

Lecture 2

Probability theory: basics

Event A : $0 \leq P(A) \leq 1$ is the prob. of A

$$P(\bar{A}) = 1 - P(A)$$

Joint prob. : $P(A \wedge B) = P(A, B)$
"A and B"

If events A & B are indep.,

$$P(A, B) = P(A) P(B)$$

$$P(A \vee B) = P(A) + P(B) - P(A, B)$$

"A or B"

Mutually exclusive events: $P(A, B) = 0$,

$$P(A \vee B) = P(A) + P(B).$$

Conditional prob.:

$$P(A, B) = P(B|A) P(A), \text{ or}$$

$$P(B|A) = \frac{P(A, B)}{P(A)}.$$

Conditional independence: $P(A, B|C) = P(A|C) P(B|C)$
 $A \perp B | C$

Discrete random variables:

$x \in X$ finite or countably infinite

$$p(x) = P(X=x), \quad 0 \leq p(x) \leq 1,$$

↑
discrete rv

$$\sum_{x \in X} p(x) = 1.$$

Continuous random variables:

Cumulative distribution function (cdf):

$$P(x) = P(X \leq x) \Rightarrow P(a < X \leq b) = P(b) - P(a)$$

Probability density function (pdf):

$$p(x) = \frac{dP(x)}{dx}.$$

$$\text{Then } P(a < X \leq b) = P(b) - P(a) = \int_a^b dx p(x).$$

$$P(x < X < x + dx) \approx p(x) dx.$$

Marginal distribution:

$$p(x) = \sum_y p(x, y) \text{ or}$$

$$p(x) = \int dy p(x, y)$$

Product rule: $p(x, y) = p(x)p(y|x)$

D variables: $p(x_1, \dots, x_D) = p(x_1)p(x_2|x_1) \times p(x_3|x_1, x_2) \times \dots \times p(x_D|x_1, \dots, x_{D-1})$

Moments:

mean

$$E[X] = \sum_x x p(x) \quad \text{or} \quad E[X] = \int dx x p(x)$$

$$E[aX + b] = aE[X] + b$$

\uparrow \uparrow
const

$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i], \text{ etc.}$$

$E(x)$

Variance

$$V[X] = E[(X - \mu)^2] = E[X^2] - \mu^2.$$

$\stackrel{\text{"}}{=} \sigma^2$

$$\text{std}[X] = \sqrt{V[X]} = \sigma.$$

$$V[aX + b] = a^2 V[X];$$

for indep. rv we have

$$V\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n V[X_i]$$

Mode

$$\tilde{x}^* = \underset{\tilde{x}}{\operatorname{argmax}} p(\tilde{x})$$

\uparrow may not be unique (multimodal distributions)

Conditional expectation:

$$E[X|Y]$$

Note that $E_Y [E_X[X|Y]] =$

$$= E_Y \left[\int dx x p(x|y) \right] \ominus \text{(dead ends)}$$

$$\begin{aligned} \textcircled{=} \int dy p(y) [\int dx \underset{p(x|y)}{\approx} p(x|y)] &= \\ = \int dx x \underbrace{\int dy p(x|y) p(y)}_{\substack{p(x,y) \\ p(x)}} &= E[x] \end{aligned}$$

Bayesian probabilities

Consider $\begin{cases} H = \text{hidden rv (e.g., a model param)} \\ Y = \text{observed rv} \end{cases}$

Then $p(h, y) = p(h|y)p(y) = p(y|h)p(h)$, or

$$p(h|y) = \frac{p(y|h)p(h)}{p(y)}, \text{ where}$$

$p(y) = \sum_h p(h, y)$ is the marginal likelihood (aka model evidence)

$$\begin{cases} p(y|h) = \text{likelihood} \\ p(h) = \text{prior} \\ p(h|y) = \text{posterior} \end{cases}$$

Ex. The Monty Hall problem

Doors: 1, 2, 3

a single prize is hidden behind one of the doors.

You select one door (but do not open it). Then the host will open one of the other 2 doors, without the prize behind it. Then you ~~can~~ can choose again: stick with your choice or switch to the remaining closed door.

(only other
door)

Then you get the prize behind your final choice of door.

For example, you initially choose door 1, then the host opens door 3 (with nothing behind it). Should you stick w/door 1 or switch to door 2? Does it make any difference?

Define H_i = hypothesis that prize is behind door $i = 1, 2, 3$

choose $P(H_1) = P(H_2) = P(H_3) = \frac{1}{3}$ for the priors

$\underbrace{\text{Conditional likelihoods:}}_{\substack{\text{after choosing door } i \\ \text{host selects at random}}}$

door opened by host $\left\{ \begin{array}{l} P(Y=2|H_1) = \frac{1}{2}, \quad P(Y=3|H_1) = \frac{1}{2} \\ P(Y=2|H_2) = 0, \quad P(Y=3|H_2) = 1 \\ P(Y=2|H_3) = 1, \quad P(Y=3|H_3) = 0 \end{array} \right.$

Similar for the initial choices of door 2 or 3.

Now, consider

$$P(H_i | Y=3) = \frac{P(Y=3 | H_i) P(H_i)}{P(Y=3)}$$

$$\begin{aligned} P(Y=3) &= \sum_{i=1}^3 P(Y=3 | H_i) P(H_i) \\ &= \frac{1}{3} \left[\frac{1}{2} + 1 + 0 \right] = \frac{1}{2}. \end{aligned}$$

$$\left\{ \begin{array}{l} P(H_1 | Y=3) = \frac{\frac{1}{2} \times \frac{1}{3}}{\frac{1}{2}} = \frac{1}{3}, \\ P(H_2 | Y=3) = \frac{1 \times \frac{1}{3}}{\frac{1}{2}} = \frac{2}{3}, \\ P(H_3 | Y=3) = \frac{0 \times \frac{1}{3}}{\frac{1}{2}} = 0. \leftarrow \text{no car behind door 3} \end{array} \right.$$

Thus, it's twice as good to switch to door 2 than to stay with door 1 (!)

what if $Y=2$?

$$P(Y=2) = \frac{1}{3} \left[\frac{1}{2} + 0 + 1 \right] = \frac{1}{2},$$

$$\left\{ \begin{array}{l} P(H_1 | Y=2) = \frac{\frac{1}{2} \times \frac{1}{3}}{\frac{1}{2}} = \frac{1}{3}, \\ P(H_2 | Y=2) = \frac{0 \times \frac{1}{3}}{\frac{1}{2}} = 0, \\ P(H_3 | Y=2) = \frac{1 \times \frac{1}{3}}{\frac{1}{2}} = \frac{2}{3}. \end{array} \right.$$

Now it's twice as good to switch to door 3 rather than stay w/door 1.

all other combinations can be obtained by permutation of door indices.

Standard distributions

Bernoulli

$$\begin{cases} p(Y=1) = \theta, \\ p(Y=0) = 1-\theta. \end{cases} \Rightarrow p(y|\theta) = \theta^y(1-\theta)^{1-y}$$

$$y \in \{0, 1\}$$

$$0 \leq \theta \leq 1$$

Binomial

observe N Bernoulli trials,

define $s = \sum_{n=1}^N \mathbb{I}(Y_n = 1)$ - the total # successes

Then $p(s|N, \theta) = \binom{N}{s} \theta^s (1-\theta)^{N-s}$, where

$$\binom{N}{s} = \frac{N!}{(N-s)! s!}$$

Sigmoid

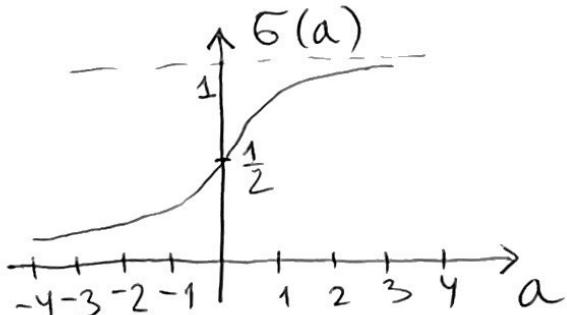
$\zeta=2$ classifier:

$$p(y|\vec{x}, \vec{\theta}) = \underbrace{\text{Ber}(y | \xi(f(\vec{x}; \vec{\theta})))}_{\text{Bernoulli}}$$

$\xi(a)$ sigmoid f^n

$$\xi(a) = \frac{1}{1+e^{-a}} = \frac{e^a}{1+e^a}$$

$$\begin{cases} \frac{d}{da} \xi(a) = \xi(a)[1-\xi(a)], \\ 1-\xi(a) = \xi(-a) \end{cases}$$



More explicitly,

$$p(y=0) = 1 - \frac{1}{1+e^{-f(\vec{x}; \vec{\theta})}} = \frac{1}{1+e^{f(\vec{x}; \vec{\theta})}},$$

$$p(y=1) = \frac{1}{1 + e^{-f(\vec{x}; \vec{\theta})}}$$

$f(\vec{x}; \vec{\theta})$ can be any function $\in \mathbb{R}$.

For example, in logistic regression

$$f(\vec{x}; \vec{\theta}) = \vec{\theta}^T \cdot \vec{x} + b.$$

The inverse of $p = \sigma(a)$ is

$$a = \sigma^{-1}(p) = \underbrace{\log\left(\frac{p}{1-p}\right)}_{\text{"log odds"}}$$

$$\text{Indeed, } \log \frac{e^a / (1 + e^a)}{1 - \frac{e^a}{1 + e^a}} = \log e^a = a.$$

$$\mathbb{I}(a) = \begin{cases} 1, & a = \text{True} \\ 0, & a = \text{False} \end{cases}$$

Categorical $\text{Cat}(y | \vec{\theta}) = \prod_{c=1}^C \theta_c \mathbb{I}(y=c)$

$$\vec{\theta} = (\theta_1, \theta_2, \dots, \theta_C)$$

$$\sum_{c=1}^C \theta_c = 1$$

$$\text{Cat}(y=c | \vec{\theta}) = \theta_c$$

One-hot encoding: $y=c \Rightarrow \vec{y} = (0, 0, \dots, \underset{\uparrow}{1}, 0, \dots)$

@ position c ,
all other entries
 \emptyset

Then $\text{Cat}(\vec{y} | \vec{\theta}) = \prod_{c=1}^C \theta_c^{y_c}$

Multinomial

$$M(\vec{n}_c | N, \vec{\theta}) = \binom{N}{n_1 \dots n_c} \prod_{c=1}^C \theta_c^{n_c},$$

where

$$N = \sum_{c=1}^C n_c,$$

$$\vec{n}_c = (n_1, n_2, \dots, n_c)$$

$$\binom{N}{n_1 \dots n_c} = \cancel{\frac{N!}{n_1! n_2! \dots n_c!}} = \frac{N!}{n_1! n_2! \dots n_c!}$$

Softmax $C \geq 2$ classifier:

$$\begin{aligned} p(y=c | \vec{x}, \vec{\theta}) &= \text{Cat}(y=c | \text{softmax}(\vec{f}(\vec{x}; \vec{\theta}))) = \\ &= \text{softmax}(\vec{f}(\vec{x}; \vec{\theta}))_c. \end{aligned}$$

Here, $\text{softmax}(\vec{a}) = \left(\frac{e^{a_1}}{Z}, \dots, \frac{e^{a_C}}{Z} \right)$,
 $\vec{a} = (a_1, a_2, \dots, a_C)$ where
 $Z = \sum_{c=1}^C e^{a_c}.$

Note that $0 \leq \text{softmax}(\vec{a})_c \leq 1$, and

$$\sum_{c=1}^C \text{softmax}(\vec{a})_c = 1.$$

For ex., in multiclass logistic regression

$$\vec{f}(\vec{x}; \vec{\theta}) = \vec{W}_{C \times D} \vec{x}_D + \vec{b}_C, \text{ s.t.}$$

$$\begin{aligned} p(y=c | \vec{x}; \vec{\theta}) &= \text{Cat}(y=c | \text{softmax}(\vec{W} \vec{x} + \vec{b})) = \\ &= \frac{e^{f_c(\vec{x}; \vec{\theta})}}{Z}. \end{aligned}$$