

Naive bayes:

Model (C - klasa, X - atribut):

$$P(C_k|\mathbf{X}) = \frac{P(\mathbf{X}|C_k) \cdot P(C_k)}{P(\mathbf{X})} = \frac{\prod_i P(X_i|C_k) \cdot P(C_k)}{\sum_k \prod_i P(X_i|C_k) \cdot P(C_k)}$$

	X ₁	X ₂	X ₃	C
Slučaj 1	Da	A	[0-10]	DA
Slučaj 2	Ne	C	[10-20]	NE
Slučaj 3	Da	A	[10-20]	NE
...
Slučaj N	Ne	B	[20-30]	DA

predviđanje klase se onda radi kao:

$$k^* = \arg \max_k \prod_i P(X_i|C_k) \cdot P(C_k)$$

$$P(a|b) = \frac{\text{count}(a \wedge b)}{\text{count}(b)}, P(a) = \frac{\text{count}(a)}{N}$$

Underflow

Zbog proizvoda velikog broja verovatnoća, može doći do underflow-a u računarskoj reprezentaciji brojeva. Zbog toga se u implementacijama najčešće $P(C_k|\mathbf{X})$ menja sa monotonom transformacijom $\log(P(C_k|\mathbf{X}))$, kako bi se proizvod prebacio u sumu koja je numerički stabilnija.

Smoothing

Empirijske procene verovatnoća nisu stabilne na malim uzorcima (tj. mnogo zavise od uzorka kada je uzorak mali), pa se stoga radi *smoothing* procenjenih verovatnoća (k je broj kategorija atributa).

$$P(x) = \frac{\text{count}(x)+1}{N+k}, \text{ što se zove Laplace smoothing}$$

ili generalno:

$$P(x) = \frac{\text{count}(x)+\alpha}{N+k \cdot \alpha}, \text{ što se zove Additive smoothing, a } \alpha \text{ se ponekad zove pseudocount}$$

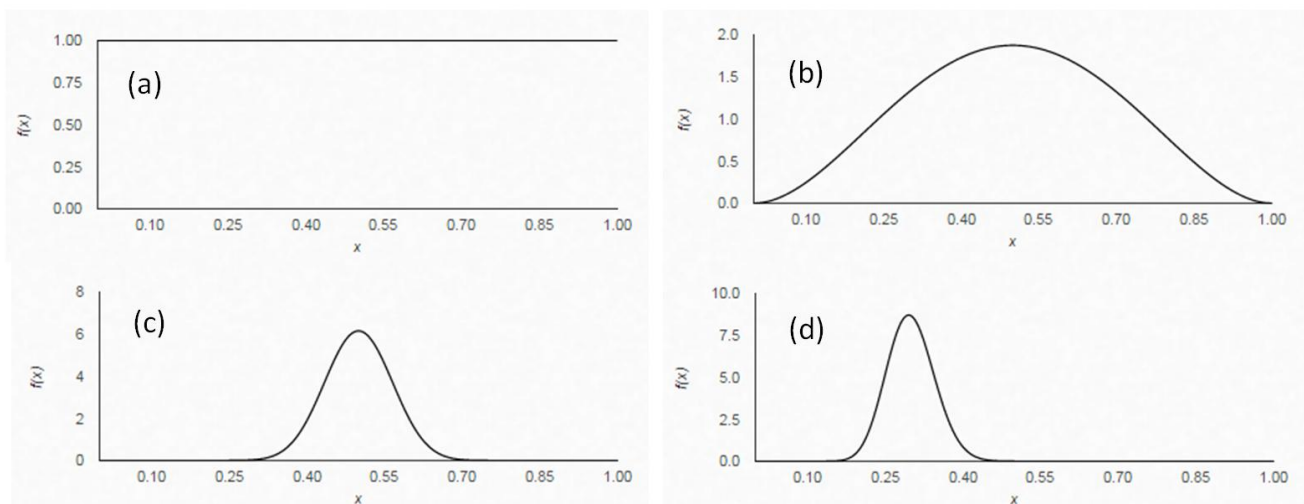
Dodavanje *pseudocount*-a se tumači kao da za $P(x)$ postoji apriori verovanje da $P(x)$ može da uzme bilo koju vrednost (kada je $\alpha = 1$) ili da imamo predznanje da $P(x)$ treba da bude u određenom opsegu vrednosti (vidi Sliku 1). Veće vrednosti *pseudocount*-a sugerišu naše jaše apriori verovanje i zahtevaju više podataka kako bi se to verovanje promenilo. α može biti različito za različite klase, a ukoliko je isto, to daje apriori jednaku težinu svim klasama.

U slučaju binarne promenljive x , ova apriori verovatnoća se može modelovati *Beta* distribucijom (Slika 1), a u slučaju viševrednosne promenljive, ona odgovara *Dirihle* distribuciji.

Ovakva rezonovanja (korišćenje apriori verovanja) su osnov za Bajesov pristup u Mašinskom učenju, i pojavljuju se u velikom broju modela¹.

Smoothing tehnika se uglavnom koristi kada imamo manje podataka, ali ne škodi i kada se koristi sa puno podataka. Najzad, vrlo je korisna kada imamo puno podataka, ali je puno šuma u njima, tj. ne može im se potpuno verovati.

¹ ovo je takođe protivno klasičnoj (frekvencionističkoj) statistici, koja ne prihvata takvo subjektivno "verovanje", kao ni da parametre treba tretirati kao slučajne promenljive



Slika 1. Apriori verovanje o verovatnoći $P(x)$ prikazano beta-distribucijom. (a) Neinformativna raspodela sa $\alpha=1$; (b) Slabo apriori verovanje u jednakost klasa $\alpha=3$; (c) Jako apriori verovanje u jednakost klasa $\alpha=30$; (d) Jako apriori verovanje u odnos klasa $.3/.7$, sa $\alpha_1=30$ i $\alpha_2=70$

Numeričke promenljive

U slučaju kada su promenljive X_i numeričke, verovatnoća $P(X_i|C_k)$ se može modelovati nekom kontinualnom distribucijom, a najčešće *Normalnom* distribucijom²:

$$P(X_i|C_k) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i-\mu)^2}{2\sigma^2}}$$

gde su parametri koje treba proceniti iz podataka (uzorka) μ i σ , i to za svaku klasu posebno.

Naive Bayes model, iako za cilj ima da proceni verovatnoću $P(C|X)$ radi predviđanja klase, on to radi računajući (modelujući) verovatnoće $(X|C)$, što ga čini *generativnom modelom*, nasuprot *diskriminacionim modelima* koje ćemo kasnije obraditi. To takođe znači da se sa *Naive Bayes* modelom mogu generisati slučajni (random, virtuelni) uzorci iz svake klase, što ponekad može biti korisno.

Korisne komande

Matlab	Python	R
histc	pandas.crosstab	?
HashMap (java)	dict	?
normpdf	scipy.stats.norm.pdf	?
...

² Normalna distribucija je u skladu sa pretpostavkom *Naive bayes* modela, koja kaže da su atributi X_i međusobno nezavisni ukoliko poznajemo klasu (uslovno nezavisni). Normalna raspodela (uslovljena klasom: $P(X_i | C_k)$) takođe pretpostavlja da na X_i najviše utiče klasa, a ostatak čine puno "sitnih" faktora koji konstituišu normalnu rasporedu (po centralnoj graničnoj teoremi).