

Pre-reg  
Matrices  $\rightarrow A_{m \times n}$

Determinants.

Defn : If  $x \in \mathbb{R}^n$  is a vector such that,  $Ax = \lambda x$ , ( $\lambda$  is a scalar), then

- i)  $\lambda \rightarrow$  eigenvalue of  $A$ .
- ii)  $x \rightarrow$  eigen vector corr' to  $A$

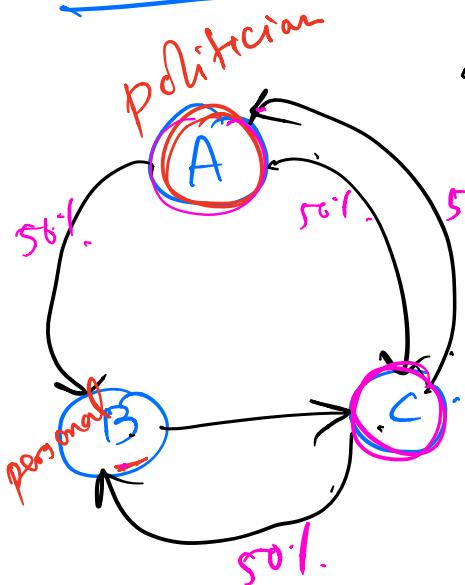
$\det(A - \lambda I) = 0$  (Caley-Hamilton equation)

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & -1 \end{bmatrix}, \quad \begin{vmatrix} 1-\lambda & 2 & 0 \\ 0 & -\lambda & 1 \\ 0 & 0 & -1-\lambda \end{vmatrix} = 0$$

$$\Rightarrow \lambda = \underbrace{3 \text{ values}}$$

$n \times n \rightarrow n$  eigen values ( $\in \mathbb{R} / \mathbb{C}$ )

# Page-Rank (Page & Sergei B.).



1 Billion web pages

so! How to rank these pages?  
Ranks are recursive.

If a node has  $m$  neighbours, each of them get  $\frac{1}{m}$  traffic

Ref. → [Markov chains]

$R_A$  is higher     $R_A = R_C \times \frac{1}{2}$

$$R_B = R_A \times \frac{1}{2} + R_C \times \frac{1}{2}$$

$$R_C = R_B \times 1 + R_A \times \frac{1}{2}$$

Flow Balance

$$\chi = \begin{pmatrix} R_A \\ R_B \\ R_C \end{pmatrix}$$

$$A \rightarrow \begin{bmatrix} 0 & 0 & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & 1 & 0 \end{bmatrix} \begin{pmatrix} R_A \\ R_B \\ R_C \end{pmatrix} = 1 \begin{pmatrix} R_A \\ R_B \\ R_C \end{pmatrix}$$

all column sums are 1

$A\chi - \chi \Rightarrow \chi$  is an eigen vector of  $A$

# Network analytics, stochastics

$$Ax = \alpha$$

$$(A - I)x = 0$$

Finding exact solution is time taking if  $A$  is large.

If  $A$  has size say  $1000 \times 1000$

① start with arbitrary  $v$ .

$v, A^1 v, A^2 v, A^3 v, A^4 v, \dots$

②  $v, A^1 v, A^2 v, A^3 v, A^4 v, \dots$  [PF thm.] in MC.

$$A^n v \Rightarrow x$$

$$A^{n+1} v = A \cdot A^n v \approx x$$

$$\Rightarrow Ax \approx x$$

Iteratively getting approx soln  
approximation  
(Machine learning)

Thm 1. If  $A$  is a real, symmetric matrix, then the eigen-vectors of  $A$ , form an orthogonal basis. #  $\{v_1, v_2, \dots, v_n\} \rightarrow$  Basis for  $C(A)$

$$Av_i = \lambda_i v_i$$

$$v_i \perp v_j$$

Thm 2 (Spectral / Eigen decomposition)  
If  $A$  is symmetric, then you can express  $A$  as  $A = Q \Lambda Q^+$ , where

$$Q = \begin{bmatrix} | & & | \\ v_1 & v_2 & \dots & v_n \\ | & & | \end{bmatrix}, \quad A = \begin{pmatrix} \lambda_1^n & & & \\ & \lambda_2^n & & \\ & & \ddots & \\ 0 & & & \lambda_n^n \end{pmatrix}$$

$$A = Q \Lambda Q^+$$

orthogonal

$A^n$   $A^{100}$   
 $A$  is  $10000 \times 10000$  |  $A = Q \Lambda Q^+$

$$QQ^+ = I, Q^+Q = I$$

$$\begin{aligned} A &= Q \Lambda Q^+ \\ A^2 &= Q \Lambda Q^+ Q \Lambda Q^+ \\ &= Q \Lambda^2 Q^+ \end{aligned}$$

## SVD (Singular Value Decomposition)

How to find the SVD?

$$A = U \Sigma V^+$$

$$A^+ = V \Sigma U^+$$

$$AA^t = U \Sigma V^+ \cancel{V} \Sigma U^+ \sim U \Sigma^2 U^+$$

$$\Rightarrow (AA^t)u = U \Sigma^2 \cancel{U^+ u} I$$

$$= U \begin{bmatrix} \sigma_1^2 & & \\ \sigma_2^2 & \ddots & \\ & & 0 \end{bmatrix}$$

$$(AA^t)u_i = \sigma_i^2 u_i$$

$\sigma_i^2$ 's are the eigenvalues of  $AA^t$   
 $u = [u_1 \ u_2 \ \dots]$ ,  $u_i$ 's are eigen-vectors

What if A is a general matrix?  
 → Rectangular

$$A = U \Sigma V^+$$

U is orthogonal  
 $V$  is orthogonal  
 $\sigma_i$ 's singular values

$$\Sigma = \begin{bmatrix} \sigma_1 \sigma_2 \dots \sigma_r & \\ & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} 2 & 3 & 1 \\ -1 & 2 & 0 \end{bmatrix} \quad AA^T$$

QR Decomposition

$$A = QR \rightarrow \text{upper triangular.}$$

Orthogonal

Regression

$$y = X\beta + \varepsilon, \quad \hat{\beta} = (X^T X)^{-1} X^T y$$

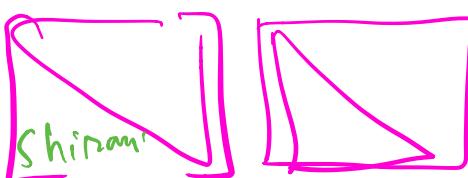
$$X_{n \times p} \xrightarrow[\text{data size}]{\text{dimension}} X = QR, \quad = R^T (Q^T Q)^{-1} R$$

$$\text{Big Data} \quad p = 1 \text{ million} \quad = R^T R$$

$X_{n \times 10000}$

$X^T X$   $\rightarrow$  inversion of large matrix is hard

Ref: Elements of stat learning Tib



Sparse

## Principal Comp. Analysis

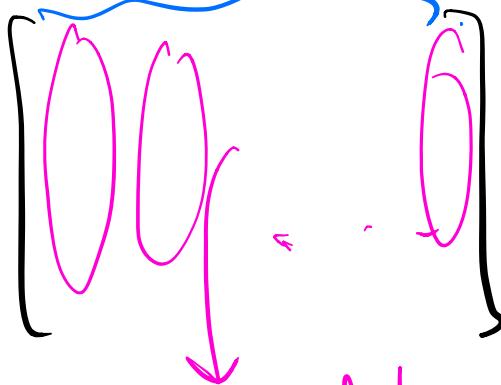
$X = \begin{bmatrix} & & \\ & & \end{bmatrix}$  }  $n$  rows,  
 $p$  features.

$n \rightarrow$   
 $p \rightarrow$  large

Assume  $X$  to  
be centred.  
[otherwise, do  $X \leftarrow (I - \frac{1}{n}J)X$ ]

genetics (Health +  $\mathbb{I}$ ).  
 $p = 100k$  features.

$n = 100$   
patients



$n \ll p$ .

Traditional Stat.

$n > p$

many would be correlated

Var-Cov. Matrix (sample)

$$\sum_{(pxp)} = \frac{1}{n-1} \sum (x_i - \bar{x})(x_i - \bar{x})^T$$

$$p = 100k, \quad \Sigma_2 \quad \boxed{\begin{aligned} \Sigma_{ij} &= \text{cov}(i^{\text{th cov}}, j^{\text{th cov}}) \\ &\text{symmetric} \end{aligned}}$$

SVD (Dimension Red.)

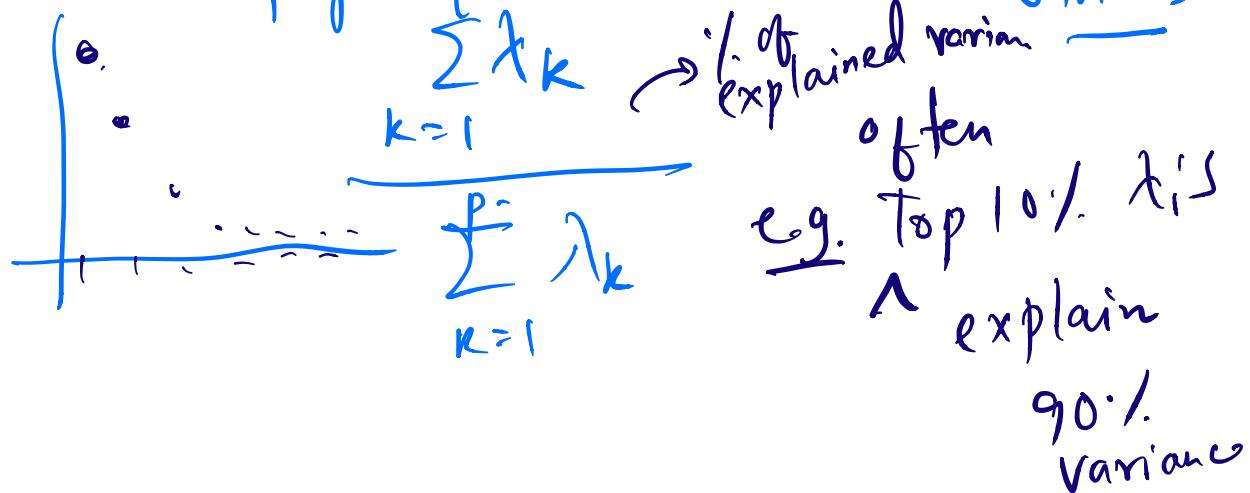
$$\Sigma = Q \Lambda Q^+$$

Thumb Rule: Reduce Size  
 $\Rightarrow$  Do SVD

$$\Lambda = \begin{bmatrix} \lambda_1, \lambda_2, \dots, \lambda_p \end{bmatrix}$$

Usually, a few top  $\lambda$ 's are very large

screeplot If you take  $i$  components, compared to others



X<sub>nxp.</sub>

$\alpha_{pxp}$

$$Q = \left[ \begin{matrix} \omega_1 & -\omega_K \\ \vdots & \vdots \\ \omega_1 & -\omega_K \end{matrix} \right] - \left[ \begin{matrix} \omega_p \\ \vdots \\ \omega_p \end{matrix} \right]$$

$$Y = X Q_k \quad \text{linear dim. reduction} \quad Q_k \rightarrow p \times k.$$

$X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$   $n \times 1$  lae-  
fat matrix.

$y =$   skinny matrix  
 $n \times k$



Choose  $k$  so the 90% variance is exp.  
 i.e.  $\frac{\sum_{i=1}^k \lambda_i}{\sum \lambda_i} > 0.9.$

$\Sigma$  is PSD.  
 always  
 " (positive semi-definite).

$\Sigma$  is symmetric

$$\Sigma_{ij} = \text{cov}(i^{\text{th}} \text{ vol}, j^{\text{th}} \text{ vol}) \\ = \Sigma_{ji}$$

$\Sigma$  is symmetric (Thm 2).

$$\rightarrow \Sigma = Q \Lambda Q^{-1} \\ = Q \Lambda Q^{-1}$$

(so, diagonalizable)

Manifold learning.  
 Spectral embedding  
 $t$ -SNE } (dimension reduction)

References are  
 in several online  
 resources.

# Probability Distributions

## ① Continuous Distr.

- 1) Gaussian.
- 2) Cauchy.
- 3) exponential
- 4) Gamma
- 5) Beta

## ② Discrete

- 1) Binomial
- 2) Bernoulli
- 3) Poisson
- 4) Geometric.

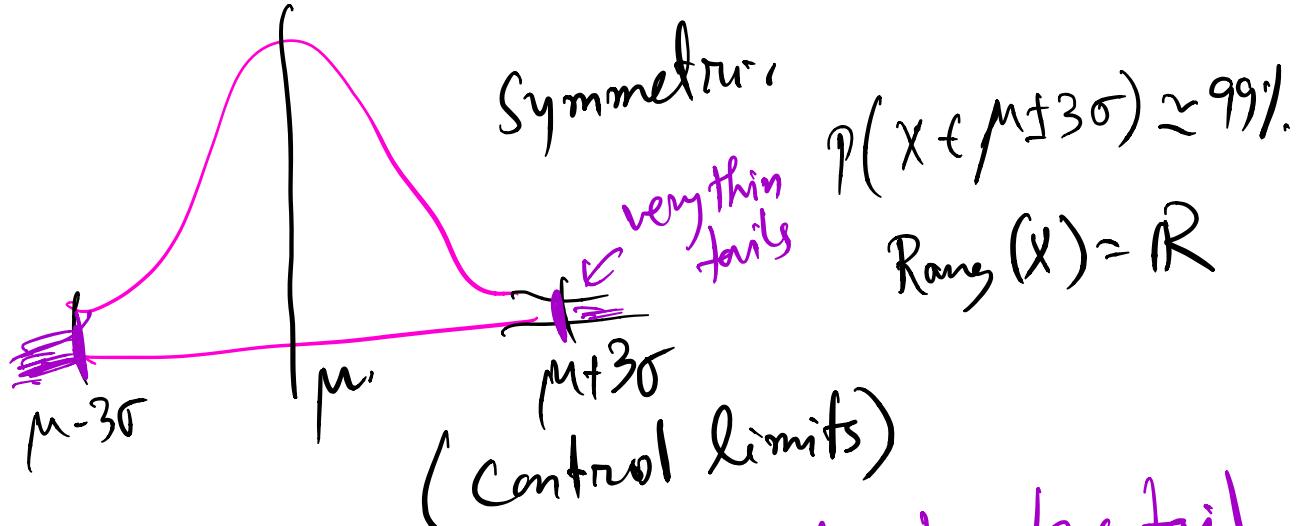
\* There are more complex stuff, but I'm going to ignore them for now.

Continuous [  $f(x) \rightarrow \text{PDF}$  ].

$$P(a < X < b) = \int_a^b f(x) dx \quad \& \quad f \geq 0.$$

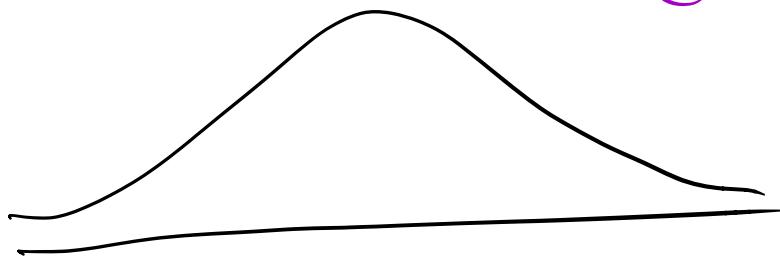
① normal  $(\mu, \sigma^2)$ .  $- \frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}$

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



Cauchy has long tail.

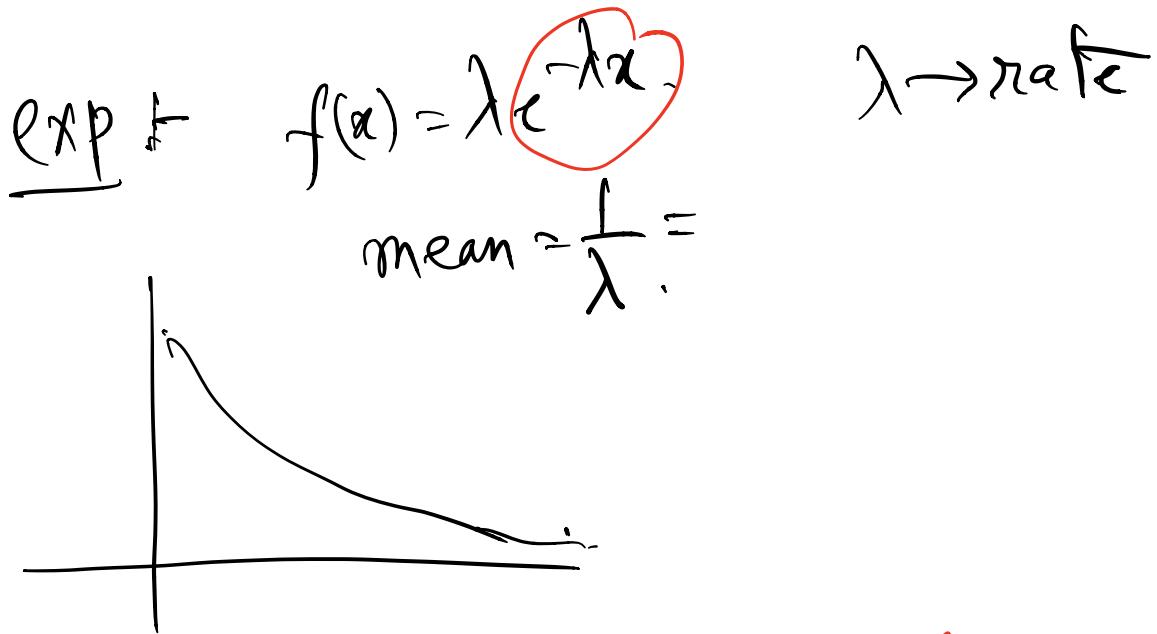
$$f(x) = \frac{1}{\pi(1+x^2)}$$



$$f(x) = \frac{1}{\pi H(x-\mu)^2}$$

Cauchy's mean is undefined

[Feller 1-].



Gamma:  $f(x) = e^{-\lambda x} x^{\alpha-1} \frac{\lambda^\alpha}{\Gamma(\alpha)}$

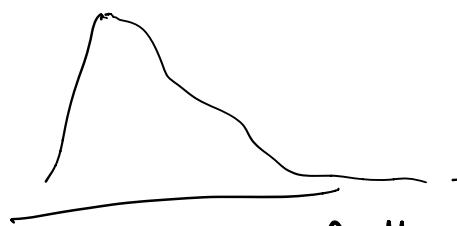
(Poisson process)  $\xrightarrow{x=1} \text{exp.}$

Binomial →  $n$  trials,  $p = P(\text{succ.})$   
 $P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$

Bernoulli →  $n=1$ ,  $\text{Bin}(1, p)$   
 $P(X=1) = p, P(X=0) = 1-p$   
 $X \in \{0, 1\}$  H/T

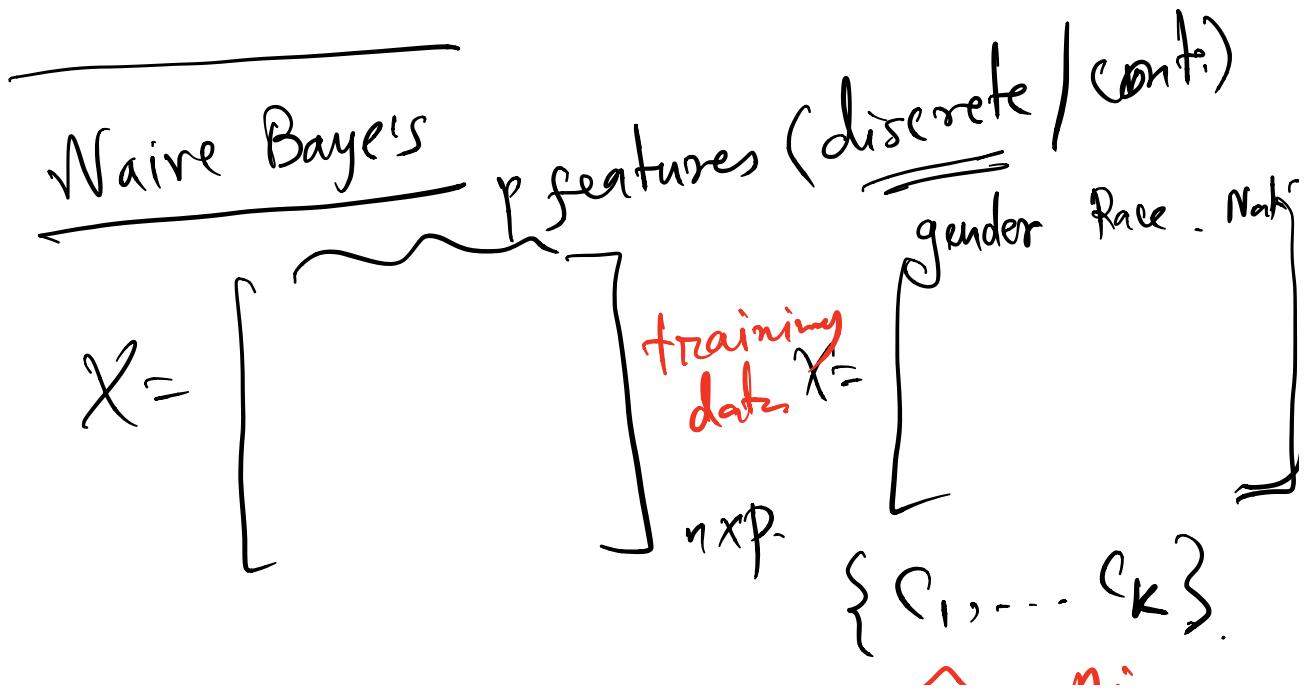
Poisson  $\rightarrow P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$

$x \geq 0, x \in \mathbb{Z}$ .



If falls very very quickly.

Modelling Rare events  
(accidents on Hwy  
Khos).



Posterior

$$P(C_i | \alpha) = \frac{P(\alpha | C_i) P(C_i)}{\sum P(\alpha | C_i) P(C_i)}$$

$\hat{\pi}_i \approx \frac{n_i}{n}$  [JEE days]  
 $C_1, C_2, C_3$

Test sample  
 $\chi = (\chi_1, \chi_2, \dots, \chi_p)$        $P(C_1 | \alpha), P(C_2 | \alpha)$   
 $P(C_3 | \alpha)$

Naive assumption: Given  $C_i, \chi_1, \chi_2, \dots, \chi_p$   
 are indep.      UCI, IRIS

$$P(\alpha | C_i) = \prod_{j=1}^p P(\chi_j | C_i)$$

[works  
fantastic  
in many  
cases]

$$\hat{P}(\chi_j | C_i)$$

$$P(C_1 | \alpha) = 42\%$$

$$P(C_2 | \alpha) = 27\%$$

$$P(C_3 | \alpha) = 31\%$$

$\chi_j \rightarrow$  nationality → [Indian  
Chinese  
Sri Lankan.]  
 assign  $\underline{\chi}$  to  $C_1$

$$\hat{P}(x_i \in \text{SriLankan} | C_i) \Rightarrow \text{training}$$

# {Srilankans in  $C_i$  in  
training  
data})

$$\frac{\# \{C_i\}}{\# \{C_i\}}.$$

Assign  $\alpha$  to  $k$ , where  
 $k = \operatorname{argmax}_i \hat{P}(C_i | x)$ .

Some answers to questions/doubts :

$$\begin{bmatrix} x_0 & x_1 & x_2 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \text{cov}(x_1, x_1) & & \\ & \text{var}(x_2) & \\ & & \text{var}(x_3) \end{bmatrix} \xrightarrow{i,j} \text{cov}(x_i, x_j)$$

orthogonal  $\{v_1, \dots, v_n\}$   
 $v_i \perp v_j \quad \forall i \neq j$

orthonormal : orthogonal  
 $\|v_i\|=1$   
 $\forall i$

## References:

- ① Tib Shirani, ESL (Stanford)
- ② Wiki pages on Naive Bayes
- ③ Bhima-Sankaran & Rao, Linear algebra.
- ④ Gilbert Strang, (MIT), linear algebra

errata: ① In the page ranking example,  
I'd referred  $\alpha$  as eigen-value  
I meant  $\alpha$  to be the eigenvector

② In the naive baye's case, I was,

during the presentation, trying to mean.  
 $P(C_1|x) \propto a$ ,  $P(C_2|x) \propto b$ ,  $P(C_3|x) \propto c$ .  
and so on. [ignoring the common denominator  
~nator]. I <sup>mistakenly</sup> referred to them as exact equalities.

The actual probabilities would be.  
resp.

$$\frac{a}{a+b+c}, \frac{b}{a+b+c}, \frac{c}{a+b+c}$$

Feel free to get in touch if you  
want to discuss anything about ML.

