

# Metode de diagnoză

## Cursul 12

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>

15 Ianuarie 2024

Ne vom concentra asupra *sampler*-urilor Metropolis și NUTS. Deoarece în urma inferenței se generează un eşantion finit din distribuția a posteriori, este important să vedem dacă acesta este valid – altfel analiza va fi eronată.

- Există mai multe teste pe care le putem efectua; unele sunt vizuale, altele sunt cantitative.
- Aceste teste sunt create pentru a identifica probleme legate de eşantioane, dar nu pot dovedi că ele aproximează distribuția corectă.
- Ele ne spun doar dacă eşantioanele noastre arată rezonabil.

Soluții în caz că găsim probleme cu eșantionul:

- Creșterea mărimii eșantionului;
- Ștergerea unui număr de valori din eșantion de la începutul acestuia. Acest procedeu se mai numește *burn-in*. Faza de *tuning* a PyMC reduce nevoia de a face burn-in.
- Modificarea parametrilor sampler-ului, de exemplu creșterea lungimii fazei de tuning sau creșterea parametrului target-accept pentru sampler-ul NUTS. Câteodată, PyMC va oferi sugestii pentru aceste schimbări.
- Re-parametrizarea modelului, adică exprimarea modelului într-un mod diferit, dar echivalent.
- Transformarea datelor. Am văzut deja că centrarea datelor îmbunătățește eșantionarea în modelele liniare.

- Vom analiza un model concret, anume un model ierarhic *minimalist*, cu doi parametri:
  - ▶ un parametru *global*,  $a$ , și
  - ▶ un parametru *local*  $b$  (parametrul intra-grup).
- Va fi un model fără date, deci fără verosimilitate.
- Vom discuta două parametrizări alternative:

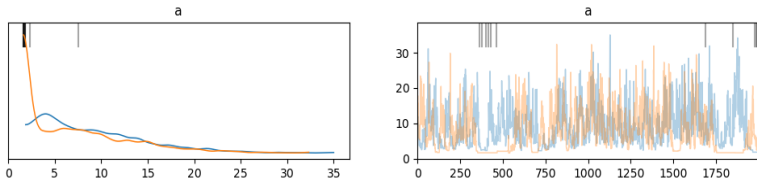
```
with pm.Model() as centered_model:
    a = pm.HalfNormal('a', 10)
    b = pm.Normal('b', 0, a, shape=10)
    idata_cm = pm.sample(2000, target_accept=0.9,
                        return_inferencedata=True, chains=2)
```

```
with pm.Model() as non_centered_model:
    a = pm.HalfNormal('a', 10)
    b_offset = pm.Normal('b_offset', mu=0, sigma=1, shape=10)
    b = pm.Deterministic('b', 0 + b_offset * a)
    idata_ncm = pm.sample(2000, target_accept=0.9,
                        return_inferencedata=True, chains=2)
```

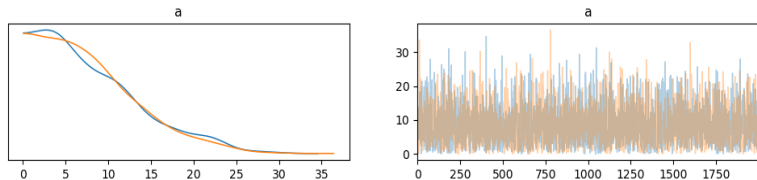
Pentru un sampler MCMC ca NUTS sau Metropolis, poate să dureze ceva timp până când converge.

- Din punct de vedere teoretic, metodele MCMC vin cu garanții de convergență în condiții foarte generale, dar pentru un eșantion infinit.
- În practic, nu putem avea decât un eșantion finit, așa că va trebui să ne bazăm pe teste empirice care în cel mai bun caz ne spun (prin indicații sau/și atenționări) că ceva rău s-ar putea întâmpla când ele eșuează, dar nu garantează că totul este ok dacă nu eșuează.
- O primă metodă este să efectuăm o verificare vizuală cu funcția `plot_trace` din ArviZ:

```
az.plot_trace(idata_cm, var_names=['a'], divergences='top', compact=False)
```



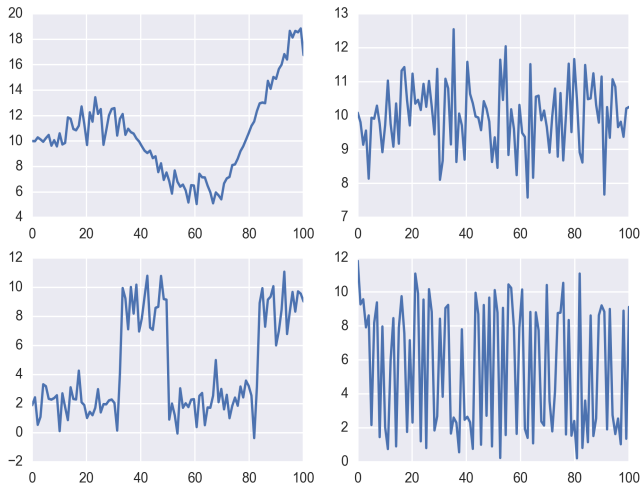
```
az.plot_trace(idata_ncm, var_names=['a'], compact=False)
```



Observăm probleme legate de primul model:

- *KDE* (figura din stânga) pentru cel de-al doilea este mai neted decât primul: un KDE neted este OK, în timp ce unul mai puțin neted poate indica o problemă, de exemplu nevoia de a mări eșantionul.
- *Urma (trace)* din dreapta ar trebui să arate ca un zgomot alb (*white noise*), adică n-ar trebui să recunoaștem nici un tipar, ca în modelul 2: spunem că avem un *mixaj bun* în cazul acesta.
- Pentru primul model, suprapunerea este mai mică; de asemenea, un lanț (*chain*) rămâne blocat între “extragerile” 350–500, iar apoi 1750–1900, ceea ce înseamnă că sampler-ul respinge toate noile propuneri cu excepția celor foarte apropiate.
- Astfel, eșantionarea este foarte încetă și nu foarte eficientă; dacă eșantionul ar fi infinit, lucrul acesta n-ar fi o problemă, dar pentru unul finit, acesta introduce o *deplasare (bias)* a rezultatului.
- Pentru compensarea acestui lucru, am putea mări eșantionul, dar dacă deplasarea este mare, mărimea necesară a acestuia va crește foarte repede.

## Exemple de *trace*: bun (dreapta) și rău (stânga):





- De asemenea, un bun trace MCMC are proprietatea de *auto-similaritate*: orice două intervale de “extrageri” trebuie să arate similar.
- Dacă prima parte arată diferit, aceasta este o indicație că sampler-ul are nevoie de un burn-in sau un eșantion mai mare.
- Dar de cele mai multe ori, dacă două părți ale trace-ului arată diferit, probabil avem nevoie de o reparametrizare.
- Pentru modelele dificile, se poate încerca o combinație de aceste metode.

- Putem de asemenea compara lanțurile (numărul lor depinde de nr. de procesoare și poate fi specificat cu argumentul `chains`): ele ar trebui să arate similar.
- Pentru a le compara cantitativ vom folosi statistica Rhat: ea calculează varianța inter-lanțuri utilizând varianța intra-lanțuri.
- În mod ideal, aceasta ar trebui să fie egală cu 1; în practică, o valoare bună este situată sub 1.1 – o valoare mai mare semnalează lipsa convergenței.
- Ea poate fi calculată cu funcția `az.r_hat`; este de asemenea returnată de funcția `az.summary` (by default) sau de `az.plot_forest` (opțional):

```
summaries = pd.concat([az.summary(idata_cm, var_names=['a']),  
                      az.summary(idata_ncm, var_names=['a'])])  
summaries.index = ['centered', 'non_centered']  
summaries
```

	mean	sd	hdi_3%	hdi_97%	mcse_mean	mcse_sd	ess_bulk	ess_tail	r_hat
centered	8.351	5.740	1.610	18.899	0.515	0.365	52.0	31.0	1.03
non_centered	8.003	5.975	0.036	19.238	0.089	0.063	3082.0	1727.0	1.00

Una din cantitățile returnate de `summary` este `mcse`.

- Aceasta este o estimare a erorii introduse de metoda de eșantionare.
- Ea ia în considerare faptul că “extragerile” nu sunt cu adevărat independente.
- `mcse_mean` este eroarea standard a mediei  $\bar{x}$  a  $n$  blocuri, unde fiecare bloc este o porțiune din eșantion:

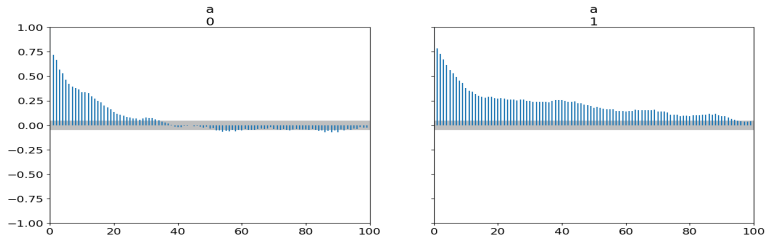
$$\text{mcse\_mean} = \frac{\sigma(\bar{x})}{\sqrt{n}}.$$

- Analog, pentru `mcse_sd` (eroarea standard a deviației standard).
- Aceste erori ar trebui să fie mai mici decât precizia pe care ne-o dorim în rezultate.

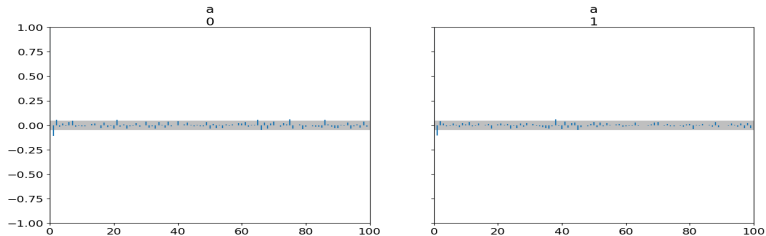
Un eșantion ideal dintr-o distribuție ar trebui să aibă o autocorelație egală cu 0:

- Un eșantion este autocorelat când o valoare la o iterație dată nu este independentă de valorile obținute la alte iterații.
- În practică, eșantioanele generate cu metoda MCMC vor fi autocorelate, în special Metropolis-Hastings, și mai puțin NUTS sau SMC.
- ArviZ are o funcție utilă pentru vizualizarea autocorelației:

```
az.plot_autocorr(idata_cm, var_names=['a'])
```



```
az.plot_autocorr(idata_ncm, var_names=['a'])
```



# Mărimea efectivă a eșantionului

- Un eșantion cu autocorelație posedă mai puțină informație decât un eșantion la fel de mare, dar fără autocorelație.
- De fapt, putem folosi autocorelația pentru a estima mărimea unui eșantion cu informație echivalentă, dar fără autocorelație.
- Aceasta se numește *mărimea efectivă a eșantionului* (*effective sample size: ess*) și poate fi calculată cu funcția `az.ess`; este de asemenea returnată de funcția `az.summary` (by default) sau de `az.plot_forest` (opțional)
- În mod ideal, `ess` ar trebui să fie apropiată de mărimea eșantionului.

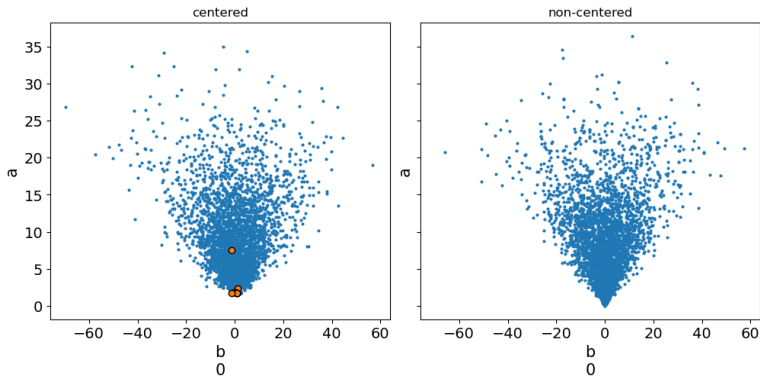
- Un avantaj al NUTS asupra Metropolis este că `ess` în NUTS este în general mult mai mare decât `ess` în Metropolis, și de aceea, folosind NUTS vom avea nevoie de eșantioane mai mici.
- PyMC ne va avertiza dacă `ess` este mai mică de 200 pentru unul din parametri (100 este de obicei îndeajuns pentru a calcula media distribuției).
- Pentru majoritatea problemelor, `ess_bulk` între 1000 și 2000 este mai mult decât de ajuns.
- Dacă dorim precizie mai mare pentru cantități ce depind de capetele (*tails*) distribuției sau evenimente ce sunt foarte rare, va trebui să considerăm același criteriu, dar pentru `ess_tail`.

Vom explora teste ce sunt specifice metodei NUTS, și care nu exprimă o proprietate a eșantioanelor generate. Acestea sunt bazate pe așa-zisele *divergențe* (PyMC afișează mesaje despre apariția acestora).

- Divergențele pot indica faptul că NUTS a întâlnit o zonă de curbură foarte mare în distribuția a posteriori pe care nu o poate explora eficient.
- Acest lucru ne spune că sampler-ul poate eventual să rateze o zonă din spațiul parametrilor, și astfel rezultatele vor fi alterate.
- Divergențele tind să apară lângă zonele problematice din spațiul parametrilor și astfel putem identifica unde ar putea fi probleme.
- O posibilitate de a vizualiza divergențele este cu funcția `az.plot_pair` cu argumentul `divergences=True`:

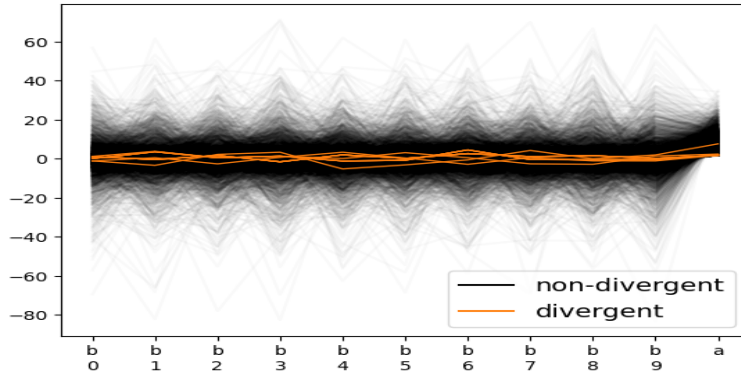


```
_, ax = plt.subplots(1, 2, sharey=True, sharex=True, figsize=(10, 5),
                    constrained_layout=True)
for idx, tr in enumerate([idata_cm, idata_ncm]):
    az.plot_pair(tr, var_names=['b', 'a'], coords={'b_dim_0':[0]}, kind='scatter',
                divergences=True, divergences_kwargs={'color':'C1'},
                ax=ax[idx])
ax[idx].set_title(['centered', 'non-centered'][idx])
```

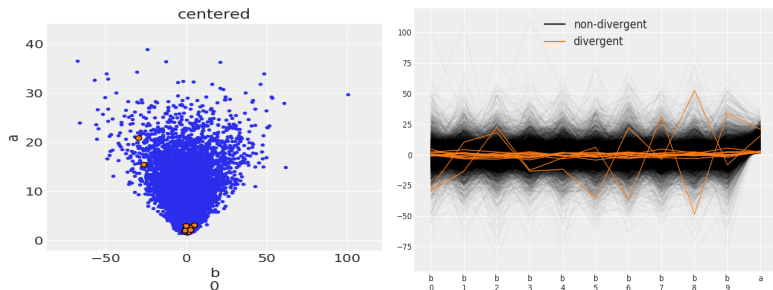


- Punctele mici și albastre reprezintă valorile eșantionului, pe când cele mai mari sunt divergențele.
- Observăm că ele apar doar în modelul 1 și sunt în principal concentrate în vârful norului de puncte, ceea ce ne semnalează o zonă cu probleme.
- În `plot_trace`, divergențele sunt semnalate de simbolul “|”.
- O altă metodă de a vizualiza divergențele este dată de funcția `az.plot_parallel`:

```
az.plot_parallel(idata_cm)
```



- Aici se vede că divergențele sunt concentrate în jurul lui 0, și pentru  $b$ , și pentru  $a$ .
- Vizualizarile de tipul precedent sunt importante, pentru că ne arată care părți ale spațiului parametrilor sunt problematice.
- PyMC folosește o regulă simplificată (*heuristic*) pentru a eticheta divergențele, iar aceasta poate spune că avem o divergență atunci când în realitate ea nu există.
- Ca regulă de bază, atunci când divergențele sunt împrăștiate în spațiul parametrilor, avem probabil un fals pozitiv; dacă ele sunt concentrate, atunci probabil există o problemă acolo.



Atunci când obținem divergențe, există în general 3 reguli de a scăpa de ele:

- 1 să creștem numărul de pași de tuning, de exemplu `pm.sample(tuning=1000)`;
- 2 să creștem valoarea `target_accept` de la valoarea standard de 0.8 la valori mai mari (maxim 1);
- 3 să reparametrizăm modelul: după cum am văzut în exemplu, modelul 2 este o reparametrizare a modelului 1.