

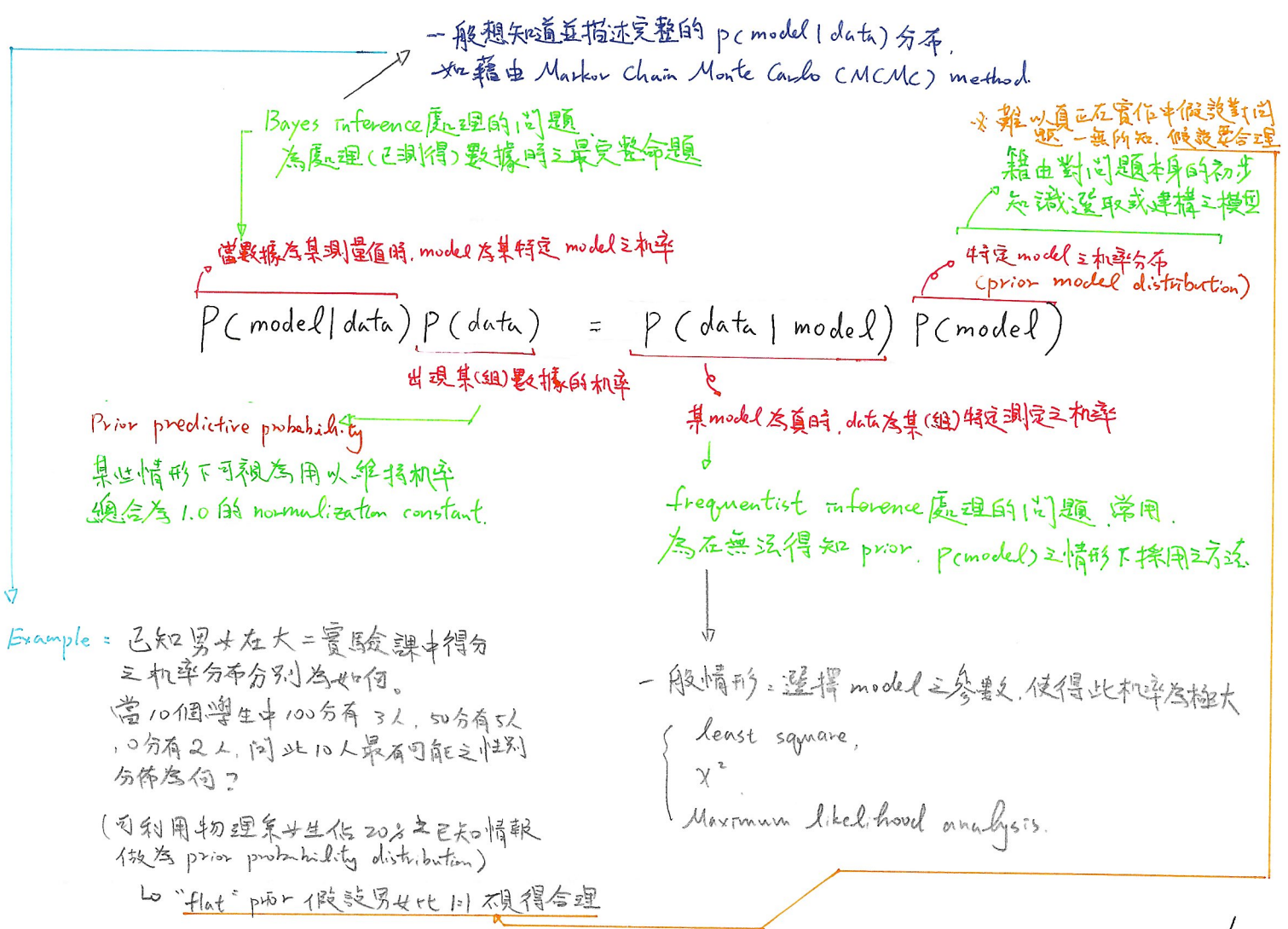
I. Purpose of modeling data:

- 1. Simplifying the description (e.g. describing as Gaussian distributions with certain mean values and standard deviations)
- 2. Constraining physical parameters and understanding the underlying physical principles. (e.g. comparing the power spectrum of CMB to inflation models)

II. Good (necessary) practice when dealing with data:

- 1. knowing what we are doing.
- 2. Being able to tell people why we are doing so. (Not good to say: 老師叫我們這樣做) Ideal to loose everybody at hello ----
- 3. Reproducibility
- 4. Portability (other people should easily understand our data product, and can utilize them. When you define quantities in your own way, or when you invent a new data format, make them transparent and acceptable)

III. Concept: from Bayesian to frequentist (類神經網路及AI為其它泛時之經驗提昇)
↓
由基本的機率公理出發



IV. Goal of frequentist inference

1. describe data with certain model parameters.
2. knowing the uncertainty of parameters.
3. assessing goodness-of-fitting. (若 goodness-of-fitting 不好, 則前兩步得知之信息不具任何參考價值)

A. Least Square analysis (frequentist)

(i) model: 測量 y 為變數 x 的函數, 並且此函數有 M 個可變的參數 $a_j, j=0, \dots, M-1$

$$y(x) = y(x | a_0, \dots, a_{M-1})$$

(ii) 如有 N 組 (y_i, x_i) 的測量, ($i=0, \dots, N-1$), 且僅 y 有測量誤差, 並且測量誤差可以用標準差為 σ 之 normal distribution 描述.

則

$$\underbrace{P(\text{data} | \text{model})}_{\text{likelihood}} \propto \prod_{i=0}^{N-1} \left\{ \exp \left[-\frac{1}{2} \left(\frac{\overbrace{y_i}^{\text{measurement}} - \underbrace{y(x_i)}_{\text{model}}}{\underbrace{\sigma}_{\text{measurement error}}} \right)^2 \right] \Delta y \right\}$$

frequentist 推論時工作為調整參數使得 $P(\text{data} | \text{model})$ 為極大, 亦即極小化

$$\left[\sum_{i=0}^{N-1} \frac{[y_i - y(x_i)]^2}{\sigma^2} \right]$$

Bayesian 推論則討論 $P(\text{model} | \text{data}) \propto P(\text{data} | \text{model}) P(\text{model})$ 之完整分布

χ^2 (chi-square) fitting 假設每個測量可能有不同之誤差 σ_i , 而極小化

$$\chi^2 \equiv \sum_{i=0}^{N-1} \left(\frac{y_i - y(x_i | a_0 \dots a_{M-1})}{\sigma_i} \right)^2$$

也就是找 $\sum_{i=0}^{N-1} \chi^2$ 為 0 之位置.

可藉由解析推導 (線性問題) 或 Monte-Carlo simulation 推估 χ^2 大於某特定值之機率. 依此可定義 goodness-of-fitting. 例如, 若 χ^2 大於 1000 的機率為 10^{-10} , 而我們的 best-fit 給出之極小的 χ^2 為 2000, 則可斷定 model 不恰當, 在此情形推導出之參數及其 uncertainties 皆不具意義.

* 機率太高也有問題. 一般代表 σ 被高估, 為常見之數據擬合作弊行為.

best-fit 合理的 χ^2 值應接近 $N-M$

(i.e., 每個獨立的測量貢獻 ~ 1.0 到 χ^2)