

# Motoare de inferență

## Cursul 11

Programare și modelare probabilistă - anul III

Facultatea de Informatică, UAIC

*e-mail:* [adrian.zalinescu@uaic.ro](mailto:adrian.zalinescu@uaic.ro)

*web:* <https://sites.google.com/view/fiicoursepmp/home>

8 Ianuarie 2024

Până acum, ne-am concentrat atenția asupra construcției de modele, interpretarea rezultatelor și critica modelelor.

- Ne-am bazat astfel doar pe *butonul de inferență*, interpretat în PyMC de funcția `pm.sample`, pentru a calcula distribuția a posteriori.
- În acest curs ne vom concentra pe unele detalii privind motorul de inferență din spatele acestei funcții.
- Scopul programării probabiliste este de a face eșantionarea din distribuția a posteriori, nu de a înțelege cum se face aceasta.
- Totuși, acest lucru este important pentru a înțelege procesul de inferență și ar putea să ne ajute atunci când aceste metode eșuează, pentru a vedea ce se poate face.

- 1 Motoare de inferență
- 2 Metode non-Markoviene
  - Grid computing
  - Metoda pătratică
  - Metode variaționale
- 3 Metode Markoviene
  - Metropolis-Hastings
  - Hamiltonian Monte Carlo
  - Monte Carlo secvențial

# Motoare de inferență

Deși conceptual sunt simple, metodele Bayesiene se pot dovedi dificile din punct de vedere matematic sau numeric.

- Motivul este că verosimilitatea marginală, numitorul din formula lui Bayes

$$p(\theta | y) = \frac{p(y | \theta) \cdot p(\theta)}{p(y)},$$

ia de obicei forma unei integrale costisitoare din punct de vedere computațional,

$$p(y) = \int p(y | \theta) \cdot p(\theta) d\theta.$$

- Din acest motiv, distribuția a posteriori este de obicei estimată numeric folosind algoritmi din familia *Markov Chain Monte Carlo* (MCMC) sau, mai recent, *algoritmi variaționali*.
- Existența acestor metode, numite motoare de inferență, au motivat dezvoltarea limbajelor de programare probabilistă precum PyMC.

- Scopul limbajelor de programare probabilistă este de a separa procesul construcției modelelor de procesul de inferență pentru a facilita pașii iterativi ai *construcției modelului, evaluării și modificării/extinderii* acestuia.
- Prin tratarea procesului de inferență ca o cutie neagră (black box), utilizatorii limbajelor de programare probabilistă sunt liberi să se concentreze pe problemele lor specifice, lăsându-le pe acestea să gestioneze detaliile computaționale.
- Dar înainte de apariția acestora, cei care proiectau modelele probabiliste își scriau singuri metodele de inferență, în general croite după modelul considerat sau simplificau modelele pentru a le face potrivite pentru anumite aproximări matematice.
- Acest lucru prezintă avantajul eleganței și chiar unei maniere mai eficiente de a calcula distribuția a posteriori, dar poate fi supusă erorilor și este consumatoare de timp.

## Grid computing

Acesta este o metodă de *forță brută*:

Să presupunem că vrem să calculăm distribuția a posteriori pentru un model cu un singur parametru. Atunci aproximarea de tip *grid* se face în modul următor:

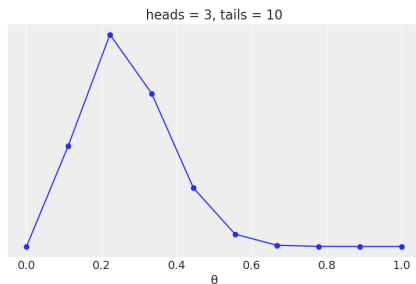
- Se definește un interval rezonabil pentru paramtrul  $\theta$  (distribuția a priori ne poate da un indiciu).
- Se partiționează intervalul (în general în sub-intervale de lungime egală).
- Pentru fiecare punct  $\theta$  din *diviziune* (*grid*), se multiplică verosimilitatea cu probabilitatea a priori:  $p(y \mid \theta) \cdot p(\theta)$ .

Codul de mai jos implementează această metodă pentru a calcula distribuția a posteriori pentru modelul aruncării cu banul:

```
def posterior_grid(grid_points=50, heads=6, tails=9):  
    """  
    A grid implementation for the coin-flipping problem  
    """  
    grid = np.linspace(0, 1, grid_points)  
    prior = np.repeat(1/grid_points, grid_points) # uniform prior  
    likelihood = stats.binom.pmf(heads+tails, grid)  
    posterior = likelihood * prior  
    posterior /= posterior.sum()  
    return grid, posterior
```

Să presupunem că am aruncat de 13 ori o monedă și am observat trei steme:

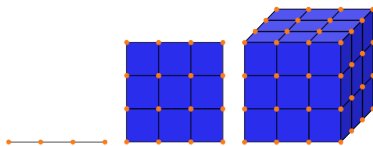
```
data = np.repeat([0, 1], (10, 3))  
points = 10  
h = data.sum()  
t = len(data) - h  
grid, posterior = posterior_grid(points, h, t)  
plt.plot(grid, posterior, 'o-')  
plt.title(f'heads = {h}, tails = {t}')  
plt.yticks([])  
plt.xlabel('θ');
```



Se poate observa că un număr mai mare de puncte (echivalent: o diviziune mai fină) duce la o aproximare mai bună:

- distribuția a posteriori este obținută ca limită a aproximării atunci când numărul de puncte crește către  $+\infty$ .
- Pe de altă parte, acest proces crește costul computațional.
- Cel mai mare dezavantaj al acestei metode este că se comportă din ce în ce mai rău pe măsură ce numărul de parametri crește (absența *scalabilității*).





Pe lângă faptul că numărul de puncte din diviziune crește, se manifestă un alt fenomen:

- regiunea din spațiul parametrilor unde este concentrată distribuția a posteriori este din ce în ce mai mică față de volumul eșantionat.
- Acest fenomen întâlnit în statistică și machine learning se numește *blestemul dimensionalității* sau *concentrarea măsurilor*:
- De exemplu, într-un hypercub, marea parte a volumului este lângă frontiera acestuia, nu în centru.
- Din punctul nostru de vedere, problema este că mare parte a timpului este petrecută alegând valori care au o contribuție aproape nulă pentru distribuția a posteriori, risipind în acest fel resurse valoroase.

# Metoda pătratică

Cunoscută și ca *metoda Laplace* sau *aproximarea normală*, *metoda pătratică* constă în aproximarea distribuției a posteriori, pe care o notăm aici  $p(x)$ , cu o distribuție Gaussiană,  $q(x)$ .

Metoda are doi pași:

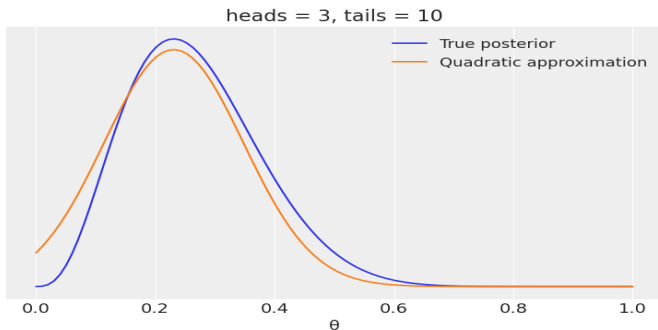
- Găsim *modul* distribuției a posteriori, adică  $\mu = \operatorname{argmax} p(x)$ . Acesta va fi media lui  $q(x)$ .
- Calculăm matricea Hessiană a lui  $p(x)$  în  $x = \mu$ . Vom calcula deviația standard a lui  $q(x)$  ca inversa rădăcinii pătrate a Hessianului (ce dă curbura lui  $p(x)$  în  $\mu$ ).

```
with pm.Model() as normal_aproximation:
    p = pm.Beta('p', 1., 1.)
    w = pm.Binomial('w', n=1, p=p, observed=data)
    mean_q = pm.find_MAP()
    std_q = ((1/pm.find_hessian(mean_q, vars=[p]))**0.5)[0]
mean_q['p'], std_q
```

Să vedem cum arată aproximarea pătratică pentru modelul beta-binomial:

```
# analytical calculation
x = np.linspace(0, 1, 100)
plt.plot(x, stats.beta.pdf(x , h+1, t+1),
         label='True posterior')

# quadratic approximation
plt.plot(x, stats.norm.pdf(x, mean_q['p'], std_q), label='Quadratic approximation')
plt.legend(loc=0, fontsize=13)
plt.title(f'heads = {h}, tails = {t}')
plt.xlabel('θ', fontsize=14)
plt.yticks([])
```



## Observații:

- Putem aplica metoda Laplace variabilelor nemărginite, adică variabilelor care au suportul (trăiesc) în  $\mathbb{R}^N$ , din cauză că
- distribuția Gaussiană este o distribuție nemărginită, așa că
- dacă încercăm să modelăm cu ajutorul ei o distribuție mărginită, vom ajunge să estimăm o densitate pozitivă acolo unde ea ar trebui să fie nulă.
- Putem totuși aplica această metodă dacă transformăm mai întâi variabila mărginită într-una nemărginită.
- De exemplu, folosim de obicei o distribuție HalfNormal pentru a modela deviația standard din cauză că aceasta ia valori pozitive; putem să o facem nemărginită dacă extragem logaritmul din aceasta.

Acestea sunt alternative la metodele Markoviene care ar putea fi o alegere mai bună în cazurile în care:

- avem seturi de date foarte mari (big data) și/sau
- pentru distribuții a posteriori foarte costisitoare din punct de vedere computațional.

Ideea este de a aproxima distribuția a posteriori cu distribuții mai simple (cum am făcut cu metoda Laplace, dar într-un mod mai elaborat).

- Vom găsi această distribuție mai simplă prin rezolvarea unei probleme de optimizare, găsind *cea mai apropiată* distribuție de cea a posteriori.
- Trebuie deci să definim ce înseamnă *apropiere*.

O modalitate de a măsura această apropiere între distribuții este prin folosirea *divergenței Kullback-Leibler (KL)*:

$$D_{KL}(q(\theta) \parallel p(\theta|y)) := \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta,$$

unde  $q$  este distribuția mai simplă folosită pentru aproximarea celei a posteriori,  $p(\theta|y)$ . Cum  $p(\theta|y)$  este necunoscută, căutăm moduri alternative de a găsi soluții ale problemei variaționale

$$\operatorname{argmin}_{q(\theta)} D_{KL}(q(\theta) \parallel p(\theta|y)).$$

Avem:

$$\begin{aligned} D_{KL}(q(\theta) \parallel p(\theta|y)) &= \int q(\theta) \log \frac{q(\theta)}{p(\theta, y)} p(y) d\theta \\ &= \int q(\theta) \left[ \log \frac{q(\theta)}{p(\theta, y)} + \log p(y) \right] d\theta \\ &= \int q(\theta) \log \frac{q(\theta)}{p(\theta, y)} d\theta + \int q(\theta) \log p(y) d\theta \\ &= \int q(\theta) \log \frac{q(\theta)}{p(\theta, y)} d\theta + \log p(y), \end{aligned}$$

deoarece  $\int q(\theta) d\theta = 1$ .

Deci

$$D_{KL}(q(\theta) \parallel p(\theta|y)) = - \underbrace{\int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta}_{\text{evidence lower bound (ELBO)}} + \log p(y)$$

- Cum  $D_{KL}$  este pozitivă, avem  $\log p(y) \geq ELBO$ , adică distribuția marginală (evidența) este întotdeauna mai mare decât  $ELBO$ , de unde numele.
- Deoarece  $\log p(y)$  este constantă, problema revine la a maximiza  $ELBO$ .
- $q(\theta)$  va fi aleasă dintr-o familie de distribuții cât mai simple:
  - ▶ o soluție este de a alege  $q(\theta)$  ca fiind descrisă de distribuții 1-dimensionale independente:

$$q(\theta) = \prod_j q_j(\theta_j),$$

- ▶ iar  $q_j$  alese din *familia exponențială* de distribuții, ce cuprinde distribuțiile normală, exponențială, Beta, Dirichlet, Gamma, Poisson, categorială (deci și Bernoulli).



- Din păcate, metoda de mai sus vine cu un algoritm specific pentru fiecare model, deci nu avem o metodă universală de inferență, (pentru toate modelele).
- Din fericire, au fost propuse soluții pentru automatizarea metodelor variaționale.
- O metodă recent propusă este *Automatic Differentiation Variational Inference (ADVI)*, în <http://arxiv.org/abs/1603.00788>, ai cărei pași sunt:
  - ▶ transformarea tuturor distribuțiilor mărginite în distribuții nemărginite (așa cum a fost discutat pentru metoda Laplace);
  - ▶ aproximarea parametrilor nemărginiți cu o distribuție Gaussiană (acei  $q_j$ ); este de remarcat faptul că o distribuție Gaussiană pe spațiul transformat al parametrilor este non-Gaussiană pe spațiul original.
  - ▶ folosirea diferențierii automate pentru a maximiza *ELBO*.
- Documentația PyMC oferă multe exemple privind folosirea inferenței variaționale în PyMC.

Constituie o familie, cunoscute drept metode *MCMC* (*Markov Chain Monte Carlo*), ce dau posibilitatea de a eșantiona din distribuția a posteriori, atât timp cât putem calcula punctual verosimilitatea și distribuția a priori.

- Metodele MCMC sunt capabile de a lua mai multe eșantioane din regiunile de probabilitate mai mare decât cele de probabilitate mai mică, în felul acesta depășind în performanță grid computing.
- La nivel fundamental, esențial în statistică este calculul mediilor de tipul următor:

$$\mathbb{E}[f] = \int p(\theta)f(\theta)d\theta.$$

- Cu MCMC, aproximăm cantitatea de mai sus cu

$$\mathbb{E}_\pi[f] = \frac{1}{N} \sum_{n=1}^N f(\theta_n),$$

unde  $\theta_1, \dots, \theta_N$  reprezintă un eșantion finit conform distribuției  $p$ , deoarece

$$\lim_{N \rightarrow \infty} \mathbb{E}_\pi[f] = \mathbb{E}[f],$$

conform LNM (legea numerelor mari).

- A fi siguri că un eșantion particular aproximează bine cantitatea căutată nu este o problemă ușoară.
- În practică, se fac teste empirice pentru a fi siguri că avem o bună aproximare MCMC.

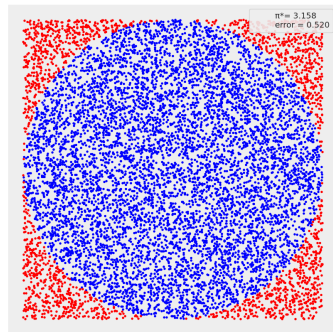
În continuare vom aminti câteva aspecte ale celor două părți componente ale MCMC:  
*Monte Carlo și lanțuri Markov.*

- Metodele *Monte Carlo* reprezintă o clasă largă de algoritmi care folosesc eșantioane aleatoare pentru a calcula sau simula un anumit proces.
- Unul dintre dezvoltatorii acestei metode, Stanislaw Ulam, a avut ideea că pentru multe probleme ce sunt dificil de rezolvat sau de formulat de o manieră precisă, ele pot fi studiate prin prelevarea de eșantioane din acestea.
- De exemplu, calculul probabilității de a obține o anumită configurație la jocul de *Solitaire*.
- Una din primele aplicații a fost rezolvarea unei probleme de fizică nucleară.
- În zilele noastre, chiar și PC-urile sunt suficient de puternice pentru a rezolva multe probleme interesante folosind metoda Monte Carlo.

Un exemplu clasic este determinarea valorii lui  $\pi$ , folosind următoarea procedură:

- Aruncați  $N$  puncte la întâmplare în interiorul unui pătrat de latură  $2R$ ;
- Desenați un cerc de rază  $R$  înscris în cerc și numărați câte puncte  $N_{int}$  sunt în interiorul acelui cerc.
- Estimați  $\pi$  ca raportul  $4N_{int}/N$ .

```
N = 10000
x, y = np.random.uniform(-1, 1, size=(2, N))
inside = (x**2 + y**2) <= 1
pi = inside.sum()*4/N
error = abs((pi - np.pi) / pi) * 100
outside = np.invert(inside)
plt.figure(figsize=(8, 8))
plt.plot(x[inside], y[inside], 'b.')
plt.plot(x[outside], y[outside], 'r.')
plt.plot(0, 0, label=f' $\pi$  = {pi:4.3f}\n'
        error = {error:4.3f}', alpha=0)
plt.axis('square')
plt.xticks([])
plt.yticks([])
plt.legend(loc=1, frameon=True, framealpha=0.9)
```



- Un *lanț Markov* este un obiect matematic care consistă într-un șir de stări și un set de probabilități de tranziție care descriu mișcarea între stări.
- Având la dispoziție un astfel de lanț, putem efectua o *plimbare aleatoare* (*random walk*) alegând un punct de plecare și apoi trecând prin diverse stări conform probabilităților de tranziție.
- Dacă am putea găsi un lanț Markov cu probabilitățile de tranziție proporționale cu distribuția din care dorim să eșantionăm (în cazul nostru, distribuția a posteriori), eșantionarea devine doar o plimbare aleatoare printre valorile căutate (în cazul nostru, prin spațiul de parametri).
- O cerință des întâlnită este ca lanțul Markov să fie *reversibil*, adică probabilitatea de a trece din starea  $i$  în starea  $j$  să fie aceeași ca cea de a trece din starea  $j$  în starea  $i$ .

Este probabil cea mai populară metodă de tip MCMC.

- Pentru unele distribuții, ca cea Gaussiană, există algoritmi eficienți de eșantionare, dar pentru altele nu este cazul.
- *Metropolis-Hastings* ne permite să eșantionăm din orice distribuție de probabilitate,  $p(x)$ , atât timp cât putem calcula măcar o valoare proporțională cu acesta, ignorând astfel factorul de normalizare, care de obicei este greu de estimat (în cazul nostru, acesta este verosimilitatea marginală).

Pentru a înțelege acest concept, vom folosi următoarea analogie:

- Să presupunem că suntem interesați în a găsi volumul unui lac și de asemenea punctul în care lacul are cea mai mare adâncime.
- De asemenea, lacul este mâlos și foarte întins, așa că o aproximare de tip grid nu e o idee prea bună.
- Strategia de eșantionare necesită doar o barcă și un băț foarte lung:

Pașii sunt următorii:

- ① Alegem un loc aleatoriu din lac și mutăm barca acolo.
- ② Folosim bățul pentru a măsura adâncimea lacului.
- ③ Mutăm barca într-un alt loc și măsurăm din nou.
- ④ Comparăm cele două măsurători astfel:
  - ▶ Dacă locul nou este mai adânc decât primul, notăm adâncimea noului loc și repetăm pasul 2.
  - ▶ Dacă este mai puțin adânc, avem două opțiuni: să acceptăm sau să respingem. Acceptare înseamnă să notăm adâncimea noului loc și să repetăm pasul 2. Respingere înseamnă să mergem înapoi în locul anterior și să notăm (din nou!) adâncimea acestuia.

Regula pentru a decide dacă să acceptăm sau să respingem este cunoscută drept *criteriul Metropolis-Hastings*:

- în principiu, aceasta stipulează că vom accepta noul loc cu o probabilitate care este proporțională cu raportul între adâncimea noului loc și a celui vechi.
- Folosind această procedură, vom găsi nu numai volumul total și adâncimea maximă a lacului, dar și forma fundului acestuia.



Evident, fundul lacului reprezintă densitatea distribuției pe care o căutăm (cu semn schimbat), adâncimea lacului este modul distribuției, iar volumul lacului este factorul de normalizare.

Astfel, algoritmul Metropolis-Hastings are următorii pași:

- 1 Alegem o valoare inițială pentru parametru,  $x_i$  (la întâmplare sau conform intuiției).
- 2 Alegem un nou parametru,  $x_{i+1}$ , prin simularea dintr-o distribuție ușor de eșantionat, precum distribuția Gaussiană sau uniformă,  $q(x_{i+1}|x_i)$ .
- 3 Calculăm probabilitatea de a accepta noul parametru prin utilizarea criteriului Metropolis-Hastings:

$$p_a(x_{i+1}|x_i) = \min \left( 1, \frac{p(x_{i+1})q(x_i|x_{i+1})}{p(x_i)q(x_{i+1}|x_i)} \right).$$

- 4 Dacă probabilitatea de la pasul 3 este mai mare decât valoarea rezultată dintr-o distribuție uniformă pe  $[0, 1]$ , atunci acceptăm noua stare; dacă nu, rămânem în starea anterioară.
- 5 Iterăm pasul 2 până obținem un eșantion suficient de mare.

## Observații:

- Dacă distribuția  $q(x_{i+1}|x_i)$  este simetrică, obținem criteriul *Metropolis*:

$$p_a(x_{i+1}|x_i) = \min \left( 1, \frac{p(x_{i+1})}{p(x_i)} \right).$$

- Pașii 3 și 4 asigură că eșantionarea va fi făcută predominant în zonele cu valori ale lui  $p$  mai mari (mai adânci ale lacului).
- Distribuția țintă (distribuția a posteriori în cazul nostru) este aproximată de lista parametrilor din eșantion. În caz de acceptare, adăugăm listei noua valoare  $x_{i+1}$ ; dacă respingem, adăugăm listei valoarea lui  $x_i$  (chiar dacă valoarea se repetă).
- În final, distribuția este doar un array de valori, simplu de manipulat.

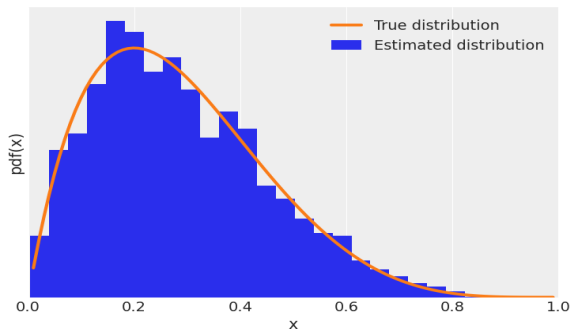
## Exemplu de implementare a algoritmului Metropolis

- pentru o distribuție oarecare:

```
def metropolis(func, draws=10000):  
    """A very simple Metropolis implementation"""  
    trace = np.zeros(draws)  
    old_x = 0.5 # func.mean()  
    old_prob = func.pdf(old_x)  
    delta = np.random.normal(0, 0.5, draws)  
    for i in range(draws):  
        new_x = old_x + delta[i]  
        new_prob = func.pdf(new_x)  
        acceptance = new_prob / old_prob  
        if acceptance >= np.random.random():  
            trace[i] = new_x  
            old_x = new_x  
            old_prob = new_prob  
        else:  
            trace[i] = old_x  
    return trace
```

- pentru cazul particular al unei distribuții Beta:

```
func = stats.beta(2, 5)
trace = metropolis(func=func)
x = np.linspace(0.01, .99, 100)
y = func.pdf(x)
plt.xlim(0, 1)
plt.plot(x, y, 'C1-', lw=3, label='True distribution')
plt.hist(trace[trace > 0], bins=25, density=True, label='Estimated distribution')
plt.xlabel('x')
plt.ylabel('pdf(x)')
plt.yticks([])
plt.legend()
```



- Eficiența algoritmului depinde de distribuția  $q$ , numită *distribuția de propunere*:
- dacă starea propusă este foarte departe de cea actuală, atunci probabilitatea de respingere este foarte mare;
- dacă starea propusă este foarte aproape, explorarea spațiului de parametri se face foarte încet.
- În general, distribuția de propunere este o distribuție Gaussiană multivariată, a cărei matrice de covarianță este determinată în faza de *tuning* (*calibrare*).
- PyMC calibrează covarianța folosind următoarea regulă generală:
  - ▶ rata de acceptare ideală este între 50% pentru o distribuție Gaussiană 1-dimensională și 23% pentru o distribuție Gaussiană  $n$ -dimensională.

- Metodele MCMC vin cu garanția că dacă eșantionul generat este suficient de mare, vom găsi o aproximare precisă a distribuției căutate.
- În practică, s-ar putea ca acest lucru să ia mai mult timp decât avem la dispoziție.
- Din această cauză, au fost propuse alternative la Metropolis-Hastings.
- Multe din aceste alternative, ca de altfel și Metropolis-Hastings, au fost concepute pentru a rezolva probleme din *mecanica statistică*, o ramură a fizicii care studiază proprietățile sistemelor atomice și moleculare.
- O astfel de modificare este *Hamiltonian Monte Carlo* (HMC): în termeni simpli, un *Hamiltonian* descrie energia totală a unui sistem fizic.

Metoda HMC este practic aceeași ca Metropolis-Hasting, cu *excepția* faptului că propunerea unei noi poziții *nu este aleatoare*:

- în loc de a muta barca în mod aleatoriu, folosim curbura fundului lacului:
- imprimăm unei bile perfecte (fără frecare) o mișcare pe fundul lacului pornind de la poziția curentă pentru un scurt moment până o oprim;
- acceptăm sau respingem noua poziție folosind criteriul Metropolis (unde nu apărea  $q$ ).
- repetăm procedura *de multe ori*.

Metoda aceasta evită unul din principalele dezavantaje ale algoritmului Metropolis-Hastings: o explorare eficientă a spațiului de parametri presupune respingerea multor pași propuși.

Pe de altă parte, acest lucru vine cu un preț: trebuie estimat *gradientul* funcției țintă (bila se va mișca în direcția acestuia, cu o viteză proporțională cu magnitudinea acestuia). Astfel:

- fiecare pas al HMC este mai costisitor decât unul al Metropolis-Hastings, dar
- probabilitatea de acceptare a acestuia este mai mare decât pentru Metropolis-Hastings.

Pentru a înclina balanța în favoarea HMC, trebuie din nou *calibrați* (*tuned*) unii parametri ai HMC:

- procedura cere mai multă experiență decât pentru Metropolis-Hastings (ceea ce face HMC mai puțin popular), dar
- PyMC vine cu un *sampler* relativ nou, cunoscut ca *No-U-Turn Sampler* (NUTS).
- Dezavantajul acestuia este că merge doar cu distribuții continue (nu putem calcula gradienti pentru distribuții discrete).

Animații pentru vizualizarea a diverși algoritmi MCMC:

<https://chi-feng.github.io/mcmc-demo/>



Una dintre problemele cu care se confruntă Metropolis-Hasting și NUTS (și alte variante HMC) este că dacă distribuția a posteriori are vârfuri multiple (separate de regiuni de foarte mică probabilitate), aceste metode pot să rămână blocate în zone uni-modale și să le rateze pe altele.

- Multe soluții dezvoltate pentru a trece de aceste probleme legate de minime locale multiple sunt legate de ideea de a introduce *temperatura* sistemului.
- Această idee vine, din nou, din mecanica statistică:
  - ▶ dacă temperatura sistemului este  $0^\circ\text{K}$  (minim absolut), acesta este blocat (înghețat) într-o singură stare;
  - ▶ dacă temperatura sistemului este infinită, atunci stările sistemului sunt echiprobabile;
  - ▶ noi suntem interesați de stări intermediare, între aceste extreme.

- Pentru modelele Bayesiene, această idee se adaptează printr-o variantă a formulei lui Bayes:

$$p(\theta|y)_\beta = p(y|\theta)^\beta p(\theta),$$

unde diferența este legată de specificarea parametrului  $\beta$ , cunoscut drept *temperatura inversă*:

- ▶ dacă  $\beta = 0$ , atunci *posteriorul temperat*  $p(\theta|y)_\beta$  este distribuția a priori;
  - ▶ dacă  $\beta = 1$ , atunci  $p(\theta|y)_\beta$  este distribuția a posteriori (până la factorul de normalizare).
- Cum eșantionarea din distribuția a priori este în general mai ușoară decât cea din distribuția a posteriori, vom începe eșantionarea din distribuția mai simplă (cu  $\beta = 0$ ) și transformarea încetă în distribuția complexă pe care o avem drept țintă.

O metodă care exploatează această idee este *Monte Carlo secvențial* (SMC). În PyMC ea este implementată în felul următor:

- 1 Inițializăm  $\beta = 0$ .
- 2 Generăm un eșantion de  $N$  valori,  $S_\beta$ , din distribuția temperată.
- 3 Creștem puțin pe  $\beta$ .
- 4 Calculăm un set de  $N$  ponderi  $W$ , conform noii distribuții temperate.
- 5 Obținem  $S_W$  prin reeșantionarea  $S_\beta$  conform  $W$ .
- 6 Pornim  $N$  lanțuri Metropolis, pornind pe fiecare valoare din  $S_W$ .
- 7 Ne întoarcem la pasul 3 până când  $\beta \geq 1$ .

Pasul de reeșantionare funcționează prin ștergerea valorilor de probabilitate mică și înlocuirea acestora cu valori de probabilitate mare.

Eficiența acestei metode depinde mult de valorile intermediare ale lui  $\beta$ , cunoscute drept *schema de răcire*.

- Cu cât este mai mică diferența între două valori succesive ale lui  $\beta$ , cu atât cele două distribuții temperate vor fi mai apropiate, și astfel tranziția de la o etapă la alta va fi mai ușoară.
- Pe de altă parte, dacă pașii sunt prea mici, vom irosi multe resurse computaționale, fără a îmbunătăți semnificativ acuratețea rezultatelor.
- Din fericire, SMC poate calcula automat valorile intermediare ale lui  $\beta$ ; schema de răcire va fi adaptată în funcție de dificultatea problemei: distribuțiile mai complicate vor necesita mai multe etape.

- În figura de mai jos, primul grafic arată un eșantion de 5 valori (în portocaliu) la o anumită etapă.
- Al doilea grafic arată cum aceste valori sunt re-ponderate conform distribuției a posteriori ponderate (în albastru este densitatea acesteia).
- Al treilea grafic ne arată rezultatul unui anumit număr de pași Metropolis, inițializați din eșantionul re-ponderat din al doilea grafic. Remarcăm cum 2 valori (de pondere mică) nu sunt folosite pentru noile lanțuri Markov.

