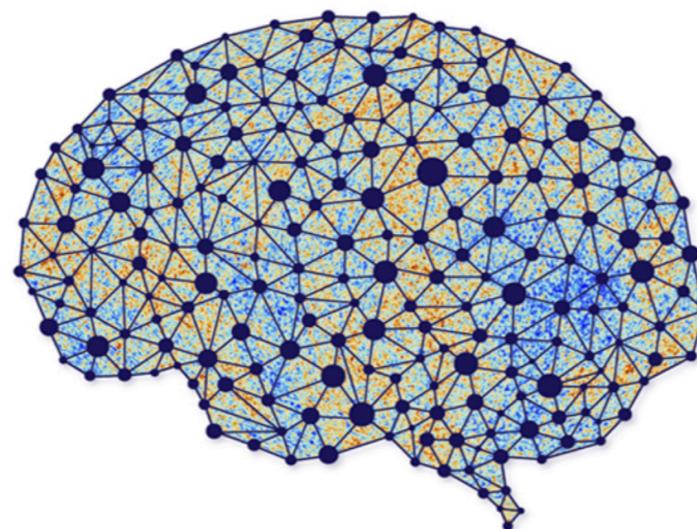


# Physics 361 - Machine Learning in Physics

## Lecture 3 – Background

Jan. 28<sup>th</sup> 2025



AI  
∩  
Universe

$$\begin{aligned} \mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i \bar{\psi} \gamma^\mu \psi + h.c. \\ & + \bar{\phi}_i \gamma_{ij} \phi_j + h.c. \\ & + D_\mu \phi^\dagger - V(\phi) \end{aligned}$$

Moritz Münchmeyer

# Unit 1: Background

## 1.2 Classical Statistics and Data Analysis Background

Sources:

- Cowan - Statistical data analysis
- also mostly covered in [deeplearningbook.org](http://deeplearningbook.org)

# Estimators

- A (point) estimator makes a prediction for some quantity of interest  $\lambda$ .  
We write  $\hat{\lambda}$ , where "hat" means estimator.  
E.g. estimator of the mean height  $\bar{h}$   
$$\hat{h} = \frac{1}{N} \sum h_i$$
- Given a dataset  $\{x^{\text{obs}}\}$  drawn / sampled from a PDF  $P(x)$ , an estimator is some function  
$$\hat{\lambda} = \mathbb{E}[\{x^{\text{obs}}\}]$$

$\nwarrow$  some function, "guessed" or derived
- An optimal estimator is
  - unbiased :  $\langle \hat{\lambda} \rangle = \lambda$  "correct on average"
  - minimum variance :  $\text{Var}[\hat{\lambda}]$  is as small as possible (smallest error)

## Some common estimators

notation:  $\bar{x}$   
„bar“ means mean

(sample) mean:  $\hat{x} = \frac{1}{n} \sum_{i=1}^n x_i^{\text{obs}}$

variance:  $\hat{V} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{x})^2$

covariance:  $\hat{\text{cov}} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{x})(y_i - \hat{y})$

These estimators have a variance, e.g.

$$V[\hat{x}] = \frac{\sigma^2}{n}$$

where  $\sigma^2$  is the variance of  $x$

## Likelihood, Posterior, Prior

- The Likelihood is the probability of measuring data  $\vec{d}$  given a model  $M$  with parameters  $\vec{\lambda}$ .

$f(\vec{d} | M, \vec{\lambda})$  often not explicitly written  
often used as the Loss function  
in machine Learning.

- The Posterior is the probability of parameters  $\vec{\lambda}$  given observed data  $\vec{d}$ .

$$P(\vec{\lambda}, M | \vec{d})$$

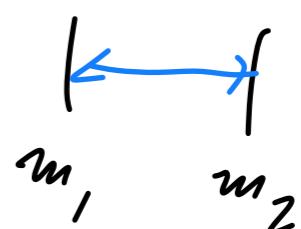
e.g.  $M$ : standard model  
 $\vec{\lambda}$ : Higgs mass  
 $\vec{d}$ : collision data

- We can use Bayes theorem to get P from L

$$P(\vec{\lambda} | \vec{d}) = \frac{f(\vec{d} | \vec{\lambda}) P(\vec{\lambda})}{P(\vec{d})}$$

- The prior  $P(\vec{\lambda})$  is the probability of  $\vec{\lambda}$  before we perform the measurement.

E.g. flat in some interval



- Gaussian around some prior measurement.

Prior is important if the data is not very constraining.

- Finally we have the evidence

$P(\vec{d})$  is the probability of the data under model  $M$  for ANY parameters  $\vec{\lambda}$ .

$$P(\vec{d}) = \int L(d|\vec{\lambda}) P(\vec{\lambda}) d\vec{\lambda}$$

Often used in model comparison.

# Maximum Likelihood and Maximum Aposteriori

- Combining the concepts of estimators and likelihoods we define the

## Maximum Likelihood estimator

$$\hat{\lambda}_{ML} = \underset{\lambda}{\operatorname{argmax}} \ell(\vec{d} | \vec{\lambda})$$

Sometimes we can do it analytically:

$$\frac{\partial \ell(\vec{d} | \vec{\lambda})}{\partial \lambda} \Big|_{\lambda = \hat{\lambda}} = 0$$

If not analytically tractable use an numerical optimizer.

- We define the maximum a posteriori estimator

$$\hat{\lambda}_{MAP} = \underset{\lambda}{\operatorname{argmax}} P(\vec{\lambda} | \vec{d}) \quad (\text{MAP})$$

## Gaussian Likelihood

Gaussian Likelihoods are often good approximations of data in physics.

Example: We measure a person's weight  $w$ . To get an uncertainty we measure 100 times.

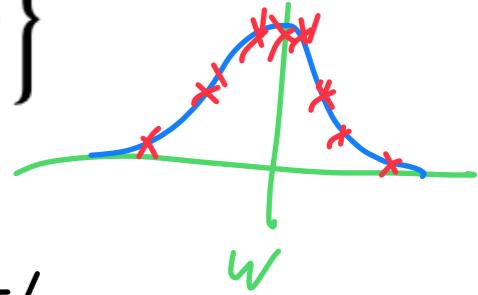
We assume that the data is described by  $d_i = w + n_i$

$\xrightarrow{\text{truth}}$        $\uparrow$   
 $n_i$  Gaussian noise with variance  $\sigma_w^2$

We want to measure the parameters  $w$  and  $\sigma_w^2$ .

For a single observation the Likelihood is:

$$\mathcal{L}(d|w, \sigma_w) \equiv P(d|w, \sigma_w) = \frac{1}{\sqrt{2\pi\sigma_w^2}} \exp\left\{-\frac{(d-w)^2}{2\sigma_w^2}\right\}$$



For  $m$  independent measurements  $\mathcal{L}$  is the product:

$$\mathcal{L}(\{d_i\}_{i=1}^m | w, \sigma_w) = \frac{1}{(2\pi\sigma_w^2)^{m/2}} \exp\left\{-\frac{\sum_{i=1}^m (d_i - w)^2}{2\sigma_w^2}\right\}$$

We get the posterior from Bayes theorem:

$$P(w, \sigma_w | \{d_i\}) = \frac{\mathcal{L}(\{d_i\} | w, \sigma_w) P(\lambda)}{P(\{d_i\})}$$

If the prior is flat we can work with the maximum Likelihood estimator.

The max. likelihood estimator for  $w$  is:

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\sum_{j=1}^m (d_j - w)}{\sigma_w^2 (2\pi\sigma_w^2)^{m/2}} \exp \left\{ -\frac{\sum_{i=1}^m (d_i - w)^2}{2\sigma_w^2} \right\}$$

$$\frac{\partial \mathcal{L}}{\partial w} = 0 \Leftrightarrow \sum_{j=1}^m (d_j - w) = 0$$

solve for  $w$ :

$$w = \hat{w} = \frac{1}{m} \sum_{i=1}^m d_i$$

This is the mean, as expected.

We could also do  $\frac{\partial \mathcal{L}}{\partial \sigma_w}$  to find  $\hat{\sigma}_w$ .

While here the result is „obvious”, this is a general procedure to find estimators given a model (= Likelihood).

## Sampling the posterior: MCMC

- We have Learned how to get the posterior for parameters  $\vec{\lambda}$  given data  $\vec{d}$ .
- It is given by  $P(\vec{\lambda} | \vec{d}) = \frac{P(\vec{d} | \vec{\lambda}) P(\vec{\lambda})}{P(\vec{d})}$   
Often we only need the unnormalized posterior  
 $P(\vec{\lambda} | \vec{d}) \propto P(\vec{d} | \vec{\lambda}) P(\vec{\lambda})$   
The evidence does not depend on  $\vec{\lambda}$  and can thus be ignored when finding  $\vec{\lambda}$ .
- There is still a problem to deal with: computational difficulty.

Example: Assume  $P(\vec{\lambda} | \vec{d})$  depends on 20 parameters.

How do we work with such a high dimensional posterior?

If we evaluate it on a grid with 100 points per dimension we'd need  $(100)^{20}$  evaluations.

This cannot be done.

Instead one works with such posteriors by sampling from them using a method called Markov Chain Monte Carlo (MCMC).

---

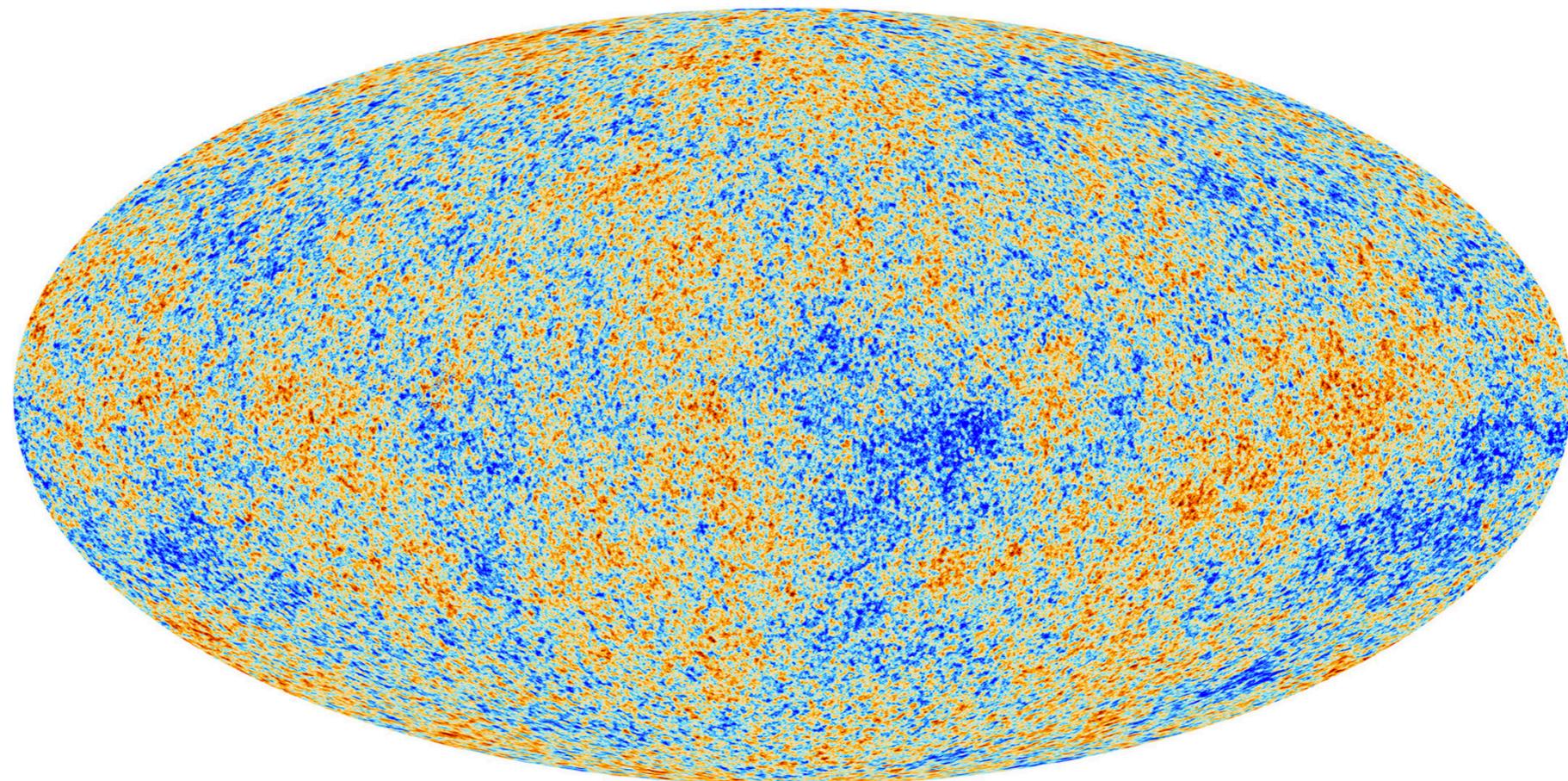
Perhaps covered later in this course.

# Unit 1: Background

**1.3 An example of data analysis in physics that uses these concepts**

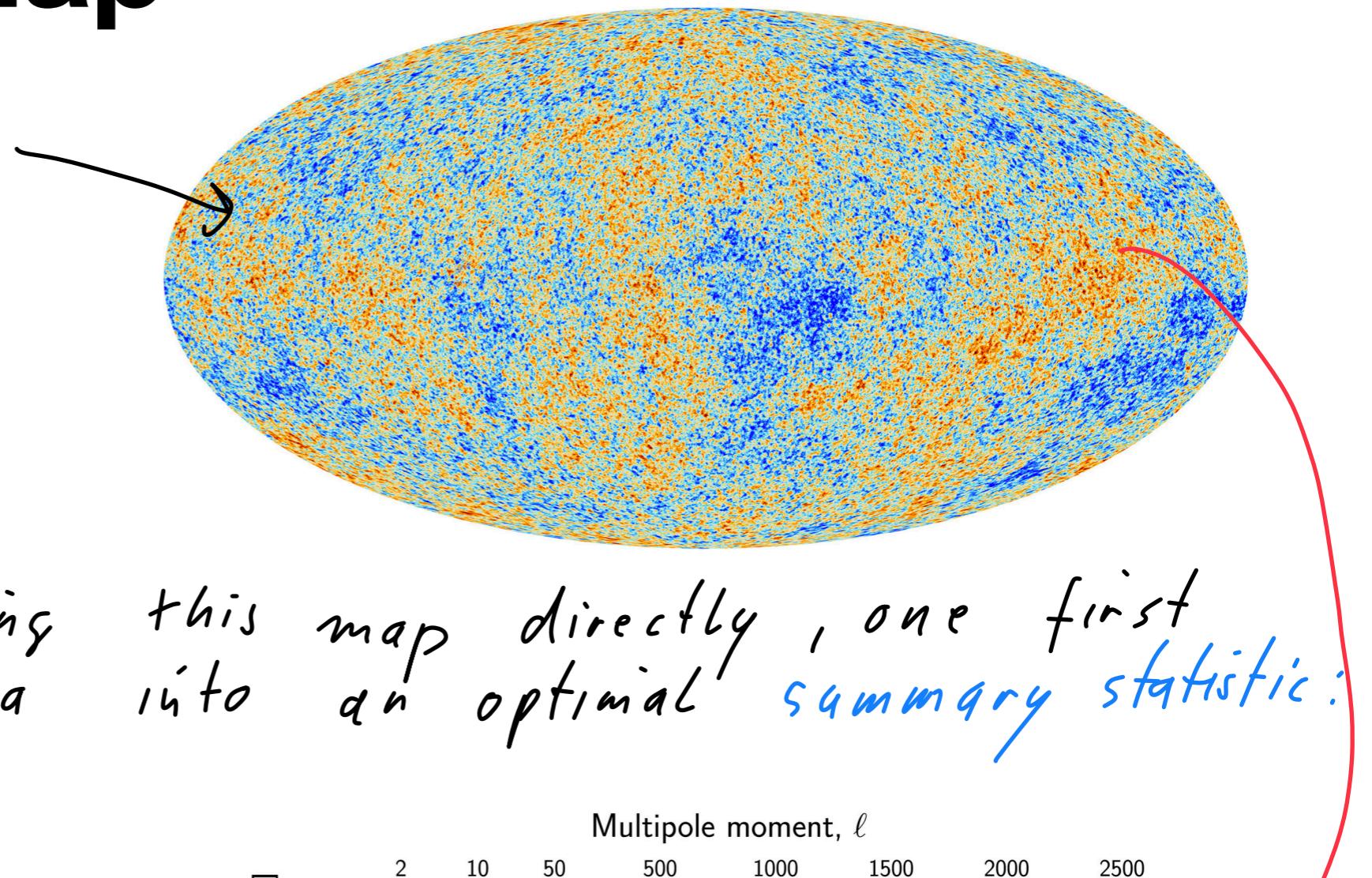
# A physics application of these concepts: The CMB power spectrum

- I want to show you an application of these ideas to real physics. The CMB power spectrum analysis is one of the jewels of physics, telling us much of what we know about the history of the universe.
- Of course this is a complicated topic and I can only give you a brief idea.
- The Cosmic Microwave Background is a radiation that permeates the universe. It has a temperature of about 3K and has been measured extremely precisely.



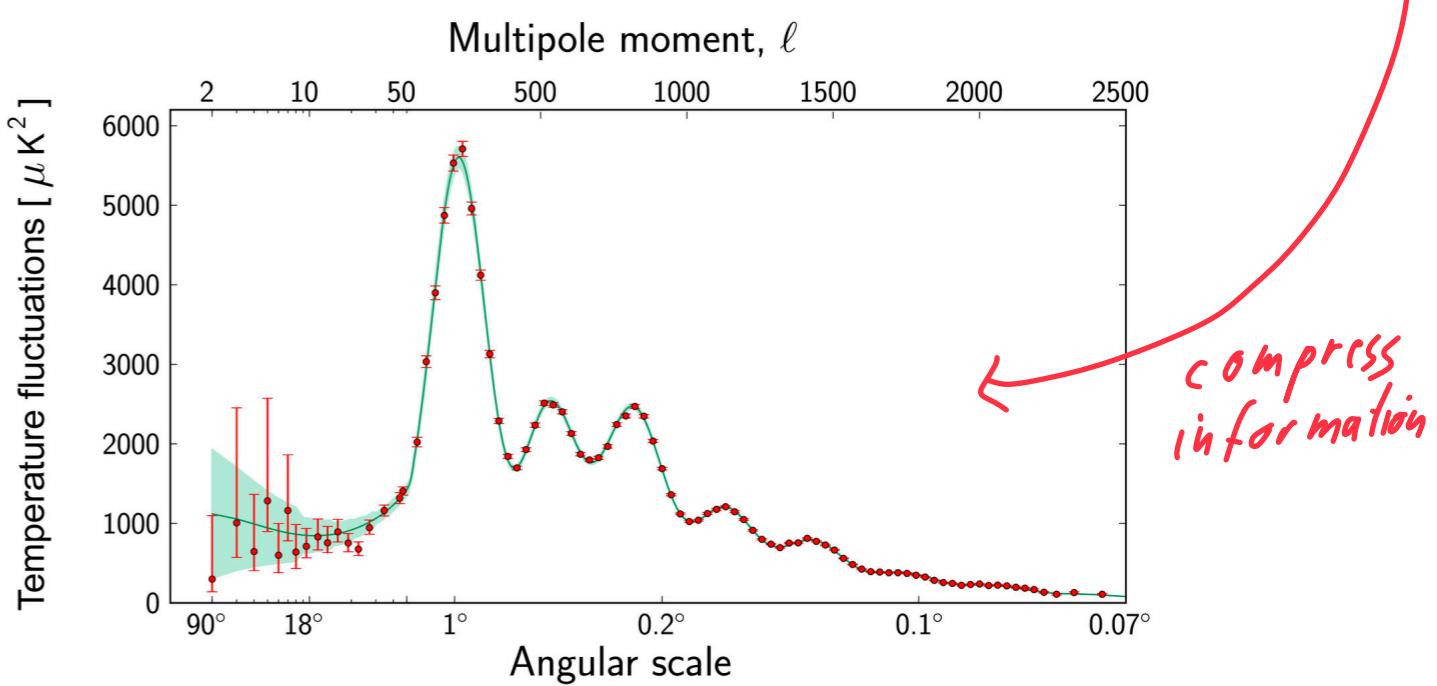
# The CMB map

CMB temperature  
field  $T^{CMB}(\vec{u})$



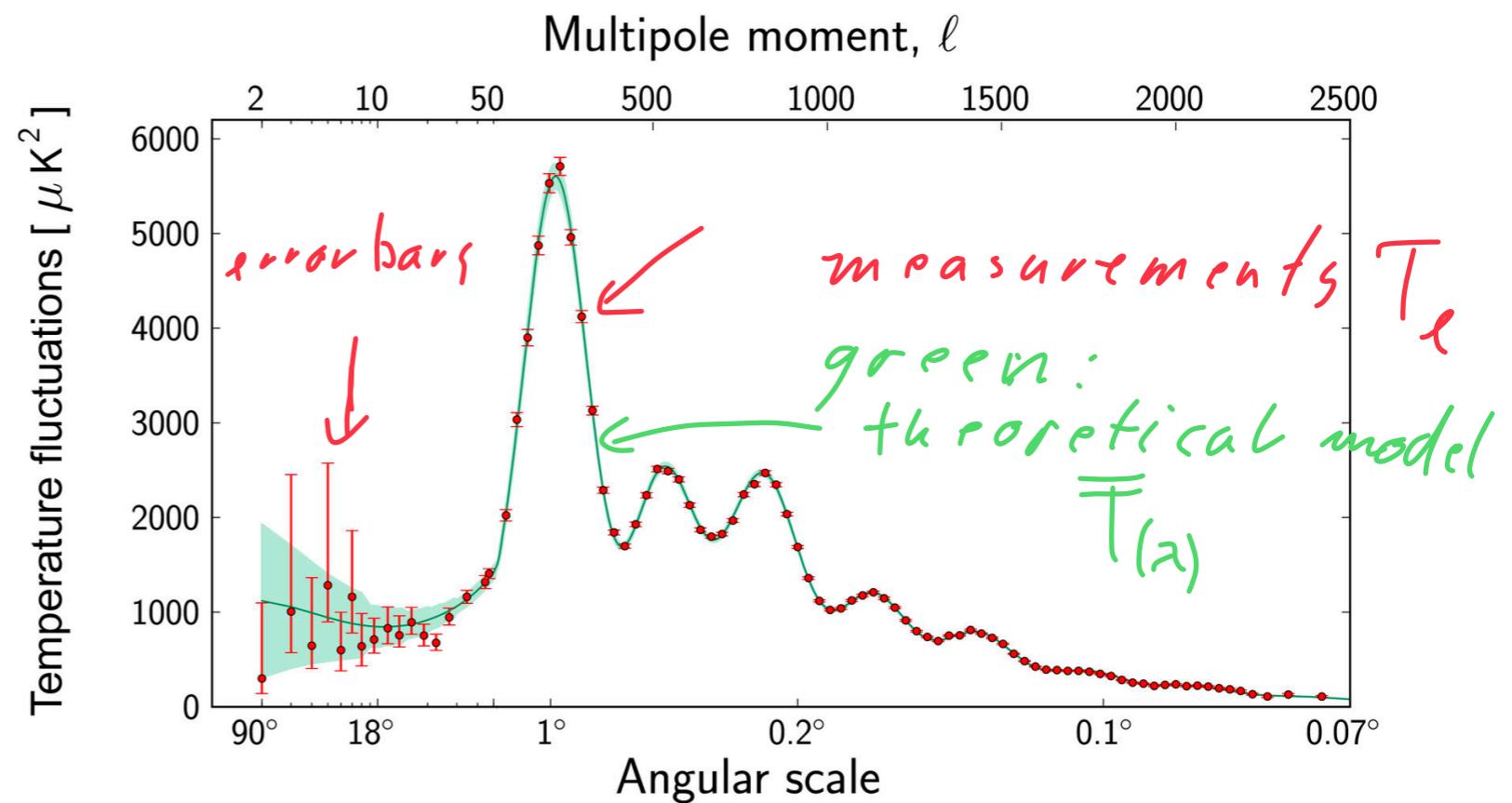
Rather than analyzing this map directly, one first compresses the data into an optimal summary statistic:

CMB power spectrum



# The dataset

- Our measurements will be the power spectrum as a function of scale. It can be extracted from the previous map by taking its Fourier transform and squaring the mode amplitudes.
- Our data points are:



data set:  $\{T_\ell\}$  with some error bars  
 $T$   
Fluctuation amplitudes of CMB temperature.

# The likelihood

- We will make a Gaussian likelihood.
- Our theory model is that the mean curve (green) is a known function of cosmological parameters, which we want to measure.

Theoretical model (green curve):

$\bar{T}_e(\vec{\lambda})$  depends on cosmological parameters, e.g.:

$\mathcal{R}_c$ : Dark matter density

$\mathcal{R}_b$ : baryon density

It can be calculated from the Laws of nature (e.g. General Relativity)

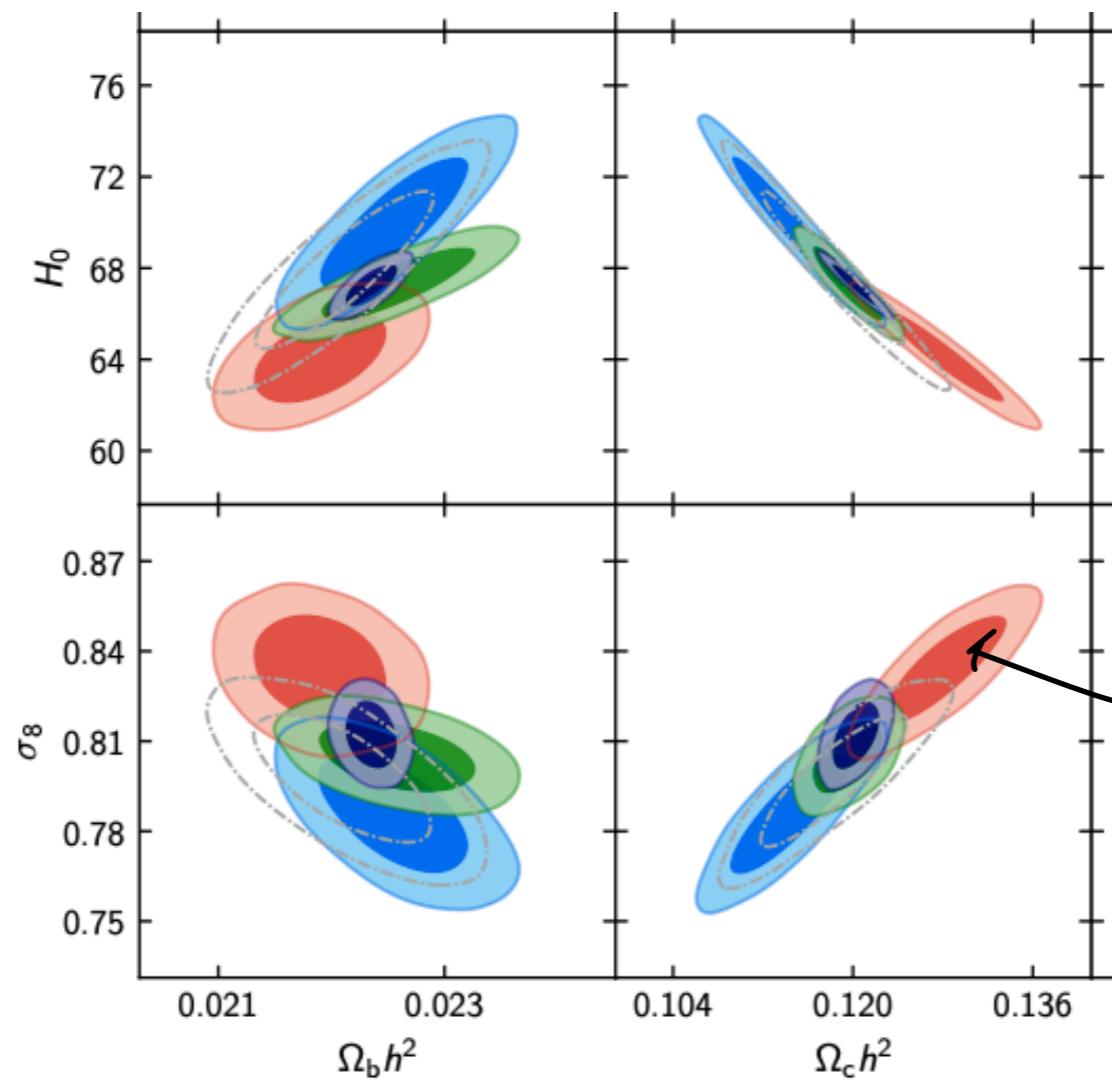
$\Rightarrow$  Gaussian Likelihood

$$\ln \mathcal{L}(\{\mathcal{T}_e\} | \vec{\lambda}) \propto \frac{1}{2} \sum_{\ell=1}^n \frac{(\mathcal{T}_e - \bar{T}_e(\vec{\lambda}))^2}{2 \sigma_e^2}$$

$\Rightarrow$  Maximize Likelihood to find theory parameters  $\vec{\lambda}$  that fit the data best.

# Sampling the posterior

- Since we have a likelihood, we can now use Bayes theorem and get the posterior.
- Then, one can sample from the posterior, to get the likely values of the desired cosmological parameters.
- The result looks are Monte Carlo plots like this:



$H_0$  : Hubble parameter  
 $\Omega_c$  : dark matter  
 $\Omega_b$  : baryons  
out and two  
sigma regions  
of the posterior  
 $P(\chi | \text{data})$

# Unit 1: Background

## 1.4 Information theory background

## What is information theory?

- Quantify how much information is in a given signal.
  - Optimal compression
- The information in a signal depends on the probability of the signal.
  - A certain signal/event contains no information.
  - A highly unlikely event contains a lot of information.
- In machine Learning, information theory is used to characterize PDFs and their similarity.

# Entropy

- In statistical physics (thermodynamics) the entropy is given by

$$S = (k_B) \log \Omega$$

↑ here we set it to 1

$\Omega$ : number of microstates (equally likely)

- We can re-write this as

$$S = - \log P \quad \text{where} \quad p = \frac{1}{\Omega}$$

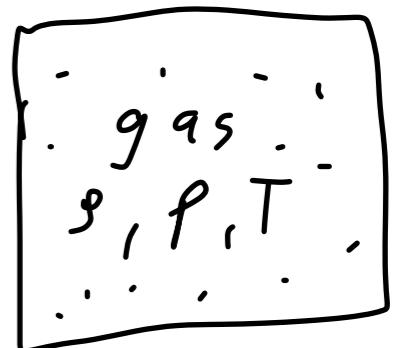
is the probability of each m.state

- The general definition of entropy

Shannon entropy

$$S = - \sum_i p_i \log p_i$$

- log here is base e  $\rightarrow$  entropy is in „nats“  
(base 2  $\rightarrow$  „ „ in bits)



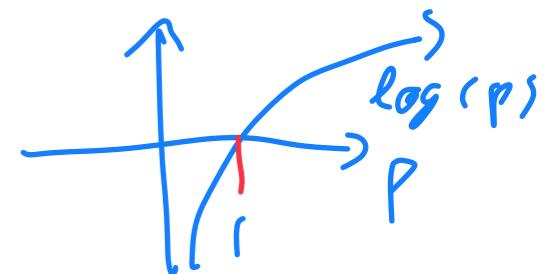
# Properties of the entropy

- discrete PDF

$$S = - \sum_i p_i \underbrace{\log(p_i)}$$

-  $\log(p_i)$  is the „self-information“ of event  $i$

$$= - \mathbb{E}_{x \sim P(x)} [\log(P(x))]$$



A unlikely event has a large selfinfo.

$$\text{self information } I(x) = -\log(P(x))$$

A certain event has self inform. 0 maximum possible entropy is

- continuous PDF

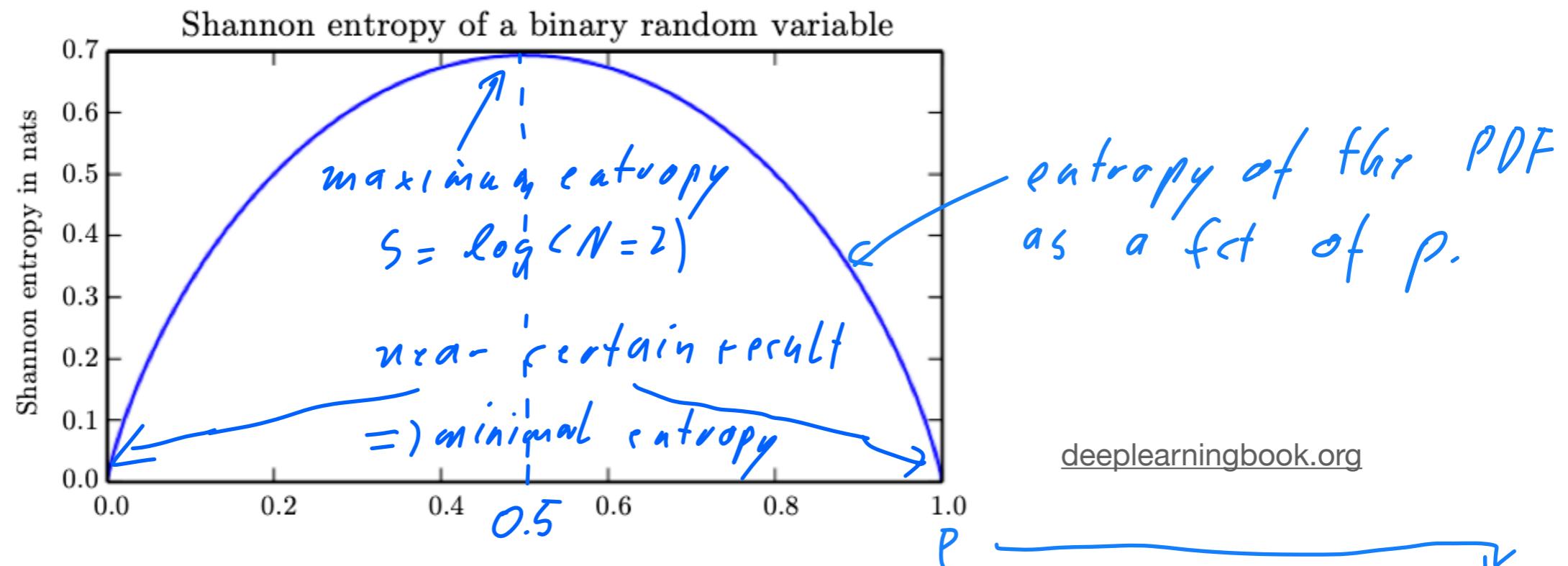
$$S = - \int dx p(x) \log(p(x))$$

$$= - \mathbb{E}_{x \sim P(x)} [\log(P(x))]$$

More uniform distributions have a higher entropy.  
(spread out)

$$S = \log(N)$$

# Example : Shannon entropy of a binary variable



Example : binary random variable  $X : \{1 : p, 0 : 1-p\}$   
calculate the entropy

$$S = - \sum_i p_i \log(p_i)$$

$$= - (1-p) \log(1-p) - p \log(p)$$

# Course logistics

- **Reading for this lecture:**
  - **For example:** Deeplearningbook.org chapter 3 and 5.
- **Problem set:** First problem set this week.