

# Curs 14

Cristian Niculescu

## 1 A priori conjugate: Beta și normală

### 1.1 Scopurile învățării

1. Să înțeleagă beneficiile a priori conjugate.
2. Să poată să actualizeze o a priori beta dată fiind o verosimilitate Bernoulli, binomială sau geometrică.
3. Să înțeleagă și să poată folosi formula pentru actualizarea unei a priori normale fiind dată o verosimilitate normală cu dispersie cunoscută.

### 1.2 Introducere și definiție

Cu o a priori conjugată, a posteriori este de același tip, de exemplu pentru verosimilitate binomială, a priori beta devine o a posteriori beta. A priori conjugate sunt utile deoarece ele reduc actualizarea Bayesiană la modificarea parametrilor repartiției a priori (așa numiții hiperparametri) în locul calculului de integrale.

Ne vom concentra pe 2 exemple importante de a priori conjugate: beta și normală. O listă mult mai cuprinzătoare este în tabelele din [http://en.wikipedia.org/wiki/Conjugate\\_prior\\_distribution](http://en.wikipedia.org/wiki/Conjugate_prior_distribution).

**Definiție.** Presupunem că avem date cu funcția de verosimilitate  $f(x|\theta)$  depinzând de un parametru  $\theta$ . Mai presupunem că repartiția a priori pentru  $\theta$  este una dintr-o familie de repartiții parametrizate. Dacă repartiția a posteriori pentru  $\theta$  este în aceeași familie ca repartiția a priori, spunem că a priori este o [a priori conjugată](#) pentru verosimilitate.

### 1.3 Repartiția beta

Arătăm că repartiția beta este o a priori conjugată pentru verosimilități binomiale, Bernoulli și geometrice.

### 1.3.1 Verosimilitate binomială

Am văzut că [repartiția beta este o a priori conjugată pentru repartiția binomială](#). Aceasta înseamnă că dacă funcția de verosimilitate este binomială și repartiția a priori este beta, atunci repartiția a posteriori este tot beta.

Mai concret, presupunem că funcția de verosimilitate are o repartiție binomială( $N, \theta$ ), unde  $N$  este cunoscut și  $\theta$  este parametrul (necunoscut) de interes. Avem de asemenea că data  $x$  este un întreg între 0 și  $N$ . Atunci, pentru o a priori beta avem următorul tabel:

hypothesis	data	prior	likelihood	posterior
$\theta$	$x$	$\text{beta}(a, b)$	$\text{binomial}(N, \theta)$	$\text{beta}(a + x, b + N - x)$
$\theta$	$x$	$c_1 \theta^{a-1} (1 - \theta)^{b-1}$	$c_2 \theta^x (1 - \theta)^{N-x}$	$c_3 \theta^{a+x-1} (1 - \theta)^{b+N-x-1}$

Tabelul este simplificat scriind coeficienții de normalizare ca  $c_1, c_2$  și  $c_3$ .

$$c_1 = \frac{(a + b - 1)!}{(a - 1)!(b - 1)!}, \quad c_2 = C_N^x = \frac{N!}{x!(N - x)!}, \quad c_3 = \frac{(a + b + N - 1)!}{(a + x - 1)!(b + N - x - 1)!}.$$

### 1.3.2 Verosimilitate Bernoulli

[Repartiția beta este o a priori conjugată pentru repartiția Bernoulli](#). Acesta este de fapt un caz special al repartiției binomiale, deoarece  $\text{Bernoulli}(\theta)$  este aceeași ca  $\text{binomial}(1, \theta)$ . În tabelul de mai jos, arătăm actualizările corespunzând succesului ( $x = 1$ ) și eșecului ( $x = 0$ ) pe linii separate.

hypothesis	data	prior	likelihood	posterior
$\theta$	$x$	$\text{beta}(a, b)$	$\text{Bernoulli}(\theta)$	$\text{beta}(a + 1, b)$ or $\text{beta}(a, b + 1)$
$\theta$	$x = 1$	$c_1 \theta^{a-1} (1 - \theta)^{b-1}$	$\theta$	$c_3 \theta^a (1 - \theta)^{b-1}$
$\theta$	$x = 0$	$c_1 \theta^{a-1} (1 - \theta)^{b-1}$	$1 - \theta$	$c_3 \theta^{a-1} (1 - \theta)^b$

Constantele  $c_1$  și  $c_3$  au aceleași formule ca în cazul precedent (al verosimilității binomiale) cu  $N = 1$ .

### 1.3.3 Verosimilitate geometrică

Repartiția geometrică( $\theta$ ) descrie probabilitatea a  $x$  eșecuri înaintea primului succes, unde probabilitatea succesului în fiecare încercare independentă este  $\theta$ . Pmf corespunzătoare este  $p(x) = \theta(1 - \theta)^x$ .

Acum presupunem că avem o dată  $x$  și ipoteza noastră  $\theta$  este că  $x$  este extrasă dintr-o repartiție geometrică( $\theta$ ). Din tabel vedem că [repartiția beta este o a priori conjugată pentru o verosimilitate geometrică](#):

ipoteza	data	a priori	verosimilitatea	a posteriori
$\theta$	$x$	$\text{beta}(a, b)$	$\text{geometrică}(\theta)$	$\text{beta}(a + 1, b + x)$
$\theta$	$x$	$c_1 \theta^{a-1} (1 - \theta)^{b-1}$	$\theta(1 - \theta)^x$	$c_3 \theta^a (1 - \theta)^{b+x-1}$

**Exemplul 1.** În timp ce călătoreau prin Regatul Ciupercilor, Mario și Luigi

au găsit niște monede neobișnuite. Ei au căzut de acord asupra unei a priori  $f(\theta) \sim \text{beta}(5, 5)$  pentru probabilitatea aversului, dar nu au fost de acord ce experiment să facă pentru a investiga  $\theta$ .

a) Mario decide să arunce o monedă de 5 ori. El obține un avers în 5 aruncări.

b) Luigi decide să arunce o monedă până la primul avers. El obține 4 reversuri înaintea primului avers.

Arătați că Mario și Luigi ajung la aceeași a posteriori pentru  $\theta$  și calculați această a posteriori.

**Răspuns.** Tabelul lui Mario:

ipoteza	data	a priori	verosimilitatea	a posteriori
$\theta$	$x = 1$	$\text{beta}(5, 5)$	$\text{binomială}(\theta)$	???
$\theta$	$x = 1$	$c_1\theta^4(1 - \theta)^4$	$C_5^1\theta(1 - \theta)^4$	$c_3\theta^5(1 - \theta)^8$

Tabelul lui Luigi:

ipoteza	data	a priori	verosimilitatea	a posteriori
$\theta$	$x = 4$	$\text{beta}(5, 5)$	$\text{geometrică}(\theta)$	???
$\theta$	$x = 4$	$c_1\theta^4(1 - \theta)^4$	$\theta(1 - \theta)^4$	$c_3\theta^5(1 - \theta)^8$

Atât a posteriori a lui Mario cât și a lui Luigi au forma unei repartiții  $\text{beta}(6, 9)$ . Factorul de normalizare este același în ambele cazuri deoarece este determinat cerând ca probabilitatea totală să fie 1.

## 1.4 Normala generează normală

**Repartiția normală este a priori conjugată cu ea însăși.** În particular, dacă funcția de verosimilitate este normală cu dispersie cunoscută, atunci o a priori normală dă o a posteriori normală. Acum atât ipotezele și datele sunt continue.

Presupunem că avem o măsurare  $x \sim N(\theta, \sigma^2)$ , unde dispersia  $\sigma^2$  este cunoscută. Adică, media  $\theta$  este parametrul nostru necunoscut de interes și știm că verosimilitatea vine dintr-o repartiție normală cu dispersia  $\sigma^2$ . Dacă alegem o pdf a priori normală

$$f(\theta) \sim N(\mu_{\text{prior}}, \sigma_{\text{prior}}^2),$$

atunci pdf a posteriori este de asemenea normală:  $f(\theta|x) \sim N(\mu_{\text{post}}, \sigma_{\text{post}}^2)$ , unde

$$\frac{\mu_{\text{post}}}{\sigma_{\text{post}}^2} = \frac{\mu_{\text{prior}}}{\sigma_{\text{prior}}^2} + \frac{x}{\sigma^2}, \quad \frac{1}{\sigma_{\text{post}}^2} = \frac{1}{\sigma_{\text{prior}}^2} + \frac{1}{\sigma^2}. \quad (1)$$

Următoarea formă a acestor formule este mai ușor de citit și arată că  $\mu_{\text{post}}$  este o medie ponderată între  $\mu_{\text{prior}}$  și data  $x$ .

$$a = \frac{1}{\sigma_{\text{prior}}^2}, \quad b = \frac{1}{\sigma^2}, \quad \mu_{\text{post}} = \frac{a\mu_{\text{prior}} + bx}{a + b}, \quad \sigma_{\text{post}}^2 = \frac{1}{a + b}. \quad (2)$$

Cu aceste formule în minte, putem exprima actualizarea prin tabelul:

hypothesis	data	prior	likelihood	posterior
$\theta$	$x$	$f(\theta) \sim N(\mu_{\text{prior}}, \sigma_{\text{prior}}^2)$	$f(x \theta) \sim N(\theta, \sigma^2)$	$f(\theta x) \sim N(\mu_{\text{post}}, \sigma_{\text{post}}^2)$
$\theta$	$x$	$c_1 \exp\left(\frac{-(\theta - \mu_{\text{prior}})^2}{2\sigma_{\text{prior}}^2}\right)$	$c_2 \exp\left(\frac{-(x - \theta)^2}{2\sigma^2}\right)$	$c_3 \exp\left(\frac{-(\theta - \mu_{\text{post}})^2}{2\sigma_{\text{post}}^2}\right)$

Demonstrația formulelor generale se face analog ca în următorul exemplu numeric.

**Exemplul 2.** Presupunem că avem a priori  $\theta \sim N(4, 8)$  și verosimilitatea  $x \sim N(\theta, 5)$ . Presupunem de asemenea că avem o măsurare  $x_1 = 3$ . Arătați că repartiția a posteriori este normală.

**Răspuns.**

a priori:  $f(\theta) = c_1 e^{-(\theta-4)^2/16}$ ; verosimilitatea:  $f(x_1|\theta) = c_2 e^{-(x_1-\theta)^2/10} = c_2 e^{-(3-\theta)^2/10}$ .

Înmulțim a priori cu verosimilitatea pentru a obține a posteriori:

$$f(\theta|x_1) = c_1 c_2 e^{-(\theta-4)^2/16} e^{-(3-\theta)^2/10} = c_1 c_2 \exp\left(-\frac{(\theta-4)^2}{16} - \frac{(3-\theta)^2}{10}\right).$$

Completăm pătratul din exponent

$$\begin{aligned} -\frac{(\theta-4)^2}{16} - \frac{(3-\theta)^2}{10} &= -\frac{5(\theta-4)^2 + 8(3-\theta)^2}{80} \\ &= -\frac{13\theta^2 - 88\theta + 152}{80} \\ &= -\frac{\theta^2 - \frac{88}{13}\theta + \frac{152}{13}}{80/13} \\ &= -\frac{(\theta - 44/13)^2 + 152/13 - (44/13)^2}{80/13}. \end{aligned}$$

De aceea, a posteriori este

$$f(\theta|x_1) = c_1 c_2 e^{-\frac{(\theta-44/13)^2 + 152/13 - (44/13)^2}{80/13}} = c_3 e^{-\frac{(\theta-44/13)^2}{80/13}}.$$

Aceasta are forma pdf pentru  $N(44/13, 40/13)$ , q.e.d.

Verificăm aceasta cu formulele (2).

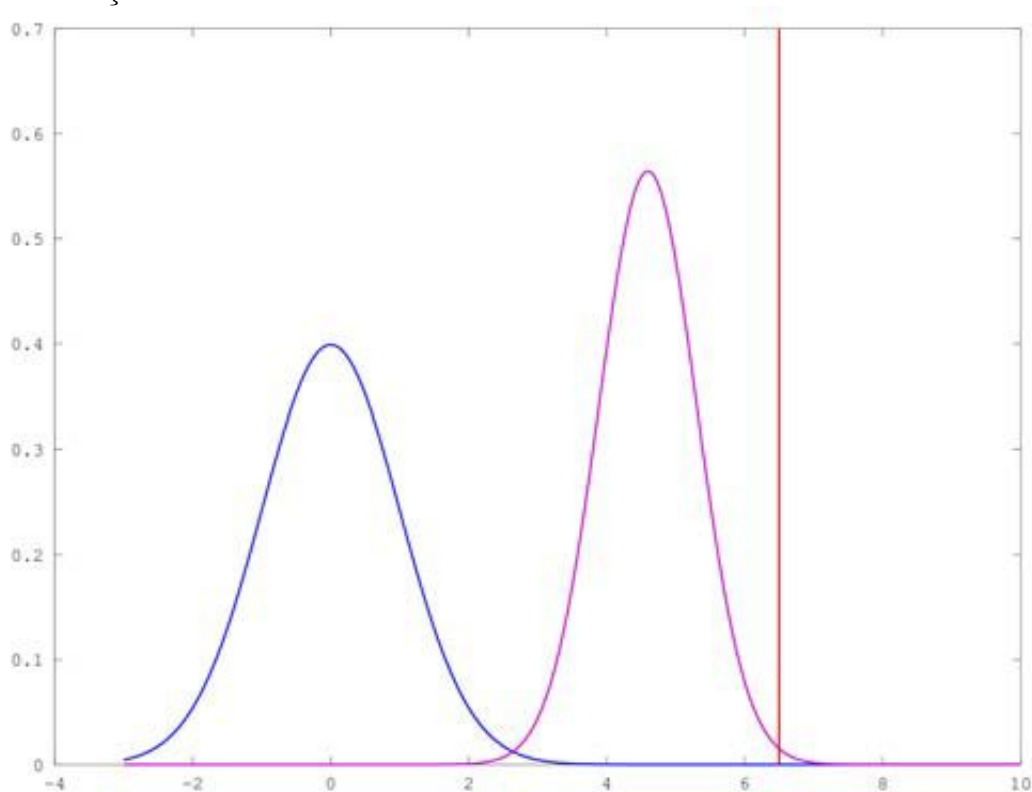
$$\mu_{\text{prior}} = 4, \sigma_{\text{prior}}^2 = 8, \sigma^2 = 5 \implies a = \frac{1}{8}, b = \frac{1}{5}.$$

De aceea

$$\begin{aligned} \mu_{\text{post}} &= \frac{a\mu_{\text{prior}} + bx}{a+b} = \frac{\frac{1}{8} \cdot 4 + \frac{1}{5} \cdot 3}{\frac{1}{8} + \frac{1}{5}} = \frac{44}{13} \approx 3.38, \\ \sigma_{\text{post}}^2 &= \frac{1}{a+b} = \frac{1}{\frac{1}{8} + \frac{1}{5}} = \frac{40}{13} \approx 3.08. \end{aligned}$$

**Exemplul 3.** Presupunem că știm datele  $x \sim N(\theta, \sigma^2)$  și avem a priori  $N(0, 1)$ . Obținem o valoare a datelor  $x = 6.5$ . Descrieți schimbările pdf pentru  $\theta$  în actualizarea de la a priori la a posteriori.

**Răspuns.** Iată graficul pdf-urilor a priori, a posteriori cu data marcată cu o line roșie.



A priori în albastru, a posteriori în mov, data în roșu.

Media a posteriori va fi o medie ponderată dintre media a priori și dată. Vârful pdf a posteriori va fi între vârful a priori și linia roșie. Avem

$$\sigma_{\text{post}}^2 = \frac{1}{1/\sigma_{\text{prior}}^2 + 1/\sigma^2} = \sigma_{\text{prior}}^2 \cdot \frac{\sigma^2}{\sigma_{\text{prior}}^2 + \sigma^2} < \sigma_{\text{prior}}^2.$$

Adică a posteriori are dispersie mai mică decât a priori, i.e. data ne face mai siguri despre unde este  $\theta$  în domeniul său.

#### 1.4.1 Mai mult de o dată

**Exemplul 4.** Presupunem că avem datele  $x_1, x_2, x_3$ . Folosiți formulele (1) pentru a actualiza succesiv.

**Răspuns.** Notăm media și dispersia a priori cu  $\mu_0$ , respectiv  $\sigma_0^2$ . Mediile și

dispersiile actualizate vor fi  $\mu_i$ , respectiv  $\sigma_i$ . Avem succesiv

$$\begin{aligned}\frac{1}{\sigma_1^2} &= \frac{1}{\sigma_0^2} + \frac{1}{\sigma^2}; \quad \frac{\mu_1}{\sigma_1^2} = \frac{\mu_0}{\sigma_0^2} + \frac{x_1}{\sigma^2} \\ \frac{1}{\sigma_2^2} &= \frac{1}{\sigma_1^2} + \frac{1}{\sigma^2} = \frac{1}{\sigma_0^2} + \frac{2}{\sigma^2}; \quad \frac{\mu_2}{\sigma_2^2} = \frac{\mu_1}{\sigma_1^2} + \frac{x_2}{\sigma^2} = \frac{\mu_0}{\sigma_0^2} + \frac{x_1 + x_2}{\sigma^2} \\ \frac{1}{\sigma_3^2} &= \frac{1}{\sigma_2^2} + \frac{1}{\sigma^2} = \frac{1}{\sigma_0^2} + \frac{3}{\sigma^2}; \quad \frac{\mu_3}{\sigma_3^2} = \frac{\mu_2}{\sigma_2^2} + \frac{x_3}{\sigma^2} = \frac{\mu_0}{\sigma_0^2} + \frac{x_1 + x_2 + x_3}{\sigma^2}.\end{aligned}$$

Exemplul se generalizează la  $n$  valori ale datelor  $x_1, \dots, x_n$ :

**Formule de actualizare normală-normală pentru  $n$  date**

$$\frac{\mu_{\text{post}}}{\sigma_{\text{post}}^2} = \frac{\mu_{\text{prior}}}{\sigma_{\text{prior}}^2} + \frac{n\bar{x}}{\sigma^2}, \quad \frac{1}{\sigma_{\text{post}}^2} = \frac{1}{\sigma_{\text{prior}}^2} + \frac{n}{\sigma^2}, \quad \bar{x} = \frac{x_1 + \dots + x_n}{n}. \quad (3)$$

Din nou dăm o formă mai simplă de citit, arătând că  $\mu_{\text{post}}$  este o medie ponderată între  $\mu_{\text{prior}}$  și media de selecție  $\bar{x}$ :

$$a = \frac{1}{\sigma_{\text{prior}}^2}, \quad b = \frac{n}{\sigma^2}, \quad \mu_{\text{post}} = \frac{a\mu_{\text{prior}} + b\bar{x}}{a + b}, \quad \sigma_{\text{post}}^2 = \frac{1}{a + b}. \quad (4)$$

**Interpretare:**  $\mu_{\text{post}}$  este o medie ponderată între  $\mu_{\text{prior}}$  și  $\bar{x}$ . Dacă numărul datelor este mare, atunci ponderea  $b$  este mare și  $\bar{x}$  va avea o puternică influență asupra a posteriori. Dacă  $\sigma_{\text{prior}}^2$  este mică, atunci ponderea  $a$  este mare și  $\mu_{\text{prior}}$  va avea o puternică influență asupra a posteriori. Pentru a rezuma:

1. Multe date au o mare influență asupra a posteriori.
2. Siguranța mare (dispersia mică) în a priori are o mare influență asupra a posteriori.

## 2 Alegerea a priori

### 2.1 Scopurile învățării

1. Să învețe că alegerea a priori afectează a posteriori.
2. Să vadă că o a priori prea rigidă poate face dificilă folosirea datelor.
3. Să vadă că mai multe date scad dependența a posteriori de a priori.
4. Să poată face o alegere rezonabilă a a priori, bazată pe înțelegerea a priori a sistemului considerat.

### 2.2 Introducere

Până acum ni s-a dat totdeauna o pdf sau pmf a priori. În acest caz, deducția statistică din date este în esență o aplicație a teoremei lui Bayes. Când a

priori este cunoscută, nu sunt controverse despre cum trebuie procedat. Artă statisticii începe când a priori nu este cunoscută cu siguranță. Sunt 2 școli principale despre cum să procedăm în acest caz: [Bayesiană](#) și [frecvenționistă](#). Acum urmăm abordarea Bayesiană. Vom învăța și abordarea frecvenționistă. Reamintim că fiind cunoscute datele  $D$  și ipoteza  $H$  am folosit teorema lui Bayes pentru a scrie

$$P(H|D) = \frac{P(D|H) \cdot P(H)}{P(D)}$$

a posteriori  $\propto$  verosimilitate  $\cdot$  a priori

**Bayesiană:** Bayesianii fac deducții folosind a posteriori  $P(H|D)$  și de aceea au nevoie totdeauna de o a priori  $P(H)$ . Dacă a priori nu este cunoscută cu certitudine, Bayesianul trebuie să încerce să facă o alegere rezonabilă. Sunt multe feluri de a face asta și oameni rezonabili pot face alegeri diferite. În general este o practică bună justificarea alegerii și explorarea unui domeniu de a priori pentru a vedea dacă toate indică aceeași concluzie.

**Frecvenționistă:** Foarte scurt, frecvenționiștii nu încearcă să creeze o a priori. În schimb, ei fac deducții folosind verosimilitatea  $P(D|H)$ .

2 beneficii ale abordării Bayesiene:

1. Probabilitatea a posteriori  $P(H|D)$  pentru ipoteză date fiind dovezile este de obicei exact ce am vrea să știm. Bayesianul poate spune ceva ca "parametrul de interes are probabilitatea 0.95 de a fi între 0.49 și 0.51".
2. Presupunerile care se iau în considerare pentru alegerea a priori pot fi clar precizate.

**Mai multe date bune:** Totdeauna [mai multe date bune](#) permit concluzii mai puternice și scad influența a priori. Accentul ar trebui să fie atât pe date bune (calitate) cât și pe mai multe date (cantitate).

## 2.3 Exemplu: zaruri

Presupunem că avem un sertar plin de zaruri, fiecare dintre acestea având 4, 6, 8, 12 sau 20 de fețe. De această dată nu știm câte zaruri de fiecare tip sunt în sertar. Un zar este ales la întâmplare din sertar și aruncat de 5 ori. Rezultatele sunt în ordine 4, 2, 4, 7 și 5.

### 2.3.1 A priori uniformă

Presupunem că nu avem idee care poate fi repartiția zarurilor din sertar. În acest caz este rezonabil să folosim o a priori plată. Iată tabelul de actualizare pentru probabilitățile a posteriori care rezultă din actualizarea după fiecare

aruncare. Pentru a încăpea toate coloanele, am eliminat a posteriori nenormalizate.

hyp.	prior	lik <sub>1</sub>	post <sub>1</sub>	lik <sub>2</sub>	post <sub>2</sub>	lik <sub>3</sub>	post <sub>3</sub>	lik <sub>4</sub>	post <sub>4</sub>	lik <sub>5</sub>	post <sub>5</sub>
$H_4$	1/5	1/4	0.370	1/4	0.542	1/4	0.682	0	0.000	0	0.000
$H_6$	1/5	1/6	0.247	1/6	0.241	1/6	0.202	0	0.000	1/6	0.000
$H_8$	1/5	1/8	0.185	1/8	0.135	1/8	0.085	1/8	0.818	1/8	0.876
$H_{12}$	1/5	1/12	0.123	1/12	0.060	1/12	0.025	1/12	0.161	1/12	0.115
$H_{20}$	1/5	1/20	0.074	1/20	0.022	1/20	0.005	1/20	0.021	1/20	0.009

Cunoscând datele, a posteriori finală este puternic ponderată spre ipoteza  $H_8$  că a fost ales un zar cu 8 fețe.

### 2.3.2 Alte a priori

Pentru a vedea cât de mult a posteriori de mai sus depinde de alegerea noastră a a priori, încercăm alte a priori. Presupunem că avem un motiv de a crede că sunt de 10 ori mai multe zaruri cu 20 de fețe în sertar decât de fiecare alt tip. Tabelul devine:

hyp.	prior	lik <sub>1</sub>	post <sub>1</sub>	lik <sub>2</sub>	post <sub>2</sub>	lik <sub>3</sub>	post <sub>3</sub>	lik <sub>4</sub>	post <sub>4</sub>	lik <sub>5</sub>	post <sub>5</sub>
$H_4$	0.071	1/4	0.222	1/4	0.453	1/4	0.650	0	0.000	0	0.000
$H_6$	0.071	1/6	0.148	1/6	0.202	1/6	0.193	0	0.000	1/6	0.000
$H_8$	0.071	1/8	0.111	1/8	0.113	1/8	0.081	1/8	0.688	1/8	0.810
$H_{12}$	0.071	1/12	0.074	1/12	0.050	1/12	0.024	1/12	0.136	1/12	0.107
$H_{20}$	0.714	1/20	0.444	1/20	0.181	1/20	0.052	1/20	0.176	1/20	0.083

Chiar și aici a posteriori finală este puternic ponderată spre ipoteza  $H_8$ .

Dar dacă zarurile cu 20 de fețe sunt de 100 de ori mai probabile decât fiecare din celelalte?

hyp.	prior	lik <sub>1</sub>	post <sub>1</sub>	lik <sub>2</sub>	post <sub>2</sub>	lik <sub>3</sub>	post <sub>3</sub>	lik <sub>4</sub>	post <sub>4</sub>	lik <sub>5</sub>	post <sub>5</sub>
$H_4$	0.0096	1/4	0.044	1/4	0.172	1/4	0.443	0	0.000	0	0.000
$H_6$	0.0096	1/6	0.030	1/6	0.077	1/6	0.131	0	0.000	1/6	0.000
$H_8$	0.0096	1/8	0.022	1/8	0.043	1/8	0.055	1/8	0.266	1/8	0.464
$H_{12}$	0.0096	1/12	0.015	1/12	0.019	1/12	0.016	1/12	0.053	1/12	0.061
$H_{20}$	0.9615	1/20	0.889	1/20	0.689	1/20	0.354	1/20	0.681	1/20	0.475

Cu o astfel de convingere a priori puternică în zarurile cu 20 de fețe, a posteriori finală dă o mare pondere teoriei că datele sunt dintr-un zar cu 20 de fețe, chiar dacă este extrem de improbabil ca un zar cu 20 de fețe să producă un maxim de 7 în 5 aruncări. A posteriori dă acum șanse aproximativ egale ca să fi fost ales un zar cu 8 fețe versus un zar cu 20 de fețe.

### 2.3.3 A priori rigide

**Disonanță cognitivă ușoară.** O convingere a priori prea rigidă poate copleși orice cantitate de date. Presupunem că suntem convinși că zarul trebuie să fie cu 20 de fețe. Deci punem a priori a noastră  $P(H_{20}) = 1$  cu



celelalte 4 ipoteze având probabilitatea 0. Iată ce se întâmplă în tabelul de actualizare.

hyp.	prior	lik <sub>1</sub>	post <sub>1</sub>	lik <sub>2</sub>	post <sub>2</sub>	lik <sub>3</sub>	post <sub>3</sub>	lik <sub>4</sub>	post <sub>4</sub>	lik <sub>5</sub>	post <sub>5</sub>
$H_4$	0	1/4	0	1/4	0	1/4	0	0	0	0	0
$H_6$	0	1/6	0	1/6	0	1/6	0	0	0	1/6	0
$H_8$	0	1/8	0	1/8	0	1/8	0	1/8	0	1/8	0
$H_{12}$	0	1/12	0	1/12	0	1/12	0	1/12	0	1/12	0
$H_{20}$	1	1/20	1	1/20	1	1/20	1	1/20	1	1/20	1

Indiferent care sunt datele, o ipoteză cu probabilitatea a priori 0 va avea probabilitatea a posteriori 0. În acest caz nu vom scăpa niciodată de ipoteza  $H_{20}$ , cu toate că putem experimenta o ușoară disonanță cognitivă.

**Disonanță cognitivă severă.** A priori rigide pot de asemenea duce la absurdități. Presupunem că suntem convinși că zarul trebuie să fie cu 4 fețe. Deci punem  $P(H_4) = 1$  și celelalte probabilități a priori 0. Cu datele cunoscute, la a 4-a aruncare intrăm în impas. 7 nu poate veni de la un zar cu 4 fețe. Totuși, aceasta este singura ipoteză pe care o permitem. A posteriori a noastră nenormalizată este o coloană de zerouri care nu poate fi normalizată.

hyp.	prior	lik <sub>1</sub>	post <sub>1</sub>	lik <sub>2</sub>	post <sub>2</sub>	lik <sub>3</sub>	post <sub>3</sub>	lik <sub>4</sub>	unnorm.	post <sub>4</sub>	post <sub>4</sub>
$H_4$	1	1/4	1	1/4	1	1/4	1	0	0	0	???
$H_6$	0	1/6	0	1/6	0	1/6	0	0	0	0	???
$H_8$	0	1/8	0	1/8	0	1/8	0	1/8	0	0	???
$H_{12}$	0	1/12	0	1/12	0	1/12	0	1/12	0	0	???
$H_{20}$	0	1/20	0	1/20	0	1/20	0	1/20	0	0	???

Trebuie să ne ajustăm convingerile despre ce este posibil sau, mai probabil, suspectăm o greșeală accidentală sau deliberată a datelor.

## 2.4 Exemplu: malaria

Iată un exemplu real adaptat din *Statistics, A Bayesian Perspective* de Donald Berry:

Prin anii 1950, oamenii de știință au început să formuleze ipoteza că purtătorii genei falciforme erau mai rezistenți la malarie ca nepurtătorii. Existau probe indirecte pentru această ipoteză. Aceasta ajută și la explicarea persistenței în populație a unei gene altfel dăunătoare. Într-un experiment oamenii de știință au injectat 30 de voluntari africani cu malarie. 15 dintre voluntari purtau o copie a genei falciforme și ceilalți 15 erau nepurtători. 14 din cei 15 nepurtători și doar 2 din cei 15 purtători au făcut malarie. Susține acest mic eșantion ipoteza că gena falciformă protejează împotriva malariei?

Fie  $S$  un purtător al genei falciforme și  $N$  un nepurtător.  $D+$  indică dezvoltarea malariei și  $D-$  indică nedeveloparea malariei. Datele pot fi puse într-un tabel.

	$D+$	$D-$	
$S$	2	13	15
$N$	14	1	15
	16	14	30

Înainte să analizăm datele ar trebui să spunem câteva cuvinte despre experiment și proiectarea lui. În primul rând, este clar neetic: pentru a obține ceva informație au infectat 16 oameni cu malarie. Trebuie de asemenea să ne îngrijorăm despre deplasare. Cum au ales subiecții testului? Este posibil ca nepurtătorii să fi fost mai slabi și astfel mai susceptibili la malarie decât purtătorii? Berry arată că este rezonabil să presupunem că o injecție este similară cu o mușcătură de țânțar, dar nu este garantat. Acest ultim punct înseamnă că dacă experimentul arată o relație între celule-seceră și protecție împotriva malariei injectate, trebuie să considerăm ipoteza că protecția împotriva malariei transmisă de țânțari este mai slabă sau inexistentă. În sfârșit, vom formula ipoteza noastră ca ”celulele-seceră protejează împotriva malariei”, dar în realitate tot ce putem spera să spunem dintr-un studiu ca acesta este că ”celula-seceră este corelată cu protecția împotriva malariei”.

**Modelul.** Pentru modelul nostru, fie  $\theta_S$  probabilitatea că un purtător injectat  $S$  face malarie și, analog, fie  $\theta_N$  probabilitatea că un nepurtător injectat  $N$  face malarie. Presupunem independența între toți subiecții experimentului. Cu acest model, verosimilitatea este o funcție de  $\theta_S$  și  $\theta_N$ :

$$P(\text{date}|\theta_S, \theta_N) = c\theta_S^2(1 - \theta_S)^{13}\theta_N^{14}(1 - \theta_N).$$

Ca de obicei lăsăm factorul constant  $c$  ca o literă. (Este produsul a 2 coeficienți binomiali:  $c = C_{15}^2 \cdot C_{15}^{14}$ .)

**Ipoteze.** Fiecare ipoteză constă dintr-o pereche  $(\theta_N, \theta_S)$ . Pentru a păstra lucrurile simple vom considera doar un număr finit de valori pentru aceste probabilități. Am putea considera mult mai multe valori sau chiar un domeniu continuu pentru ipoteze. Presupunem că  $\theta_N$  și  $\theta_S$  pot avea fiecare valorile 0, 0.2, 0.4, 0.6, 0.8 sau 1. Aceasta duce la tabele 2 dimensionale.

Primul este un tabel de ipoteze. Codul de culori indică următoarele:

1. Dreptunghiurile portocalii deschis de-a lungul diagonalei sunt unde  $\theta_S = \theta_N$ , i.e. celulele-seceră nu contează în niciun fel.
2. Dreptunghiurile roz și roșii de deasupra diagonalei sunt unde  $\theta_N > \theta_S$ , i.e. celulele-seceră dau protecție împotriva malariei.
3. În dreptunghiurile roșii  $\theta_N - \theta_S \geq 0.6$ , i.e. celulele-seceră dau multă protecție.
4. Dreptunghiurile albe de sub diagonală sunt unde  $\theta_S > \theta_N$ , i.e. celulele-

seceră de fapt cresc probabilitatea de a face malarie.

$\theta_N \backslash \theta_S$	0	0.2	0.4	0.6	0.8	1
1	(0,1)	(.2,1)	(.4,1)	(.6,1)	(.8,1)	(1,1)
0.8	(0,.8)	(.2,.8)	(.4,.8)	(.6,.8)	(.8,.8)	(1,.8)
0.6	(0,.6)	(.2,.6)	(.4,.6)	(.6,.6)	(.8,.6)	(1,.6)
0.4	(0,.4)	(.2,.4)	(.4,.4)	(.6,.4)	(.8,.4)	(1,.4)
0.2	(0,.2)	(.2,.2)	(.4,.2)	(.6,.2)	(.8,.2)	(1,.2)
0	(0,0)	(.2,0)	(.4,0)	(.6,0)	(.8,0)	(1,0)

Ipotezele asupra nivelului de protecție dată de  $S$ : roșu = mare; roz = mic; portocaliu = 0; alb = negativ.

Următorul este tabelul verosimilităților. (De fapt am profitat de indiferența noastră la scalare și am scalat toate verosimilitățile cu  $100000/c$  pentru a face tabelul mai prezentabil.) Observăm că, la precizia tabelului, multe verosimilități sunt 0. Codul de culori este același ca în tabelul de ipoteze. Am evidențiat cele mai mari verosimilități cu o margine albastră.

$\theta_N \backslash \theta_S$	0	0.2	0.4	0.6	0.8	1
1	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0.8	0.00000	1.93428	0.18381	0.00213	0.00000	0.00000
0.6	0.00000	0.06893	0.00655	0.00008	0.00000	0.00000
0.4	0.00000	0.00035	0.00003	0.00000	0.00000	0.00000
0.2	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000

Verosimilitățile  $p(\text{date}|\theta_S, \theta_N)$  scalate cu  $100000/c$ .

#### 2.4.1 A priori plată

Presupunem că nu avem nicio opinie despre dacă sau în ce măsură celula-seceră protejează împotriva malariei. În acest caz este rezonabil să folosim o a priori plată. Deoarece sunt 36 de ipoteze, fiecare primește o probabilitate a priori de  $1/36$ . Aceasta apare în tabelul de mai jos. Reamintim că fiecare dreptunghi din tabel reprezintă o ipoteză. Deoarece este un tabel de probabilități includem pmf-urile marginale.

$\theta_N \backslash \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N)$
1	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.8	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.6	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.4	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.2	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0	1/36	1/36	1/36	1/36	1/36	1/36	1/6
$p(\theta_S)$	1/6	1/6	1/6	1/6	1/6	1/6	1

A priori plată  $p(\theta_S, \theta_N)$ : fiecare ipoteză (dreptunghi) are aceeași probabilitate

Pentru a calcula a posteriori înmulțim tabelul verosimilităților cu tabelul a priori și normalizăm. Normalizarea ne asigură că suma din întregul tabel este 1.

$\theta_N \backslash \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N   \text{data})$
1	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0.8	0.00000	0.88075	0.08370	0.00097	0.00000	0.00000	0.96542
0.6	0.00000	0.03139	0.00298	0.00003	0.00000	0.00000	0.03440
0.4	0.00000	0.00016	0.00002	0.00000	0.00000	0.00000	0.00018
0.2	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
$p(\theta_S   \text{data})$	0.00000	0.91230	0.08670	0.00100	0.00000	0.00000	1.00000

A posteriori la a priori plată:  $p(\theta_S, \theta_N | \text{date})$

Pentru a decide dacă  $S$  dă protecție împotriva malariei, calculăm probabilitățile a posteriori pentru "protecție" și "protecție puternică". Acestea sunt calculate adunând numerele din dreptunghiurile corespunzătoare din tabelul a posteriori.

Protecție:  $P(\theta_N > \theta_S) = \text{suma din roz și roșu} = 0.99995$

Protecție puternică:  $P(\theta_N - \theta_S \geq 0.6) = \text{suma din roșu} = 0.88075$ .

Lucrând de la a priori plată, este efectiv sigur că celula-seceră dă protecție și foarte probabil că dă protecție puternică.



### 2.4.2 A priori informată

Acest experiment nu a fost făcut fără informație a priori. Erau multe dovezi că gena falciformă oferea protecție împotriva malariei. De exemplu, era raportat un mai mare procentaj de purtători care supraviețuiau până la maturitate.

Iată un mod de a construi o a priori informată: Vom rezerva o cantitate rezonabilă de probabilitate pentru ipoteza că  $S$  nu dă protecție. Să zicem că 24% împărțit egal între cele 6 celule portocalii unde  $\theta_N = \theta_S$ . Știm că nu ar trebui să punem nicio probabilitate a priori 0, deci hai să împărțim 6% din probabilitate între cele 15 celule de sub diagonală. Aceasta lasă 70% din probabilitate pentru cele 15 dreptunghiuri roz și roșii de deasupra diagonalei.

$\theta_N \backslash \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N)$
1	0.04667	0.04667	0.04667	0.04667	0.04667	0.04000	0.27333
0.8	0.04667	0.04667	0.04667	0.04667	0.04000	0.00400	0.23067
0.6	0.04667	0.04667	0.04667	0.04000	0.00400	0.00400	0.18800
0.4	0.04667	0.04667	0.04000	0.00400	0.00400	0.00400	0.14533
0.2	0.04667	0.04000	0.00400	0.00400	0.00400	0.00400	0.10267
0	0.04000	0.00400	0.00400	0.00400	0.00400	0.00400	0.06000
$p(\theta_S)$	0.27333	0.23067	0.18800	0.14533	0.10267	0.06000	1.0

A priori informată  $p(\theta_S, \theta_N)$ : folosește informația a priori că celula seceră protejează.

Apoi completăm pmf a posteriori.

$\theta_N \backslash \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N   \text{data})$
1	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0.8	0.00000	0.88076	0.08370	0.00097	0.00000	0.00000	0.96543
0.6	0.00000	0.03139	0.00298	0.00003	0.00000	0.00000	0.03440
0.4	0.00000	0.00016	0.00001	0.00000	0.00000	0.00000	0.00017
0.2	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
$p(\theta_S   \text{data})$	0.00000	0.91231	0.08669	0.00100	0.00000	0.00000	1.00000

A posteriori la a priori informată:  $p(\theta_S, \theta_N | \text{data})$

Calculăm din nou probabilitățile a posteriori ale "protecției" și "protecției puternice".

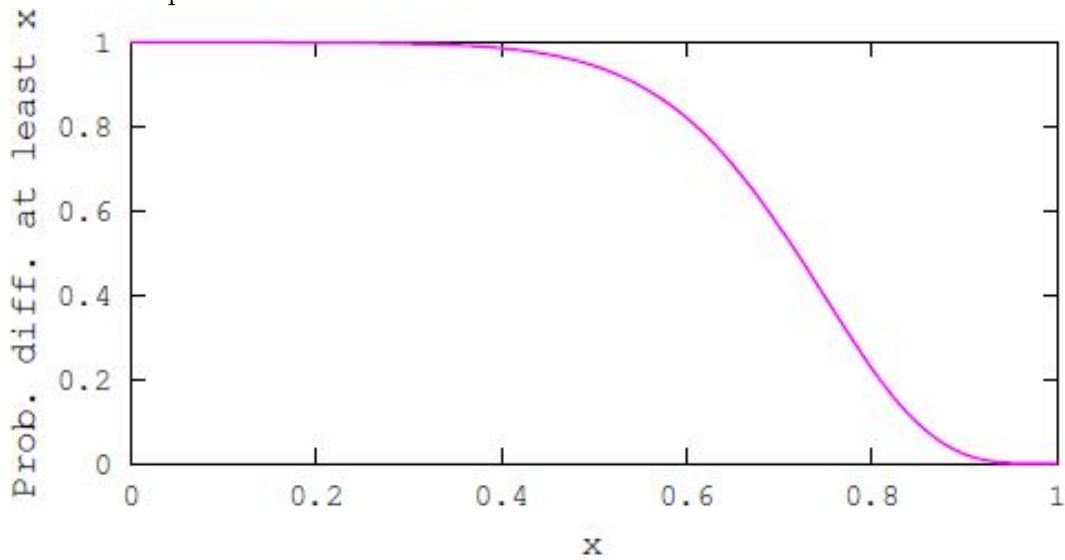
Protecție:  $P(\theta_N > \theta_S) = \text{suma din roz și roșu} = 0.99996$

Protecție puternică:  $P(\theta_N - \theta_S \geq 0.6) = \text{suma din roșu} = 0.88076$ .

Observăm că a posteriori informată este aproape identică cu a posteriori din a priori plată.

### 2.4.3 PDALX

Următoarea reprezentare este bazată pe a priori plată. Pentru fiecare  $x$ , dă probabilitatea ca  $\theta_N - \theta_S \geq x$ . Pentru a o face netedă au fost folosite mult mai multe ipoteze.



Probabilitatea că diferența  $\theta_N - \theta_S$  este cel puțin  $x$  (PDALX).

Observăm că este practic sigur că diferența este cel puțin 0.4.