

Imperial Centre for Inference & Cosmology

The github repository

<u>https://github.com/florent-leclercq/Bayes_InfoTheory</u>

< > C BB GitHub, Inc. [US] github.com/florent-leclercq/Bayes_InfoTheory					
Search or jump to	Pull requests Issues Marketp	lace Explore			
	florent-leclercq / Bayes_InfoTheory Code Stores 0 The pull request	s o IIII Projects o IIII Wiki	Unwatch -	2 🛧 Star 0 😵 Fork 0	
	Lectures on Bayesian statistics and information theory Manage topics				
	© 39 commits	្រំ 1 branch	♡ 0 releases	L contributor	
	Branch: master - New pull request		Create new file Upload files	Find File Clone or download -	
	florent-leclercq updated notebooks, corrected error in ABC discrepancy		Clone with SSH	(?) Use HTTPS	
	🖿 data	added machine learning	Use an SSF Searc	h Copy Send to My Flow	
	.gitignore	updated gitignore	git@github.com:	florent-leclercq/Bayes	
	B ABC_discrepancy_effective_likelihood.ip updated notebooks, corrected error in ABC discre		BC discrep	Download ZID	
	ABC_rejection.ipynb	updated ABC notebooks		u your ugo	
	ABC synthetic likelihood.ipynb	updated notebooks, corrected error in A	BC discrepancy	a year ago	

git clone https://github.com/florent-leclercq/Bayes_InfoTheory.git (or with SSH)

 Course website: <u>http://florent-leclercq.eu/teaching.php</u> (this lecture is actually part of a series of 3)

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.

Outline

- Probability theory and Bayesian statistics: reminders
- Ignorance priors and the maximum entropy principle
- Gaussian random fields (and a digression on non-Gaussianity)
- Bayesian signal processing and reconstruction:
 - Bayesian de-noising
 - Bayesian de-blending
- Bayesian decision theory and Bayesian experimental design
- Bayesian networks, Bayesian hierarchical models and Empirical Bayes



Reminders

- 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)



• Bayesian model comparison:

$$\mathcal{B}_{12} = \frac{p(d|\mathcal{M}_1)}{p(d|\mathcal{M}_2)}$$

Ignorance priors and the maximum entropy principle



Notebook 1: <u>https://github.com/florent-</u> <u>leclercq/Bayes_InfoTheory/blob/master/LighthouseProblem.ipynb</u> Notebook 2: <u>https://github.com/florent-</u> <u>leclercq/Bayes_InfoTheory/blob/master/MaximumEntropy.ipynb</u>

Ignorance priors, functional equations and transformation groups

Ignorance priors: impose an invariant state of knowledge according to some transformation:

p(T(x))dT(x) = p(x)dx

- Simplest case: symmetry under the exchange of two models \mathcal{M}_1 and \mathcal{M}_2 :
 - $p(\mathcal{M}_1) = p(\mathcal{M}_2)$ $p(\mathcal{M}_1) + p(\mathcal{M}_2) = 1 \qquad \implies p(\mathcal{M}_1) = p(\mathcal{M}_2) = \frac{1}{2} \qquad \mathbb{Z}_2\text{-symmetry}$
- "Location parameter": $T(x) = x + s \quad \forall s$ dT = dx $p(x) = p(x + s) \quad \forall s$ p(x) = C Flat prior U(1)-symmetry
- "Scale parameter": $T(x) = ax \quad \forall a$ $dT = a \, dx$ $p(x) = a \, p(ax) \quad \forall a \quad \Longrightarrow \quad p(x) = C/a$ Jeffreys prior U(1)-symmetry
- General case: specify a group of transformations and solve the functional equation.



Maximum ignorance for one variable is generally not the same thing as maximum ignorance for a non-linear function of that variable.

The maximum entropy principle

- Maximising the entropy = a general method to select priors while accounting for:
 - indifference about states of equal knowledge
 - relevant prior information
- What should H[p] be for a source of information producing N finite "words" with probabilities p_n ?
- Desiderata:
 - If all words are equiprobable ($\forall n, \ p_n = \frac{1}{N}$), H[p] must grow with N
 - If words are generated in two steps (1- choosing a subset of words; 2- choosing a word in this subset), then the entropy is the sum of the entropy assigned to each step

$$\Longrightarrow$$
 Theorem (Shannon): $H[p] \propto -\sum p_n \log_2 p_n$

n

Information theory

Pictures taken at the Science Museum in South Kensington...



Claude Shannon 1916–2001

Claude Shannon was a mathematician and electrical engineer. Whilst working at Bell Laboratories he discussed artificial intelligence and the far-reaching possibilities of digital computers with Alan Turing. For the first time, Shannon gave 'information' a mathematical value. He inspired a new academic field, information theory, which underpinned developments in modern communication.

The Bell System Technical Journal

Vol. XXVII July, 1948 No. 3

A Mathematical Theory of Communication By C. E. SHANNON

INTRODUCTION

THE recent development of various methods of modulation such as PCM institution of PCM which exchange bandwidth for signal-to-noise ratio has in rensified the interest in a general theory of communication. A hasis for such a theory is contained in the important papers of Nyupit's and Hartieg' on this subject. In the present paper we will extend the theory to include a number of new factors, in particular the effect of noise in the channel, and the savings possible due to the statistical structure of the original message and due to the mature of the final destination of the information.

The fundamental problem of communication is that of reproducing at one point either exactly or approximately a mesage selected at another point. Frequently the mesages have meaning, that is they refer to a rate correlated according to some system with certain physical or conceptual entities. These semantic aspects of communication are irrelevant to the engineering problem. The significant aspect is that the setual using a sem selected free as set of possible selection, not just the one which will actually be choose niscre this is unknown at the time of besign.

If the number of messages in the set is finite then this number or any monotonic function of this number can be scheme from the next, all choices formation produced when one remained on the next and the information produced when one remained on the Narrely the mest natural ductic is the hyperbary. As found in the next set of the next set being equally likely. As found in the next set of the next set message and when we have a continuous range of messages, we will in all cases use an essentially logarithmic measure.

The logarithmic measure is more convenient for various reasons: I. It is practically more useful. Parameters of engineering importance

¹ Nyrquist, H., "Certain Factors Affecting Telegraph Speed," Bill System Technical Jeuroid, April 1924, p. 324: "Certain Tepics in Telegraph Transmission Theory," A. J. E. E. Trans, v. 47, April 1928, p. 515. "Hartley, R. V. L., "Transmission of Information," Bill System Technical Journal, July 1928, p. 535.

The loaded dice

- For a fair dice, $p_n = \frac{1}{6} \quad \forall n \in [1, 6]$: the principle of indifference was enough.
- Now let's say that the average value after many trials is not
 3.5 but 4. What is the probability law in this case?
- We want to maximise H[p] given two constraints:

$$< n >_p = \sum_{n=1}^{6} n p_n = 4$$
 and $\sum_{n=1}^{6} p_n = 1$

- (1) brute force way:
 - get p_5 and p_6 as a function of p_1 , p_2 , p_3 , p_4
 - express $H[p] = \sum_{n=1}^{\infty} p_n \ln p_n$ as a function of p_1 , p_2 , p_3 , p_4
 - differentiate and solve $\frac{\partial H}{\partial p_n} = 0$ for $n \in [1, 4]$

The loaded dice

 (2) a more elegant solution which does not break the symmetry: the method of Lagrange multipliers

• Lagrangian:
$$\mathcal{L}[\{p_n\}, \lambda, \mu] = -\sum_{n=1}^{6} p_n \ln p_n - \lambda \left(\sum_{n=1}^{6} np_n - 4\right) - \mu \left(\sum_{n=1}^{6} p_n - 1\right)$$

• Our two constraints are $\frac{\partial \mathcal{L}}{\partial n} = 0$ and $\frac{\partial \mathcal{L}}{\partial n} = 0$

$$\frac{\partial \mathcal{L}}{\partial p_n} = 0 \quad \text{gives} \quad -1 - \ln p_n - \lambda n - \mu = 0$$
$$p_n = \frac{e^{-\lambda n}}{Z} \quad \text{with} \ \ln Z \equiv 1 + \mu$$

- The normalisation constraint fixes $Z = \sum_{n=1}^{\circ} e^{-\lambda n} = \frac{1 e^{-\delta \lambda}}{e^{\lambda} 1}$
- The constraint on the mean is obtained by noting that

$$-\frac{\mathrm{d}\ln Z}{\mathrm{d}\lambda} = -\frac{1}{Z}\frac{\mathrm{d}Z}{\mathrm{d}\lambda} = \sum_{n=1}^{6} n\frac{\mathrm{e}^{-\lambda n}}{Z} = \sum_{n=1}^{6} n p_n = 4$$

This gives an equation for e^{λ} : $e^{\lambda}/(e^{\lambda}-1) - 6/(e^{6\lambda}-1) = 4$

The loaded dice

Notebook 2: <u>https://github.com/florent-</u>

leclercq/Bayes_InfoTheory/blob/master/MaximumEntropy.ipynb



- This is an example of probability theory beyond Bayesian statistics: we obtained a numerical probability assignment, conditional on some observations, without using Bayes' theorem.
- Thermodynamics analogy:
 - Fair dice = microcanonical ensemble: $p_n = \frac{1}{N}$
 - Loaded dice = canonical ensemble:

 $p_n = \frac{e^{-\beta E_n}}{Z}$ $\beta \equiv \frac{1}{k_B T}$ E_n = energy of different states Z = partition function \equiv evidence in Bayesian statistics

Gaussian random fields

Notebook 3: <u>https://github.com/florent-</u> leclercq/Bayes_InfoTheory/blob/master/GRF_and_fNL.ipynb

Bayesian signal processing and reconstruction

Notebook 4: <u>https://github.com/florent-</u> <u>leclercq/Bayes_InfoTheory/blob/master/WienerFilter_denoising.ipynb</u> Notebook 4bis: <u>https://github.com/florent-</u> <u>leclercq/Bayes_InfoTheory/blob/master/WienerFilter_denoising_CMB.ipynb</u> Notebook 5: <u>https://github.com/florent-</u> <u>leclercq/Bayes_InfoTheory/blob/master/WienerFilter_deblending.ipynb</u>

Gaussian random fields

• Definition: any random vector \boldsymbol{x} with pdf

$$p(x|\mu, C) = \frac{1}{\sqrt{|2\pi C|}} \exp\left[-\frac{1}{2}(x-\mu)^{\mathsf{T}}C^{-1}(x-\mu)\right]$$

 $-2\ln p(x|\mu, C) = (x - \mu)^{\mathsf{T}} C^{-1}(x - \mu) + \ln |2\pi C|$... and that's it.

- Generating a Gaussian random field:
 - Draw a white noise vector ξ (uncorrelated unit-Gaussian variables)
 - Find the matrix square-root of $C: \sqrt{C}$ (any such matrix works)

• Compute
$$x=\sqrt{C}\xi+\mu$$

Notebook 3: <u>https://github.com/florent-</u> leclercq/Bayes_InfoTheory/blob/master/GRF_and_fNL.ipynb



Notebook 3: <u>https://github.com/florent-</u> leclercq/Bayes_InfoTheory/blob/master/GRF_and_fNL.ipynb

$$C_{ij} = \exp\left(-\frac{|i-j|}{20}\right)$$

covariance matrix





1

Notebook 3: <u>https://github.com/florent-</u> leclercq/Bayes_InfoTheory/blob/master/GRF_and_fNL.ipynb

$$C_{ij} = \frac{1}{(1+|i-j|/2)^2}$$

covariance matrix





Notebook 3: https://github.com/florentleclercq/Bayes InfoTheory/blob/master/GRF and fNL.ipynb

$$C_{ii} = \begin{cases} 1 & \text{if } i < N/2 \\ 100 & \text{otherwise} \end{cases}$$
$$C_{ij} = 0 \quad \text{for} \quad i \neq j$$
covariance matrix









Notebook 3: <u>https://github.com/florent-</u> leclercq/Bayes_InfoTheory/blob/master/GRF_and_fNL.ipynb



Histograms of Gaussian random fields are not always Gaussian!

Example of a non-Gaussian signal

Notebook 3: <u>https://github.com/florent-</u>

leclercq/Bayes_InfoTheory/blob/master/GRF_and_fNL.ipynb

$$s = \Phi + f_{
m NL} \Phi^2$$
 where Φ is a GRF.

In cosmology, this is called "local-type" non-Gaussianity signal histogram



The one-point pdf is skewed.

Gaussian random fields: marginals and conditionals

- We work with a "joint" Gaussian random field $\begin{pmatrix} x \\ y \end{pmatrix}$
- Marginals:

Mean:
$$\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}$$
 Covariance: $C = \begin{pmatrix} C_{xx} & C_{xy} \\ C_{yx} & C_{yy} \end{pmatrix}$

The marginal means and covariances are just the corresponding parts of the joint mean and covariance.

• Conditionals:

Mean: $\mu_{x|y} = \mu_x + C_{xy}C_{yy}^{-1}(y - \mu_y)$ Covariance: $C_{x|y} = C_{xx} - C_{xy}C_{yy}^{-1}C_{yx}$

- Data model: d = s + n where $\binom{s}{d}$ is jointly Gaussian.
- Solution:

$$\mu_{s|d} = \mu_s + C_{sd} C_{dd}^{-1} (d - \mu_d)$$
$$C_{s|d} = C_{ss} - C_{sd} C_{dd}^{-1} C_{ds}$$

• Notations: $C_{ss} \equiv S$ and $C_{nn} \equiv N$.

• Assumption:
$$C_{sn} = C_{ns} = 0$$
. Then
 $C_{dd} = S + N$
 $C_{sd} = C_{ss} + C_{sn} = C_{ss} = S$

• Final expressions:

$$\mu_{s|d} = \mu_s + S(S+N)^{-1}(d-\mu_d) = \mu_s + (S^{-1}+N^{-1})^{-1}N^{-1}(d-\mu_d)$$
$$C_{s|d} = S - S(S+N)^{-1}S = (S^{-1}+N^{-1})^{-1}$$

Notebook 4: <u>https://github.com/florent-</u> leclercq/Bayes InfoTheory/blob/master/WienerFilter denoising.ipynb

Setup signal and noise covariance matrices







Notebook 4: <u>https://github.com/florent-</u> leclercq/Bayes InfoTheory/blob/master/WienerFilter denoising.ipynb

Generate mock data



$$d = s + n$$

Notebook 4: <u>https://github.com/florent-</u> leclercq/Bayes InfoTheory/blob/master/WienerFilter denoising.ipynb

Perform Wiener filtering

$$C_{s|d} = (S^{-1} + N^{-1})^{-1}$$





Notebook 4: https://github.com/florentleclercq/Bayes InfoTheory/blob/master/WienerFilter denoising.ipynb

Draw constrained realisations

$$s_{\rm sim} = \mu_{s|d} + \sqrt{C_{s|d}}\,\xi$$







Notebook 5: https://github.com/florent-

leclercq/Bayes_InfoTheory/blob/master/WienerFilter_deblending.ipynb

• Data model:
$$d = \begin{pmatrix} x_1 \\ x_2 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} n_1 \\ n_2 \\ n_3 \end{pmatrix}$$

• Assumptions: $C_{x_1d} = \begin{pmatrix} C_{x_1x_1} & C_{x_1x_1} & 0 \end{pmatrix}$ $C_{x_2d} = \begin{pmatrix} 0 & C_{x_2x_2} & C_{x_2x_2} \end{pmatrix}$ $C_{nn} = \begin{pmatrix} C_{n_1n_1} & 0 & 0 \\ 0 & C_{n_2n_2} & 0 \\ 0 & 0 & C_{n_3n_3} \end{pmatrix}$

$$C_{dd} = \begin{pmatrix} C_{x_1x_1} + C_{n_1n_1} & C_{x_1x_1} & 0\\ C_{x_1x_1} & C_{x_1x_1} + C_{x_2x_2} + C_{n_2n_2} & C_{x_2x_2}\\ 0 & C_{x_2x_2} & C_{x_2x_2} + C_{n_3n_3} \end{pmatrix}$$

- Solution:
 - $\mu_{x_1|d} = C_{x_1d}C_{dd}^{-1}d \qquad \qquad \mu_{x_2|d} = C_{x_2d}C_{dd}^{-1}d$ $C_{x_1|d} = C_{x_1x_1} - C_{x_1d}C_{dd}^{-1}C_{dx_1} \qquad \qquad C_{x_2|d} = C_{x_2x_2} - C_{x_2d}C_{dd}^{-1}C_{dx_2}$

Bayesian decision theory

Notebook 6: <u>https://github.com/florent-</u> leclercq/Bayes InfoTheory/blob/master/DecisionTheory.ipynb

Bayesian experimental design

Bayesian decision theory

- Bayesian decision theory is **optimal decision making**, given a set of possible actions and uncertain beliefs, encoded in some pdf $p(\theta|I)$ (usually the posterior of Bayesian inference)
- Notations: $\{\theta\}$ = set of features (observable variables) $\{a\}$ = set of actions
- Expected utility hypothesis: given a set of gain functions $G(a|\theta)$, the optimal decision rule is to take the action that maximises the expected utility U(a|I), defined by

$$U(a|I) \equiv \langle G(a|\theta) \rangle_{p(\theta|I)} = \int G(a|\theta)p(\theta|I) \,\mathrm{d}\theta$$

Take action $a^* = \operatorname{argmax}_a U(a|I)$

Example: Bayesian alerts

- We look for an event *E*. We have access to $p(E|I) = p(\bar{E}|I) = 1 p(E|I)$
- There are two possible actions:
 - a_1 = raise the alert



Structures in the cosmic web



FL, Jasche & Wandelt 2015b, arXiv:1503.00730

A decision rule for structure classification

• Space of "input features":

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

• 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 3

$$U(a_j(\vec{x}_k)|d) = \sum_{i=0}^{o} G(a_j|\mathbf{T}_i) \mathcal{P}(\mathbf{T}_i(\vec{x}_k)|d)$$

How to write down the gain functions?

Gambling with the Universe



3.83

41.67

7.08

voids

sheets

□ filaments

Without data, the expected utility is

 $U(a_j) = 1 - \alpha \quad \text{if } j \neq 1$ "Playing the game" "Not playing the game"

• With $\alpha = 1$, it's a *fair game* \Longrightarrow always play "speculative map" of the LSS

 $U(a_{-1}) = 0$

• Values $\alpha > 1$ represent an *aversion for risk* increasingly "conservative maps" of the LSS



Bayesian networks

Bayesian hierarchical models

and Empirical Bayes

Bayesian networks



Bayesian networks are probabilistic graphical models consisting of:

- A directed acyclic graph (DAG)
- At each node, conditional probabilities distributions

Bayesian networks



p(C, M, E, G) = p(C) p(E|C) p(M|C, E) p(G|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; 0 < p(B) < 1; p(A|B) = p(A) $\widehat{p}(A|B,C) < p(A|C)$

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

Malmquist bias

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



Bayesian hierarchical models

- Simple inference: $p(\theta|d) \propto p(d|\theta) p(\theta)$ • Adaptive prior: $p(\theta|d) \propto p(d|\theta) p(\theta|\eta) p(\eta)$
- ... 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	$\Omega_{\rm m} \sim { m Uniform}(0,2)$			
Cosmological constant density parameter	$\Omega_{\Lambda} \sim \text{Uniform}(0,2)$			
Dark energy EOS	$w \sim \text{Uniform}(-2,0)$			
Hubble parameter	$H_0/\mathrm{km/s/Mpc} = 67.3$			
Covariates				
Coefficient of stretch covariate	$\alpha \sim \text{Uniform}(0,1)$			
Coefficient of color covariate	β (or β_0) ~ Uniform(0, 4)			
Coefficient of interaction of color correction and \boldsymbol{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	$M_0^\epsilon \sim \mathcal{N}(-19.3,2^2)$			
Residual scatter after corrections	$\sigma^2_{\rm res} \sim {\rm InvGamma}(0.003, 0.003)$			
Mean of absolute magnitude, low galaxy mass	$M_0^{\rm lo} \sim \mathcal{N}(-19.3, 2^2)$			
SD of absolute magnitude, low galaxy mass	${\sigma_{\mathrm{res}}^{\mathrm{lo}}}^2 \sim \mathrm{InvGamma}(0.003, 0.003)$			
Mean of absolute magnitude, high galaxy mass	$M_0^{\rm hi} \sim \mathcal{N}(-19.3, 2^2)$			
SD of absolute magnitude, high galaxy mass	${\sigma_{\mathrm{res}}^{\mathrm{hi}}}^2 \sim \mathrm{InvGamma}(0.003, 0.003)$			
Mean of stretch	$x_{1\star} \sim \mathcal{N}(0, 10^2)$			
SD of stretch	$R_{x_1} \sim \text{LogUniform}(-5,2)$			
Mean of color	$c_{\star} \sim \mathcal{N}(0, 1^2)$			
SD of color	$R_c \sim \text{LogUniform}(-5,2)$			
Mean of host galaxy mass	$M_{\rm g\star} \sim \mathcal{N}(10, 100^2)$			
SD of host galaxy mass	$R_{\rm g} \sim {\rm LogUniform}(-5,2)$			

BHM example: weak lensing

PSF, instrumental noise

cosmology

galaxy characteristics



Can include:

Mask

Intrinsic alignments Baryon feedback Shape measurement Photometric redshifts

- Empirical Bayes is a truncation of this scheme after a few steps (often just one).
- Particular case: $p(\eta|d) \approx \delta_{\mathrm{D}}(\eta \eta^{\star}(d)) \implies \underline{p(\theta|d)} \approx \frac{p(d|\theta) p(\theta|\eta^{\star})}{p(d|\eta^{\star})}$

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