

# Cosmology with Bayesian statistics and information theory

## Lecture 1: Aspects of probability theory

... a.k.a. *why am I not allowed to “change the prior” or “cut the data”?*

Florent Leclercq

Institute of Cosmology and Gravitation, University of Portsmouth

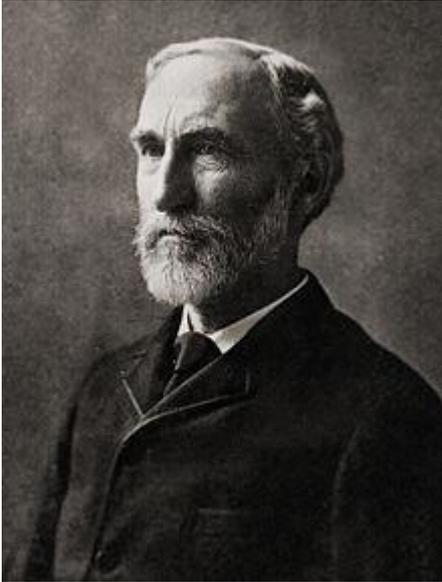
<http://icg.port.ac.uk/~leclercq/>



March 6<sup>th</sup>, 2017

# Introduction: why proper statistics matter

## An historical example: the Gibbs paradox

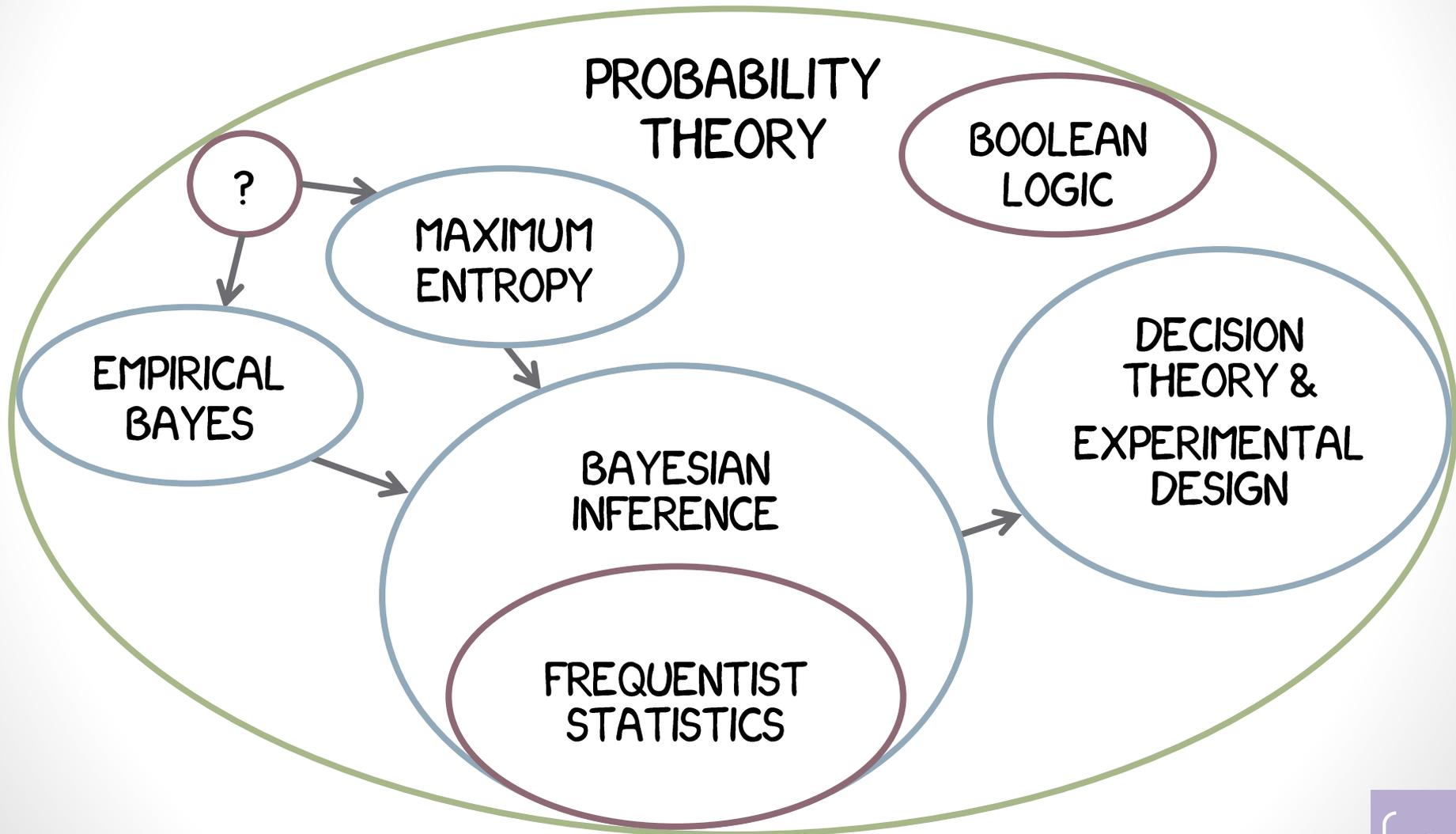


J. Willard Gibbs (1839-1903)

- Gibbs's canonical ensemble and grand canonical ensembles, derived from the maximum entropy principle, *fail to correctly predict thermodynamic properties* of real physical systems.
- The predicted entropies are always larger than the observed ones... there must exist *additional microphysical constraints*:
  - Discreteness of energy levels: radiation: Planck (1900), solids: Einstein (1907), Debye (1912), Ising (1925), individual atoms : Bohr (1913)...
  - ...Quantum mechanics: Heisenberg, Schrödinger (1927)

The first clues indicating the need for quantum physics were uncovered by seemingly “unsuccessful” application of statistics.

# Jaynes's "probability theory": an extension of ordinary logic



# Reminders

- A tribute to my PhD supervisor (Benjamin Wandelt):

Ben's summary of Bayesian statistics:

*“Whatever is uncertain gets a pdf.”*

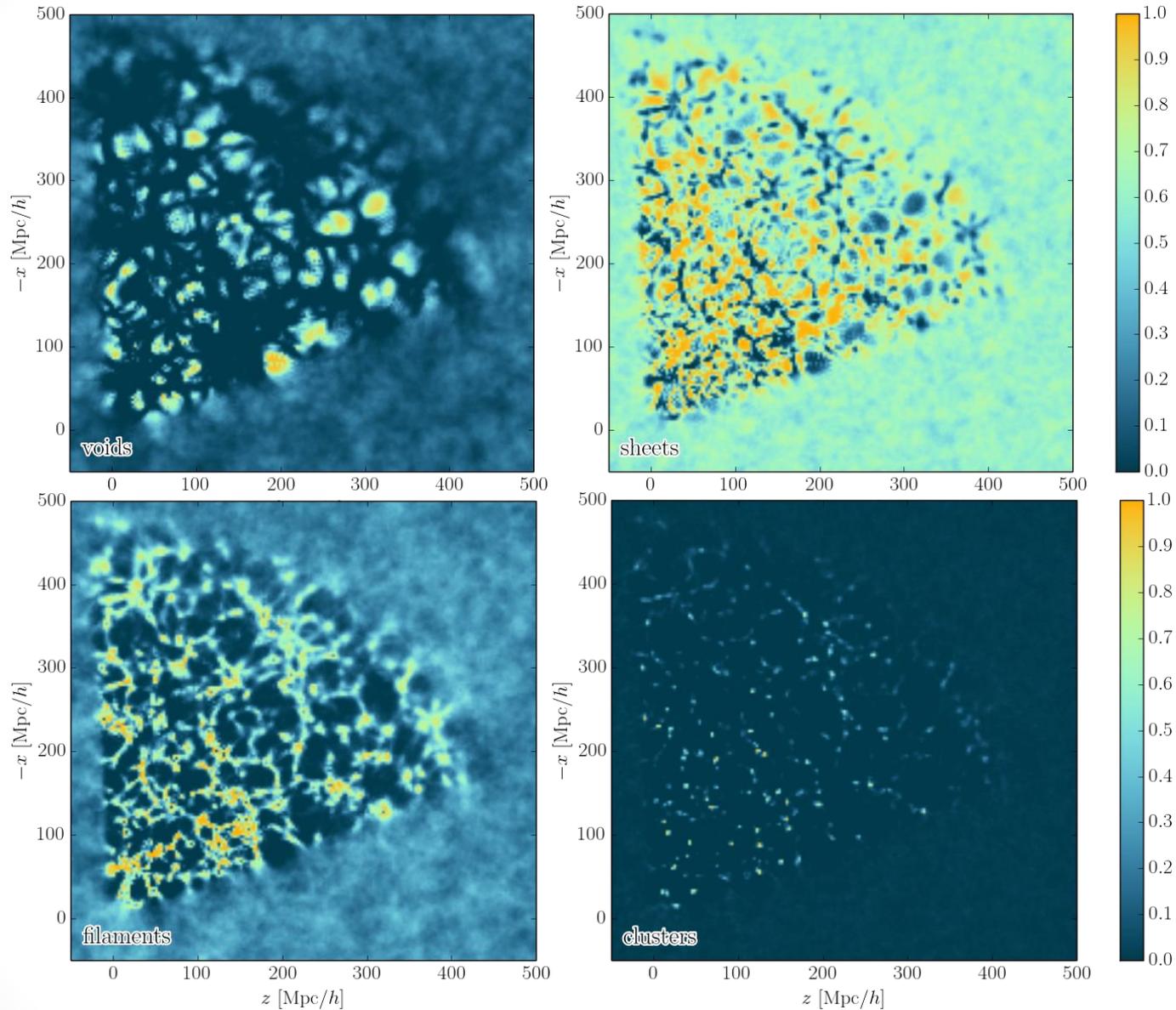
- Product rule:  $p(AB|C) = p(A|BC) p(B|C)$
- Sum rule:  $p(A + B|C) = p(A|C) + p(B|C) - p(AB|C)$

- Bayes's formula: 
$$p(s|d) = \frac{p(d|s) p(s)}{p(d)}$$

Diagram labels for Bayes's formula:  
- **posterior** points to  $p(s|d)$   
- **evidence** points to  $p(d)$   
- **likelihood** points to  $p(d|s)$   
- **prior** points to  $p(s)$

- Bayesian model comparison: 
$$\mathcal{B}_{12} = \frac{p(d|\mathcal{M}_1)}{p(d|\mathcal{M}_2)}$$

# Structures in the cosmic web



# A decision rule for structure classification

- Space of “input features”:

$$\{T_0 = \text{void}, T_1 = \text{sheet}, T_2 = \text{filament}, T_3 = \text{cluster}\}$$

- Space of “actions”:

$$\{a_0 = \text{“decide void”}, a_1 = \text{“decide sheet”}, a_2 = \text{“decide filament”}, a_3 = \text{“decide cluster”}, a_{-1} = \text{“do not decide”}\}$$

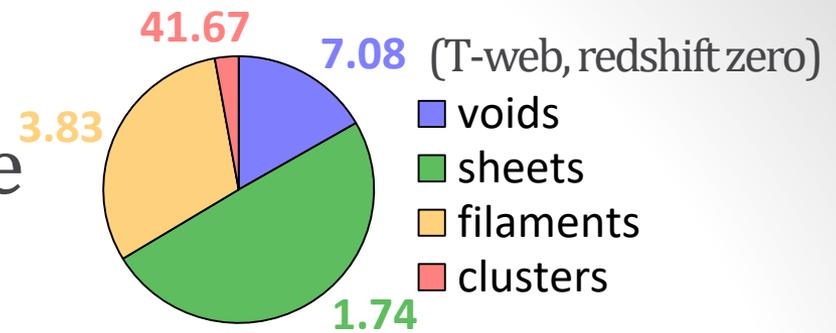
➡ A problem of **Bayesian decision theory**:

one should take the action that maximizes the utility

$$U(a_j(\vec{x}_k)|d) = \sum_{i=0}^3 G(a_j|T_i) \mathcal{P}(T_i(\vec{x}_k)|d)$$

- How to write down the gain functions?

# Gambling with the Universe



- One proposal:
 
$$G(a_j | T_i) = \begin{cases} \frac{1}{\mathcal{P}(T_i)} - \alpha & \text{if } j \in \llbracket 0, 3 \rrbracket \text{ and } i = j & \text{“Winning”} \\ -\alpha & \text{if } j \in \llbracket 0, 3 \rrbracket \text{ and } i \neq j & \text{“Losing”} \\ 0 & \text{if } j = -1. & \text{“Not playing”} \end{cases}$$

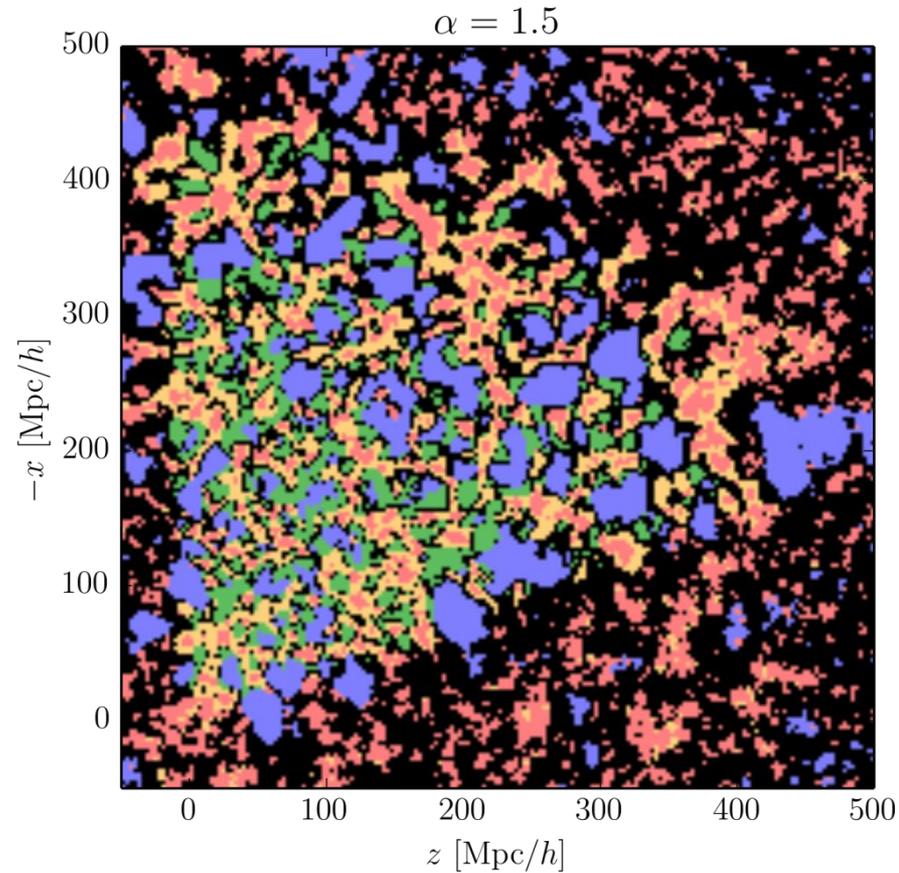
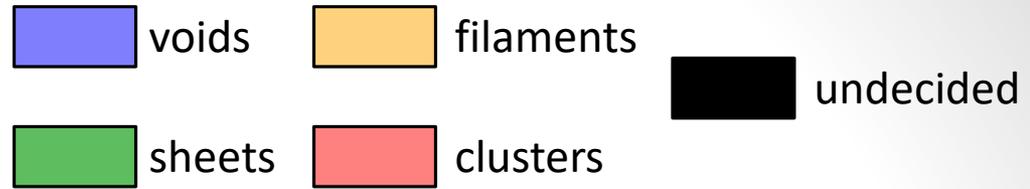
- Without data, the expected utility is

$$U(a_j) = 1 - \alpha \quad \text{if } j \neq -1 \quad \text{“Playing the game”}$$

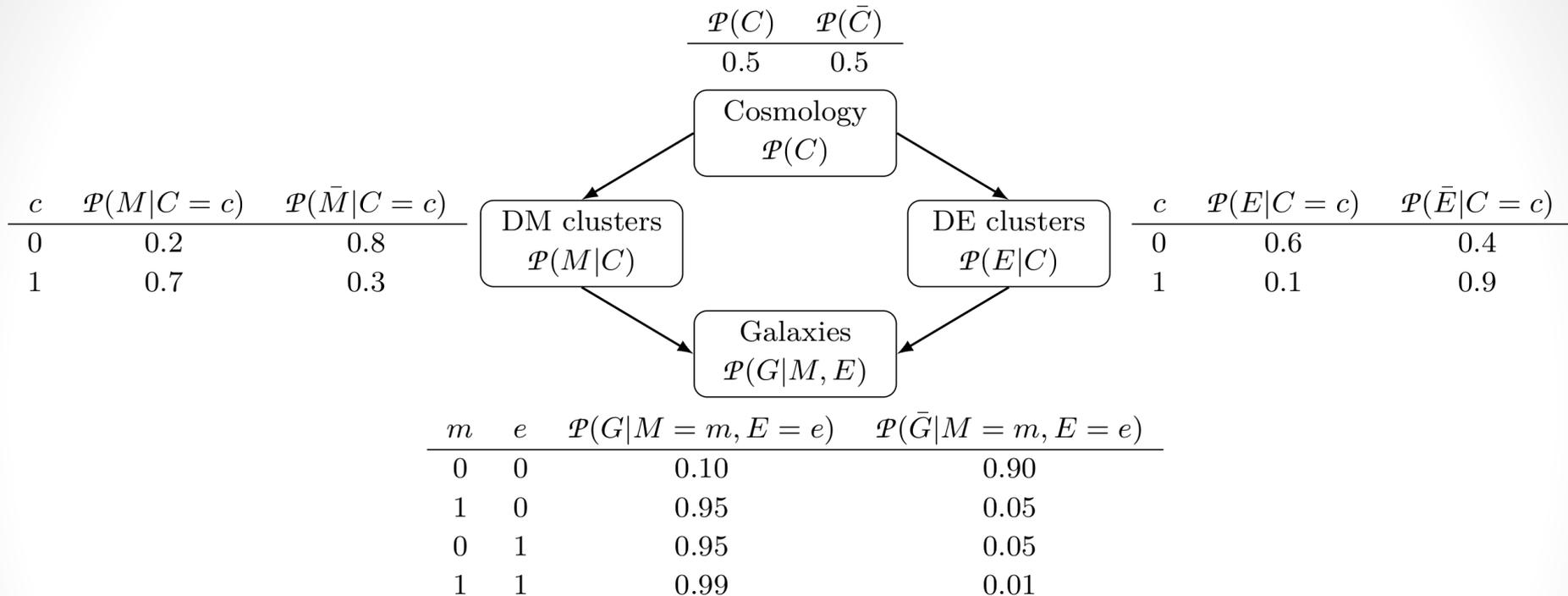
$$U(a_{-1}) = 0 \quad \text{“Not playing the game”}$$

- With  $\alpha = 1$ , it's a *fair game*  $\Rightarrow$  always play  $\Rightarrow$  “speculative map” of the LSS
- Values  $\alpha > 1$  represent an *aversion for risk*  $\Rightarrow$  increasingly “conservative maps” of the LSS

# Playing the game...



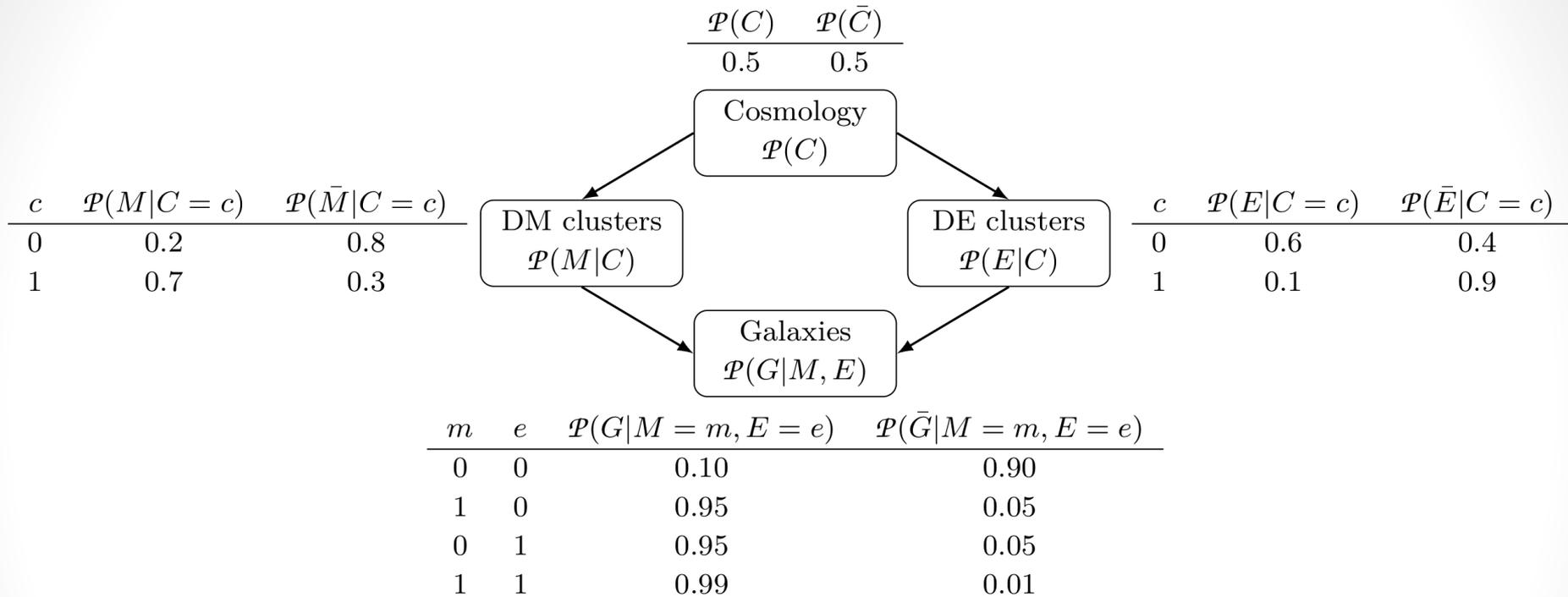
# Bayesian networks



Bayesian networks are probabilistic graphical models consisting of:

- A **directed acyclic graph**
- At each node, **conditional probabilities distributions**

# Bayesian networks



$$p(C, M, E, G) = p(C) p(E|C) p(M|C, \cancel{E}) p(G|\cancel{C}, M, E)$$

$$p(C, M, E, G) = p(C) p(E|C) p(M|C) p(G|M, E)$$

# Bayesian networks

## inference and prediction

- Inference:

$$p(M|G) = \frac{p(M,G)}{p(G)} = \frac{\sum_{c,e} p(C=c, M=1, E=e, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.4313}{0.70305} \approx 0.6135$$

$$p(E|G) = \frac{p(E,G)}{p(G)} = \frac{\sum_{c,m} p(C=c, M=m, E=1, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.3363}{0.70305} \approx 0.4783$$

$$p(\bar{M}, \bar{E}|G) = \frac{p(\bar{M}, \bar{E}, G)}{p(G)} = \frac{\sum_c p(C=c, M=0, E=0, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.0295}{0.70305} \approx 0.0420$$

- Prediction:

$$p(G|C) = \frac{p(G,C)}{p(C)} = \frac{\sum_{m,e} p(C=1, M=m, E=e, G=1)}{p(C=1)} = 0.7233$$

# Bayesian networks

the “explaining away” phenomenon

$$p(E|M, G) = \frac{p(E, M, G)}{p(M, G)} = \frac{\sum_c p(C=c, M=1, E=1, G=1)}{\sum_{c,e} p(C=c, M=1, E=e, G=1)} = \frac{0.09405}{0.4313} \approx 0.2181$$

$$p(E|G) = \frac{p(E, G)}{p(G)} = \frac{\sum_{c,m} p(C=c, M=m, E=1, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.3363}{0.70305} \approx 0.4783$$

- So we have both:

$$p(E|M) = p(E)$$

$$p(E|M, G) < p(E|G)$$

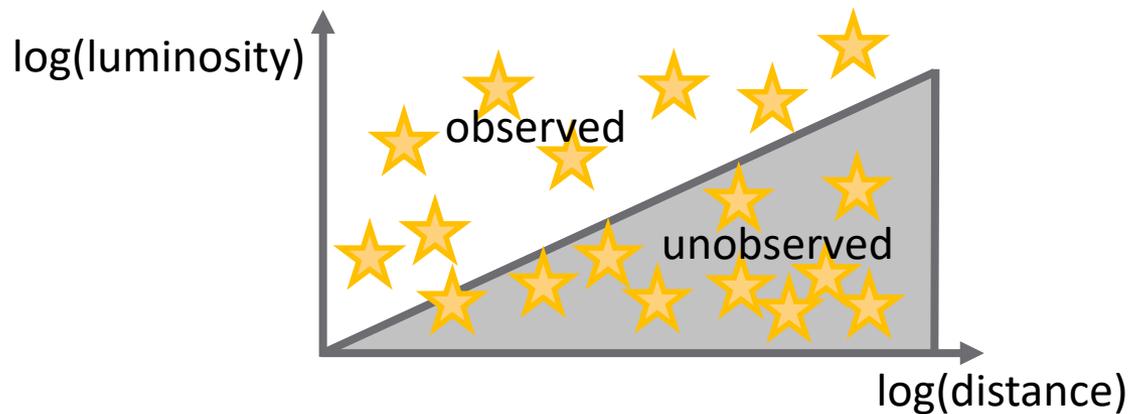
- This is “**collider bias**” or the “**explaining away**” phenomenon: two causes collide to explain the same effect.
- Particular case: “**selection bias**” or “**Berkson’s paradox**”

$$0 < p(A) < 1; \quad 0 < p(B) < 1; \quad p(A|B) = p(A)$$

$$\Rightarrow \begin{array}{l} p(A|B, C) < p(A|C) \\ p(A|\bar{B}, C) = 1 > p(A|C) \end{array} \quad C = A + B$$

# Malmquist bias

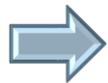
- Malmquist (1925) bias: in magnitude-limited surveys, far objects are preferentially detected if they are intrinsically bright.



$$0 < p(A) < 1; \quad 0 < p(B) < 1; \quad p(A|B) = p(A)$$

$$C = A + B$$

detected      bright      close



$$p(A|\bar{B}, C) = 1 > p(A|C)$$

# Bayesian hierarchical models

- Simple inference:  
$$p(\theta|d) \propto p(d|\theta) p(\theta)$$

prior
- Adaptive prior:  
$$p(\theta|d) \propto p(d|\theta) p(\theta|\eta) p(\eta)$$

prior      hyperprior
- ... or a full hierarchy of hyperpriors.

- Examples:

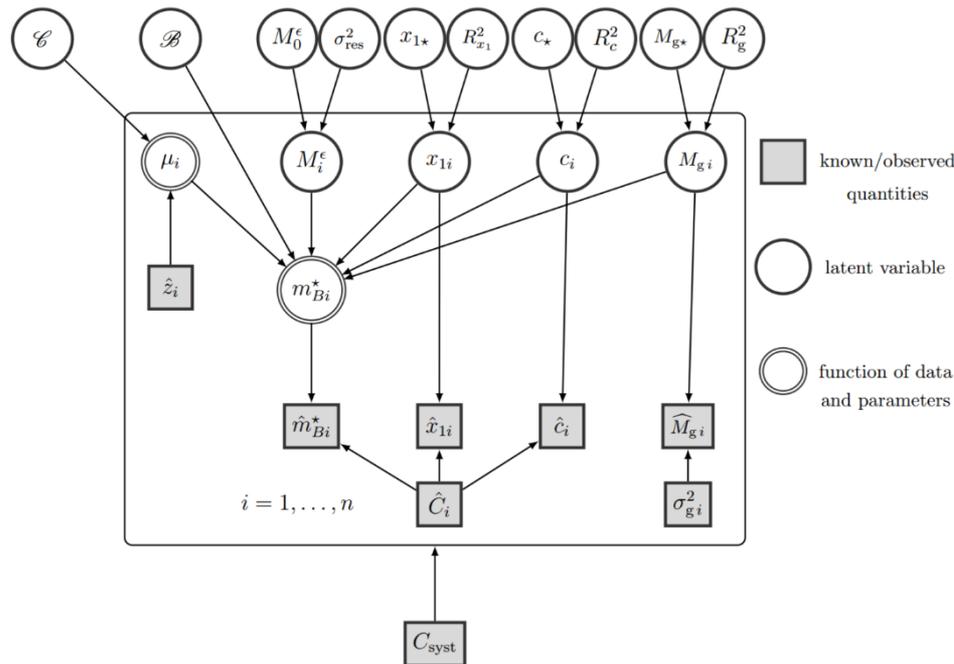
- Cosmic microwave background:

$$p(\{\Omega\}, \{C_\ell\}, s|d) \propto p(d|s) p(s|\{C_\ell\}) p(\{C_\ell\}|\{\Omega\}) p(\{\Omega\})$$

- Large-scale structure:

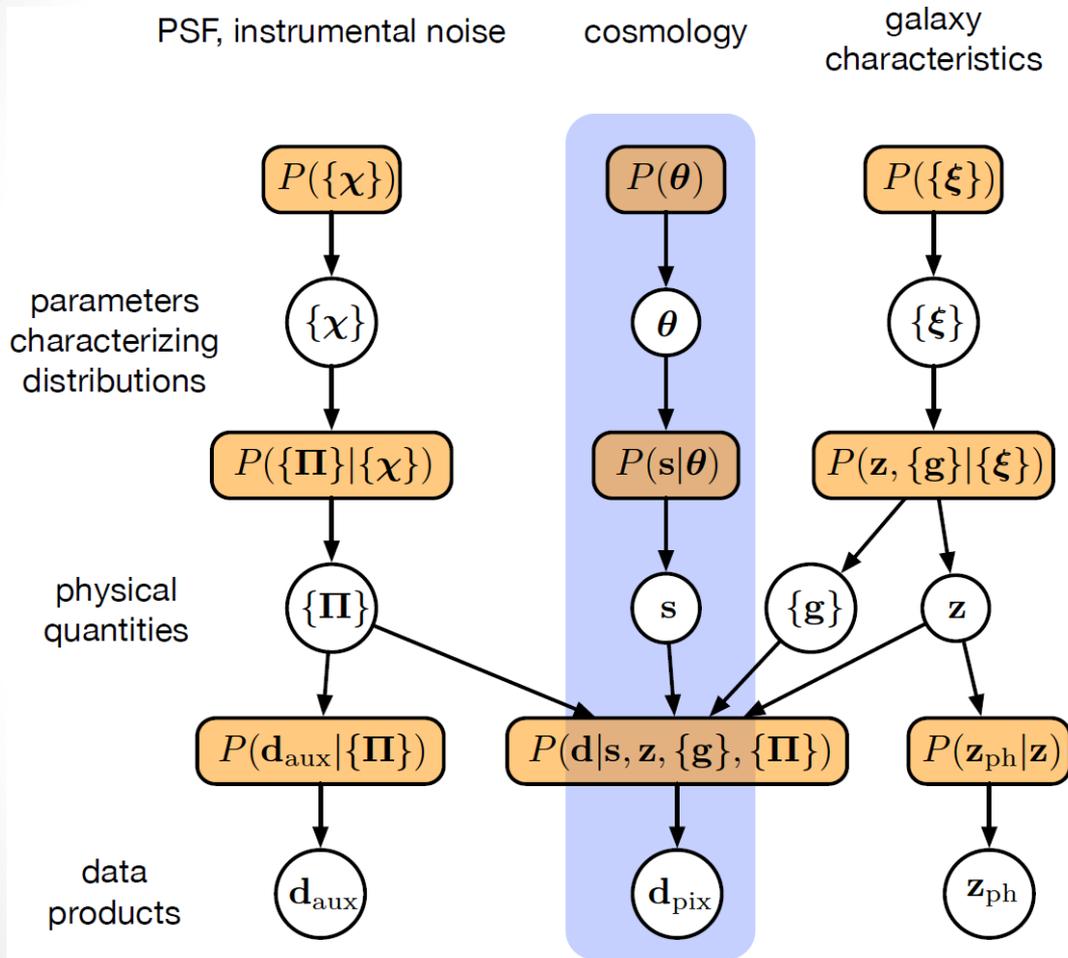
$$p(\{\Omega\}, \phi, g|d) \propto p(d|g) p(g|\phi) p(\phi|\{\Omega\}) p(\{\Omega\})$$

# BHM example: supernovae (BAHAMAS)



Parameter	Notation and Prior Distribution
Cosmological parameters	
Matter density parameter	$\Omega_m \sim \text{UNIFORM}(0, 2)$
Cosmological constant density parameter	$\Omega_\Lambda \sim \text{UNIFORM}(0, 2)$
Dark energy EOS	$w \sim \text{UNIFORM}(-2, 0)$
Hubble parameter	$H_0/\text{km/s/Mpc} = 67.3$
Covariates	
Coefficient of stretch covariate	$\alpha \sim \text{UNIFORM}(0, 1)$
Coefficient of color covariate	$\beta$ (or $\beta_0$ ) $\sim \text{UNIFORM}(0, 4)$
Coefficient of interaction of color correction and $z$	$\beta_1 \sim \text{UNIFORM}(-4, 4)$
Jump in coefficient of color covariate	$\Delta\beta \sim \text{UNIFORM}(-1.5, 1.5)$
Redshift of jump in color covariate	$z_t \sim \text{UNIFORM}(0.2, 1)$
Coefficient of host galaxy mass covariate	$\gamma \sim \text{UNIFORM}(-4, 4)$
Population-level distributions	
Mean of absolute magnitude	$M_0^c \sim \mathcal{N}(-19.3, 2^2)$
Residual scatter after corrections	$\sigma_{\text{res}}^2 \sim \text{INV GAMMA}(0.003, 0.003)$
Mean of absolute magnitude, low galaxy mass	$M_0^{\text{lo}} \sim \mathcal{N}(-19.3, 2^2)$
SD of absolute magnitude, low galaxy mass	$\sigma_{\text{res}}^{\text{lo}2} \sim \text{INV GAMMA}(0.003, 0.003)$
Mean of absolute magnitude, high galaxy mass	$M_0^{\text{hi}} \sim \mathcal{N}(-19.3, 2^2)$
SD of absolute magnitude, high galaxy mass	$\sigma_{\text{res}}^{\text{hi}2} \sim \text{INV GAMMA}(0.003, 0.003)$
Mean of stretch	$x_{1*} \sim \mathcal{N}(0, 10^2)$
SD of stretch	$R_{x_1} \sim \text{LOG UNIFORM}(-5, 2)$
Mean of color	$c_* \sim \mathcal{N}(0, 1^2)$
SD of color	$R_c \sim \text{LOG UNIFORM}(-5, 2)$
Mean of host galaxy mass	$M_{g*} \sim \mathcal{N}(10, 100^2)$
SD of host galaxy mass	$R_g \sim \text{LOG UNIFORM}(-5, 2)$

# BHM example: weak lensing



Can include:

- Mask
- Intrinsic alignments
- Baryon feedback
- Shape measurement
- Photometric redshifts

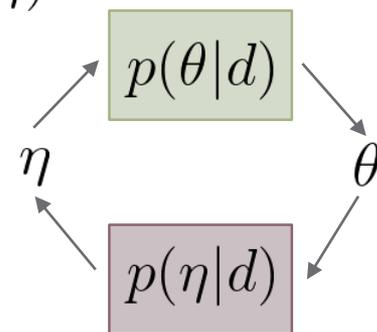
# Empirical Bayes

an alternative to maximum entropy for choosing priors

$$p(\theta|d) \propto p(d|\theta) \overset{\text{prior}}{p(\theta|\eta)} \overset{\text{hyperprior}}{p(\eta)}$$

$$\underline{p(\theta|d)} = \int p(\theta|\eta, d) p(\eta|d) d\eta = \int \frac{p(d|\theta) p(\theta|\eta)}{p(d|\eta)} \underline{p(\eta|d)} d\eta$$

$$\underline{p(\eta|d)} = \int p(\eta|\theta) \underline{p(\theta|d)} d\theta$$



➡ Iterative scheme (“Gibbs” sampler)

- **Empirical Bayes** is a truncation of this scheme after a few steps (often just one).

- Particular case:  $p(\eta|d) \approx \delta_D(\eta - \eta^*(d)) \Rightarrow \underline{p(\theta|d)} \approx \frac{p(d|\theta) p(\theta|\eta^*)}{p(d|\eta^*)}$

➡ the **Expectation-Maximization** (EM) algorithm (machine learning, data mining).