

Compararea modelelor

Cursul 9

Programare și modelare probabilistă - anul III

Facultatea de Informatică, UAIC

e-mail: adrian.zalinescu@uaic.ro

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

11 Decembrie 2023

Pâna acum, am abordat problema inferenței, adică cum să deducem parametrii din date.

În cele ce urmează ne vom concentra asupra comparației a două sau mai multe modele care sunt folosite pentru a explica aceleași date.

Aceasta nu este o problemă simplă, dar este în același timp este esențială în analiza datelor.

- 1 Verificări predictive a posteriori
- 2 Briciul lui Occam
 - Simplitatea și acuratețea
 - Overfitting și underfitting
- 3 Criterii bazate pe informație
- 4 Regularizarea distribuțiilor a priori

Verificări predictive a posteriori

Distribuția predictivă a posteriori a fost folosită pentru a verifica cât de bine explică modelele datele folosite.

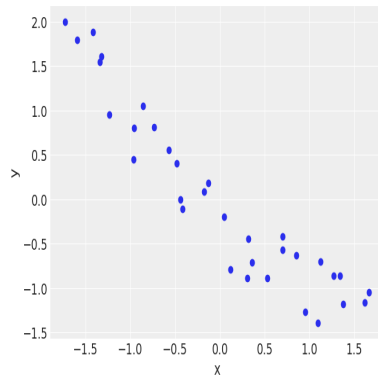
- Prin aceste verificări, sperăm să observăm limitările unui model, fie ca modalitate de a-l înțelege, fie ca o încercare de a-l îmbunătăți.
- În mod implicit, se înțelege că modelele nu vor reproduce în general toate aspectele unei probleme (la fel de bine).
 - ▶ Acest lucru nu pune probleme, deoarece modelele sunt construite cu un scop în minte.
 - ▶ O verificare predictivă a posteriori este un mod de a evalua modelul în contextul acestui scop.
 - ▶ Astfel, dacă avem mai multe modele, putem efectua acest tip de verificare pentru a le compara.

Exemplu

Generarea datelor

```
import pymc3 as pm
import numpy as np
import scipy.stats as stats
import matplotlib.pyplot as plt
import arviz as az
az.style.use('arviz-darkgrid')

dummy_data = np.loadtxt('./data/dummy.csv')
x_1 = dummy_data[:, 0]
y_1 = dummy_data[:, 1]
order = 2
x_1p = np.vstack([x_1**i for i
                  in range(1, order+1)])
x_1s = (x_1p - x_1p.mean(axis=1, keepdims=True))/
      x_1p.std(axis=1, keepdims=True)
y_1s = (y_1 - y_1.mean()) / y_1.std()
plt.scatter(x_1s[0], y_1s)
plt.xlabel('x')
plt.ylabel('y')
```



Acum, vom încerca să calibrăm aceste date cu două modele diferite (nu foarte tare), unul liniar și unul polinomial de ordinul 2 (cunoscut drept *modelul parabolic* sau *modelul pătratic*):

```
with pm.Model() as model_l:  
     $\alpha$  = pm.Normal('α', mu=0, sigma=1)  
     $\beta$  = pm.Normal('β', mu=0, sigma=10)  
     $\epsilon$  = pm.HalfNormal('ε', 5)  
     $\mu$  =  $\alpha$  +  $\beta$  * x_1s[0]  
    y_pred = pm.Normal('y_pred', mu= $\mu$ , sigma= $\epsilon$ , observed=y_1s)  
    idata_l = pm.sample(2000, return_inferencedata=True)
```

```
with pm.Model() as model_p:  
     $\alpha$  = pm.Normal('α', mu=0, sigma=1)  
     $\beta$  = pm.Normal('β', mu=0, sigma=10, shape=order)  
     $\epsilon$  = pm.HalfNormal('ε', 5)  
     $\mu$  =  $\alpha$  + pm.math.dot( $\beta$ , x_1s)  
    y_pred = pm.Normal('y_pred', mu= $\mu$ , sigma= $\epsilon$ , observed=y_1s)  
    idata_p = pm.sample(2000, return_inferencedata=True)
```

Acum, vom reprezenta grafic mediile pentru cele două modele:

```
x_new = np.linspace(x_1s[0].min(), x_1s[0].max(), 100)
```

```
 $\alpha$ _l_post = idata_l.posterior[' $\alpha$ '].mean(("chain", "draw")).values
```

```
 $\beta$ _l_post = idata_l.posterior[' $\beta$ '].mean(("chain", "draw")).values
```

```
y_l_post =  $\alpha$ _l_post +  $\beta$ _l_post * x_new
```

```
plt.plot(x_new, y_l_post, 'C1', label='linear model')
```

```
 $\alpha$ _p_post = idata_p.posterior[' $\alpha$ '].mean(("chain", "draw")).values
```

```
 $\beta$ _p_post = idata_p.posterior[' $\beta$ '].mean(("chain", "draw")).values
```

```
idx = np.argsort(x_1s[0])
```

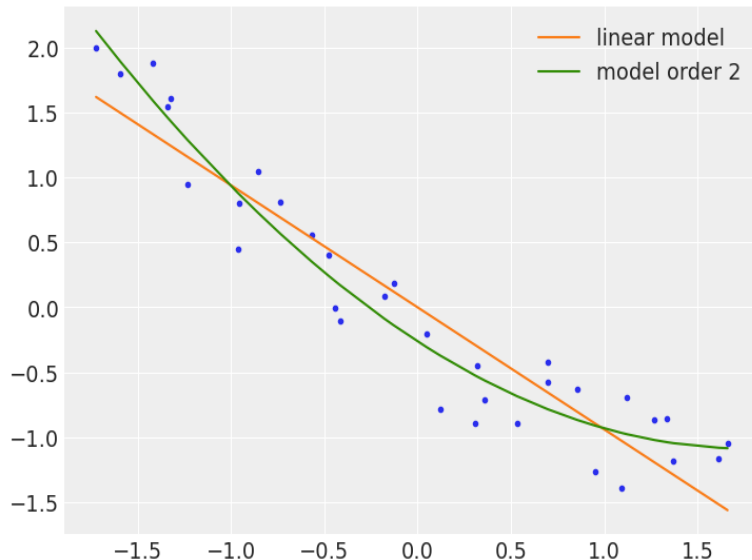
```
y_p_post =  $\alpha$ _p_post + np.dot( $\beta$ _p_post, x_1s)
```

```
plt.plot(x_1s[0][idx], y_p_post[idx], 'C2', label=f'model order {order}')
```

```
plt.scatter(x_1s[0], y_1s, c='C0', marker='.')
```

```
plt.legend()
```

Modelul pătratic pare mai bun, dar cel liniar se comportă și el bine:



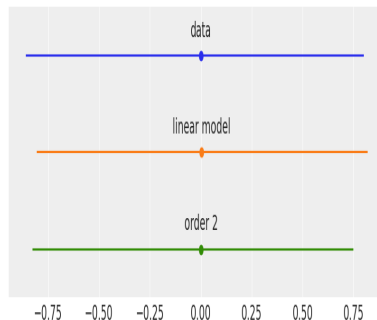
Verificările predictive a posteriori:

```
pm.sample_posterior_predictive(idata_l, model=model_l, extend_inferencedata=True)
pm.sample_posterior_predictive(idata_p, model=model_p, extend_inferencedata=True)
```

```
y_l = idata_l.posterior_predictive['y_pred']
y_p = idata_p.posterior_predictive['y_pred']

plt.figure(figsize=(8, 3))
data = [y_ls, y_l, y_p]
labels = ['data', 'linear model', 'order 2']
for i, d in enumerate(data):
    mean = d.mean()
    err = np.percentile(d, [25, 75])
    plt.errorbar(mean, -i, xerr=[-err[0], err[1]],
                 fmt='o')
plt.text(mean, -i+0.2, labels[i], ha='center',
         fontsize=14)

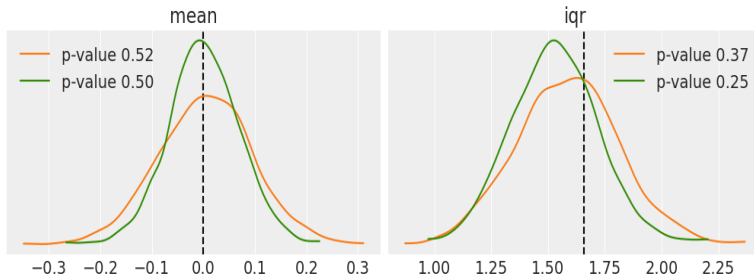
plt.ylim([-i-0.5, 0.5])
plt.yticks([])
```



- ▶ Această diagramă arată media și intervalul intercvartilic (intervalul $[Q_1, Q_3]$, unde Q_i este a i -a cvartilă) pentru date și fiecare model.
- ▶ Media și cvartilele se determină după eşantioane provenite din distribuția predictivă a posteriori pentru fiecare model.

Putem de asemenea reprezenta distribuțiile pentru medii și lungimile intervalelor intercvartilice, pentru a le compara:

```
fig, ax = plt.subplots(1, 2, figsize=(10, 3), constrained_layout=True)
def iqr(x, a=0):
    return np.subtract(*np.percentile(x, [75, 25], axis=a))
y_l_ = y_l.stack(samples=('chain', 'draw')).values.T
y_p_ = y_p.stack(samples=('chain', 'draw')).values.T
for idx, func in enumerate([np.mean, iqr]):
    T_obs = func(y_ls)
    ax[idx].axvline(T_obs, 0, 1, color='k', ls='--')
    for d_sim, c in zip([y_l_, y_p_], ['C1', 'C2']):
        T_sim = func(d_sim, 1)
        p_value = np.mean(T_sim >= T_obs)
        az.plot_kde(T_sim, plot_kwargs={'color': c},
                    label=f'p-value {p_value:.2f}', ax=ax[idx])
ax[idx].set_title(func.__name__)
ax[idx].set_yticks([])
ax[idx].legend()
```



- Linia punctată reprezintă statistica provenită din date (pentru medie și $IQR = Q_3 - Q_1$).
- Curbele (cu același cod de culoare ca pentru figura precedentă) reprezintă distribuțiile mediilor (stânga) și ale IQR (dreapta), care au fost deduse din eșantioane provenite din distribuția predictivă a posteriori.
- Valoarea p reprezintă proporția datelor simulate din statistica dată (media sau IQR) care sunt mai mari sau egale cu cea provenită din date.
- Dacă datele și simulările sunt în concordanță, ar trebui să ne așteptăm la valori p apropiate de 0.5, altfel suntem în prezența unei distribuții predictive a posteriori deplasate.

Briciul lui Occam

Simplitatea și acuratețea

Atunci când alegem dintre alternative, există un principiu de ghidare cunoscut drept *briciul lui Occam*:

- dacă există două sau mai multe explicații echivalente pentru același fenomen, o vom alege pe cea mai simplă, deoarece
- modelele simple au avantajul de a fi mai ușor de înțeles decât cele complexe.

Un alt factor pe care vrem să-l luăm în evidență este acuratețea modelului, adică cât de bine acesta se potrivește datelor. Verificarea predictivă a posteriori este bazată pe această idee de acuratețe.

- Astfel, intuitiv, vom tinde să alegem modelele care au o acuratețe mai mare și sunt în același timp simple.
- Dar dacă cel mai simplu model are cea mai mică acuratețe?
- Sau, mai general, cum putem împăca cele două tendințe: simplitatea și acuratețea?

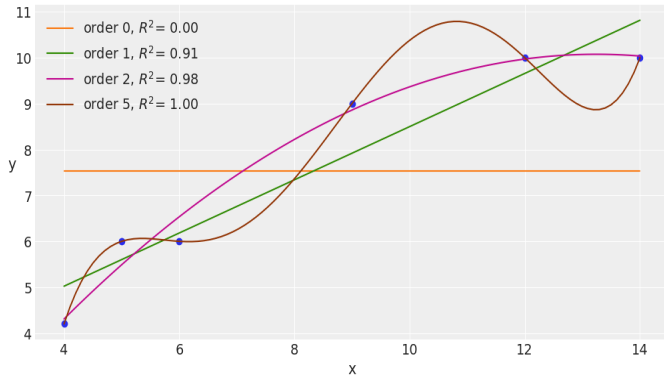
Exemplu

Vom încerca să calibrăm polinoame cu grad din ce în ce mai mare pe un set de date foarte simplu.

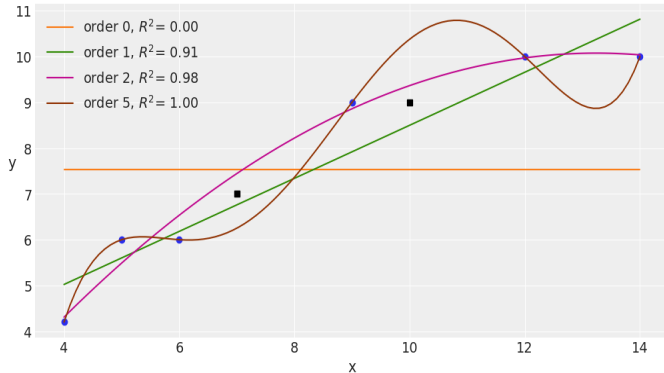
- Pentru că exemplul este doar explicativ, nu vom mai face apel la inferența Bayesiană, ci la metoda celor mai mici pătrate (care corespunde unui model Bayesian cu distribuții a priori plate):

```
x = np.array([4., 5., 6., 9., 12, 14.])
y = np.array([4.2, 6., 6., 9., 10, 10.])
plt.figure(figsize=(10, 5))
order = [0, 1, 2, 5]
plt.plot(x, y, 'o')
for i in order:
    x_n = np.linspace(x.min(), x.max(), 100)
    coeffs = np.polyfit(x, y, deg=i)
    ffit = np.polyval(coeffs, x_n)
    p = np.poly1d(coeffs)
```

```
yhat = p(x)
ybar = np.mean(y)
ssreg = np.sum((yhat-ybar)**2)
sstot = np.sum((y - ybar)**2)
r2 = ssreg / sstot
plt.plot(x_n, ffit, label=f'order{i},
 $R^2 = {r2:.2f}$ ')
plt.legend(loc=2)
plt.xlabel('x')
plt.ylabel('y', rotation=0)
```



- Observăm că acuratețea modelului, dată de coeficientul de determinare R^2 , este acompaniată de complexitatea acestuia.
- se observă, de altfel, că polinomul de grad 5 se potrivește perfect datelor!



- Observăm că acuratețea modelului, dată de coeficientul de determinare R^2 , este acompaniată de complexitatea acestuia.
- se observă, de altfel, că polinomul de grad 5 se potrivește perfect datelor!
- să presupunem că ulterior, adăugăm două puncte, $(7, 7)$ și $(10, 9)$ la setul nostru de date;
- se observă că modelul de ordin 5 nu este bun pentru a generaliza date viitoare, încă neobservate, dar potențial observabile.

Overfitting și underfitting

Când un model se potrivește prea bine cu setul de date care a fost folosit pentru a-i învăța parametrii, dar se potrivește rău cu alte seturi de date, spunem că avem *overfitting*.

- Acesta este un fenomen general întâlnit în Machine Learning și în Statistică.
- O modalitate utilă de a vizualiza overfitting-ul este de a ne imagina că seturile de date au două componente: *semnalul* și *zgomotul*.
 - ▶ semnalul este ceea ce vrem să învățăm din date;
 - ▶ presupunem că datele conțin semnal, altfel demersul ar fi inutil;
 - ▶ pe de altă parte, zgomotul este inutil și este rezultatul erorilor, limitărilor modului în care datele au fost generate sau capturate, datelor corupte, etc.
 - ▶ un model realizează overfitting dacă el învață din zgomot, ascunzând semnalul.

- Aceasta este o justificare practică pentru briciul lui Occam: prea multă complexitate duce la overfitting;
 - ▶ în principiu, putem oricând crea un model atât de complex și care să justifice totul în cel mai mic detaliu.

Underfitting

- În exemplul de mai sus, avem modelul de grad 0 (o constantă);
- Acest model captează doar media variabilei de output y , fiind complet independent de variabila de input x .
- Spunem că în cadrul acestui model se manifestă fenomenul de *underfitting*.

Balanța între simplitate și acuratețe

O variantă a briciului lui Occam:

“Lucrurile ar trebui să fie cât mai simple cu putință, dar nu mai simple” – Einstein

Ideal, ar trebui să avem un model în care nu avem nici overfitting, nici underfitting.

Așadar, va trebui să facem un compromis între simplitate și acuratețe.

Acest compromis este discutat de obicei în termeni de *varianță* și *tendențiozitate (bias)*:

- O *tendențiozitate mare* este rezultatul unui model ce are o capacitate mică de a se potrivi datelor, și astfel poate duce la underfitting.
- O *varianță mare* este rezultatul unui model ce are o sensibilitate mare la detaliile din date. Acest lucru poate cauza captarea zgomotului în date și duce astfel la overfitting.

În exemplul de mai sus, modelul de grad 0 este cel cu tendențiozitatea cea mai mare, deoarece are *tendința* de a returna o linie orizontală ce nu depinde de valoarea lui x .

Modelul de grad 5 este cel cu cea mai mare varianță: (aproape) oricum am aranja setul de date de 6 puncte, modelul se va adapta acestuia.

Măsuri pentru acuratețea predictivă

Pentru a ține cont de acuratețe, pe de o parte, și de simplitate, pe de alta, introducem mai multe concepte:

- *acuratețea în cadrul eșantionului*: este acuratețea măsurată pe *eșantion*, adică datele folosite pentru calibrarea modelului;
- *acuratețea în afara eșantionului*: este acuratețea modelului măsurată pe date care nu fac parte din eșantion – se mai numește și *acuratețea predictivă*.
 - În medie, acuratețea în cadrul eșantionului este mai mare decât cea din afara lui.
 - Folosind-o doar pe prima, am putea crede că avem un model mai bun decât este în realitate.
 - De aceea, acuratețea predictivă este preferată celei din cadrul eșantionului.

- Pe de altă parte, trebuie să ne permitem să lăsăm pentru testarea modelului, și nu calibrarea lui, o parte din datele avute la dispoziție, ceea ce reprezintă un lux pentru mulți analiști.
- Pentru a trece peste această problemă, s-au creat metode pentru a estima acuratețea predictivă folosind numai datele din eșantion. Două din aceste metode sunt:
 - *Cross-validation*: este o strategie empirică bazată pe divizarea datelor disponibile în subgrupuri folosite pentru calibrare și evaluare într-un mod alternativ.
 - *Criteriul informației*: este un termen generic pentru mai multe expresii relativ simple ce pot fi considerate ca moduri de aproximare a rezultatelor pe care le-am fi putut obține prin cross-validation.

Cross validation

Acest proces este rezumat de următoarea schemă:

- Datele se partiționează în K porțiuni, ce se încearcă a fi mai mult sau mai puțin egale (în mărime sau alte caracteristici).
- Folosim apoi $K - 1$ dintre ele pentru a antrena (calibra) modelul, iar porțiunea rămasă pentru a-l testa.
- Repetăm această procedură în mod sistematic, lăsând deoparte altă porțiune, până completăm K runde:



- Se face apoi media rezultatelor, A_k , $1 \leq k \leq K$; aceasta este cunoscută drept *K-fold cross validation*.

- În cazul în care K este mărimea setului de date, obținem *leave-one-out cross validation* (LOO-CV).
- Cross-validation este o idee simplă și folositoare, dar pentru unele modele sau pentru date foarte mari, costul computațional al acesteia poate fi prea mare.
- În cele ce urmează prezentăm alternative pentru estimarea rezultatelor obținute prin cross-validation.

Criterii bazate pe informație

Acestea sunt colecții de unelte distincte, dar interconectate, ce sunt folosite pentru a compara modelele în funcție de cum se potrivesc datelor, în timp ce țin cont de complexitatea lor printr-un factor de penalizare.

Astfel, aceste criterii formalizează intuiția dezvoltată anterior, aceea de a face un compromis între acuratețe și simplitate.

Log-verosimilitatea și devianța

Un mod intuitiv de a măsura cât de bine se potrivește datelor este de a calcula *eroarea medie pătratică* între date și predicțiile făcute de model:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \mathbb{E}[y_i \mid \theta])^2,$$

unde $\mathbb{E}[y_i \mid \theta]$ este valoarea prezisă, dat fiind parametrul estimat θ .

O măsură mai generală este *log-verosimilitatea*:

$$\sum_{i=1}^n \log p(y_i \mid \theta).$$

Atunci când verosimilitatea, $p(\cdot \mid \theta)$, este distribuită normal, aceasta este proporțională cu eroarea medie pătratică.

Din motive istorice, se mai folosește o măsură, numită *devianță*:

$$-2 \sum_{i=1}^n \log p(y_i \mid \theta).$$

- Cu cât aceasta este mai mică, cu atât este mai mare log-verosimilitatea și deci acuratețea în cadrul eșantionului;
- Așadar, modelele complexe vor avea în general o devianță mai mică decât modelele simple. Prin urmare, trebuie inclus un termen de penalizare pentru complexitate.
- În cele ce urmează, vom utiliza diverse criterii bazate pe informație, ce au în comun faptul că folosesc devianța și un termen de penalizare.

Este definit ca:

$$AIC := -2 \sum_{i=1}^n \log p(y_i \mid \hat{\theta}_{mle}) + 2p_{AIC}, \text{ unde}$$

- p_{AIC} este numărul de parametri, iar
- $\hat{\theta}_{mle}$ este *verosimilitatea maximă (MLE)*:

$$\hat{\theta}_{mle} = \operatorname{argmax}_{\theta} \prod_{i=1}^n p(y_i \mid \theta).$$

- ▶ aceasta este folosită în mod curent de non-Bayesieni și coincide cu *maximul a posteriori (MAP)* dacă se folosesc distribuții a priori plate.
- ▶ Astfel, acest criteriu este mai util în abordările non-Bayesiene, dar datorită faptului că nu folosește distribuția a posteriori, este problematic altfel.

Criteriul WAIC (widely applicable information criterion)

Este versiunea Bayesiană a criteriului AIC și, la fel, are doi termeni:

$$\begin{aligned} WAIC &= -2 \sum_{i=1}^n \log \left(\frac{1}{S} \sum_{s=1}^S p(y_i | \theta^s) \right) + 2 \sum_{i=1}^n \bigvee_{s=1}^S \log p(y_i | \theta^s) \\ &= -2lppd + 2p_{WAIC}. \end{aligned}$$

- Termenul $lppd$ (*log point-wise predictive density*) calculează verosimilitatea medie pe un eșantion de mărime S din distribuția a posteriori; este ceea ce numim *devianță* în cadrul Bayesian.
- Al doilea termen, p_{WAIC} , calculează varianța log-verosimilității pe eșantionul de mărime S ; este un factor ce penalizează complexitatea: cu cât mai mare e numărul de parametri efectivi, cu atât mai mare deviația distribuției a posteriori va fi.

Pareto smoothed importance sampling LOO-CV

Este o metodă care aproximează LOO-CV, dar fără a efectua cele K iterații.

- Ideea de bază este că este posibil să aproximăm LOO-CV prin re-ponderarea potrivită a verosimilităților.
- Acest lucru se face folosind o tehnică foarte cunoscută în statistică, numită *importance sampling*.
- Problema acesteia este că este foarte instabilă, iar
- pentru a remedia acest lucru, o nouă metodă este folosită, anume *Pareto smoothed importance sampling*, ce poate fi utilizată pentru estimări mai stabile ale LOO.
- Interpretarea este similară AIC și $WAIC$: cu cât este mai mică valoarea, cu atât mai bine.

Comparația modelelor cu PyMC

Aceasta se face simplu:

```
pm.compute_log_likelihood(idata_1,model=model_1)
waic_1 = az.waic(idata_1, scale="deviance")
waic_1
```

► Computed from 8000 posterior samples and 33 observations log-likelihood matrix.

	Estimate	SE
deviance_waic	28.85	5.36
p_waic	2.51	-

Dacă dorim să calculăm PSIS-LOO-CV în loc de WAIC, trebuie să folosim `az.loo`. PyMC prezintă 3 valori:

- o estimare punctuală a WAIC;
- eroarea standard (SE) a estimării punctuale;
- numărul efectiv de parametri, `p_waic`.

Deoarece valorile WAIC/PSIS-LOO-CV sunt interpretate în mod relativ (fiecare în modelul propriu), ArviZ oferă două funcții auxiliare pentru ușurarea comparației:

- prima este `az.compare`:

```
pm.compute_log_likelihood(idata_p,model=model_p)
cmp_df = az.compare({'model_l':idata_l, 'model_p':idata_p},
                    method='BB-pseudo-BMA', ic="waic", scale="deviance")
cmp_df
```

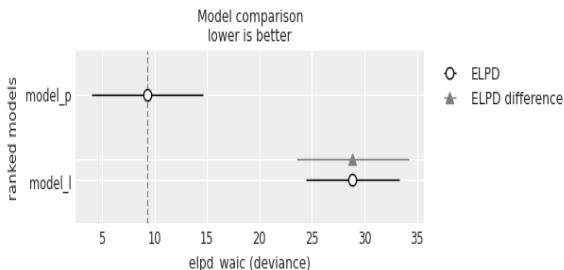
	rank	elpd_waic	p_waic	elpd_diff	weight	se	dse	warning	scale
model_p	0	9.340557	2.708534	0.000000	0.998944	5.395026	0.000000	False	deviance
model_l	1	28.853447	2.506985	19.512889	0.001056	4.421055	5.343268	False	deviance

- ▶ primele două coloane conțin estimările punctuale ale WAIC, respectiv numărul efectiv de parametri, `p_waic`;
- ▶ a treia coloană reprezintă diferența între valoarea WAIC a celui mai bun model și restul;
- ▶ câteodată, nu vrem doar să selectăm cel mai bun model, ci să facem predicții, incluzându-le pe toate conform unor ponderi, ce sunt date în a patra coloană (și reprezintă probabilitatea unui model, având în vedere datele);

- ▶ a cincea și a șasea coloană reprezintă erorile standard ale estimărilor punctuale ale WAIC, respectiv erorile standard ale diferențelor între valorile WAIC;
- ▶ a șaptea coloană arată dacă valoarea WAIC este fiabilă, caz în care False este indicat.

Putem de asemenea vizualiza grafic aceste informații, folosind:

```
az.plot_compare(cmp_df)
```



- cercurile goale reprezintă valorile WAIC, cea mai mică fiind de asemenea indicată cu o linie punctată verticală;
- liniile negre orizontale reprezintă valorile erorilor standard ale WAIC;
- pentru toate modelele cu excepția celui mai bine plasat, avem un triunghi și o linie gri orizontală care indică diferența între valoarea WAIC a celui mai bun model și restul, respectiv eroarea standard a acestei diferențe.

Regularizarea distribuțiilor a priori

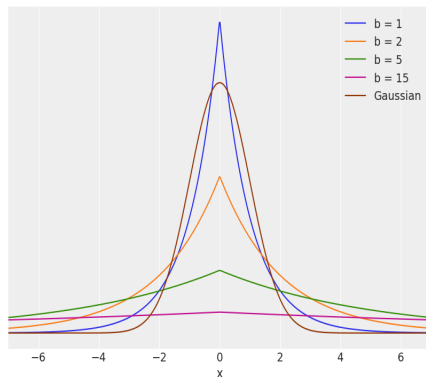
Folosirea distribuțiilor a priori informative sau slab informative este o modalitate de a introduce tendențiozitate în model.

- Dacă este făcută riguros, acest lucru este benefic, deoarece previne fenomenul de overfitting și astfel modelele vor fi capabile să facă predicții bune.
- Această idee de a adăuga tendențiozitate fără a afecta acuratețea modelului, se numește *regularizare*.
- Ea ia de obicei forma penalizării valorilor mari ale parametrilor din model;
- Astfel, se reduce informația pe care un model este capabil să o reprezinte, și în felul acesta reduc șansele ca modelul să capteze zgomotul în locul semnalului.

- În statistica non-Bayesiană, această idee ia forma a două modificări a metodei celor mai mici pătrate, numite *regresia de creastă* (*ridge regression*) și *regresia Lasso*.
- Dintr-un punct de vedere Bayesian, regresia de creastă poate fi interpretată ca folosind distribuții normale pentru coeficienții β (a unui model liniar), cu o deviație standard mică ce-i împinge pe aceștia spre 0.
- Acesta este modelul pe care l-am folosit în regresia liniară.
- Pe de altă parte, regresia Lasso poate fi interpretată ca *probabilitatea maximă a posteriori* (MAP) calculată dintr-un model care folosește distribuții a priori de tip *Laplace* pentru coeficienții β .

- Distribuția Laplace arată similar cu cea Gaussiană, dar derivata densității ei în 0 este nedefinită (arată ca un vârf ascuțit):

```
plt.figure(figsize=(8, 6))
x_values = np.linspace(-10, 10, 1000)
for df in [1, 2, 5, 15]:
    distri = stats.laplace(scale=df)
    x_pdf = distri.pdf(x_values)
    plt.plot(x_values, x_pdf,
             label=f'b = {df}')
x_pdf = stats.norm.pdf(x_values)
plt.plot(x_values, x_pdf, label='Gaussian')
plt.xlabel('x')
plt.yticks([])
plt.legend()
plt.xlim(-7, 7)
```



- Distribuția Laplace concentrează mai multă masă de probabilitate în jurul lui 0 decât distribuția normală, contribuind astfel la regularizare.