

(BRK)

Probabilities \rightarrow Much of successful machine learning algorithms
are based on probabilistic reasoning

Terms
Random Variable \rightarrow denotes something about which we are uncertain.
 \rightarrow outcome of a random experiment

Ex. $A =$ True if a randomly drawn person is female

$P(A) =$ the fraction of possible worlds in which A is true.

Random variable A is a function of sample space S
defined as $A : S \rightarrow \{0, 1\}$.
Set of all possible worlds.

Sample Space \rightarrow all set of possible worlds

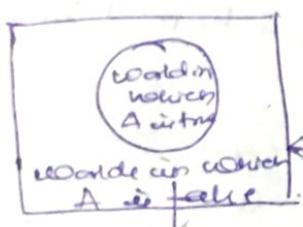
Ex. Set of students in our class

Random variable Gender: $S \rightarrow \{\text{m}, \text{f}\}$

Height: $S \rightarrow \text{Reals}$

Event is a subset of S Ex. subset of S where $G = \text{f}$

visualizing A in Venn diagram.



Sample Space of all possible worlds
Area is 1

Axioms of probability

$0 \leq P(A) \leq 1$ Probability is between 0 to 1

false = 0
true = 1

$$P(\text{True}) = 1$$

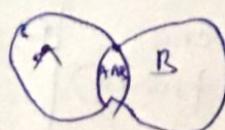
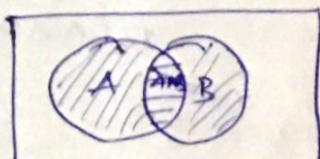
$$P(\text{false}) = 0$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

when gambling based on uncertainty formalism A,
you can be exploited by an opponent

if

Your uncertainty formalism A violates these axioms.

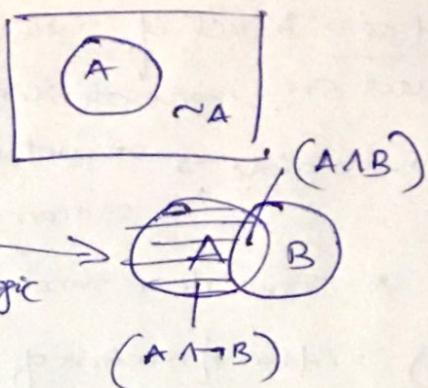


Elementary probability

$$P(\sim A) + P(A) = 1$$

marginal joint

$$P(A) = P(A \cap B) + P(A \cap \sim B)$$



If A & B are boolean, apply in boolean logic

$$A = [A \text{ and } (B \text{ or } \sim B)]$$

$$= [(A \text{ and } B) \text{ or } (A \text{ and } \sim B)]$$

Can be reexpressed as

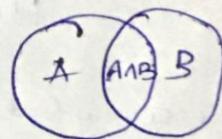
$$D(A) = D(A \text{ and } B) + D(A \text{ and } \sim B) - P((A \text{ and } B) \text{ and } (A \text{ and } \sim B))$$

$$P(A) = P(A \text{ and } B) + P(A \text{ and } \sim B) - P(A \text{ and } B \text{ and } A \text{ and } \sim B)$$

Conditional probability

$P(A|B)$ = what is the probability that A is true given B is true

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



Ex. what is the probability that you have blue eyes given that you are female.

Let say A is blue eyes
 B is female.

$$\text{what is } P(\text{blue eyes} | \text{female}) = \frac{b}{b+n}$$

Generalized The chain rule:

$$P(A \cap B) = P(A|B) P(B) = P(B|A) P(A) \xrightarrow{\text{Chain Rule}} P(A \cap B) = P(A|B) P(B)$$

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

$$= P(A|B) P(B|C) P(C)$$

Hence.

from ①

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

Bayes' rule

we call

$P(A)$ \Rightarrow Prior

$P(A|B)$ \Rightarrow Posterior

SOMK - 2

$$\text{otherwise } P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B|A) \cdot P(A) + P(B|\neg A) \cdot P(\neg A)} \quad // P(B) \text{ rewritten.}$$

$$P(A|B \wedge x) = \frac{P(B|A \wedge x) \cdot P(A \wedge x)}{P(B|x)}$$

x could be anything.

Ex. A = you have flu
 B = you just coughed.

Assume $P(A) = 0.05$

$$P(B|A) = 0.8$$

$$P(B|\neg A) = 0.2$$

What is $P(\text{flu}|\text{cough}) = \frac{P(\text{cough}|\text{flu}) \cdot P(\text{flu})}{P(\text{cough}|\text{flu}) \cdot P(\text{flu}) + P(\text{cough}|\neg \text{flu}) \cdot P(\neg \text{flu})}$

$$= \frac{0.8 \times 0.05}{0.8 \times 0.05 + 0.2 \times 0.95} = 0.17$$

\Rightarrow Posterior Probability

we are interested in function approximation, coming up with training algorithms that can learn functions $X \rightarrow Y$
 learning

$$f : X \rightarrow Y$$

email spam.

Equivalent thing is $P(Y|x)$

$$P(\text{spam}=Y| \text{email})$$

What is the probability distribution of probability values.

The single most important thing to learn is joint probability distribution. \rightarrow Assignment of probability to a collection of random variables.

Joint distribution recipe with M variables.

- ① Make a truth-table listing all combination of values.
 M boolean variable $\Rightarrow 2^M$ rows
- ② for each combination of values, say how probable it is
- ③ probabilities must sum to 1

If you know joint distribution, then there is nothing that you cannot answer on probabilities, Joint, Conditional or any subset

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

Joint probability
 $P(A, B) = P(A=1, B=1)$ (irrespective of C)
 $= 0.25$ to 1

$$P(A) = \sum_{\text{rows matching } A} P(\text{row})$$

$$P(A/B) = \frac{\sum_{\substack{\text{rows matching } A \\ \text{and } B}} P(\text{row})}{\sum_{\text{rows matching } B} P(\text{row})} = \text{Conditional probability}$$

$$= \frac{0.25 + 0.1}{0.25 + 0.1 + 0.1 + 0.05} = 0.5$$

$$P(\sim A/B) = 1 - P(A/B) \quad \text{for given Value being same.}$$

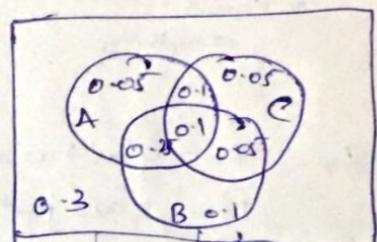
if will sum to 1

Joint probability sounds like a solution to Learning $f: X \rightarrow Y$
or $P(Y/X)$

But problems are

- There are simple boolean, but real world (2^M is huge)
- Computational values.
- Also it is not easy to learn the probability distribution.

Ex



A	B	C	Prob
0	0	0	0.3
0	0	1	0.05
0	1	0	0.1
0	1	1	0.05
1	0	0	0.05
1	0	1	0.1
1	1	0	0.25
1	1	1	0.1

Ex.

SRNK (8)

Joint dist with 100 attributes / reasonable in real world

$$\# \text{rows} = 2^{100} (\approx 1000^{10} \text{ or } 10^{30})$$

people on earth (10^9)

There are more rows but the table will have fraction of ones with
e.g. example (20.999)

Data sparsity is a problem.

What to do?

1. Be smart about how we estimate probability from sparse data
 - Maximum likelihood estimates
 - Maximum a posteriori estimates
2. Be smart about how to represent joint distribution.
 - Bayes network, graphical models.

Being Smart about Estimating Probabilities

Ex. Coinflip $H=1$ $T=0$
 d_{heads} d_{tails}

$$P(X=1) ?$$

$$P(X=0) ?$$

Algorithm for estimating the probability by flipping repeatedly

$$P(X=1) = \frac{d_1}{d_1 + d_0}$$

Estimating $\Theta = P(X=1)$

Test A: 100 flips 51 heads 49 tails \rightarrow lot of data.

$$\hat{P}(X=1) = \frac{d_1}{d_1 + d_2} = \frac{51}{100} \approx 0.51$$

Test B: 3 flips 2 heads 1 tail \rightarrow This is problem with simple universe, the true probability might be skewed.
 $\hat{P}(X=1) = \frac{2}{3} = 0.67$

Indication of problem with small data.

Case c: (Online learning)

Keep flipping, want single learning algorithm that gives reasonable estimate after each flip

$$\theta = P(X=1)$$

$d_1 = \# \text{ of observed heads } (X=1)$

$d_0 = \# \text{ of observed tails } (X=0)$

Given that there is also some prior \Rightarrow All fair-minted coins have equal probability
we can add imaginary hallucinated flips of 10H & 10T

$$\frac{d_1 + 10}{(d_1 + 10) + (d_0 + 10)}$$

\rightarrow which could reflect the prior.

At 0 coin flips, the ratio is 0.5

$$\text{At 1} \quad \text{as} \quad \frac{10}{21}$$

At million, 10 is negligible ϵ_1 and $\rightarrow \infty$ the ratio

Converges to $\frac{d_1}{d_1 + d_0}$.

But why 10, instead of 100, it would take more data to converge in

General algorithm, $B_1 = \# \text{ hallucinated } X=1s$

$B_0 = \# \text{ hallucinated } X=0s$

The formula could be replaced by

$$\frac{d_1 + B_1}{(d_1 + B_1) + (d_0 + B_0)} \quad (\text{observed + hallucination})$$

more confidence on prior, higher value of B

As the flips increase B becomes negligible \Rightarrow In case of sparse data, we could use the prior.

Principles for estimating probabilities

Principle 1 \rightarrow Maximum likelihood

Choose parameters θ that maximize $P(\text{data}|\theta)$

$$\text{Ex. } \hat{\theta}_{\text{MLE}} = \frac{\lambda_1}{\lambda_1 + \lambda_0}$$

Principle 2 \rightarrow Maximum a posteriori probability, $= \frac{P(\text{data}|\theta) \cdot P(\theta)}{P(\text{data})}$

Choose parameters θ that maximize $P(\theta|\text{data})$ (Posterior.)

$$\text{Ex. } \theta_{\text{MAP}} = \frac{\lambda_1 + \# \text{ hallucinated-1s}}{(\lambda_1 + \# \text{ hallucinated-1s}) + (\lambda_0 + \# \text{ hallucinated-0s})}$$

To maximize $P(\theta|\text{data})$, $P(\text{data})$ has normal effect
hence maximize the numerator.

Maximum likelihood Estimation

$$P(X=1) = \theta \quad P(X=0) = (1-\theta)$$

Let say the coin flips result in $\{1, 0, 0, 1, 1\}$

$$S = \{1, 0, 0, 1, 1\} \quad \theta, 1-\theta, (1-\theta), 0, 0$$

$$P(\text{data}|\theta) = \prod S_i$$

$$= \theta \cdot (1-\theta) \cdot (1-\theta) \cdot \theta \cdot \theta$$

$$= \theta^3 (1-\theta)^2$$

Generalize for Observed heads & tails.

$$P(D|\theta) = P(\lambda_1, \lambda_0|\theta) = \theta^{\lambda_1} (1-\theta)^{\lambda_0}$$

Hence maximizing $\theta, \bar{\theta}$ would be in maximizing the value of $P(D|\theta)$

Maximising the function or its log is going to be the same since log is monotonic.

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \ln(P(D|\theta)) \quad \begin{matrix} \text{data} \\ \text{likelihood} \end{matrix}$$

$$= \underset{\theta}{\operatorname{argmax}} \ln \theta^{L_1} (1-\theta)^{L_0}$$

Set derivative $\frac{d}{d\theta} \ln P(D|\theta) = 0$ $\frac{\partial \ln \theta}{\partial \theta} = \frac{1}{\theta}$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \ln [\theta^{L_1} (1-\theta)^{L_0}]$$

$$= \frac{\partial}{\partial \theta} L_1 \ln \theta + L_0 \ln (1-\theta)$$

$$L_1 \frac{1}{\theta} + L_0 \frac{\partial \ln(1-\theta)}{\partial \theta} \quad \begin{matrix} \text{using chain rule} \\ \frac{\partial \ln(1-\theta)}{\partial (1-\theta)} \cdot \frac{\partial (1-\theta)}{\partial \theta} \end{matrix}$$

$$= \frac{1}{1-\theta} - 1$$

$$\hat{\theta} = L_1 \frac{1}{\theta} - \frac{L_0}{1-\theta}$$

$$\hat{\theta} = \frac{L_1}{L_1 + L_0}$$

For a config which has a bernoulli distribution (a number)

$$P(X=1) = \theta$$

$$P(X=0) = 1-\theta$$

Each flip yields boolean value for X

$$X \sim \text{Bernoulli } P(X) = \theta^x (1-\theta)^{1-x} \quad \begin{matrix} \text{HTT HTH} \\ P(\text{HTH}) = \theta \cdot (1-\theta) \cdot (1-\theta) \cdot \theta \end{matrix}$$

$$\theta^x (1-\theta)^{1-x}$$

SENK (5)

Dataset D of independent identically distributed (iid) flips produces α_1 , ones & zeroes (Binomial)

$$P(D|\theta) = P(\alpha_1, \alpha_0 | \theta) = \theta^{\alpha_1} (1-\theta)^{\alpha_0}$$

$$\hat{\theta}^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} P(D|\theta) = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

Principle 2 \rightarrow Maximum a posteriori probability

Beta prior distribution - $P(\theta)$

$$P(\theta) = \theta^{\beta_H-1} (1-\theta)^{\beta_T-1} \sim \text{Beta}(\beta_H, \beta_T)$$

$$\text{Likelihood function: } P(D|\theta) = \theta^{\alpha_H} (1-\theta)^{\alpha_T}$$

$$\text{Posterior: } P(\theta|D) \propto P(D|\theta) P(\theta)$$

$$\hat{\theta}^{\text{MAP}} = \frac{\alpha_H + \beta_H-1}{(\alpha_H + \beta_H-1) + (\alpha_T + \beta_T-1)}$$

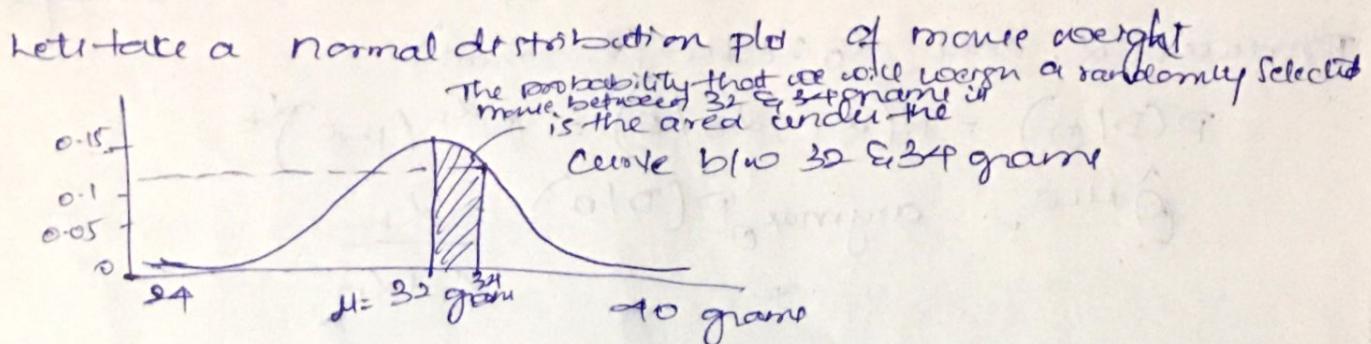
Term 2
Likelihood function: $P(\text{data}|\theta)$

Prior: $P(\theta)$

Posterior $P(\theta|\text{data})$

Conjugate prior: $P(\theta)$ is the conjugate prior for likelihood function $P(\text{data}|\theta)$ if the forms of $P(\theta)$ & $P(\theta|\text{data})$ are the same.

Probability vs Likelihood



Here if the area under curve = 0.29

\Rightarrow there's a 29% chance a randomly selected mouse will weigh between 32 & 34 grams.

Mathematically,

$$P(\text{wt } 32-34 \mid \begin{matrix} \mu = 32 \\ \text{mean} \end{matrix} \text{ & } \begin{matrix} \sigma = 2.5 \\ \text{Standard deviation} \end{matrix}) = 0.29$$

Let's say we measured a mouse with weight = 34 grams

The likelihood of weighing a 34 grams mouse is

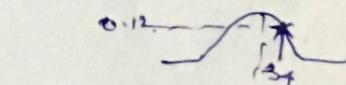
Mathematically

$$L(\text{mean} = 32 \text{ & Standard deviation: } 2.5 \mid \text{mouse weighs 34 grams}) = 0.12$$

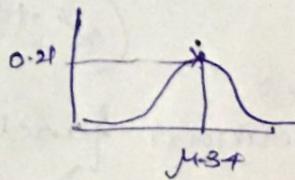
If mean is shifted to 34 grams

with right being fixed we modify the shape & location of the distribution with the left side

$P(\text{data/distribution})$: area under a fixed distribution

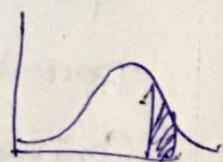


$L(\text{distribution}/\text{data})$
Distribution fixed
in Likelihood.



Likelihood are the y-axis values for fixed data points with distributions that can be moved.

$L(\text{distribution}/\text{data})$



BRMK ⑥

Bayes Rule

$$P(Y/x) = \frac{P(x/Y) P(Y)}{P(x)}$$

which is shorthand for

$$(H_{i,j}) P(Y=y_i/x=x_j) = \frac{P(X=x_j/Y=y_i) P(Y=y_i)}{P(X=x_j)}$$

Conditional Independence

$$(H_{i,j,k}) P(X=x_i/Y=y_j, Z=z_k) = P(X=x_i/Z=z_k)$$

X is conditionally independent of Y given Z if the probability distribution governing X is independent of the value of Y , given the value of Z .

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

Ex. $P(\text{Thunder}/\text{Rain, Lightning}) = P(\text{Thunder}/\text{Lightning})$

Naive Bayes makes assumption that the x_i are conditionally independent, given y ex. $P(x_1/x_2 y) = P(x_1 y)$

Given this assumption then

$$P(x_1, x_2 | y) = \frac{P(x_1/x_2 y) P(x_2/y)}{P(x_1/y) P(x_2/y)}$$

$$P(x_1, \dots, x_n | y) = P(x_1/y) \cdots P(x_n/y)$$

$$\prod P(x_i/y)$$

without conditional independence assumption $2 \cdot (2^n - 1) + 1$
with conditional $2^n + 1$

Naive Bayes classifier

Supervised learning algorithm used for classification

Based on Bayes theorem, $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$

$P(A)$ = Prior probability (Probability of event before event B)

$P(A|B)$ = Posterior Probability (Probability of event after event B occurs)

Training Dataset : [Document]

<u>Text</u>	<u>Category</u>	
"A great game"	Sports	# words in Sports 11
"The election was over"	Non-Sports	# words in non-Sports 9
"Very clean Match"	Sports	# unique words 14
"A clean but forgettable game"	Sports	
"It was a close election"	Non-Sports	

Classify ?? "A very close game" Category? Sports
NonSports

Concept

$$P(\text{Sports} / \text{"A very close game"}) = ? \quad \left\{ \begin{array}{l} \text{Whoever wins} \\ \text{is the answer} \end{array} \right.$$

$$P(\text{Non-Sports} / \text{"A very close game"}) = ?$$

Has to find probability

$$P(\text{"A very close game"}) = P(A) \times P(\text{very}) \times P(\text{close}) \times P(\text{game})$$

$$\text{Then } P(\text{"A very close game"} / \text{Sports}) = P(A / \text{Sports}) \times P(\text{very} / \text{Sports}) \times P(\text{close} / \text{Sports}) \times P(\text{game} / \text{Sports})$$

$$P(\text{"A very close game"} / \text{Non-Sports}) = P(A / \text{Non-Sports}) \times P(\text{very} / \text{Non-Sports}) \times P(\text{close} / \text{Non-Sports}) \times P(\text{game} / \text{Non-Sports})$$

Calculating probabilities $P(\text{Sport}) = 3/5$ $P(\text{Non-Sport}) = 2/5$

$P(A/\text{Sport}) = 2/11$, $P(\text{Very}/\text{Sport}) = 1/11$, $P(\text{close}/\text{Sport}) = 0 > 0/11$
and hence $P(\text{"A very close game"/Sport}) = 0 \Rightarrow$ this is incorrect

Solution

Laplace Smoothing

$$\theta_j = \frac{x_j + \alpha}{N + \alpha d}$$

$\alpha > 0$, mainly $\alpha = 1$ (always)

$d \rightarrow$ no. of unique words or distinct words.

In other words

$$P(\text{word}) = \frac{\text{word count} + 1}{\text{Total no. of words} + \text{No. of unique words.}}$$

$$\text{with Laplace smoothing } P(\text{close}/\text{Sport}) = \frac{0+1}{11+14}$$

word	$P(\text{word}/\text{Sport})$	$P(\text{word}/\text{Non-Sport})$
------	-------------------------------	-----------------------------------

a

$$\frac{2+1}{11+14}$$

$$\frac{1+1}{9+14}$$

Very

$$\frac{1+1}{11+14}$$

$$\frac{0+1}{9+14}$$

close

$$\frac{0+1}{11+14}$$

$$\frac{1+1}{9+14}$$

game

$$\frac{2+1}{11+14}$$

$$\frac{0+1}{9+14}$$

$$P(a/\text{sport}) + P(\text{Verse}/\text{sport}) + P(\text{Close}/\text{sport}) + P(\text{game}/\text{sport}) \stackrel{(b-6)}{=} 2.76 \times 10^{-5} = 0.000276$$

$$P(a/\text{non-sport}) + P(\text{Verse}/\text{non-sport}) + P(\text{Close}/\text{non-sport}) + P(\text{Not Sport})$$

$$= \frac{2}{23} \times \frac{1}{23} \times \frac{2}{23} \times \frac{1}{23} \times 0.4$$

$$= 0.572 \times 10^{-5}$$

$$P(\text{Sentence}/\text{Sport}) > P(\text{Sentence}/\text{non-sport})$$

\Rightarrow Sport Category

Naive Bayes in a nutshell

$$P(Y=y_k/x_1, \dots, x_n) = \frac{P(Y=y_k) P(x_1, \dots, x_n | Y=y_k)}{\sum_i P(Y=y_i) P(x_1, \dots, x_n | Y=y_i)}$$

Assuming Conditional independence among x_i 's

$$P(Y=y_k/x_1, \dots, x_n) = \frac{P(Y=y_k) \prod_i P(x_i | Y=y_k)}{\sum_i P(Y=y_i) \prod_i P(x_i | Y=y_i)}$$

For new data

$$Y^{\text{new}} \leftarrow \underset{y_k}{\operatorname{argmax}} P(Y=y_k) \prod_i P(x_i^{\text{new}} | Y=y_k)$$

Classifying Text Documents

Turn words in document into bag of words.

Y discrete valued Ex. Spam or not

$X = \langle x_1, x_2, \dots, x_n \rangle = \text{document}$

x_i is boolean $\begin{cases} 1 & \text{word in document} \\ 0 & \text{otherwise} \end{cases}$