

Probability Theory

Frequentist interpretation $p = \frac{f}{N}$

$$p(\text{head}) = \frac{f}{N}$$

In N tosses we have roughly $N/2$ heads

of times something happen

↓
of trials

Bayesian interpretation : p is about our beliefs.

$$p(\text{head}) = \frac{1}{2}$$

means we believe head and tail are equally likely.

If we observe something to the contrary we will change our opinion.

However the math of probability is still the same in both approaches.

Key ideas needed for Bayesian point of view

Joint probability

$$P(A, B) = P(A \cap B) = P(A|B) P(B)$$

Conditional probability

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

Probability distributions; Discrete distribution

Random variable X $x \in X$

$$p(x) = P(X=x)$$

Joint distribution $P(X=x, Y=y)$

$$P(X=x | Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

Bayes' Rule / Theorem

$$\frac{P(Y=y | X=x) P(X=x)}{\sum_{x'} P(X=x') P(Y=y | X=x')}$$

Example 1:

Test

Disease	0	1	
0	0.9	0.1	1
1	0.2	0.8	1
-	-	-	

$$P(y|x)$$

$$\times 0.996$$

$$P(x,y)$$

$$P(x) = .004$$

Test

Disease	0	1	
0	0.8664	0.0996	0.996
1	0.0008	0.0032	0.004
-	0.8672	0.1028	

$$P(x|y)$$

Test

Disease	0	1	
0	0.9991	0.9689	-
1	0.0009	0.0311	-
-	1	1	

normalize

$$\begin{aligned}
 & 0.004 \times 0.8 \\
 = & \frac{0.004 \times 0.8 + 0.996 \times 0.1}{0.0032 + 0.0996} \\
 = & \frac{0.0032}{0.0032 + 0.0996} \approx 0.031
 \end{aligned}$$

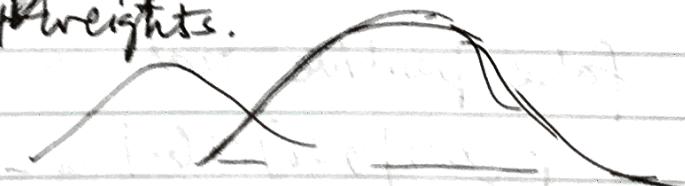
Ignoring this prior gives the wrong impression
 \Rightarrow base rate fallacy.

Example:?

Generative classifiers

$$P(y=c|x) = \frac{P(y=c) P(x|y=c)}{\sum_{c'} P(y=c') P(x|y=c')}$$

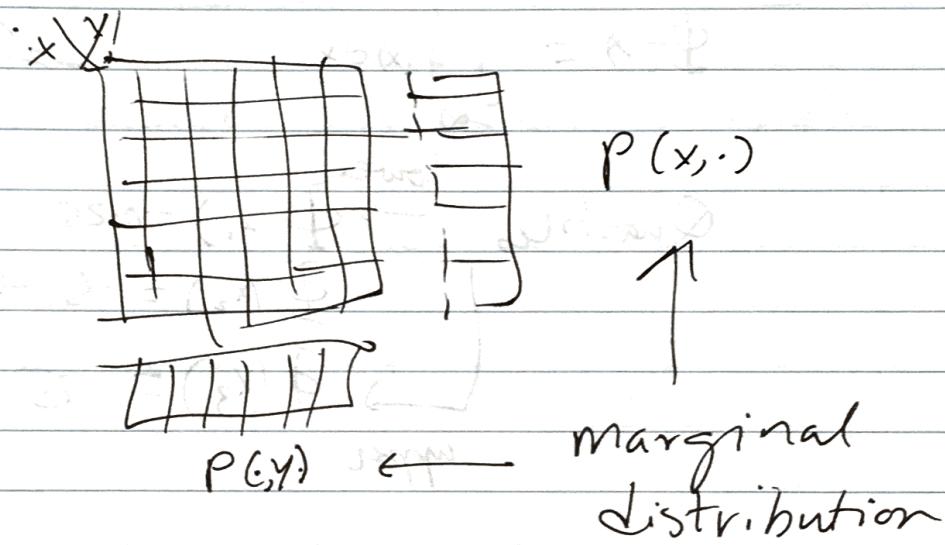
Consider data cloud coming from
 two normal distribution with
 different weights.



Independence and conditional independence

$X \perp Y$ means $P(X, Y) = P(X)P(Y)$

In that case $P(X|Y) = P(X, \cdot)$



Example: $P(X, Y) = P(X)P(Y)$

Conditional independence

$$X \perp Y | Z \quad P(X, Y | Z) = P(X | Z)P(Y | Z)$$

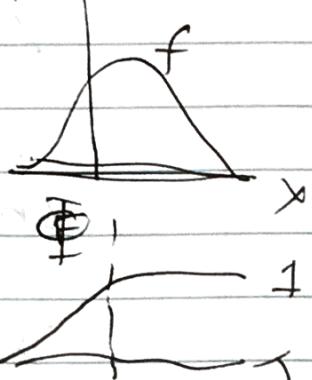
$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$

Continuous distribution

$$P(a < X \leq b) = \int_a^b f(x) dx$$

Cumulative

$$\Phi(x) = \int_{-\infty}^x f(x) dx$$



Quantiles $\rightarrow \Phi(x_1) = 0.25$

$\Phi(x_2) = 0.5 \rightarrow \text{Median}$

$\Phi(x_3) = 0.75$

upper

Mean: $E[X] = \sum_x x p(x)$ or $\int x p(x) dx$

Variance: $\text{Var}[X] = E[(X - E[X])^2]$

$$= E[X^2] - (E[X])^2$$

$$\sum_x f(x)(x - \mu)^2 = \sum_x f(x)x^2 - 2 \sum_x f(x)x\mu + \mu^2 \sum_x f(x)$$

\uparrow
 $E[X]$

$$= E[X^2] - 2\mu^2 + \mu^2 = E[X^2] - \mu^2$$

$$= E[X^2] - [E[X]]^2$$

$$\text{std}(X) = \sqrt{\text{Var}[X]}$$

~~$E[X^2] = \theta + (1-\theta)^2 E[X^2]$~~

Some common discrete distribution

Bernoulli distribution

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

$$\begin{aligned} \text{Ber}(X|\theta) &= \theta & \text{if } x=1 \\ &= 1-\theta & \text{if } x=0 \end{aligned} \quad \begin{array}{l} \text{Success} \\ \text{failure} \end{array} \quad \begin{array}{l} \text{like } p \\ q=1-p \end{array}$$

Note that $X^2 = X$

$$E[X] = \theta \quad \text{Var}[X] = ?$$

$$\text{Var}[X] = E[X^2] - E[X]^2 = E[X] - E[X]^2$$

$$= \theta - \theta^2 = \theta(1-\theta)$$

Binomial distribution ($= pg$)

$$X \sim \text{Bin}(n, \theta) \quad X \in \{0, 1, \dots, n\}$$

successes in n trials.

$$\text{Bin}(k|n,\theta) \triangleq \binom{n}{k} \theta^k (1-\theta)^{n-k}$$

$k = x_1 + \cdots + x_n$
 $\underbrace{\hspace{1cm}}_{\text{independent}} \uparrow$
 Bernoulli variables

$$E[x] = ? \quad \text{var}[x] = ?$$

Related Concept
 Prob. Generating Function

$$G(z) = E[z^x] = \sum_{x=0}^{\infty} p(x) z^x$$

Moment Generating Function

$$M_x(t) = E[e^{tx}] = \sum_x p(x) e^{tx}$$

or $f(x) e^{tx}$

Characteristic Function

$$\varphi_x(t) = E[e^{itx}] = \sum_x p(x) e^{itx}$$

or $\int f(x) f(tx) e^{itx}$

For $\text{Bin}(n, \theta)$

$$G(z) = \sum_{k=0}^n \binom{n}{k} \theta^k (1-\theta)^{n-k} z^k = [\theta z + (1-\theta)]^n$$
$$= [1 + \theta(z-1)]^n$$

$$M_x(t) = \{1 + \theta(e^t - 1)\}^n$$

etc.

$$\frac{d^n}{dt^n} M_x(t) = \frac{d^n}{dt^n} \sum_x p(x) e^{tx} = \sum_x p(x) x^n e^{tx}$$

$$\left. \frac{d^n}{dt^n} M_x(t) \right|_{t=0} = \sum_x p(x) x^n = E[X^n]$$

So, let us get ~~to~~ down to business

$$\frac{d}{dt} [1 + \theta(e^t - 1)]^n = n [1 + \theta(e^t - 1)]^{n-1} \theta e^t$$

$$\begin{aligned} \frac{d^2}{dt^2} [1 + \theta(e^t - 1)]^n &= n [1 + \theta(e^t - 1)]^{n-1} \theta e^t \\ &\quad + n(n-1) [1 + \theta(e^t - 1)]^{n-2} \theta^2 e^{2t} \end{aligned}$$

Setting $t=0$ $E[X] = n\theta$, $E[X^2] = n\theta + n(n-1)\theta^2$

$$\begin{aligned}\text{Var}(x) &= E[x^2] - E[x]^2 \\ &= n\theta + n(n-1)\theta^2 - n^2\theta^2 \\ &= n(\theta - \theta^2) = n\theta(1-\theta)\end{aligned}$$

Known results.

Multinomial distribution

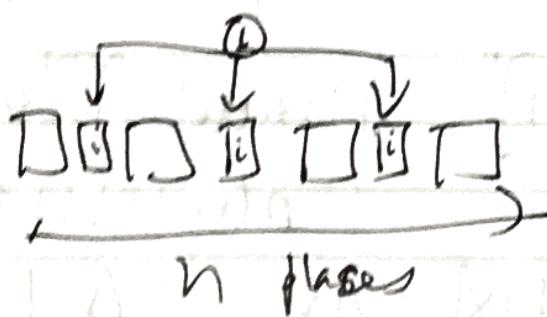
Think of k -sided die with $\theta_1, \dots, \theta_k$ as the probabilities of the different sides falling

$$\binom{n}{x_1, \dots, x_k} = \frac{n!}{x_1! \cdots x_k!}$$

$$\sum_{i=1}^k x_i = n$$

$$M_n(x | n, \theta) = \binom{n}{x_1, \dots, x_k} \prod_{j=1}^k \theta_j^{x_j}$$

vector



$$\frac{n!}{x_1! x_2! \cdots x_i! \cdots x_k!}$$

~~Example 2 DNA~~

$$n = 1$$

$$M(x|1, \theta) = \prod_{j=1}^k \theta_j \cdot I(x_j = 1)$$

$x = (0, 0, 1, 0, 0, 0)$, for example

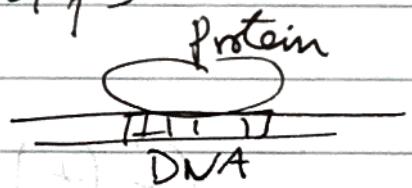
1 hot encoding

$$\text{Cat}(x|\theta) \triangleq M_n(x|1, \theta)$$

Categorical distribution

Example: DNA Seq. motifs

A, C, G, T



At different locations we have different multinomial distributions

ATTCA

ATCTA

TTTTA

ATACA

A motif is specified by
on a matrix with
entries

$$(\theta_A, \theta_C, \theta_G, \theta_T)$$

$$(0.25, 0.25, 0, 0.5)$$

θ_{ib}
locator bases

Poisson distrn

Fun stuff of taking limits

For $\text{Bin}(n, \theta)$ $G(z) = [1 + \theta(z-1)]^n$

If we take the limit $n \rightarrow \infty$

with $\theta_n = \frac{\lambda}{n}$

$$G(z) = [1 + \frac{\lambda}{n}(z-1)]^n$$

$$\rightarrow e^{\lambda(z-1)} = G_\lambda(z)$$

$$G_\lambda(z) = \sum_{k=0}^{\infty} e^{-\lambda} \frac{\lambda^k}{k!} z^k$$

This is Poisson distn

$$\text{Poi}(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$$

$$M_x(t) = e^{\lambda(e^t - 1)}$$

$$E[X] = \lambda \quad \text{Var}[X] = \lambda$$

Empirical distrⁿ

Data ~~Set~~ $\mathcal{D} = \{x_1, \dots, x_N\}$

$P_{\text{emp}}(A)$

$$= \frac{1}{N} \sum_{i=1}^N \delta_{x_i}(A) = \frac{1}{N} \sum_{i=1}^N I(x_i \in A)$$



If we assign weights $w_i \geq 0$ with $\sum w_i = 1$

$$\sum_{i=1}^N w_i \delta_{x_i}(A)$$

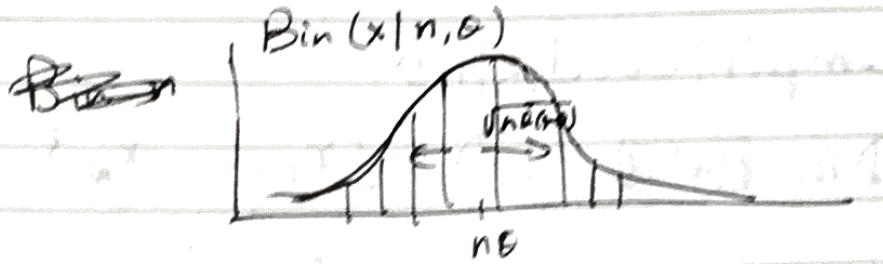
Now it is time to move to continuous distributions.

For that, consider the limit of Binomial distribution

$$z = \frac{x - n\theta}{\sqrt{n\theta(1-\theta)}}$$

What is the distribution of z , when n is large

Note that $E[z] = 0$ $\text{Var}[z] = 1$



Now note that the characteristic function has the following ~~properties~~ properties

$$\begin{aligned}\Phi_{X+a}(t) &= E[e^{it(X+a)}] = \boxed{\text{e}^{ita}} E[e^{itX}] \\ &= e^{ita} \Phi_X(t)\end{aligned}$$

Also $\Phi_{cx}(t) = E[e^{itcx}]$

$$\begin{aligned}&= \Phi_X(ct)\end{aligned}$$

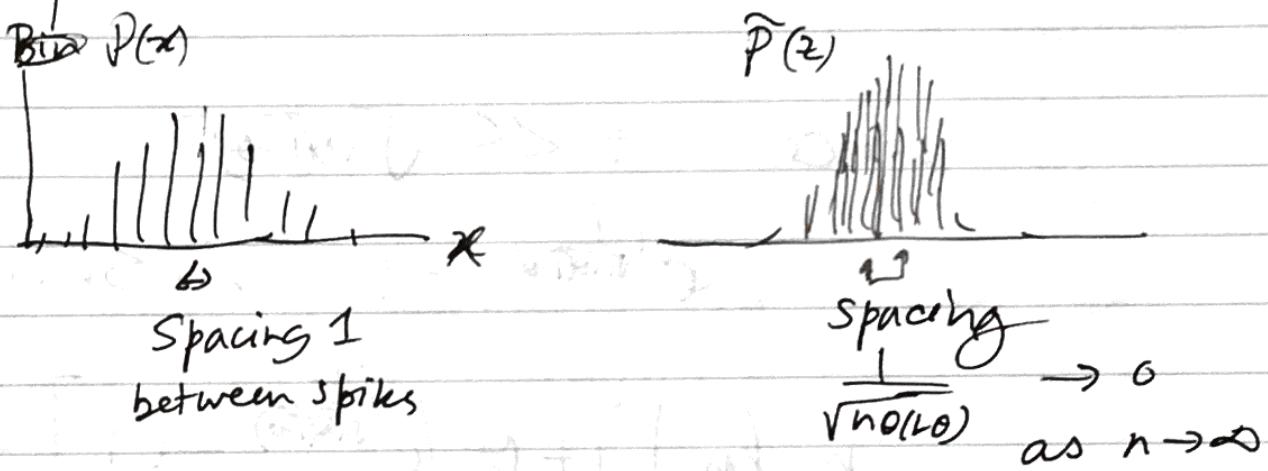
If $X \sim \text{Bin}(n, \theta)$

$$\begin{aligned}\Phi_X(t) &= E[e^{itX}] \\ &= (1 + \theta(e^{it} - 1))^n\end{aligned}$$

~~Z~~ $= \frac{X - n\theta}{\sqrt{n\theta(1-\theta)}}$

$$\begin{aligned}\varphi_z(t) &= \varphi_{X-n\theta}\left(\frac{t}{\sqrt{n\theta(1-\theta)}}\right) \\ &= \varphi_x\left(\frac{t}{\sqrt{n\theta(1-\theta)}}\right) e^{-\frac{i\pi n\theta t}{\sqrt{n\theta(1-\theta)}}}\end{aligned}$$

So far, this is exact.



Approximating a continuous distribution as $n \rightarrow \infty$.

$\varphi(t)$ is essentially the Fourier transformation of $P(z)$. The variable t are like momenta.

To produce all this small scale structure $1/\sqrt{n}$ we need $t \sim \sqrt{n}$. What if we consider $t \sim 1 \ll \sqrt{n}$.

$$c_f(t) = \left[1 + O\left(e^{\frac{it}{n\theta(1-\theta)}} - 1\right) \right]^n e^{-\frac{int}{n\theta(1-\theta)}}$$

~~$$\left[e^{it} \left(1 + O\left(\frac{t}{\sqrt{n\theta(1-\theta)}}\right) \right)^n e^{-\frac{int}{n\theta(1-\theta)}} \right]$$~~

$$= e^{n \ln \left[1 + O\left(e^{\frac{it}{n\theta(1-\theta)}} - 1\right) \right] - \frac{int}{n\theta(1-\theta)}}$$

Since $t \ll \sqrt{n\theta(1-\theta)}$

$e^{\frac{it}{n\theta(1-\theta)}} - 1$ is small

$$e^{n \ln \left[1 + O\left(e^{\frac{it}{n\theta(1-\theta)}} - 1\right) \right] - \frac{int}{n\theta(1-\theta)}}$$

$$= n O\left(e^{\frac{it}{n\theta(1-\theta)}} - 1\right) = \frac{n\theta^2}{2} \left(\frac{it}{n\theta(1-\theta)}\right)^2 - \frac{int}{n\theta(1-\theta)}$$

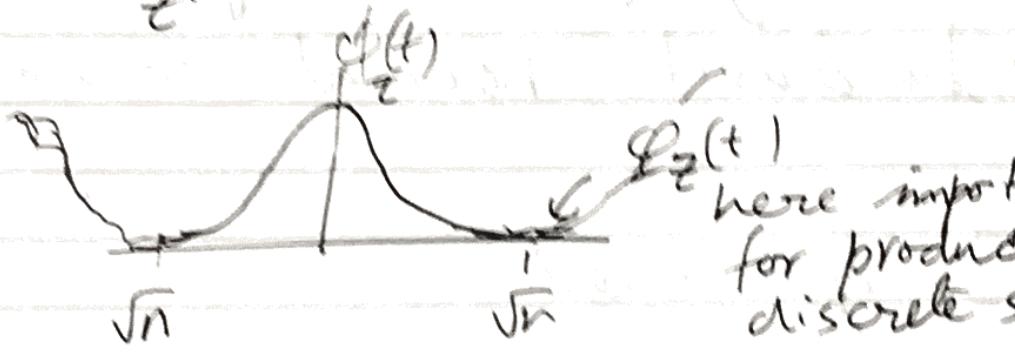
$$= n O\left\{ \frac{it}{\sqrt{n\theta(1-\theta)}} + \frac{1}{2} \left(\frac{it}{\sqrt{n\theta(1-\theta)}} \right)^2 \right\} - \frac{n\theta^2}{2} \left(\frac{it}{\sqrt{n\theta(1-\theta)}} \right)^2 - \frac{int}{\sqrt{n\theta(1-\theta)}}$$

$$= -\frac{n\theta(1-\theta)}{2} \frac{t^2}{n\theta(1-\theta)} + O\left(\frac{n t^3}{(n\theta(1-\theta))^{3/2}}\right)$$

$$= -\frac{t^2}{2} + O\left(\frac{n^{\frac{3}{2}}}{n^{1/2}}\right)$$

So for $|t| < n^{1/2}$

$$\varphi_2(t) \approx e^{-t^2/2}$$



here important
for producing
discrete spikes

If we only keep $e^{-t^2/2}$, the inverse Fourier transform is going to be smoother

$$\begin{aligned} & \int_{-\infty}^{\infty} e^{-t^2/2} e^{-itz} \frac{dt}{2\pi} \\ &= \int_{-\infty}^{\infty} e^{-\frac{(t+iz)^2}{2}} dt \frac{e^{-z^2/2}}{2\pi} \quad \boxed{\text{B}} \\ &= \frac{\sqrt{2\pi}}{\sqrt{2\pi}} e^{-z^2/2} \end{aligned}$$

That is the probability density function of the Standard normal variate

In general

Normal distribution

Pdf

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$E[X] = \mu, \text{ var}[X] = \sigma^2$$