

# Data Modeling

## Main References:

- ① Selected Topic 10
- ② " " 11
- ③ " " 12-A
- ④ " " 13-A
- ⑤ 14031001
- ⑥ 14030216
- ⑦ IPM Workshop 2023

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِیْمِ



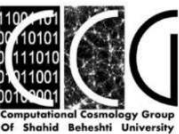
### Most relevant References: Books

## Review on Data Analysis Methods & Cosmological Simulations

Seyed Mohammad Sadegh Movahed

Physics department, Shahid Beheshti University  
Computational and Cosmological group (CCG-SBU)  
[www.smovahed.ir](http://www.smovahed.ir)

Workshop on Computational Cosmology:  
From Theory to Observation  
1-2 August 2023



- 1) Modern Cosmology, S. Donelson.
- 2) The Cosmic Microwave Background, R. Burrer
- 3) Physical foundation of Cosmology, V. Mukhanov
- 4) Cosmology, S. Weinberg
- 5) Neutrino Cosmology, J. Lesgourgues et al.
- 6) Galaxy Formation and Evolution, H. Mo, F.V. den Bosch and S. White
- 7) Statistics of Galaxy Distribution, V. Martinez, E. Saar
- 8) Cosmological Physics, J.A. Peacock
- 9) Topological Complexity of Smooth Random Functions, Adler, Robert, Taylor, Jonathan E., Springer, 2009.
- 10) Geometry, Topology and Physics, Mikio Nakahara, 1990.
- 11) Analysis and Data-Based Reconstruction of Complex Nonlinear Dynamical Systems, Using the Methods of Stochastic Processes, M. Reza Rahimi Tabar, Springer, 2019
- 12) Zomorodian, Afra. "Topological data analysis." Advances in applied and computational topology 70 (2012): 1-39.
- 13) Edelsbrunner, Herbert, and John Harer. Computational topology: an introduction. American Mathematical Soc., 2010.

### Most relevant References: Previous activities

### Most relevant References: Papers and Lectures

#### Part A: Previous workshops

- 1) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/106-data-modeling-workshop>
- 2) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/131-workshop-on-observational-data-analysis-in-cosmology-kurdistan-1398>
- 3) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/144-school-and-workshop-on-statistical-analysis-of-cosmic-fields-1400>
- 4) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/154-school-and-workshop-on-topological-based-data-analysis-1401>

#### Part B: Some lectures:

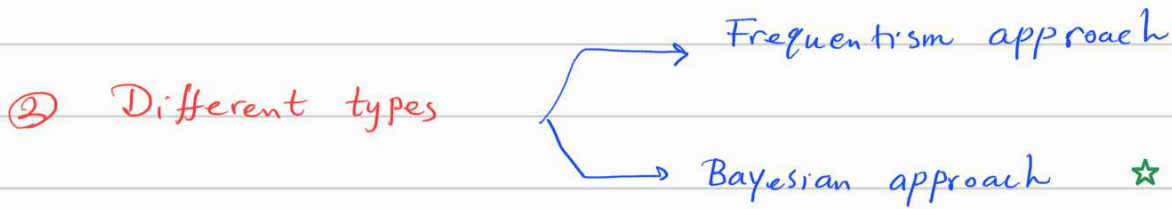
- 1) <http://facultymembers.sbu.ac.ir/movahed/index.php/courses/142-optimization-and-computational-approaches-fall-2021>
- 2) <http://facultymembers.sbu.ac.ir/movahed/index.php/courses/132-advanced-course-on-computational-physics>
- 3) <http://facultymembers.sbu.ac.ir/movahed/index.php/courses/139-stochastic-processes>

#### Part C: Some of my talks:

- 1) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/102-my-talk-at-cosmology-meeting-ipm-96>
- 2) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/148-my-talk-at-sharif-group-meeting-1401-2022>
- 3) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/145-my-talk-at-conference-on-gravity-and-cosmology-1400>
- 4) <http://facultymembers.sbu.ac.ir/movahed/index.php/talks-a-presentations/130-data-science>

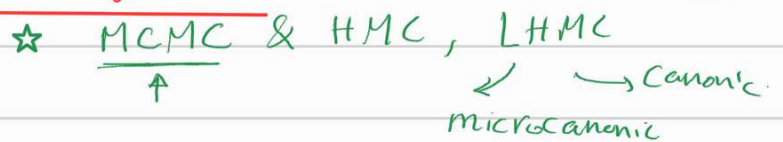
- 1) Matsubara, Takahiko. "Statistics of smoothed cosmic fields in perturbation theory. I. Formulation and useful formulae in second-order perturbation theory." The Astrophysical Journal 584.1 (2003).
- 2) Vafaei Sadr, A., and S. M. S. Movahed. "Clustering of local extrema in Planck CMB maps." Monthly Notices of the Royal Astronomical Society 503.1 (2021): 815-829.
- 3) Masoomy, H., et al. "Persistent homology of fractional Gaussian noise." Physical Review E 104.3 (2021): 034116.
- 4) Masoomy, H., S. Tajik, and S. M. S. Movahed. "Homology groups of embedded fractional Brownian motion." Physical Review E 106.6 (2022): 064115.
- 5) Lesgourgues, J. "Cosmological perturbations." Searching for New Physics at Small and Large Scales: TASI 2012. 2013. 29-97.
- 6) Mostaghel, Behrang, Hossein Moshafi, and S. M. S. Movahed. "Non-minimal derivative coupling scalar field and bulk viscous dark energy." The European Physical Journal C 77 (2017): 1-22.
- 7) Pranav, Pratyush, et al. "Topology and geometry of Gaussian random fields I: on Betti numbers, Euler characteristic, and Minkowski functionals." Monthly Notices of the Royal Astronomical Society 485.3 (2019): 4167-4208.
- 8) Pranav, Pratyush. "Topology and geometry of Gaussian random fields II: on critical points, excursion sets, and persistent homology." arXiv preprint arXiv:2109.08721 (2021).
- 9) Bardeen J. M., Bond J. R., Kaiser N., Szalay A. S., 1986, Astro-phys. J., 304, 15
- 10) Bond, J. R., et al., The Astrophysical Journal 379 (1991): 440-460.

# Different parts of Data modeling



③ Mathematical Description and foundation ☆

④ Analytical and computation Algorithms for Data modeling



⑤ Models Selection and Comparison ← Posterior ☆  
Evidence

نیل خودہ ہائرس

⑥ Goodness of fit & Confidence Interval on model's free ☆  
Parameters. کنڈیشنل ← انڈیپنڈنٹ

⑦ Bayesian Model averaging BMA

⑧ Fisher Forecast or Fisher Information Matrix ← arxiv.0906.0993  
Un-Biased Estimator  
arxiv.2409.13583

⑨ Simulation Based Inference (SBI), Field-Level Inference.

Neural Likelihood Estimation, Emulators. (Data-Based Inference)

Likelihood free Inference

$\{D\}$ : observation or Experiments  
 Synthetic data

آزمایش

- ۱) اندازه گیری (مستقیم یا شبیه سازی) (Measurement)
- ۲) برآورد خطا و ارزیابی انتشار خطا بر روی کمیت های ثانویه (Error estimation and error propagation)
- ۳) تدوین مدل با توجه به تابع مناسب (model selection) (Merit function)
- ۴) انتخاب بینش تعیین مقادیر آزاد مدل (Bayesian or Frequentist)
- ۵) استفاده از روشهای تعیین مقادیر کمیت های آزاد مدل و البته حوزه اعتبار آنها
- ۶) تعیین خوبی مدل (Goodness of fit)

$\{\theta\}$ : پارامترهای آزاد مدل

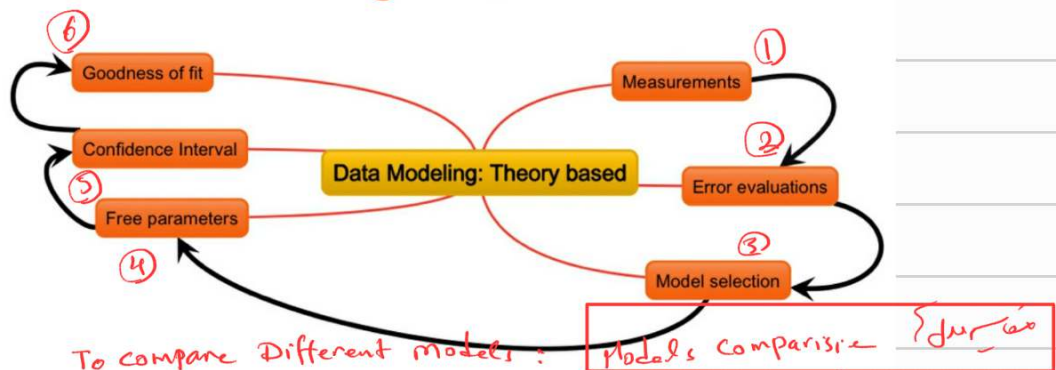
(Parameter estimation and confidence interval)

تعیین خوبی مدل (Goodness of fit)

① ☆ Theory-Based approach.

مدل خوش تعریف در نظر داریم

### General view on Theory based a challenge: a priori-that's it



Bayesian Model Averaging (BMA)

Mostaghel, Behrang, Hossein Moshafi, and S. M. S. Movahed. "Non-minimal derivative coupling scalar field and bulk viscous dark energy." The European Physical Journal C 77 (2017): 1-22.

arXiv: 2403.02120  
 arXiv: 2310.06747  
 Lecture Note: 14030230

Ex:  $U = \langle H \rangle = \frac{3K_B T}{2} + \frac{n}{2} \int d^3r u(r) g(r)$

↑  
Energy of system

↓  
Pair-Distribution  
Intra-atomic potential.

For Ideal Gas, we ignore this part

$$u(r) = 4\epsilon \left[ \left(\frac{\alpha}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right], \quad g(r) = e^{-\beta u(r)}$$

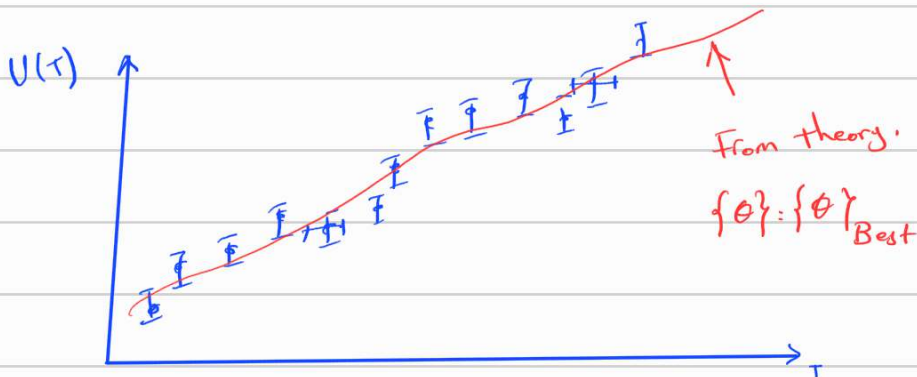
$$\{\theta\} = \{a, b, \epsilon, \alpha, \sigma, \beta\} \quad \boxed{M=6}$$

$\{\theta_i\}, i=1 \dots M \rightarrow$  # of model's Free Parameter

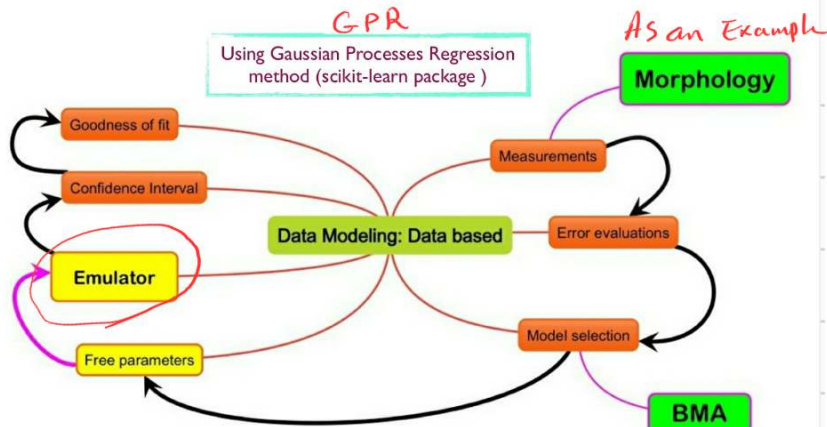
$\{\theta\} = ?$

For experimental Setup we can Record

T	U
$T_1$	$U_1$
$T_2$	$U_2$
$T_3$	$U_3$
$\vdots$	$\vdots$

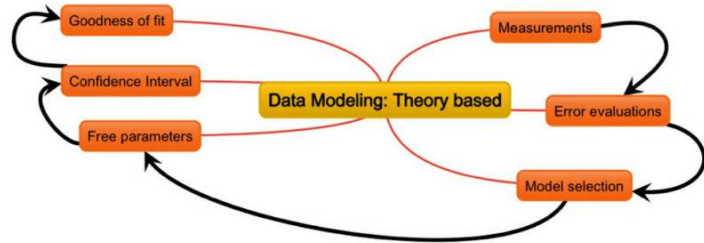


## General view on Data based



1) Pedregosa, Fabian, et al. "Scikit-learn: Machine learning in Python." the Journal of machine Learning research 12 (2011): 2825-2830.  
 2) Heydenreich, Sven, Benjamin Brück, and Joachim Harnois-Déraps. "Persistent homology in cosmic shear: constraining parameters with topological data analysis." Astronomy & Astrophysics 648 (2021): A94.

## General view on Theory based a challenge: a priori-that's it



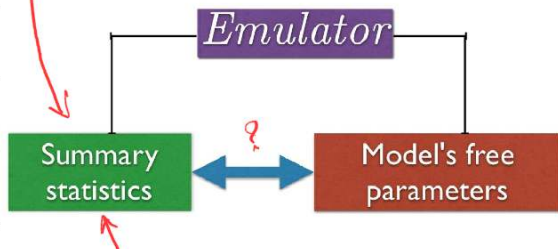
### Bayesian Model Averaging

Mostaghel, Behrang, Hossein Moshafi, and S. M. S. Movahed. "Non-minimal derivative coupling scalar field and bulk viscous dark energy." The European Physical Journal C 77 (2017): 1-22.

$U(T) = \dots ?$   
 $\{\theta\} = \{M, S, \dots\}$  : our interested parameters

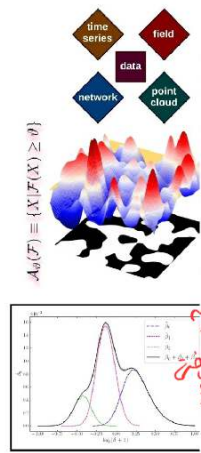
## From Morphology to Cosmological Inferences (Part 1: Emulator)

- 1) Data preparation (Acquisition, reduction, generation)
- 2) Tracer selection (field, excursion sets, critical sets, morphological measure)
- 3) Summary statistics



## From Morphology to Cosmological Inferences (Part 2: Simulation Based Inference)

- 1) Likelihood independent
- 2) No theoretical relation



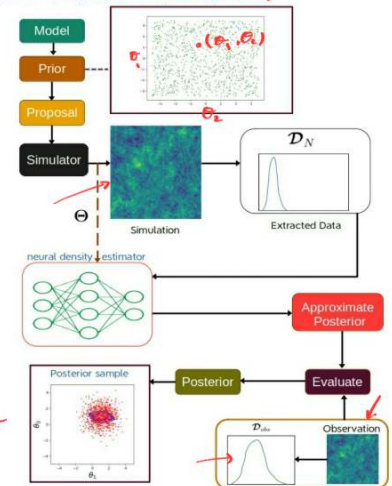
Forward modeling

According to our

physical First Principal  
and/or phenomenological  
theory Building

we are able to simulate  
our Cosmas

Credit: M.H. Jalali



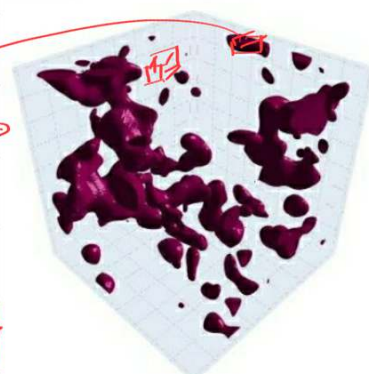
Summary Statistics  $\theta$

میدان حفر ذرات:  $\delta = \frac{n(r) - \langle n \rangle}{\langle n \rangle}$

$$\delta = \frac{n(r) - \langle n \rangle}{\langle n \rangle}$$

Density Contrast

$\delta(\vec{r}) \rightarrow \{ \text{size}, m, Q, T, P, \dots \}$   
 $\{\theta\}$



## ② Bayesian approach with Mathematical Foundation

### A history about Bayesian strategy

- By Thomas Bayes (1702-1761)
- Pierre Simon Laplace (1812)
- Fisher, Neyman, Wald,...
- Gelfand and Smith (1990)
- Now

### تعیین مقدار برای کمیتهای آزاد

~~Frequentist~~ رهیافت

(۱) در این روش داده ها و نتایج ریزحالت یک پیکربندی به حساب می آیند و داده ها تکرار پذیرند

(۲) کمیتهای مدل مجهول ولی ثابت هستند

(۳) هیچ گونه اطلاعی از مدل مورد استفاده نیست

(۴) سازوکاری برای رهایی از کمیتهای اضافی (Nuisance) ندارد

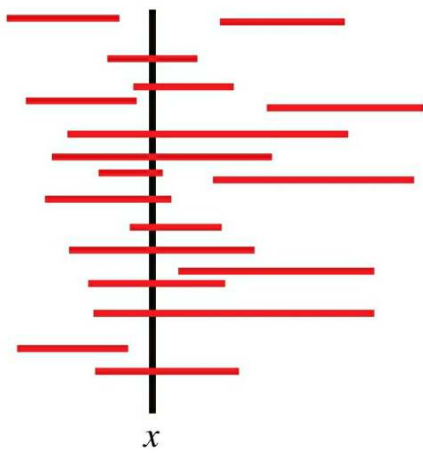
→ رهیافت Bayesian

(۱) در این روش داده ها و نتایج بخشی از یک آنسامبل هستند

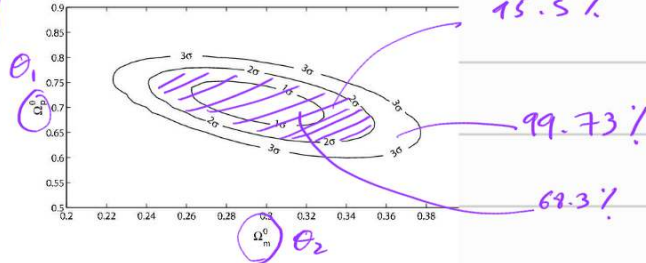
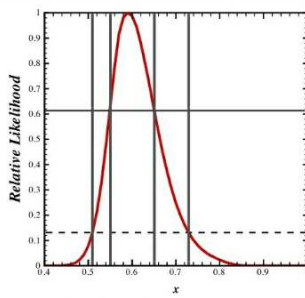
(۲) کمیتهای آزاد مدل مجهول هستند که ما تنها به صورت احتمالی می توانیم مقدار آنها را تعیین کنیم

Prior Informat  
(۳) هر گونه اطلاعات اولیه در این رهیافت قابل استفاده هستند

(۴) سازوکاری برای رهایی از کمیتهای اضافی (Nuisance) دارد



بر اساس نگرش frequentist  
برای مثال ۹۰ درصد  
مشاهدات به مقدار ثابت x  
منجر شده است



بر اساس نگرش Bayesian با  
احتمال ۹۰ درصد مقدار دقیق  
کمیت x در این بازه قرار می گیرد

Mathematics:

Bayes theorem

$$P(x|y)P(y) = P(y|x)P(x)$$

احتمال وقوع آدرس x بر سر یک نامه را با آدرس y

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

احتمال وقوع آدرس y بر سر یک نامه را با آدرس x

Suppose that

$$\{D\} = \{D_1, D_2, \dots, D_N\}$$

$$\{\theta\} = \{\theta_1, \theta_2, \dots, \theta_M\}$$

$$Y \left( \{\theta\} \right)$$

Thea

given a model.

$$P(\{\theta\}, M | \{D\}) = \frac{L(\{D\} | \{\theta\}, M) P(\{\theta\}, M)}{\int d\{\theta\} L(\{D\} | \{\theta\}, M) P(\{\theta\}, M)}$$

Posterior  
یعنی احتمال پستی

☆ Probability of having  $(\{\theta\}, M)$  given  $\{D\}$

Likelihood Distribution توزیع درست نمایی

Probability of having  $\{D\}$  given  $\{\theta\}, M$

→ Prior Distribution توزیع پیشینی

اطلاعات پیشین در مورد پارامترها

$\{\text{mass}, Q\}$   
↑

If we have no prior information regarding  $\{\theta\}$ ,

therefore we set  $P(\{\theta\}, M) = \text{cts}$

$$P(\{\theta\}, M | \{D\}) \propto L(\{D\} | \{\theta\}, M)$$

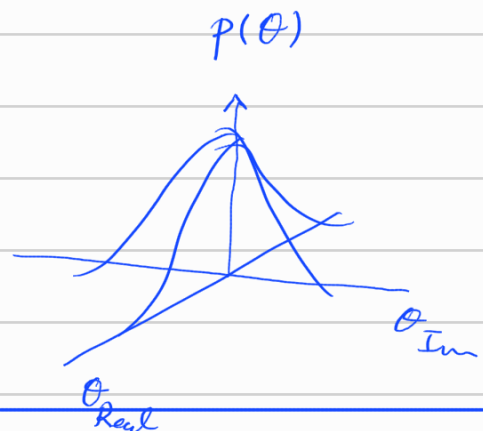
$$P(\{\theta\} = \{\theta\}_{\text{Best}}, M | \{D\}) \equiv \text{To be maximized}$$

Products

$$L(\{D\} | \{\theta\} = \{\theta\}_{\text{Best}}) \equiv \text{To be maximized}$$

$M_1$ ,  $\theta = \text{imaginary variable}$   
 $P(\theta)$

$$\theta = \theta_{\text{Real}} + i\theta_{\text{Im}}$$





# Mathematical Form of Posterior and Likelihood ?

$$Y_{\text{theory}}(\{x\}) \checkmark$$

$$\mathcal{L}(\{D\} | \{\theta\}) = ?$$

## پیشینه سازی تابع Posterior

$$D: \{x_i, y_i\} \quad i = 1, \dots, N$$

$$\Theta: \{\theta_j\} \quad Y(x; \{\Theta\}) = \sum_{k=1}^M \theta_k f_k(x)$$

$$p(\Theta | D) = \frac{p(D | \Theta)p(\Theta)}{\int d\Theta p(D | \Theta)p(\Theta)}$$

← Bayes theorem

if  $p(\Theta) = \text{cts}$

در صورتی که تابع prior وجود نداشته باشد در نتیجه تحلیل posterior به تحلیل Likelihood تبدیل می شود

$$p(\Theta | D): p(D | \Theta)$$

$$\mathcal{L}(D | \Theta) = \prod_{i=1}^N P_i = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{[y_i - Y(x_i; \{\Theta\})]^2}{2\sigma_i^2}\right)$$

(۱) با فرض اینکه اندازه گیری ها از یکدیگر مستقل باشند

$$\chi^2(\{\Theta\}) = \sum_{i=1}^N \frac{[y_i - Y(x_i; \{\Theta\})]^2}{\sigma_i^2} = \sum_{i=1}^N \frac{\left[y_i - \sum_{k=1}^M \theta_k f_k(x_i)\right]^2}{\sigma_i^2}$$

(۲) با توجه به قضیه حد مرکزی می توان در نظر گرفت که هر نقطه اندازه گیری حول مقدار بهینه اش به صورت گوسی توزیع شده است

$$\mathcal{L}(D | \Theta): e^{-\frac{\chi^2(\{\Theta\})}{2}}$$

پیشینه شدن Likelihood معادل با کمینه شدن  $\chi^2$  است

12

$$\mathcal{L}(\{D\} | \{\theta\}_M) = \frac{1}{\sqrt{(2\pi)^N \text{Det Cov}_D}}$$

$$e^{-\frac{\Delta^T \text{Cov}_D^{-1} \cdot \Delta}{2}}$$

Multi-variate Gaussian Function

$$\Delta^T = (y - Y_{th})^T_{1 \times N}$$

$$\Delta = (y - Y_{th})_{N \times 1}$$

$$\text{Cov}_D = \langle \Delta^T \Delta \rangle =$$

$$\begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} & \dots & \sigma_{1N} \\ \sigma_{21} & \sigma_2^2 & & & \\ \vdots & & \ddots & & \\ \sigma_{N1} & & & & \sigma_N^2 \end{bmatrix}_{N \times N}$$

$$\chi^2 = \Delta^T \text{Cov}_D^{-1} \cdot \Delta$$

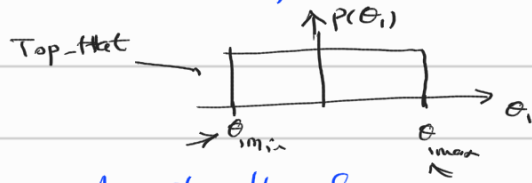
For Diagonal Cov  $\chi^2 = \sum_{i=1}^N \frac{[y_i - Y_{th}(\{\theta\}, x_i)]^2}{\sigma_i^2}$

ماتریس کوواریانس

③ Computational Method to minimize  $\chi^2$  or maximize

Prior = cts

Likelihood

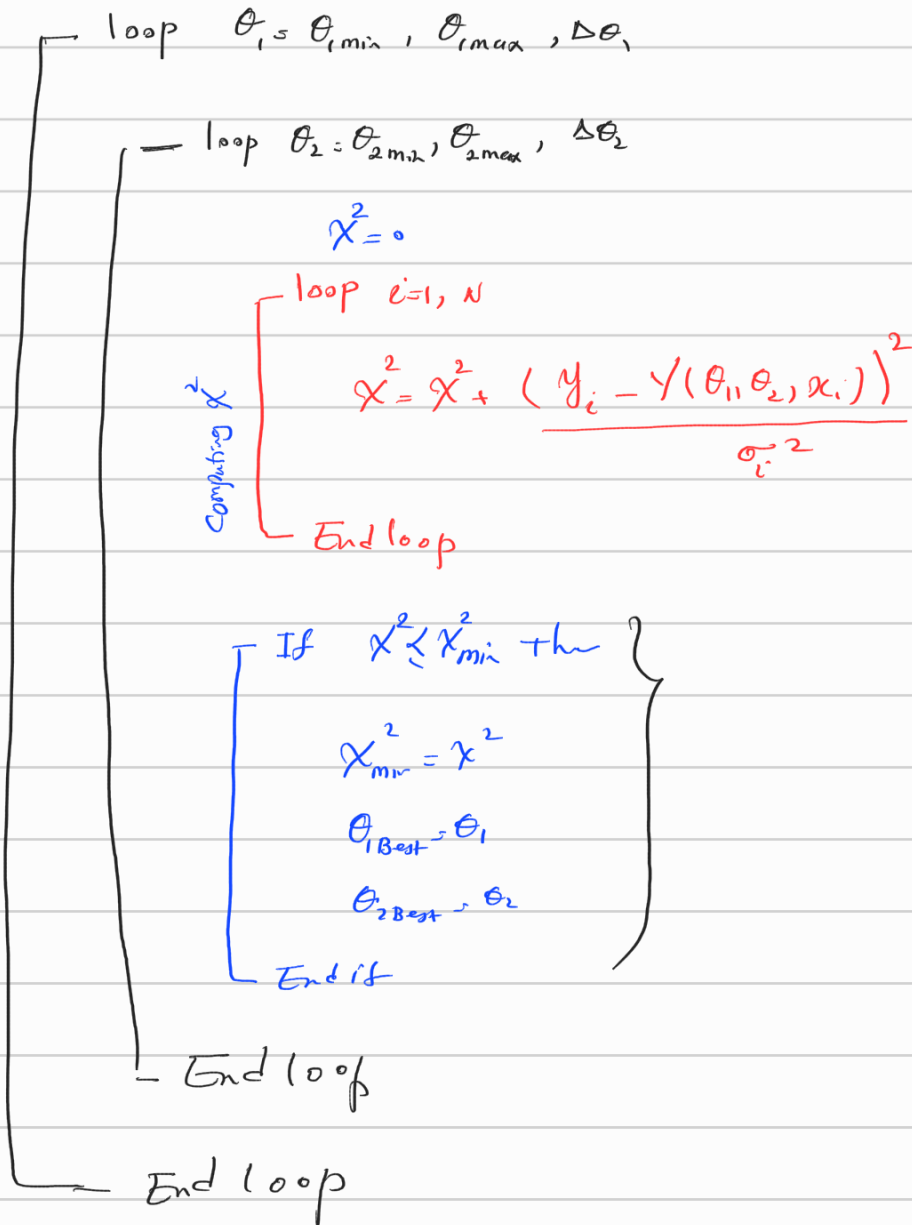


An algorithm for  
Deterministic scanning.

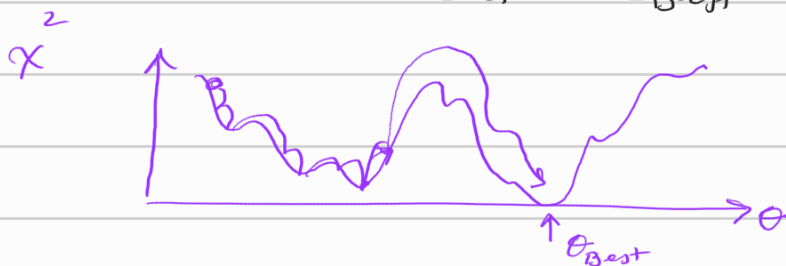
Ex 1:  $\{\theta\} = \{\theta_1, \theta_2\}$   $M=2$

Import  $\{D\} = \{(x_i, y_i, \sigma_i)\}$   $i=1, \dots, N$

$$\chi_{min}^2 = 1e10$$

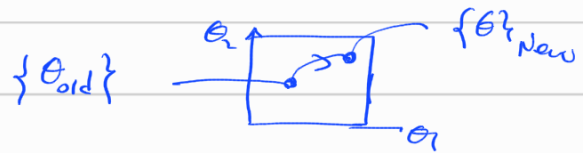


Write  $\theta_{1Best}$ ,  $\theta_{2Best}$ ,  $\chi_{min}^2$



Ex 2

# MCMC Algorithm



Import Data  $\{(x_i, y_i, \sigma_i)\}$   $i=1, \dots, N$

Select  $\{\theta\}_{old}$  Randomly

$$\theta_1 = \theta_{1min} + R_1 (\theta_{1max} - \theta_{1min})$$

$$\theta_2 = \theta_{2min} + R_2 (\theta_{2max} - \theta_{2min})$$

loop  $i=1, N$

$$\chi^2_{old} = \chi^2_{old} + \frac{[y_i - Y_{the}(\{\theta\}_{old}, x_i)]^2}{\sigma_i^2}$$

End loop

[M]

loop on MCMC according to Metropolis-Hasting Algorithm

Select  $\{\theta\}_{New}$  according to  $\{\theta\}_{old}$



$$\{\theta\}_{New} = \{\theta\}_{old} + \{\Delta\theta\}$$

Random Number

→  $R_1, R_2$

$\epsilon_1$  Box-Muller algorithm

$$\Delta\theta_1 = \sqrt{-2 \ln R_1} \cos(2\pi R_2)$$

$$\Delta\theta_2 = \sqrt{-2 \ln R_1} \sin(2\pi R_2)$$

$\chi^2$  for New locate

$\chi^2_{old}$

loop  $i=1, N$

$$\chi^2 = \chi^2 + \frac{[y_i - Y_{the}(\{\theta\}_{New}, x_i)]^2}{\sigma_i^2}$$

End loop

$\chi^2_{New} = \chi^2$

$\Delta\chi^2 = \chi^2_{New} - \chi^2_{old}$

$$AR \equiv \min \left\{ 1, \frac{L_{New}}{L_{old}} \right\} = \min \left\{ 1, e^{-\frac{\Delta \chi^2}{2}} \right\}$$

$$e^{-\frac{(\mathcal{H}_{New} - \mathcal{H}_{old})}{k_B T}}$$

Algorithm →

$\xi \equiv$  call Random no.

If  $\xi \leq e^{-\frac{\Delta \chi^2}{2}}$  then

Metropolis-Hastings method

$$\chi_{old}^2 = \chi_{New}^2$$

$$\{\theta\}_{old} = \{\theta\}_{New}$$

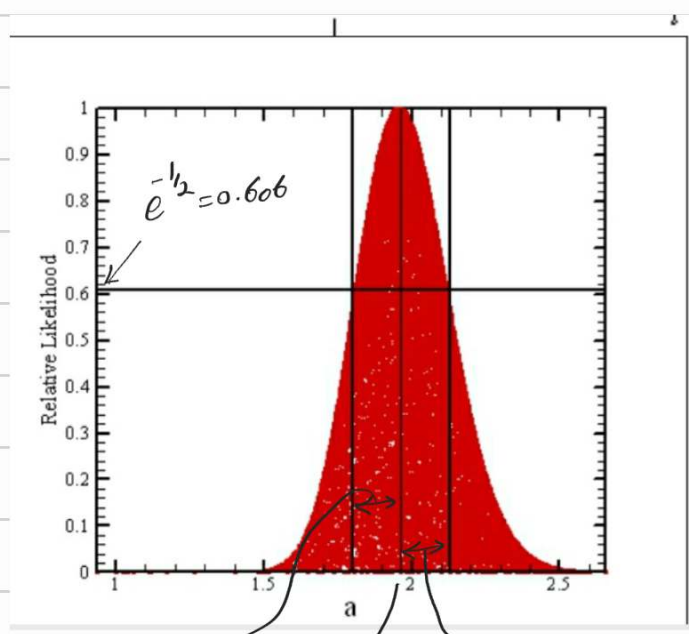
End if

Write  $\{\theta\}_{old}, e^{-\frac{(\chi^2 - \chi_{min}^2)}{2}}$

Write  $\{\theta\}_{old}, \chi_{old}^2$

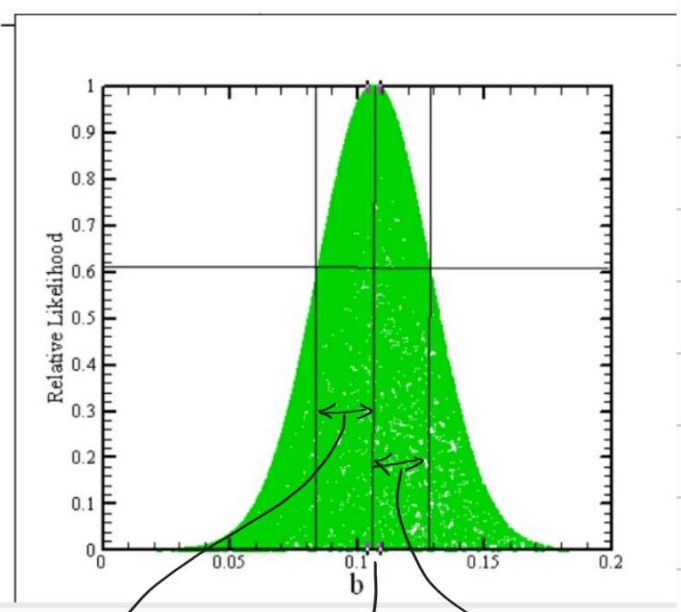
$\{\chi_{min}^2, \{\theta\}_{Best}\}$

End loop MCMC



$\sigma_a^-$        $a_{Best}$        $\sigma_a^+$

40



$\sigma_b^-$        $b_{Best}$        $\sigma_b^+$