① A symetric. $\forall v \quad v^T A v \geq 0$.
spectral decomposition: $A = Q D Q$
Q contain all of A's eigenvectors
D diagonal and contain corresponding
eigenvalues, $\lambda_i$.

$\Longrightarrow \quad 0 \leq v^T A v = v^T Q D(\lambda) Q^T v = \sum_i \lambda_i |Q^T v|_i^2$

$\Longrightarrow \quad \lambda_i \geq 0 \quad$ for all $i$ (it previous true for $\forall v$).

$\Longrightarrow \quad A = Q D(\lambda) Q^T = Q D(\sqrt{\lambda}) D(\sqrt{\lambda}) Q^T$
$\quad \equiv X X^T \quad (X = Q D(\sqrt{\lambda}))$

Conversly, it $A = X X^T$, Then

$v^T A v = v^T X X^T v \equiv \omega^T \omega = \langle \omega, \omega \rangle \geq 0$ //

② A, B are PSD.
define $C = \theta A + (1-\theta) B$ for $\theta \in [0,1]$
then; for $\forall v$;

$v^T C v = v^T (\theta A + (1-\theta) B) v =$
$\underbrace{\theta v^T A v}_{\geq 0} + \underbrace{(1-\theta)}_{\geq 0} \underbrace{v^T B v}_{\geq 0} \geq 0 \Longrightarrow C \text{ is PSD}$ //

# Calculus & Probability

(1) $\quad x^T A x = x_i A_{ij} x_j \qquad$ (Einstien notation of summation)

$$\frac{\partial}{\partial x_k}(x_i A_{ij} x_j) = \frac{\partial x_i}{\partial x_k} A_{ij} x_j + x_i \frac{\partial A_{ij}}{\partial x_k} x_j$$

$$= \delta_{ik} A_{ij} x_j + x_i A_{ij} \frac{\partial x_j}{\partial x_k} =$$

$$= \delta_{ik} A_{ij} x_j + x_i A_{ij} \delta_{jk} = A_{kj} x_j + x_i A_{ik}$$

$$= A_{kj} x_j + A_{ik} x_i = (Ax + A^T x)_k$$

(2) $\quad p = (p_1 \cdots p_n) \qquad \sum_i p_i = 1 \qquad p_i \geq 0$

$$H(p) = - \sum_i p_i \log p_i$$

$$G(p, \lambda) = H(p) - \lambda \sum_i p_i$$

(we would like to find $p_0$ that minimize the entropy given the constraint

$\text{I} \quad \dfrac{\partial G}{\partial p_i} = -1 - \ln p_i - \lambda \overset{!}{=} 0$

$\text{II} \quad \dfrac{\partial G}{\partial \lambda} = - \sum_i p_i = 1$

I: $\quad \ln p_i = -1 - \lambda \implies p_i = e^{-1-\lambda}$

II: $\quad 1 = - \sum_i e^{-1-\lambda} = N e^{-1-\lambda} \implies$

$$e^{-1-\lambda} = \frac{1}{N} \implies \lambda = -1 - \ln \frac{1}{N}$$

I: $\quad -1 - \ln p_i + 1 + \ln \frac{1}{N} = 0$

$$\implies p_i = \frac{1}{N} \qquad \text{uniform dist.}$$

③ (a) $P[x_0 \geq \max(x_1, \ldots x_{n-1})] =$ Total prob.

$x_i \in [0,\infty)$

$$= \int_0^\infty P(x_0 \geq \max(x_1 \ldots x_{n-1}) | x_0 = a) f_{x_0}(a) da$$

$$\int_0^\infty P(a \geq \max(x_1 \ldots x_{n-1})) f_{x_0}(a) da =$$

$$\int_0^\infty P\left(\bigwedge_{i=1}^{n-1} x_i \leq a\right) f_x(a) da \underset{\text{independent}}{=} \int_0^\infty \prod_{i=1}^{n-1} P(x_i \leq a) f_{x_0}(a) da$$

$$\underset{\text{identical}}{=} \int_0^\infty (P(x_0 \leq a))^{n-1} f_x(a) da = \int_0^\infty (F_{x_0}(a))^{n-1} f_x(a) da$$

③

$$F_{x_0}(a) = \int_0^a f_x(\tilde{a}) d\tilde{a} \quad , \quad f_x(a) = \left.\frac{\partial F_x}{\partial x}\right|_{x=a} \equiv F_x'(a)$$

$$\int_0^\infty (F_{x_0}(a))^{n-1} F_x'(a) da = \int_0^\infty \left(\frac{1}{n} [F_x(a)]^n\right)' da$$

$$= \frac{1}{n} F_x^n(a) \Big|_0^\infty = \frac{1}{n}(1 - 0) = \frac{1}{n}$$

$$\implies P(x_0 \geq \max(x_1 \ldots x_{n-1})) = \frac{1}{n}$$

# Decision Rules & Concentration Bounds

① $\quad L(h) = \sum\limits_{x,y} P(X=x, Y=y)\, \Delta_{z_0}(y, h(x))$

$\sum\limits_{x} P(X=x) \left[ \sum\limits_{y=0}^{L} P(Y=y \mid X=x)\, \Delta_{z_0}(y, h(x)) \right]$

$\Longrightarrow$ for a given $x$

$L(h) \propto \sum\limits_{y=0}^{L} P(Y=y \mid X=x)\, \Delta_{z_0}(y, h(x)) \quad \stackrel{h(x)=\hat{y}}{=}$

$\sum\limits_{y=0}^{L} P(Y=y \mid X=\hat{x}) - P(Y=\hat{y} \mid X=x) \qquad h(x)=\hat{y}$

$\Longrightarrow$ minimizing $L(h) \Longleftrightarrow$ maximizing
$$P(Y=\hat{y} \mid X=x)$$

$\Longrightarrow$ we find $y$ for wich $P(Y=y \mid X=x)$
is the largest

$\Longrightarrow \quad h(x) = \arg\max\limits_{y \in [0, \ldots L]} P(Y=y \mid X=x)$

(1) $f(\vec{x}, \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)$

(a) $P(y=1|x) > P(y=0|x) \implies y=1$

Bayes: $\dfrac{f_x(x|Y=1)P(Y=1)}{f_x(x)} > \dfrac{f_x(x|Y=0)P(Y=0)}{f_x(x)}$

$\implies f_x(x|Y=1)P(Y=1) > f_x(x|Y=0)P(Y=0)$

$p \cdot \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)\right) > (1-p)\exp\left(-\frac{1}{2}(x-\mu_0)^T\Sigma^{-1}(x-\mu_0)\right)$

$\dfrac{p}{1-p} > \exp\left(-\frac{1}{2}\left[(x-\mu_0)^T\Sigma^{-1}(x-\mu_0) - (x-\mu_1)^T\Sigma^{-1}(x-\mu_1)^T\right]\right)$

$= \exp\left(-\frac{1}{2}\left[x^T\Sigma^{-1}(\mu_1-\mu_0) + (\mu_1-\mu_0)^T\Sigma^{-1}x \right.\right.$
$\left.\left. + \mu_0^T\Sigma^{-1}\mu_0 - \mu_1\Sigma^{-1}\mu_1\right]\right)$

$\implies \ln\left(\frac{1-p}{p}\right) < x^T\Sigma^{-1}(\mu_1-\mu_0) + (\mu_1-\mu_0)^T\Sigma^{-1}x$
$\qquad\qquad\qquad + \mu_0^T\Sigma^{-1}\mu_0 - \mu_1\Sigma^{-1}\mu_1$

$\left(\begin{array}{l} \text{for sym } A. \\ A^{-1}A = I \implies (A^{-1}A)^T = A^T(A^{-1})^T = A(A^{-1})^T = I \\ \qquad \implies A^{-1} = (A^{-1})^T, \quad A^{-1} \text{ symmetric} \end{array}\right)$

$\Sigma$ is cov matrix, Thus sym, $\implies \Sigma^{-1}$ also sym

$\implies x^T\Sigma^{-1}(\mu_1-\mu_0) = (\mu_1-\mu_0)^T\Sigma^{-1}x$

$\implies \boxed{(\mu_1-\mu_0)^T\Sigma^{-1}x > \ln\left(\frac{1-p}{p}\right) + \dfrac{\mu_1^T\Sigma^{-1}\mu_1 - \mu_0^T\Sigma^{-1}\mu_0}{2}}$

(c)(3) for $d$ Dimensions the decision boundary should be a $d-1$ surface on which $P(y=1|x) = P(y=0|x)$. we got a $d-1$ degrees of freedom in the equation in section a

(3) $S = \sum_{i=1}^{n} x_i$ , $x_i \sim U(-3,5)$

$S = n\bar{X}$

$P(S > n^2 + 0.2n) = P(\bar{X} > n + 0.2) = \left( \mu = \frac{5-3}{2} = 1 \right)$ Uniform dist

$= P(\bar{X} - \mu \geq n - 0.8) \leq P(|\bar{X} - \mu| \geq n - 0.8)$

$\leq 2\exp\left(\frac{-2n(n-0.8)^2}{8^2}\right) \leq 0.1$

$\Rightarrow n(n-0.8)^2 > 32 \ln(0.05)$

$n(n-0.8)^2 + 32\ln(0.05) > 0$

real root by wolfram$\alpha$ is $5.12$

so we choose $n \geq 6$ //

(4) (a) $E(R_{ij}) = \frac{1}{m} L_j + \frac{m-1}{m} \cdot 0 = \frac{1}{m} L_j$ (uniform dist) (for servers ass)

$E(R_i) = E\left(\sum_j R_{ij}\right) = \sum_j E(R_{ij}) = \sum_j \frac{1}{m} L_j = \frac{L}{m}$

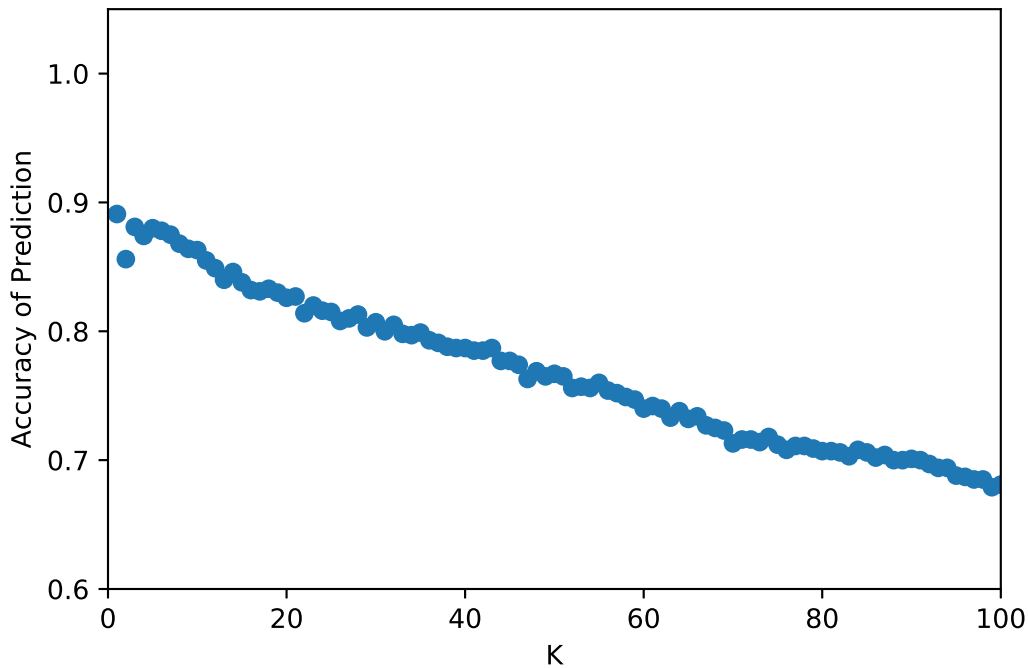(b) $P(R_i \geq (1+\sigma)E(R_i)) = \left( \frac{e^{\sigma}}{(1+\sigma)^{(1+\sigma)}} \right)^{\frac{L}{m}} \equiv C$

(c) $P\left(\bigvee_{i=1}^{m} R_i \geq (1+\sigma)E(R_i)\right) = P\left(\bigvee_{i=1}^{m} R_i \geq (1+\sigma)\frac{L}{m}\right)$

$= \sum_i P\left(R_i \geq (1+\sigma)\frac{L}{m}\right) = \sum_i C = mC$ //

# Programming Assignment (Discussion)

(b) Running knn for K=10 and a train data of size 1000 we got accuracy of ~0.87. meaning 87% of all test images were classified correctly. The expected accuracy of a random predictor should be ~0.1, There are 10 possible labels and each label has an equal prob to be chosen.

(c) we got the best accuracy for K=1, and it is easy to see that the accuracy decreases with K. it is possible that for this dataset, the nearest neighbor for most images will carry the same label, but still, the images can't be perfectly clustered in an euclidian space so we are likely to find many images with different labeles in the neighborhood of a given image.

(d) we can see that the accuracy increases with N. meaning the more train data we have, The more likely we are to include relativly close neighbors with labels identical to that of a given image.

Accuracy of Prediction in KNN Algorithm with Respect to K

Accuracy of Prediction in KNN Algorithm (K = 1)
with Respect to Train Data Size