

Random Variables & Probability Distributions

Variable : whose value can change.

Deterministic
[governed by deterministic equations]

Random
[Value is not pre-determined]

can change in unpredictable ways but has a probability distribution.

Discrete

Discrete R.V:

Coin Toss

$$x = \begin{cases} 1 & \text{Heads} \\ 0 & \text{Tails} \end{cases}$$

$$P(X=0) = q = 1 - p$$

Bernoulli Distribution $\leftarrow \begin{cases} P(X=1) = p \\ 0 \leq p \leq 1 \end{cases}$

Dice Roll

Binomial Distribution : N coin tosses

WRITE CODE
to plot this
for large 'n'.

$$P_r(X=k) = f(k, n, p) = {}^n C_k p^k (1-p)^{n-k}$$

Poisson Distribution : $f(k; \lambda) =$

$$\frac{\lambda^k e^{-\lambda}}{k!} = P_r(X=k)$$

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

Number of events in a fixed time interval

$$\Omega = \{x_1, x_2, x_3, \dots, x_n\}$$

$$0 \leq P(X=x_i) \leq 1$$

$$\sum_{i=1}^n P(X=x_i) = 1$$

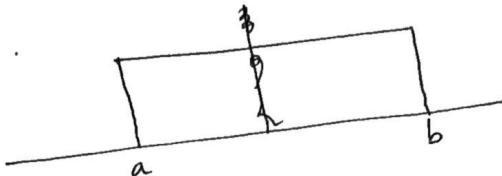
Continuous R.V.

Any measuring device has finite precision & range

when outcome can be any real number [uncountable].

probability of X taking any one value is zero.
Practical: when no. of possibilities is very large

Uniform Distribution:



Distribution $f(x) = \frac{1}{b-a}$

$$P(X=x) = 0$$

$$P(x \leq X \leq x + \Delta x) = f(x) \Delta x, \text{ if } \Delta x \approx 0$$

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$f(x) \geq 0$$

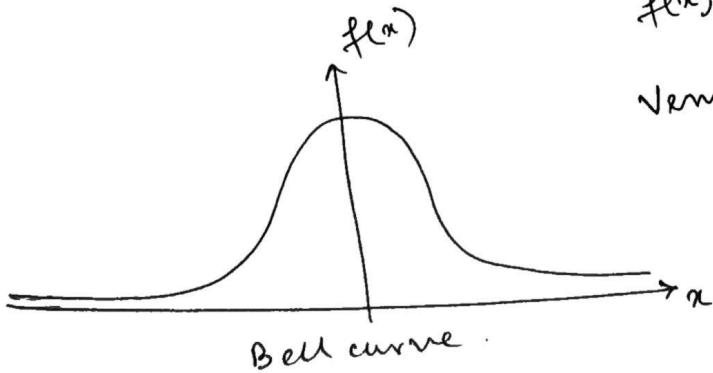
$\Rightarrow f(x) \geq 0$ $\Rightarrow f(x) \leq 1$ $\Rightarrow f(x) \leq 1$
 $f(x)$ can be greater than 1 over a finite range of x values.

Gaussian:

Velocity of gas molecules
Noise in communication channels.
Height of a large population.

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right]$$

Very crucial in ML too!!



Expectation, Moments & CLT

Expectation: $E[g(x)] = \int_{-\infty}^{\infty} g(x) f(x) dx$

Mean: $E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

Variance: $\sigma^2 = E[(x - \mu)^2] = E[x^2] - E[x]^2$

Moments: $m_n = E[x^n] = \int_{-\infty}^{\infty} x^n f(x) dx$

Characteristic function: $\Phi_x(\omega) = \int_{-\infty}^{\infty} f(x) e^{j\omega x} dx = E[e^{j\omega x}]$

Generating function: $M_x(s) = \int_{-\infty}^{\infty} f(x) e^{sx} dx = E[e^{sx}]$

Moment Generating Function (MGF)

$$M_x^{(n)}(0) = E[x^n] = m_n$$

If all moments are finite & the following series converges absolutely, then

$$M_x(s) = \sum_{n=0}^{\infty} \frac{m_n}{n!} s^n \quad \begin{array}{l} \text{[does not work} \\ \text{for log-normal} \\ \text{distribution]} \end{array}$$

If X & Y are independent random variables,

$$\Phi_{X+Y}(\omega) = E[e^{j\omega(X+Y)}]$$

$$= E[e^{j\omega X}] E[e^{j\omega Y}] = \Phi_X(\omega) \Phi_Y(\omega)$$

MGF of Normal Distribution

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad \int_{-\infty}^{\infty} f(x) dx = 1$$

$$\begin{aligned} M(s) &= \mathbb{E}[e^{sx}] = \int_{-\infty}^{\infty} e^{sx} \cdot f(x) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{sx} e^{-x^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(x-s)^2/2} e^{s^2/2} dx \\ &= e^{s^2/2} \cdot \left[\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(x-s)^2/2} dx \right] = e^{s^2/2} \end{aligned}$$

Central Limit Theorem

$$\bar{x}_n = \frac{x_1 + x_2 + \dots + x_n}{n} \quad \{x_i\} \text{ i.i.d. rv}$$

$$\text{Let } y_i = x_i - \mu \quad \& \quad z_n = \frac{y_1 + y_2 + \dots + y_n}{\sqrt{n}\sigma^2}$$

$$M_y(s) = 1 + m_1 s + \frac{m_2}{2!} s^2 + \dots = 1 + \frac{\sigma^2}{2} s^2 + \dots$$

$$\begin{aligned} M_{z_n}(s) &= \left[M_y\left(\frac{s}{\sqrt{n}\sigma^2}\right) \right]^n = \left[1 + \frac{\sigma^2}{2} \left(\frac{s}{\sqrt{n}\sigma^2} \right)^2 + \dots \right]^n \\ &= \left[1 + \frac{s^2}{2n} + \dots \right]^n \sim e^{s^2/2} \quad \text{for large } n \end{aligned}$$

BAYES' THEOREM

$V = [A_1, A_2, \dots, A_n]$ is a partition of S , meaning $A_i \neq A_j$
 are mutually exclusive
 & $\sum_i p(A_i) = 1$

Let B be an arbitrary event:

$$\begin{aligned} B &= B \cap S = B \cap [A_1 \cup A_2 \cup A_3 \dots \cup A_n] \\ &= (B \cap A_1) \cup (B \cap A_2) \dots \cup (B \cap A_n) \end{aligned}$$

$B \cap A_i$ & $B \cap A_j$ are
 mutually exclusive $\therefore P(B) = P(B \cap A_1) + P(B \cap A_2) + \dots + P(B \cap A_n)$

Total probability theorem: $P(B \cap A_i) = P(B|A_i) P(A_i)$
 $= P(A_i|B) P(B)$

$$\Rightarrow P(A_i|B) = \frac{P(B|A_i) P(A_i)}{P(B)}$$

$$= \frac{P(B|A_i) P(A_i)}{\sum_i P(B|A_i) P(A_i)}$$

BAYES' THEOREM

Example: A test for ~~Cancer~~ is known to be ~~95%~~ 95% accurate.
 A person submits to test & is found positive. Let the person
 belong to a population of 1L people out of which 2000
 actually have the disease. What's the prob. of the
 person being actually infected?

$$P(T_+|c) = 0.95 \quad P(T_-|H) = 0.05$$

$$P(H) = 0.98 \quad P(c) = 0.02$$

$$P(c|T_+) = \frac{P(T_+|c) P(c)}{P(T_+)} = \underline{\underline{0.278}}$$

$$P(T_+) = P(T_+|c) P(c) + P(T_+|H) P(H)$$

Example : $P(D|B_1) = \frac{100}{2000} = 0.05$ $P(D|B_2) = \frac{200}{500} = 0.4$

Papoulis
page 3A

$$P(D|B_3) = \frac{100}{1000} = 0.1 \quad P(D|B_4) = \frac{100}{1000} = 0.1$$

$$P(B_1) = P(B_2) = P(B_3) = P(B_4) = 1/4$$

$$P(D) = 0.1625$$

$$P(B_2|D) = \frac{0.4 \times 0.25}{0.1625} = 0.615$$

Example
Papoulis
page 104

in coin tosses k heads
Prob of head in $(n+1)$ th toss?

$$\text{Ans: } \frac{k+1}{n+2}$$

Naive Bayes & Gaussian Naive Bayes

Bayer's Theorem

$\cup = [A_1, A_2, \dots, A_n]$ is a partition of S
 & B is an arbitrary event.

$$P(A_i/B) = \frac{P(B/A_i) P(A_i)}{\sum_i P(B/A_i) P(A_i)} = \frac{P(B/A_i) P(A_i)}{P(B)}$$

Dataset of 1L emails with 70% Normal & 30% Spam.
 find total count of each word occurring in each category
 [BAG OF WORDS model
 $w_1 = \text{work}$

words Two ways to calculate $P(w_i/N) = \frac{22}{100} = 0.22$
 # docs to calculate $P(w_i/S) = \frac{11}{100} = 0.11$

if $P(w_i/N) = 0$
 & $P(w_i/S) \neq 0$
 or reverse,
 add $\frac{1}{N}$ to all
 word counts

$$\begin{aligned} P(w_2/N) &= \frac{26}{100} = 0.26 \\ P(w_3/N) &= \frac{19}{100} = 0.19 \\ P(w_4/N) &= \frac{13}{100} = 0.13 \\ P(w_5/N) &= \frac{8}{100} = 0.08 \\ P(w_6/N) &= \frac{12}{100} = 0.12 \end{aligned}$$

$$\begin{aligned} P(w_1/S) &= \frac{11}{100} = 0.11 \\ P(w_2/S) &= \frac{4}{100} = 0.04 \\ P(w_3/S) &= \frac{14}{100} = 0.14 \\ P(w_4/S) &= \frac{31}{100} = 0.31 \\ P(w_5/S) &= \frac{21}{100} = 0.21 \\ P(w_6/S) &= \frac{19}{100} = 0.19 \end{aligned}$$

$$\begin{aligned} w_2 &= \text{deadline} \\ w_3 &= \text{hand} \\ w_4 &= \text{bonus} \\ w_5 &= \text{luxury} \\ w_6 &= \text{Money} \end{aligned}$$

Total 100 words in each category

New email received with words w_3, w_5 & w_6 [in any order]
Naive \equiv [we also assume word occurrence to be independent events]

$$P(N/w_3 \wedge w_5 \wedge w_6)$$

High bias - ignores word order
 Low variance
 \hookrightarrow works well

$$\frac{P(w_3 \wedge w_5 \wedge w_6/N) P(N)}{P(w_3 \wedge w_5 \wedge w_6)}$$

$$= \frac{P(w_3/N) P(w_5/N) P(w_6/N) P(N)}{P(w_3 \wedge w_5 \wedge w_6)}$$

$$\propto 0.19 \times 0.08 \times 0.12 \times \frac{0.7}{0.7} = \frac{0.0012768}{0.000912}$$

$$P(S/w_3 \wedge w_5 \wedge w_6) = \frac{P(w_3/S) P(w_5/S) P(w_6/S) P(S)}{P(w_3 \wedge w_5 \wedge w_6)}$$

$$\propto 0.14 \times 0.21 \times 0.19 \times \frac{0.3}{0.3} = \frac{0.0016758}{0.000793}$$

[Take log since numbers get very small underflow]

Gaussian Naive Bayes

Estimating $p(x_i | c_j)$ using Gaussian Distribution
used when ' x_i ' values can be any real number
in a range.

Eg: Boy vs. Girl classification using

height, weight & feet size.

Use mean & variance of available data
to fit Gaussian.

~~fit a Gaussian
using mean & SD.~~

Bernoulli Naive Bayes : use Bernoulli distribution

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

Categorical
Naive Bayes

word present or absent

Multinomial Naive Bayes : $x_i \in \{0, 1, 2, \dots, K\}$
word count taken into account.

Bayes Optimal Classifier

Model parameter dependencies leading to the
most general classifier. Abstract notion of hard
to implement in practice.

BAYESIAN PARAMETER ESTIMATION

Event A : K Heads in n tosses

p : probability of getting Head in one toss

Event A': Head in $(n+1)$ th toss

Since we do not know the value of ' p ' & are trying to estimate it, we need to consider it as a random variable.

Before event A, let's say we have no information about ' p '

$$\text{So, assume } f(p) = \begin{cases} 1, & 0 \leq p \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$$P(A|p) = p^k (1-p)^{n-k}$$

$$f(p|A) = \frac{P(A|p) f(p)}{\int_0^1 P(A|p) f(p) dp} = \frac{P(A|p) f(p)}{\int_0^1 p^k (1-p)^{n-k} \cdot \frac{(n+1)!}{(n-k)! k!} dp}$$

$$= p^k (1-p)^{n-k} \cdot \frac{(n+1)!}{(n-k)! k!} \quad 0 \leq p \leq 1$$

$\sim \beta(\alpha, \beta)$
Beta-distributed
 $\alpha = n$
 $\beta = k$ here.

$$P(A'|A) = \int_0^1 \underbrace{P(A'|p)}_{p} f(p|A) dp$$

$$= \frac{k+1}{n+2}$$

This is generally very complicated to estimate in most problems. Hence, we instead estimate the value of ' ϕ ' that best explains the observed data.

There are two ways of doing it.

MLE & MAP

MLE: Assume that you have no prior info about the parameters [uniform distribution]

MAP: Assume availability of prior info based on reasonable considerations.

also leads to "regularisation"
which helps in preventing overfitting.

(x, y) Training Data $(\hat{x}, \hat{y})^+$

\hat{x} new input data

\hat{y} predicted output data

$$p(\hat{y} | x, y, \hat{x}) = \underbrace{\int p(\hat{y} | \phi, \hat{x})}_{\text{model}} \underbrace{p(\phi | x, y) d\phi}$$

$$p(\phi | x, y) = \frac{p(y | x, \phi) p(\phi)}{p(y | x)}$$

$$= \frac{p(y | x, \phi) p(\phi)}{\int p(y | x, \phi) p(\phi) d\phi}$$

BAYESIAN PARAMETER ESTIMATION.

$$\text{MLE \& MAP : } p(\hat{y} | x, y, \hat{x}) = p(\hat{y} | \phi_E, \hat{x})$$

$$\phi_E = \arg \max_{\phi} p(\phi | x, y)$$

MLE & Linear Regression

Let $y_i \text{ follows } y_i = g(x_i, \theta) + \epsilon_i$ curve fitting

$$\therefore \theta_{MLE} = \arg \max \sum_i \log p(y_i | x_i, \theta)$$

$$= \arg \max \sum_i \log p(\epsilon_i)$$

$$= \arg \max \sum_i \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\epsilon_i^2/2\sigma^2} : \begin{matrix} \text{Assuming} \\ \text{Gaussian} \\ \text{Error/} \\ \text{Noise} \end{matrix}$$

$$= \arg \max \left[\sum_i \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_i \epsilon_i^2 \right]$$

$$= \arg \max \arg \min \sum_i \epsilon_i^2$$

$$= \arg \min \sum_i [y_i - g(x_i, \theta)]^2$$

Least Squares Error
Cost Function

MLE & Classification

Predict output $y \in \{0, 1\} \rightarrow$

using a probabilistic function

$$g(x, \theta) \text{ or } g(\theta) = P[Y|g(x) = 1]$$

Disease / Healthy
Bank default / pay
Fake news / factual
could be Logistic Regression
or ANN, etc.

$$\theta_{MLE} = \arg \max p(y_i | x_i, \theta)$$

$$= \arg \max \prod_i g_i^{y_i} (1-g_i)^{1-y_i}$$

$$= \arg \max \sum_i [y_i \log g_i + (1-y_i) \log (1-g_i)]$$

$$= \arg \min \left(\sum_i y_i \log g_i + (1-y_i) \log (1-g_i) \right)$$

Binary Cross-Entropy

used in
Logistic
Regression

MLE & MAP Maximum Likelihood Estimate

Naïve Bayes assumes independence of parameters which is generally not true in real world.
 So, need to build models & estimate parameters.

$$f_{\theta}(\theta|x, y) = \frac{p(y|x, \theta)f_{\theta}(\theta)}{p(y|x)} \xrightarrow{\text{Likelihood}} \text{prior } p(\theta|x) = p(\theta)$$

posterior estimate

MLE
 Assume $f_{\theta}(\theta)$ to have uniform distribution
 i.e. no prior knowledge of parameters available.

$$\begin{aligned}\therefore \theta_{MLE} &= \arg \max f_{\theta}(\theta|x, y) \\ &= \arg \max p(y|x, \theta) \xrightarrow{\text{L}} \\ &= \arg \max \prod_i p(y_i|x_i, \theta) \\ &= \arg \max \sum_i \log p(y_i|x_i, \theta) \xrightarrow{\text{Log likelihood}}\end{aligned}$$

Coin toss experiment
 p = prob. of heads (1) tails = 0
 $p(y_i|p) = p^{y_i}(1-p)^{1-y_i} \quad y_i \in \{0, 1\}$
 \hookrightarrow Bernoulli dist.

$$\begin{aligned}LL &= \log p(y|p) \\ &= \sum_{i=1}^n \log p(y_i|p) = \sum_{i=1}^n \log p^{y_i}(1-p)^{1-y_i} \\ &= \log p \sum y_i + \log(1-p) \sum (1-y_i) \\ &= n^{(1)} \log p + n^{(0)} \log(1-p).\end{aligned}$$

$$\frac{\partial LL}{\partial p} = 0 \Rightarrow \frac{n^{(1)}}{p} - \frac{n^{(0)}}{1-p} = 0 \Rightarrow \frac{1}{p} - 1 = \frac{n^{(0)}}{n^{(1)}} \Rightarrow p_{MLE} = \frac{n^{(1)}}{n^{(1)} + n^{(0)}}$$

[Intuitive Estimate in MLE]

MAP Estimation - Maximum a posteriori

$$f_p(\theta | y, x) = \frac{p(y|x, \theta) f(\theta|x)}{p(y|x)}$$

$$\begin{aligned}\theta_{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} \, f_p(\theta | y, x) \\ &= \underset{\theta}{\operatorname{argmax}} \, p(y|x, \theta) f(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \left[\sum_i \log p(y_i|x_i, \theta) + f(\theta) \right]\end{aligned}$$

Usually clear info of $p(\theta)$ is not available.
 So use $f_p(\theta)$ to enforce certain constraints.
 [regularisation]
 Eg: choose $f(\theta)$ such that there is low probability
 for θ taking high values.

~~Now let's say we wish to make a prediction
 for new data point~~

$$p(\tilde{y}|\tilde{x}, \theta) = p(\tilde{y}|\tilde{x}, y, x) = p(\tilde{y}|\tilde{x}, \theta) f(\theta|y, x) d\theta$$

~~but actually, $p(\tilde{y}|\tilde{x}, y, x) = \int p(\tilde{y}|\tilde{x}, \theta) f(\theta|y, x) d\theta$~~
 usually we estimate this just using θ_{ML} or θ_{MAP}
 which is not same as RHS.
 But estimating RHS is very hard since
 it requires $p(x)$ too.

MLE for binomial

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log p$$

MAP for Coin Toss

$$\hat{p}_{\text{MAP}} = \underset{\hat{p}}{\operatorname{argmax}} \sum_{i=1}^N \log p(y_i | \hat{p}) + \log p(\hat{p})$$

$$f = p(y_i = 1)$$

$$\begin{aligned} n^{(1)} &= \text{No. of } 1s \\ n^{(0)} &= \text{No. of } 0s \\ n^{(1)} + n^{(0)} &= n \end{aligned}$$

$$\begin{aligned} &= \underset{\hat{p}}{\operatorname{argmax}} \left[n^{(1)} \log \hat{p} + n^{(0)} \log (1 - \hat{p}) \right. \\ &\quad \left. + \log \frac{\hat{p}^{\alpha-1} (1-\hat{p})^{\beta-1}}{B(\alpha, \beta)} \right] \\ &\quad \text{Beta-distribution} \\ &\quad \alpha, \beta > 0 \end{aligned}$$

$$= \underset{\hat{p}}{\operatorname{argmax}} \left[n^{(1)} \log \hat{p} + n^{(0)} \log (1 - \hat{p}) \right]$$

$$+ n \log \alpha (\alpha-1) \log \hat{p} + \beta (\beta-1) \log (1 - \hat{p}) \right]$$

$$= \underset{\hat{p}}{\operatorname{argmax}} \left[\frac{n^{(1)} + \alpha - 1}{n + \alpha + \beta - 2} \log \hat{p} + \frac{(n^{(0)} + \beta - 1)}{n + \alpha + \beta - 2} \log (1 - \hat{p}) \right]$$

$$\alpha = 1, \beta = 1$$

$$\Rightarrow \hat{p}_{\text{MAP}} = \hat{p}_{\text{MLE}}$$

$$\hat{p} = \frac{n^{(1)} + \alpha - 1}{n + \alpha + \beta - 2}$$

$$\begin{aligned} \frac{\partial L}{\partial \hat{p}} &= 0 \Rightarrow \frac{n^{(1)} + \alpha - 1}{\hat{p}} - \frac{n^{(0)} + \beta - 1}{1 - \hat{p}} = 0 \\ &\Rightarrow \hat{p}^{-1} = \frac{n^{(0)} + \beta - 1}{n^{(1)} + \alpha - 1} \\ &\Rightarrow \hat{p}_{\text{MAP}} = \frac{n^{(0)} + \beta - 1}{n^{(1)} + \alpha - 1} \end{aligned}$$

$$\hat{p}(y_i | y) = \frac{p(y_i | \hat{p}) p(\hat{p})}{\int_0^1 p(y_i | \hat{p}) p(\hat{p}) d\hat{p}}$$

$$\hat{p}(y_i | y) = \frac{p(y_i | \hat{p}) p(\hat{p})}{\int_0^1 p(y_i | \hat{p}) p(\hat{p}) d\hat{p}}$$

$$\begin{aligned} &= \frac{\left[\prod_i p(y_i | \hat{p}) \right] p(\hat{p})}{\int_0^1 \left[\prod_i p(y_i | \hat{p}) \right] p(\hat{p}) d\hat{p}} \\ &\hookrightarrow \text{Assume Beta distribution} \\ &= \frac{\hat{p}^{n^{(1)}} (1 - \hat{p})^{n^{(0)}}}{B(n^{(1)} + \alpha, n^{(0)} + \beta)} \hat{p}^{\alpha-1} (1 - \hat{p})^{\beta-1} \end{aligned}$$

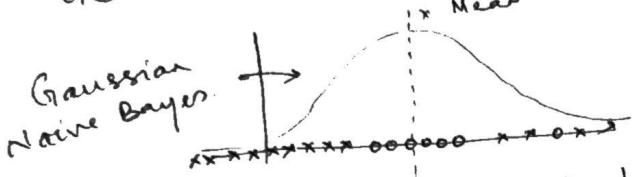
~~Bayesian Estimation~~

GENERATIVE VS. DISCRIMINATIVE MODELS

Bayesian Reasoning
(BPE)

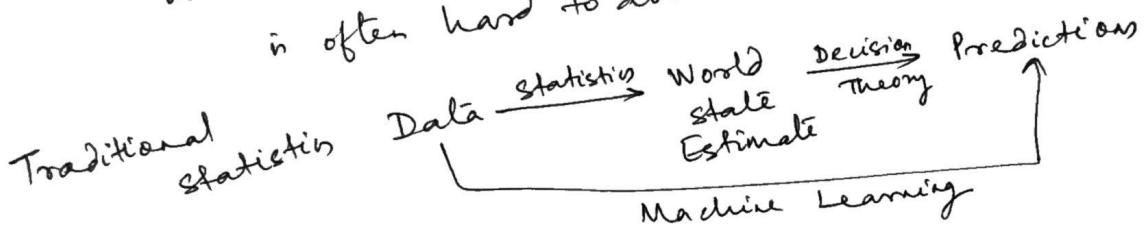
function Approximation
(MLE & MAP)
ML

Generative models can go wrong when there are outliers.



SVM, for example, can handle this with ease.

Generative may work here if we remove outliers or make better choice of distribution, but that is often hard to do.



Generative

Goal: To explain data
(concept learning)

More accurate when they use correct models

Explainable

Generate realistic data samples
(artwork, super-resolution, colorization, etc.)
simulation & planning
for reinforcement learning

Generative Adversarial Networks
(GANs)

Discriminative

Goal: To make predictions
(shortcut coaching)

More robust against outliers & bad models

Blackbox as we use more complex models

GANs

Generators

Generate new data (fake) to fool the discriminator into believing that samples ~~data~~ generated is real

vs. Discriminators

Identify fake (new generated) samples & force the generator to improve.

Applications: Art AI

Imaginary fashion models

Model dark matter distribution & predict gravitational lensing

DeepFakes!!

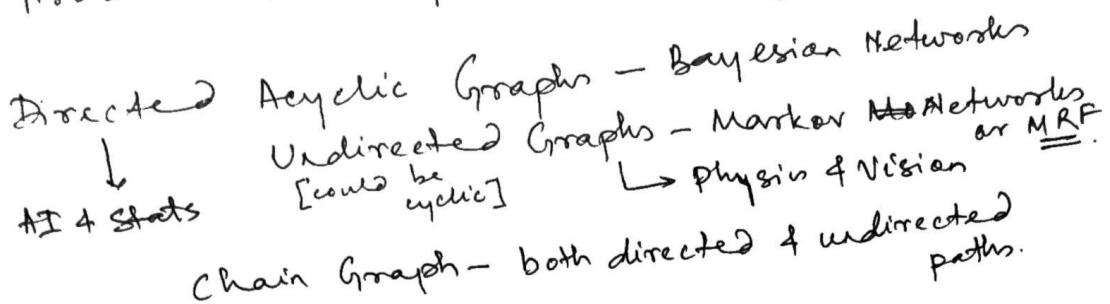
MNIST handwritten digits

Bayesian Networks - Help in identifying how causes ~~for~~ generate outcomes / events.
[work both ways bottom up & top down]

PROBABILISTIC GRAPH MODELS

[Bayesian Networks & Markov Network]

Probability Theory + Graph Theory

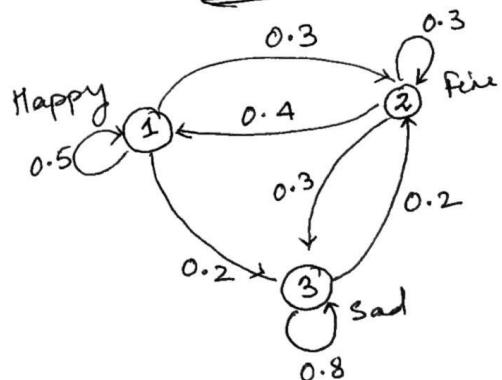


DAGs - Encode causal relationships
Easier to learn from data
But cannot represent cyclic dependencies

$\text{PGM} \neq \text{Markov chain}$
↓
Nodes are elements of state space
& edges denote transition probabilities
Nodes are random variables & edges denote conditional dependencies.

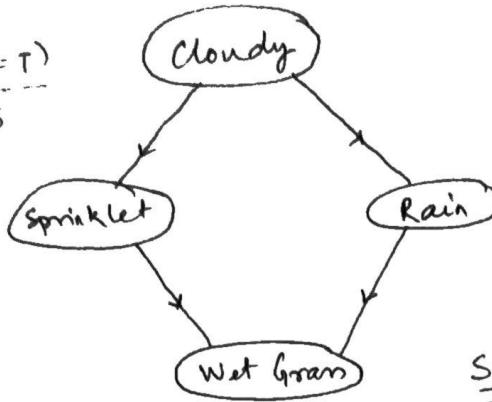
Markov: Future is independent of the past given the present.

Markov chain



	$P(C=F)$	$P(C=T)$
0.5		0.5

C	$P(S=F)$	$P(S=T)$		$P(R=F)$	$P(R=T)$
F	0.5	0.5		0.8	0.2
T	0.9	0.1		0.2	0.8



C	$P(R=F)$	$P(R=T)$
F	0.8	0.2
T	0.2	0.8

pfs,

$$P(C, S, R, W) = P(C)P(S|C)P(R|C, S)P(W|R, S)$$

$$= P(C)P(S|C)P(R|C)P(W|R, S)$$

[Simplification of conditional probabilities]

Most common task: probabilistic inference.

$$P(S=T|W=T) = \frac{P(S=T \wedge W=T)}{P(W=T)}$$

$$\stackrel{\text{Bottom Up}}{=} \sum_{R, C \in T/F} P(S=T \wedge W=T \wedge R=F \wedge C=F)$$

$$\quad\quad\quad \curvearrowleft P(W=T)$$

$$\sum_{R, S, C \in T/F} P(S=T/F \wedge W=T \wedge R=T/F \wedge C=T/F)$$

$$\stackrel{\text{Top Down}}{\curvearrowleft} P(W=T | C=T)$$

$$P(W=T | C=T)$$

$$= \frac{0.2781}{0.6471} = 0.430$$

∴ if grass is wet, it is more likely to be due to rain!

$$P(R=T|W=T) = \frac{0.4581}{0.6471} = 0.708$$

"Explaining Away"

$$P(S=T|W=T \wedge R=F) = 0.1945$$

Berkson's Paradox

S & R actually are independent
Also called "selection bias"

EXPECTATION MAXIMIZATION

Bayesian Networks

Structure	Observability	Method
Known	Full	MLE
Known	Partial	Expectation Maximization (or Gradient Descent)
Unknown	Full	Search through Model space [NP-Hard]
Unknown	Partial	EM + Search Model space

Two coins with biases θ_A & θ_B

Pick a coin at random & toss it 10 times
Repeat this 5 times

Record $X = (x_1, x_2, \dots, x_5)$ & $(z_1, z_2, z_3, z_4, z_5) = Z$

$x_i \in \{0, 1, 2, \dots, 10\}$ = No. of heads observed

$z_i \in \{A, B\}$ = Coin chosen

Estimating θ_A & θ_B is easy if this full info is available
But what if z_i are hidden variables??

Simple Approach -

Assume some initial values for $\theta_A^{(0)}, \theta_B^{(0)}$
Then estimate $\{z_i\}$ based on which coin is more likely to generate the corresponding $\{x_i\}$
Estimate $\theta_A^{(1)}, \theta_B^{(1)}$ using this completed info
Repeat the above till convergence reached.

EM Approach - Compute probability of each possible completion of table.
(create weighted training set of all these possibilities)
Then use a modified form of MLE.

$$x = \{5, 9, 8, 4, 7\} \quad ?$$

Simple Approach [MLE]

$$\text{Assume } \theta_A^{(0)} = 0.60 \quad \theta_B^{(0)} = \cancel{0.40} \quad 0.5$$

$$P(x_1|A) = \frac{10}{9} \cdot 0.6^5 \cdot 0.4^5 = 0.0007776 \cdot 0.2 \quad \left| \begin{array}{l} P(x_1|B) = 0.25 \\ P(x_2|B) = 0.5098 \end{array} \right.$$

$$P(x_3/A) = \frac{C_9 \cdot 6^9 \cdot 0.4}{C_{12} \cdot 6^{12} \cdot 0.12} = \frac{0.000483}{0.2684} = 0.0018$$

$$P(x_3/A) = \frac{10}{8} \cdot 0.6 \cdot 0.4 = 0.2604$$

$$P(A_5|B) = 0.117$$

$$P(x_5/A) = C_f \cdot 6^7 \cdot 0.4^3 = 0.22$$

$$P(x_5/A) = \frac{5+4}{10} = 0.45$$

$$\theta_A^{(1)} = \frac{9+8+7}{30} = 0.8 \quad \theta_B = \frac{1}{20} = 0.05$$

Expectation Maximization

another way to estimate MLE / MAP

x_i	play
5	0.2
9	0.04

maximization		$p(x_i B)$	w_{Ai}	w_{Bi}	A	B
x_i	$p(x_i A)$					
5	0.2	0.25	0.45	0.55		
9	0.04	0.0098	0.80	0.20		
8	0.12	0.0044	0.73	0.27		
4	0.11	0.205	0.35	0.65		
7	0.22	0.123	0.65	0.35		
					\downarrow	
		$p(x_i A)$				$1 - w_A$
			$\frac{p(x_i A)}{p(x_i A) + p(x_i B)}$			

$$\hat{\theta}_A^{(1)} = \frac{\sum x_i w_{Ai}}{10 \sum w_{Ai}} = 0.712$$

$$\theta_B^{(1)} = \frac{\sum x_i w_{Bi}}{10 \sum w_{Bi}} = 0.583$$

$$\frac{\sum x_i (1 - w_{Ai})}{10 \sum (1 - w_{Ai})} = \frac{\sum x_i - \sum w_{Ai}}{10(5 - \sum w_{Ai})}$$

After 10 steps converges to

$$\theta^{(10)} \approx 0.80$$

$$\theta_a^{(10)} \approx 0.52$$

May converge to local ~~maximum~~^{maxima} for non-concave functions.
but convergence is guaranteed.

GAUSSIAN MIXTURE MODELS

Feature Extraction from speech data

Multi-object tracking

essentially, used when data is multi-modal, i.e. there are many "peaks" in the distribution.

$$P(\vec{x}) = \sum_{i=1}^K \phi_i \mathcal{N}(\vec{x} / \vec{\mu}_i, \Sigma_i)$$

Mean vector
Covariance Matrix
Multi-variate Gaussian Distribution.

$$\mathcal{N}(\vec{x} / \vec{\mu}_i, \Sigma_i) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_i|}} \exp \left[-\frac{1}{2} (\vec{x}_i - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{x}_i - \vec{\mu}_i) \right]$$

$$\sum \phi_i = 1$$

weights
[ratio of points
in each
category]

Parameters estimated

using Expectation Maximization
since analytic solution is
usually not possible.

EM is guaranteed to converge

EM algo.

E-step : using $\phi_i, \vec{\mu}_i$ & Σ_i

estimate $p(c_i | x_j)$ for each data point

↳ K classes/components

M-step : using assigned c_i for each x_j , [weighted]

update $\phi_i, \vec{\mu}_i$ & Σ_i

Initialisation :

Randomly pick $\vec{\mu}_i$ from the given unlabelled data, x_j

Set variance estimator to variance of whole data
OR, use K-means to find initial parameter values

set $\phi_i = \frac{1}{K} \cdot 1$

(speeds up considerably)

MNIST - K=10 components

Each image of size $N \times N$ is a vector of size N^2 of Bernoulli distributions (one per pixel) - PRML book

E-step

$$\gamma_{ji} = \frac{\phi_i N(\vec{x}_j | \vec{\mu}_i, \Sigma_i)}{\sum_{i=1}^K \phi_i N(\vec{x}_j | \vec{\mu}_i, \Sigma_i)} = p(c_i | \vec{x}_j)$$

M-step

$$\vec{\mu}_i = \frac{1}{N_i} \sum_{j=1}^N \gamma_{ji} \vec{x}_j$$

$$\vec{\Sigma}_i = \frac{1}{N_i} \sum_{j=1}^N \gamma_{ji} (\vec{x}_j - \vec{\mu}_i) (\vec{x}_j - \vec{\mu}_i)^T$$

updated value

$$\phi_i = \frac{N_i}{N}, \quad N_i = \sum_{j=1}^N \gamma_{ji}$$

: Effective no. of points in category i

Log-Likelihood

$$\ln p(x | \mu, \Sigma, \phi) = \sum_{j=1}^N \ln \left\{ \sum_{i=1}^K \phi_i N(\vec{x}_j | \vec{\mu}_i, \Sigma_i) \right\}$$

Repeat the process till LL or parameters converge.

EM required since analytical solution of parameters which maximise LL not available.

K-means	GMM
hard assignment of \vec{x}_j to c_i , $\gamma_{ji} \in \{0, 1\}$ fast convergence even for high dimensional data Assumes clusters to be spherical Easier to interpret	Soft assignment $0 \leq \gamma_{ji} \leq 1$ EM is slow can handle arbitrary shapes more complex Difficult to interpret