \frac{8\pi G}{3H_0^2} \frac{H_0^2}{2} \rho (1+3W(a)) \Leftrightarrow \frac{\ddot{a}}{a} = - \frac{H_0^2}{2} \left[ \Omega_m a^{-3} + \Omega_{DE} (1+3W(a)) \right]$$

(for a general dark energy fluid)

For the case of a cosmological constant:

$$\Rightarrow \frac{\ddot{a}}{a} \Big|_{t_0} = - \frac{H_0^2}{2} \left[ \Omega_m + \Omega_{DE} (1+3W) \right]$$

$$\text{or, for } \Omega_{DE} = \Omega_\Lambda \Rightarrow \frac{\ddot{a}}{a} \Big|_{t_0} = - \frac{H_0^2}{2} (\Omega_m - 2\Omega_\Lambda)$$

and so

$$q_0 = \frac{1}{2}\Omega_m - \Omega_\Lambda$$

(note that  $q_0$  is independent of  $H_0$ , which is consistent, since they are directly different orders of the Taylor expansion  $\rightarrow$  acceleration is independent of velocity)

So, models with the same acceleration (same  $q_0$ ) all lie in a line

$$y = ax + b$$

where  $y = \Omega_\Lambda$ ,  $x = \Omega_m$ ,  $a = 1/2$ ,  $b = -q_0$  (which is  $>0$  for an accelerated model)

**This line defines the direction of the contour** (with some width due to the uncertainty on the measured acceleration)

This shows that  $\Omega_m$  and  $\Omega_\Lambda$  are **correlated** in the acceleration they produce. The two parameters define a straight line along which all models have exactly the same acceleration and will have exactly the same likelihood values  $\rightarrow$  a **degeneracy direction**

Consider a model 1. Now, if a model 2 has a higher  $\Omega_\Lambda$  with respect to model 1, then by also increasing its  $\Omega_m$  value the acceleration produced by model 2 will be the same as for model 1  $\rightarrow$  they are **correlated (positively correlated)**  $\rightarrow$  contours from bottom-left to top-right.

If to keep the acceleration constant when one of the parameters increase, the other would need to decrease, then they would be  $\rightarrow$  **anti-correlated (negatively correlated)**  $\rightarrow$  contours from top-left to bottom-right

Two different ways of increasing luminosity distance:

- 1) Increase  $\Omega_\Lambda$
- 2) Decrease  $\Omega_m$

This causes the degeneracy between  $\Omega_\Lambda$  and  $\Omega_m$

**It is then impossible to distinguish those 2 models (or any model along the degeneracy direction) with SN measurements (or any other DL based method).**



In general, **cosmological probes are very good in constraining degeneracy directions** (i.e. combinations of cosmological parameters) but not so good in constraining individual parameters.

In our case, SN measurements are good in constraining the orthogonal direction to the degeneracy direction i.e., the deviation from the acceleration line (or the width of the contours).

Note that a parameter defined along the width of the contours would be highly constrained - this parameter corresponds to the last principal components in a PCA analysis of the parameter space covariance matrix.

Notice that  $(0.5 \Omega_m - \Omega_\Lambda = \text{constant})$  is a perfect degeneracy. Why then do the contours close and do not extend infinitely along the degeneracy direction?

This is because the linear dependence of  $D_L$  on  $q_0$  is only a good approximation at second-order of the  $a(t)$  expansion. **In reality, there are other terms and degeneracy is not perfect → the contours close and show a preference for  $\Omega_m < 1$  (and  $\Omega_\Lambda > 0$ )**

## Could we get better estimates for individual parameters?

The Pantheon results are:

Analysis	Model	$\Omega_m$	$\Omega_\Lambda$
SN-stat	$\Lambda$ CDM	$0.284 \pm 0.012$	$0.716 \pm 0.01$
SN-stat	$o$ CDM	$0.348 \pm 0.040$	$0.827 \pm 0.06$
SN	$\Lambda$ CDM	$0.298 \pm 0.022$	$0.702 \pm 0.02$
SN	$o$ CDM	$0.319 \pm 0.070$	$0.733 \pm 0.11$

Notice the constraints are looser (worse) if:

- systematics are included in the data error budget
- curvature is left free (one more free parameter to add to the general degeneracies)

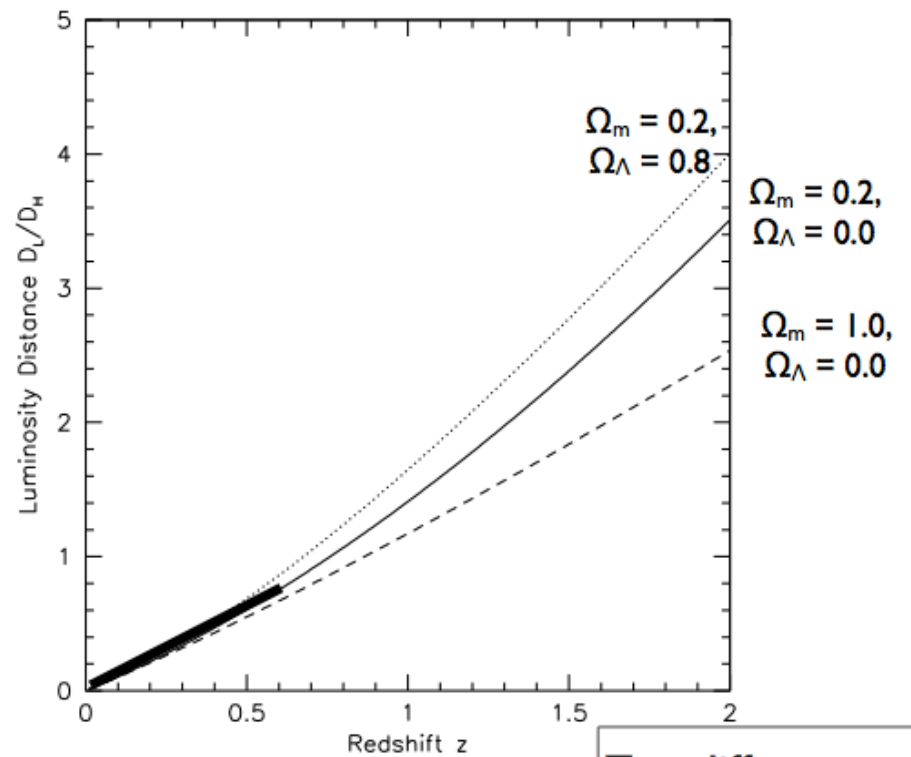
Notice that in the (flat)  $\Lambda$ CDM case, the result for  $\Omega_\Lambda$  is just  $1 - \Omega_m$

We can improve the constraints by **combining various cosmological probes** such as to break the degeneracy.

For example, consider an observable that would depend directly on the curvature of the Universe. In the  $(\Omega_m, \Omega_\Lambda)$  plane we see that lines of constant curvature are more or less orthogonal (i.e. **complementary**) to lines of constant acceleration.

The joint likelihood analysis of those two datasets would produce contours in the intersection of the two directions  $\rightarrow$  i.e. potentially small round contours  $\rightarrow$  constraining simultaneously the two parameters  $\Omega_m$  and  $\Omega_\Lambda$ .

## Do the cosmological observations prove the existence of dark energy?



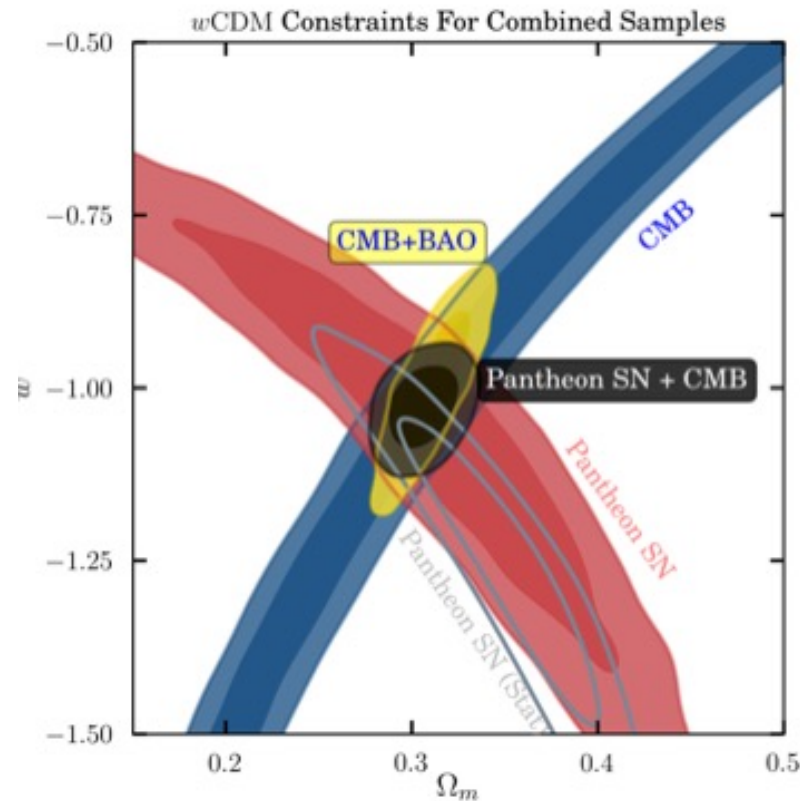
The evidence from the data is for acceleration (based on the shape of the  $D_L(z)$  function).

The “**evidence for dark energy**” is a **model-dependent** conclusion (i.e. based on the assumption of an underlying cosmology) and therefore less robust than the evidence for acceleration.

**Let us consider now the  $w$ CDM scenario**, where dark energy has a constant equation of state, but not necessarily equal to -1 (which would be  $\Lambda$ CDM).

**$w$ CDM**  $\rightarrow$  there are 4 independent background cosmological parameters:  $H_0$ ,  $\Omega_m$ ,  $\Omega_{DE}$ ,  $w$  (or alternatively  $H_0$ ,  $\Omega_m$ ,  $\Omega_K$ ,  $w$ ), or only 3:  $H_0$ ,  $\Omega_m$ ,  $w$ , if flatness is also assumed ( $\Omega_K = 0$  and  $\Omega_{DE} = 1 - \Omega_m$ )

***Constraints on the  $(\Omega_m, w)$  plane***  
(after marginalizing over the other parameters)



The SN-Pantheon contours (red) are in a very **different direction** than the contours in the  $(\Omega_m, \Omega_\Lambda)$  plane that we saw previously.

This is because that (as before) they are determined by the acceleration parameter  $q_0$ , which now (from Raychadhuri's eq.) is,

$$q_0 = \frac{1}{2} \Omega_m + \Omega_{DE} \frac{(1+3w)}{2}$$

i.e.,  $\Omega_m$  and  $w$  add, instead of subtracting (contrary to the relation between  $\Omega_m$  and  $\Omega_\Lambda$ ), and so they are anti-correlated. (Note that this is just an effect of  $w$  being negative)

Moreover, **the contour is no longer an ellipse** (it is curved). This is because the line of constant luminosity distance (which in our  $O(z^2)$  approximation is the line of constant acceleration) is no longer a straight line in the  $(\Omega_m, w)$  plane.

When we move along a straight line in this plane, a change on  $\Omega_m$  induces a change on  $\Omega_{DE} \rightarrow$  the dependence of  $q_0$  on the parameters is no longer linear.

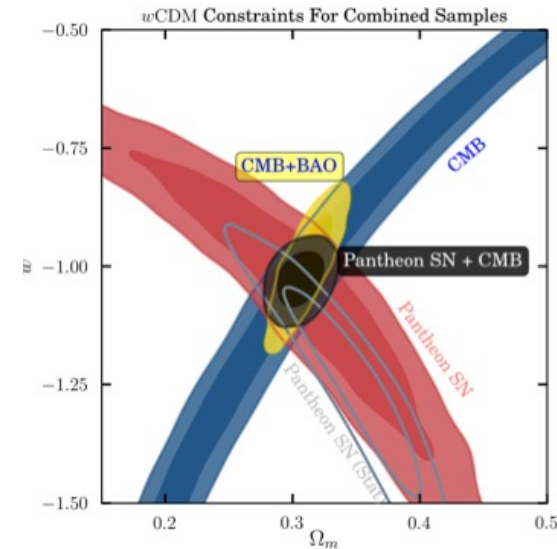
Indeed, if we replace  $\Omega_\Lambda = 1 - \Omega_m$  in the expression for  $q_0$ , we get

$$q_0 = \frac{1}{2} + \frac{3}{2}w(1 + \Omega_m)$$

i.e.,  $y = a(1-x)^{-1}$

where  $y = w$ ,  $x = \Omega_m$ ,  $a = 2/3 (q_0 - 1/2)$

corresponding to the red curved contour.



**We recover the result that only in the case the cosmological function (in this case, distance modulus, luminosity distance, acceleration) depends linearly on the parameters, is the posterior distribution in the parameters space a Gaussian (leading to elliptical contours).**

The figure also shows:

- Contours from CMB-Planck measurements

orthogonal to the SN ones (they do not measure the luminosity distance or acceleration but different observables, like the angular size of the sound horizon at recombination) → they are complementary probes, and the joint contours are much reduced.

- Baryonic Acoustic Oscillations (BAO) measurements

similar to the SN ones (they measure the angular diameter distance) and also complementary to CMB.

Sample	$w$
CMB+BAO	$-0.991 \pm 0.074$
CMB+H0	$-1.188 \pm 0.062$
CMB+BAO+H0	$-1.119 \pm 0.068$
SN+CMB	$-1.026 \pm 0.041$
SN+CMB+BAO	$-1.014 \pm 0.040$
SN+CMB+H0	$-1.056 \pm 0.038$
SN+CMB+BAO+H0	$-1.047 \pm 0.038$

SNe Ia distances combined with CMB and/or BAO remain the best probe to constraint the DE equation of state :

- a **5%** measure of a constant DE EoS,  $w$ , is achievable
- currently little sensitivity to  $w(z)$

Including systematics and combined with BAO and CMB :  $w$  (cte) =  **$-1.018 \pm 0.057$**  (~6%) compatible with a cosmological constant

Let us consider now the  $w(z)$ CDM scenario, where dark energy has an evolving equation of state

$w(z)$ CDM  $\rightarrow$  there are now 5 independent cosmological parameters:  $H_0$ ,  $\Omega_m$ ,  $\Omega_{DE}$ ,  $w_0$ ,  $w_a$  (or only 4 if flatness is assumed)

The evolution of the dark energy equation of state is parameterized as  $w(z) = w_0 + w_a (1 - a)$  which is a first-order Taylor expansion in the scale factor:  $w_0 + w_a (1 - a)$

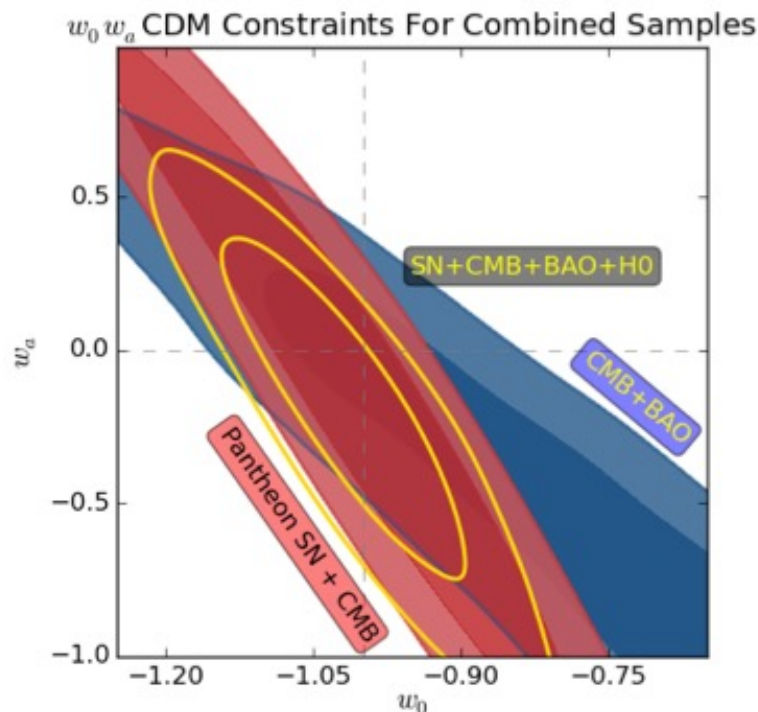
### Constraints in the $(w_0, w_a)$ plane (after marginalizing over the other parameters)

#### Some notes:

- The effect of  $w_a$  on the geometric observables is very weak  $\rightarrow$  probes of structure are more useful, since the evolution of dark energy affects structure formation

- Due to the weak constraints the figure only shows combined contours: SN+CMB, CMB+BAO, BAO+CMB, SN+CMB+BAO+H0\_prior

-  $\Lambda$ CDM is a point in this plane ( $w_0 = -1$ ,  $w_a = 0$ ) and is inside all the contours





Sample	$w_0$	$w_a$	$\Omega_m$	$H_0$	FoM
CMB+BAO	$-0.616 \pm 0.262$	$-1.108 \pm 0.771$	$0.343 \pm 0.025$	$64.614 \pm 2.447$	14.5
CMB+H0	$-1.024 \pm 0.347$	$-0.789 \pm 1.338$	$0.265 \pm 0.015$	$73.397 \pm 1.961$	9.1
CMB+BAO+H0	$-0.619 \pm 0.270$	$-1.098 \pm 0.781$	$0.343 \pm 0.026$	$64.666 \pm 2.526$	14.5
SN+CMB	$-1.009 \pm 0.159$	$-0.129 \pm 0.755$	$0.308 \pm 0.018$	$68.188 \pm 1.768$	31.4
SN+CMB+BAO	$-0.993 \pm 0.087$	$-0.126 \pm 0.384$	$0.308 \pm 0.008$	$68.076 \pm 0.858$	65.0
SN+CMB+H0	$-0.905 \pm 0.101$	$-0.742 \pm 0.465$	$0.287 \pm 0.011$	$70.393 \pm 1.079$	54.2
SN+CMB+BAO+H0	$-1.007 \pm 0.089$	$-0.222 \pm 0.407$	$0.300 \pm 0.008$	$69.057 \pm 0.796$	63.2

The **dark energy figure-of-merit** (FOM) is defined as the inverse of the area of the 1-sigma contour - or more precisely, it is the area of an ellipse that fits the contour (since it is defined from the Fisher matrix approach the contour is necessarily an ellipse).

The larger the FoM  $\rightarrow$  the smaller the contour  $\rightarrow$  the stronger the constraint.

The most powerful combination in the table is SN+CMB+BAO.

## Model selection (Goodness-of-fit)

The estimation of the cosmological parameter credible intervals (mean values and uncertainties) is not the last step of the cosmological data analysis process.

Using the SN - Pantheon example, let us look at its results:

### Nuisance parameters

(notice the uncertainties are much larger if only low-z SN are used)

Survey	$\alpha$	$\beta$	$\gamma$
Pantheon	$0.154 \pm 0.006$	$3.02 \pm 0.06$	$0.053 \pm 0.009$
Low-z	$0.154 \pm 0.011$	$2.99 \pm 0.15$	$0.076 \pm 0.030$

### Cosmological parameters

- constraints are worse if the full (stat+sys) errors are used (more realistic)

Analysis	Model	$w$	$\Omega_m$	$\Omega_\Lambda$
SN-stat	$\Lambda$ CDM		$0.284 \pm 0.012$	$0.716 \pm 0.012$
SN-stat	$o$ CDM		$0.348 \pm 0.040$	$0.827 \pm 0.066$
SN-stat	$w$ CDM	$-1.251 \pm 0.144$	$0.350 \pm 0.035$	
SN	$\Lambda$ CDM		$0.298 \pm 0.022$	$0.702 \pm 0.022$
SN	$o$ CDM		$0.319 \pm 0.070$	$0.733 \pm 0.111$
SN	$w$ CDM	$-1.090 \pm 0.220$	$0.316 \pm 0.072$	

From the table, it is clear that the results depend on the scenario assumed:

- $\Lambda$ CDM ( $\Omega_m$ ) - with few free parameters, the constraints are tighter
- oCDM ( $\Omega_m$   $\Omega_\Lambda$ ) - not only parameter uncertainties are larger, but the central values can change a lot (central values for  $\Lambda$ CDM are not even contained in the oCDM  $1\sigma$  confidence intervals)
- wCDM ( $\Omega_m$   $w$   $\Omega_K$ ) - constraints closer to the oCDM ones

**So, what is the final result?** What is our finding, is it  $\Omega_m$  0.30 or 0.32?

This is a question of **goodness-of-fit**. Among the various best-fits which one is the best?

We turn again to Bayesian inference to answer this question by performing **model comparison** tests.

There are different ways to evaluate the goodness-of-fit. The classic way is to look at the **chi-square**, while the most rigorous way is to use the **evidence**.

## Chi-square

Criteria based on the chi-square values are standard in determining the best model in all branches of physics.

The most usual quantity is the **reduced chi-square** of the best-fit, i.e., the chi-square normalised by the **number of degrees-of-freedom**,

$N_{\text{dof}} = N_d - N_p$  (where  $N_d$  is the number of datapoints - for example the number of redshift bins in the SN data - and  $N_p$  is the number of parameters in the model )

In this criterium, **the best model** (i.e., the favoured one) **is the one where the best-fit has the lowest reduced chi-square**,

$$\chi^2_{\text{red}} = \chi^2 / N_{\text{dof}}$$

## Evidence

It is the integral of the likelihood on the parameters space of a given cosmological model → it indicates the ‘average likelihood of a model’.

It may happen that a certain set of parameter values are a very good fit to the data (high likelihood values in that region of the parameter space), but overall this model can have a worse evidence than another one (for example because of having a larger number of parameters, or a large region of small likelihood values).

The evidence is thus a global way to characterize the goodness-of-fit of a model, beyond the simple assessment of finding which model has the “best best-fit”.

The evidence is a good number to show the balance between **best-fit vs. model complexity**.

**In this approach, the best model is the one with the highest Bayes factor**, computed from the evidences of the 2 models under comparison:

$$B = (\text{Evidence}_1 * \text{Prior}_1) / (\text{Evidence}_2 * \text{Prior}_2)$$

The **Jeffrey's scale** classifies the degree of preference for a model over another, based on the values of  $\ln B$ :

$< 1 \rightarrow$  inconclusive

$1 - 2.5 \rightarrow$  substantial evidence for one of the models

$2.5 - 5 \rightarrow$  strong evidence

$> 5 \rightarrow$  decisive evidence

The evidence is difficult to compute in practice with high precision, since it is a **multi-dimensional integral** of a possibly complex posterior distribution function.

Moreover, by sampling the posterior with a grid or an MCMC method, we only know a rough sample of it, which may be good enough to find the parameter constraints, but not precise enough to compute the total integral.

By design, MCMC algorithms only sample with high resolution the region near the maximum of likelihood. The tails of the distribution are usually badly sampled because they are not needed for parameter inference.

So the sample obtained with MCMC is not complete enough to compute the evidence. We need other Monte Carlo sampling methods to solve the multi-dimensional integral.

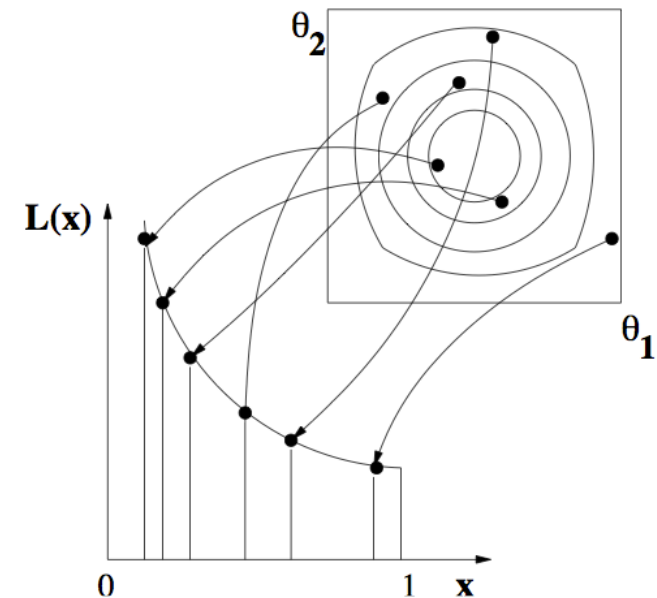
A popular algorithm for this is the **Nested Sampling**:

1. Sample  $N$  points randomly from within the prior, and evaluate their likelihoods. Initially we will have the full prior range available, i.e.,  $(0, X_0 = 1)$ .
2. Select the point with the lowest likelihood ( $L_j$ ). The prior volume corresponding to this point,  $X_j$ , can be estimated probabilistically. The average volume decrease is given as  $X_j/X_{j-1} = t$ , where  $t$  is the expectation value of the largest of  $N$  random numbers from uniform  $(0, 1)$ , which is  $N/(N + 1)$ .
3. Increment the evidence by  $E_j = L_j(X_{j-1} - X_{j+1})/2$ .
4. Discard the lowest likelihood point and replace it with a new point, which is uniformly distributed within the remaining prior volume  $(0, X_j)$ . The new point must satisfy the hard constraint on likelihood of  $L > L_j$ .
5. Repeat steps 2–4, until the evidence has been estimated to some desired accuracy.

This means: find iso-regions of likelihood. If they are ‘nested’ the integrand is monotonic  $\rightarrow$  the integral reduces to 1-dimension.

For each layer  $\rightarrow E_j = \frac{L_j}{2} (X_{j-1} - X_{j+1})$ .

The total evidence is  $\rightarrow E = \sum_{j=1}^m E_j$ ,



## Information criteria

Besides the evidence, there are alternative approximate methods, much simpler to compute, that can also be used for model selection and quantify the balance of best-fit vs. model complexity. Some popular of these **information criteria** are:

**Akaike information criterion:**  $AIC = -2 \ln L_{\text{bestfit}} + 2 n_p = \chi^2_{\text{bestfit}} + 2n_p$   
(this formula is the result of a minimisation of entropy criterium)

**Bayesian information criterion:**  $BIC = -2 \ln L_{\text{bestfit}} + n_p \ln(n_d)$   
(based on an approximation of the evidence)  
BIC penalizes more the complexity than AIC does.

**Deviance information criterion:**  $DIC = 2 \chi^2_{\text{mean}} - \chi^2_{\text{bestfit}}$   
(it is like an effective  $\chi^2$ , sensitive to the difference between the best-fit and the full distribution).

For all information criteria, **the best model is the one with the lowest value.**



## Results of model selection

In this example, SN data was used to test two very different scenarios:  $\Lambda$ CDM and UDM (model where DM and DE are one single fluid).

This model has one density parameter less, but two new additional parameters - so one parameter more than  $\Lambda$ CDM in total).

Two different UDM models were tested and (like  $\Lambda$ CDM) both are able to produce  $D_L(z)$  functions that allow for good fits to the SN data for certain values of their parameters.

Various model selection criteria were computed:

**The question is, is there enough evidence to select UDM over  $\Lambda$ CDM?**

	UDM	$\Lambda$ CDM	UDM <sub>ph</sub>
$\chi^2_{\min}$	552.59	552.77	552.75
$\chi^2_{\text{red}}$	0.9478	0.9449	0.9481
$\ln B_{\text{UA}}$	-0.0196	0	0.6850
BIC	584.485	571.902	584.644
DIC	553.250	552.770	552.814

- The first UDM model is the one with the smallest best-fit  $\chi^2$  , i.e., it contains a vector of parameter values that produced the closest fit to the data.

However, since this model has more cosmological parameters than  $\Lambda$ CDM it is penalized, and the lowest reduced chi-square turns out to be the one of  $\Lambda$ CDM. The complexity of the model (having more free parameters) is always penalized in these criteria. This is because increasing the number of parameters naturally helps in finding a closer fit (in a potentially artificial way).

- UDM\_ph is the model with largest evidence. Indeed, the Bayes factor of the second UDM model with respect to  $\Lambda$ CDM is positive, although smaller than one → the analysis shows a very slight inconclusive preference for this model UDM\_ph

- BIC shows a reasonable preference for  $\Lambda$ CDM.

- DIC shows a slight preference for  $\Lambda$ CDM.

**The analysis does not show a conclusive preference for any of the models**

(but given the close results, it shows that it is useful to compute all the criteria).