

V. Markov Chain Monte Carlo method.

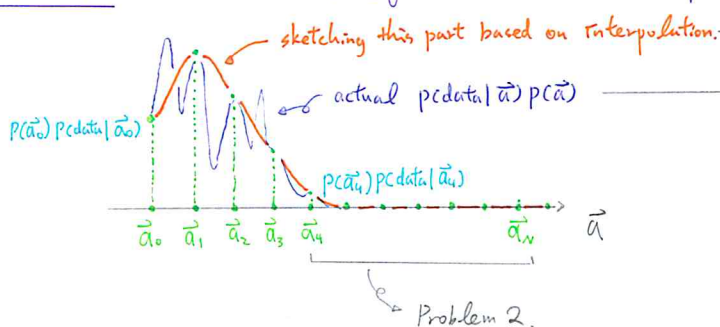
what we want to know (Bayes' theorem):

posterior probability of the model parameter \vec{a} : $P(\vec{a} | \text{data})$

$$\underbrace{P(\vec{a} | \text{data})}_{\text{posterior probability dist.}} \propto \underbrace{P(\text{data} | \vec{a})}_{\text{likelihood}} \underbrace{P(\vec{a})}_{\text{prior probability}} \equiv \underbrace{\pi(\vec{a})}_{\text{unnormalized probability distribution}}$$

We want a numerical method to draw $\pi(\vec{a})$, in particular, when $P(\text{data} | \vec{a})$ either has a very complicated analytical form or no analytical form.

approach (i) = (Monte Carlo) (quasi) uniform sampling.



Problem 1.

We usually cannot know what is the appropriate density of sampling points \vec{a}_i , such that the $P(\text{data} | \vec{a})P(\vec{a})$ can be represented with necessary details.

类比: 在原理 χ^2 minimization 时, 也有在接近 χ^2_{\min} 处 \vec{a} 才是重要的

Usually this low-probability region is not important. But in high-dimensional problem (i.e., problems with lots of free parameters), these low-probability samples can cost most of your CPU hours.

approach (ii) = Markov Chain Monte Carlo method.

Using this magical method to visit the \vec{a} space iteratively, with a probability that is proportional to $\pi(\vec{a})$

requirement for this method.

- (i) ergodic: no region is not unreachable.
- (ii) convergence: given any initial distribution of random sampling, we should be able to eventually arrive the desired distribution function $\pi(\vec{a})$
- (iii) practical: it should not be terribly difficult to construct this method, and it should not cost unreasonably large computing power

此方法的好处之一为集中运算资源在描述 $\pi(\vec{a})$ 较大之区域
在空间中

Markov chain = a sequence of random variable \vec{x}_i where the probability distribution $P(\vec{x}_i)$ only depends on \vec{x}_{i-1} , i.e., the variable \vec{x}_i has no memory about what happened "earlier" than $i-1$

Nicholas Constantine Metropolis, 美籍. 芝加哥大學物理學家, 蒙特卡洛計劃最早期參加者之一

A. Metropolis' Principle (earlier 1950s)

If we try to sample $\pi(\vec{a})$ with a Markov chain (i.e., a sequence of random variable $\vec{a}_0, \vec{a}_1, \vec{a}_2, \dots$ that is "locally correlated"), we will eventually visit everywhere in the \vec{a} space with a density of sampling that is proportional to $\pi(\vec{a})$, as long as the probability to go from one sampling point to the next, $P(\vec{a}_i | \vec{a}_{i-1})$, by construction, satisfy the **detailed balance equation**:

$$\pi(\vec{a}_1) P(\vec{a}_2 | \vec{a}_1) = \pi(\vec{a}_2) P(\vec{a}_1 | \vec{a}_2) \quad \text{Eq. (4)}$$

for any \vec{a}_1 and \vec{a}_2 .

物理詮釋: 位置 \vec{a}_1 之機率密度乘以 \vec{a}_1 到 \vec{a}_2 之躍遷機率等於位置 \vec{a}_2 之機率密度乘以 \vec{a}_2 到 \vec{a}_1 之躍遷機率

(i) A loose argument for the ergodicity (嚴格之數學證明存在, 但必須對機率論中之某些定理具有一定的認識)

integrating Eq. (4) with respect to \vec{a}_1

$$\begin{aligned} \Rightarrow \int \pi(\vec{a}_1) P(\vec{a}_2 | \vec{a}_1) d\vec{a}_1 &= \int \pi(\vec{a}_2) P(\vec{a}_1 | \vec{a}_2) d\vec{a}_1 \\ &= \pi(\vec{a}_2) \int \underbrace{P(\vec{a}_1 | \vec{a}_2) d\vec{a}_1}_{\substack{\text{無窮多} \\ \text{1.0, due to} \\ \text{probability conservation}}} = \pi(\vec{a}_2) \quad \text{Eq. (5)} \end{aligned}$$

若初始時於 \vec{a} 空間放入許多個 sampling points, 並且其密度正比於 $\pi(\vec{a})$, 且若根據此初始 sampling points 各別位置, 隨機地依 $P(\vec{a}_i | \vec{a}_{i-1})$ 機率函式選取下一批 (無窮多個) sampling points 的位置, 則下一批 sampling points 在 \vec{a} 空間中之密度分布亦等同 $\pi(\vec{a})$ 。

由此反觀之, \vec{a} 空間中任一點無論 $\pi(\vec{a})$ 有多低, 都必可由 \vec{a} space 中另一位置之 sampling point 依 $P(\vec{a}_i | \vec{a}_{i-1})$ 隨機採樣被尋訪到, 不存在不可能被尋訪到之位置。[然必須適當地選取 $P(\vec{a}_i | \vec{a}_{i-1})$ 之函數型式, 使得遍歷整個 \vec{a} 空間不困難 (i.e., 不需要無窮多個 iterations)]

(ii) A loose argument for the convergence (類似量子力學, 先看 variable 為不連續之情形之定理, 再 loosely argue 在 variable 為連續之情形定理亦成立)

由於 P_{ij} 為機率, 必然滿足 $\begin{cases} 0 \leq P_{ij} \leq 1 \\ \sum_j P_{ij} = 1 \end{cases}$

Defining $P(\vec{x}_j | \vec{x}_i) \equiv \underline{P}_{ij}$
transition matrix

則由一組 sampling points 之密度分布過渡到下一組 sampling points 之密度分布由 $\vec{p}^T \underline{P}$ 決定, \vec{p} 代表 sampling points 之密度分布 (可想像為 histogram)

定理: \vec{p}^T

- ① must has at least one unity eigenvalue
(observing Eq. (5) can see that $\pi(\vec{a})$ is an eigenvector of \underline{P}^T that has unity eigenvalue)
- ② The absolute values of \underline{P}^T eigenvalues are less than 1.0

(otherwise, after many operation of \underline{P}^T , the number distribution of sampling points will diverge, which is unacceptable.)

Introduction to Data Analysis 2023 Mar. 07

若以 \hat{P}^T 的 eigenfunctions decompose 初始之任意 sampling points 分布 \bar{u} ,
 由 page 9 最底下兩定理可預期經過無數多次 \hat{P}^T 的作用之後, 僅能留下與 π 平行的 component
 至於要經過多少次 \hat{P}^T 的作用才能達到好的 convergence 由 \hat{P}^T 第二大之 eigenvalue 之大小決定。
 若 \hat{P}^T 第二大之 eigenvalue 很小, 稱 \hat{P} 為 rapid mixing.

B. How to construct the transition matrix: Metropolis-Hasting Algorithm

(實作見課程網頁 jupyter notebook) 不同的 package 可能採用不同的 $q(\bar{a}_2|\bar{a}_1)$ 型式, 例如: 以 \bar{a}_1 為中心之多維高斯分布

- (i) pick a (arbitrary, to some extent...) proposal distribution $q(\bar{a}_2|\bar{a}_1)$
- (ii) starting at \bar{a}_1 , generate a candidate point \bar{a}_{2c} by randomly drawing from the proposal distribution
- (iii) evaluate the acceptance probability $\alpha(\bar{a}_1, \bar{a}_{2c})$ that is

$$\alpha(\bar{a}_1, \bar{a}_{2c}) = \min\left(1, \frac{\pi(\bar{a}_{2c}) q(\bar{a}_1|\bar{a}_{2c})}{\pi(\bar{a}_1) q(\bar{a}_{2c}|\bar{a}_1)}\right) \quad \text{--- Eq. (6)}$$

- (iv) with the probability $\alpha(\bar{a}_1, \bar{a}_{2c})$, accept the candidate point and set $\bar{a}_2 = \bar{a}_{2c}$.
 Otherwise, reject \bar{a}_{2c} and set $\bar{a}_2 = \bar{a}_1$.

注意若 samplers 未被 advanced, 即在步驟 (iv) 中 \bar{a}_{2c} 不斷地被 reject, 可能代表選取的 $q(\bar{a}_2|\bar{a}_1)$ 型式不理想

此演算法等致地給出 transition probability: $p(\bar{a}_2|\bar{a}_1) = q(\bar{a}_2|\bar{a}_1) \alpha(\bar{a}_1, \bar{a}_2)$, ($\bar{a}_2 \neq \bar{a}_1$)

Eq. (6) 左右同乘 $\pi(\bar{a}_1) q(\bar{a}_2|\bar{a}_1)$

$$\begin{aligned} \Rightarrow \pi(\bar{a}_1) q(\bar{a}_2|\bar{a}_1) \alpha(\bar{a}_1, \bar{a}_2) &= \min[\pi(\bar{a}_1) q(\bar{a}_2|\bar{a}_1), \pi(\bar{a}_2) q(\bar{a}_2|\bar{a}_1)] \\ &\quad \downarrow \text{左右兩項位置互換} \\ &= \min[\pi(\bar{a}_2) q(\bar{a}_1|\bar{a}_2), \pi(\bar{a}_1) q(\bar{a}_2|\bar{a}_1)] \end{aligned}$$

observing that

$$\pi(\bar{a}_2) q(\bar{a}_1|\bar{a}_2) \alpha(\bar{a}_2, \bar{a}_1) = \min[\pi(\bar{a}_2) q(\bar{a}_1|\bar{a}_2), \pi(\bar{a}_1) q(\bar{a}_2|\bar{a}_1)]$$

$$\Rightarrow \pi(\bar{a}_1) q(\bar{a}_2|\bar{a}_1) \alpha(\bar{a}_1, \bar{a}_2) = \pi(\bar{a}_2) q(\bar{a}_1|\bar{a}_2) \alpha(\bar{a}_2, \bar{a}_1)$$

依上述方式定義 transition matrix, 可滿足

$$\pi(\bar{a}_1) p(\bar{a}_2|\bar{a}_1) = \pi(\bar{a}_2) p(\bar{a}_1|\bar{a}_2)$$

ergodicity 及 convergence 理論上不成問題

C. Practical Aspect (以 Python emcee package 實作, 見課程網頁 Jupyter notebook)

可同時用多個 samplers, 天然地適合平行運算. 接近 converge 到 equilibrium state 之前
 會有過多的 samplers 處於低 $\pi(\bar{a})$ 位置, 故必須放棄前 n 步取到的 samplers. n 的大小
 無法由理論推算, 僅能靠對 samplers 的行為之觀察得知. 稱此前 n 個 iterations 為
 burn-in steps. 把這些 samplers 的分布做成 histogram 即可得到 $p(\bar{a}|\text{data})$, 依此說明
 最有可能之 \bar{a} 為何 (即 best-fit) 與 \bar{a} 之不同 components 的 uncertainties.