

Introduction to Machine Learning

Section 1

Linear Algebra

1.a

\Rightarrow symmetric matrix A is PSD such that $v^t A v = (v^t u) \text{diag}(\lambda)(u^t v) = \sum_i \lambda_i (V u^t)^2 \geq 0$ where λ is the Eigenvalue of A .
and matrix A can be decomposed as:

$$A = Q D Q^t = Q * \text{diag}(\lambda_1, \lambda_2 \dots \lambda_n) * Q^t = \\ Q * \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2} \dots \sqrt{\lambda_n}) * \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2} \dots \sqrt{\lambda_n}) * Q^t = X X^t$$

\Leftarrow A can be written as $v^t X X^t v$ we get:

$$v^t A v = v^t X X^t v = (X^t v)^t (X^t v) = \| X_v^t \|^2 \geq 0$$

1.b

for a given PSD matrix A and $\alpha \in R$

(*) $v^t(\alpha A) \geq 0 \Rightarrow v^t(u A) \geq 0$ when $u, A \geq 0$

then for PSD matrix's A, B when $A, B \geq 0$, $A + B \geq 0$

now let's apply (*) on $(A+B)$ we will get (**)

$$v^t(A + B)v = v^t A v + v^t B v \geq 0$$

then from both (*) and (**) immediately get $\alpha A + \beta B \geq 0$

the set of all $n \times n$ PSD matrices over R is not a vector space over R because its not apply the closures to multiplication in scalar property, for $\lambda < 0$ and

$$A \geq 0 \rightarrow \lambda A < 0 \rightarrow \lambda A \notin \{PSD\}$$

Calculus and Probability

1.a

for $x_1, x_2 \dots x_n$ i.i.d $U([0, 1])$ continuous random variables, lets write the Order Statistics such as $\overline{x}_1, \overline{x}_2 \dots \overline{x}_n$ when $\forall i, \overline{x}_i \leq \overline{x}_{i+1}$ first lets find the CDF of $Y = MAX\{x_1, x_2 \dots x_n\} = \overline{x}_n$:

$$F_y(x) = F_{\overline{x}_n} = \Pr(\overline{x}_n \leq k) = \Pr(\overline{x}_1 \leq k, \overline{x}_2 \leq k \dots \overline{x}_n \leq k)$$

because they i.i.d

$$\Pr(\overline{x}_1 \leq k) \Pr(\overline{x}_2 \leq k) \dots \Pr(\overline{x}_n \leq k) = [\Pr(\overline{x}_i \leq k)]^n = [F_x(k)]^n$$

$$(*)F(x_i) = \begin{cases} 0 & \text{for } x < 0 \\ x/1 & \text{for } x \in \{0, 1\} \\ 1 & \text{for } x > 1 \end{cases} \quad f(x) = \begin{cases} 1 & \text{for } x \in \{0, 1\} \\ 0 & \text{for } x \notin \{0, 1\} \end{cases}$$

we get:

$$f_y(k) = f_{\overline{x}_n}(k) = \frac{d}{dk}(F_{\overline{x}_n}(k))^n = n(F(k))^{n-1}f(k)$$

now lets set values in (*) and get :

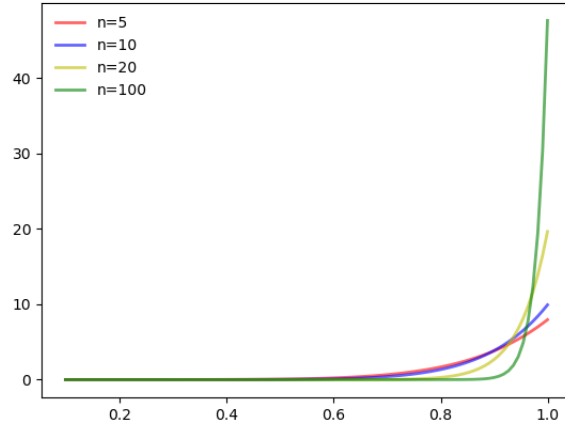
$$f_y(k) = nk^{n-1}f(k) = nk^{n-1}f(k) = nk^{n-1}I_{(1,0)} \sim Beta(n, 1)$$

therefore :

$$\lim(E[y])_{n \rightarrow \inf} = \lim\left(\frac{n}{n+1}\right)_{n \rightarrow \infty} \longrightarrow 1$$

and :

$$(var[y])_{n \rightarrow \inf} = \lim\left(\frac{n}{(n+1)^2(n+2)}\right)_{n \rightarrow \inf} \rightarrow 0$$



2

$$E[|x - \alpha|] = \int_{-\infty}^{+\infty} |x - \alpha| f(x) dx = \int_{-\infty}^{\alpha} |x - \alpha| f(x) dx + \int_{\alpha}^{+\infty} |x - \alpha| f(x) dx$$

when $\alpha \in \text{argmin}$:

$$\underbrace{(\alpha - x)f(x)}_{\rightarrow 0} + \int_{-\infty}^{\alpha} f(x) dx + \underbrace{(x - \alpha)f(x)}_{\rightarrow 0} - \int_{\alpha}^{+\infty} f(x) dx$$

$$\Rightarrow \int_{-\infty}^{\alpha} f(x) dx = \int_{\alpha}^{+\infty} f(x) dx \Rightarrow \Pr(x \leq \alpha) = \Pr(x \geq \alpha) \Leftrightarrow$$

.

$$\Pr(x \leq \alpha) = 1/2$$

Optimal Classifiers and Decision Rules

1.a

Let X and Y be random variables where Y can take values in $Y = \{1, \dots, L\}$, and Let ℓ be the 0-1 loss function defined in class , hence:

$$E[\Delta(y, f(x))] = \sum_{k=1}^L Pr(X = \hat{x}, Y = k) \Delta(k, f(x))$$

using bayes :

$$\sum_{k=1}^L Pr(X = \hat{x}) Pr(Y = k|X = \hat{x}) \Delta(k, f(k)) = Pr(X = \hat{x}) \sum_{k=1}^L Pr(y = k|X = \hat{x}) \Delta(k, f(k))$$

$$L(h) = Arg \min_{f: X \rightarrow Y} \{Pr(x = \hat{x}) \sum_{y \neq k, y \in \{1 \dots L\}} Pr(y = k|X = \hat{x})\} = f(\hat{x}) = h(x) = k$$

$$\Rightarrow h(\hat{x}) = Arg \max \{Pr(x = \hat{x})(1 - Pr(y = k|X = \hat{x})) : h(\hat{x}) = k$$

$$h(\hat{x}) = Arg \max_{y \in Y} Pr(y = i|x = \hat{x})$$

Optimal Classifiers and Decision Rules

1.b

To find decision rule for:

$$Pr[y = 1|X] > Pr[y = 0|X]$$

lets apply bayes rule on both sides. we get:

$$\frac{f_{X|Y=1}(x) Pr[Y = y] f_X(X)}{f_x} > \frac{f_{X|Y=0}(x) Pr[Y = y] f_X(X)}{f_x}$$

$$p f_1(x, \mu_1, \Sigma) > (1 - p) f_0(x, \mu_0, \Sigma)$$

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

$$(x - \mu_0)^T \sum^{-1} (x - \mu_0) - (x - \mu_1)^T \sum^{-1} (x - \mu_1) > 2 \ln\left(\frac{1-p}{p}\right)$$

where $(x - \mu)^T \sum^{-1} (x - \mu)$ is the square Mahalanobis Distance between x and μ

so our simpler Decision rule will be

$$h(x) = \begin{cases} 1 & \text{for } d^2_{\mathbf{m}}(x, \mu_0) > d^2_{\mathbf{m}}(x, \mu_1) + 2 \ln\left(\frac{1-p}{p}\right) \\ 0 & \text{otherwise} \end{cases}$$

1.c

when $d=1$ the general Matrix \sum size will be size $d \times d$, so the shape of the decision shape boundary will be just dot, in the same way when $d=2$ we will have a line, and for general d it might be d -dimensional shape...

1.d

For $d = 1, \mu_0 = \mu_1 = \mu$ and $\sigma_1 \neq \sigma_0$ we looking for equation in the decision rule formula we go had above:

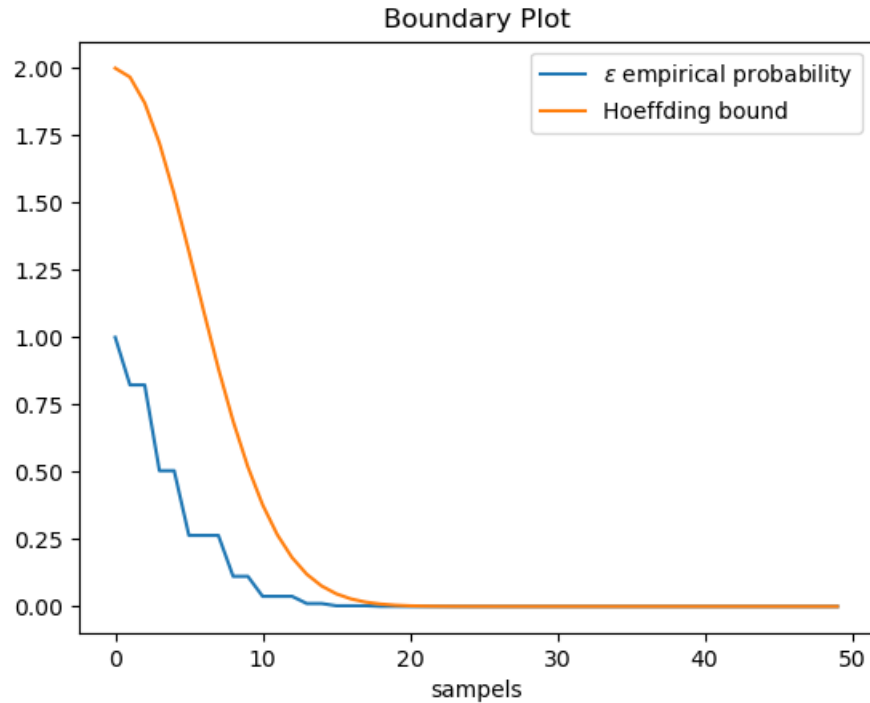
$$d^2_{\mathbf{m}}(x, \mu_0) - d^2_{\mathbf{m}}(x, \mu_1) = 2 \ln\left(\frac{1-p}{p}\right)$$

$$(x - \mu)^2 \left(\frac{1}{\sigma_0^2} - \frac{1}{\sigma_1^2} \right) = 2 \ln\left(\frac{1-p}{p}\right) \Rightarrow (x - \mu)^2 = (\sigma_0^2 - \sigma_1^2) 2 \ln\left(\frac{1-p}{p}\right)$$

$$(x - \mu) = \pm \sqrt{(\sigma_0^2 - \sigma_1^2) 2 \ln\left(\frac{1-p}{p}\right)} \Rightarrow x = \mu \pm \sqrt{(\sigma_0^2 - \sigma_1^2) 2 \ln\left(\frac{1-p}{p}\right)}$$

Programming Assignment

Visualizing the Hoeffding bound:



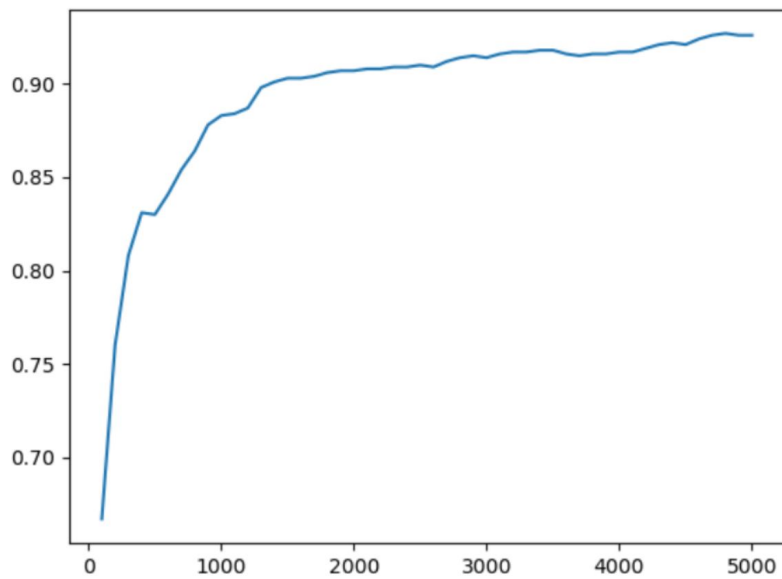
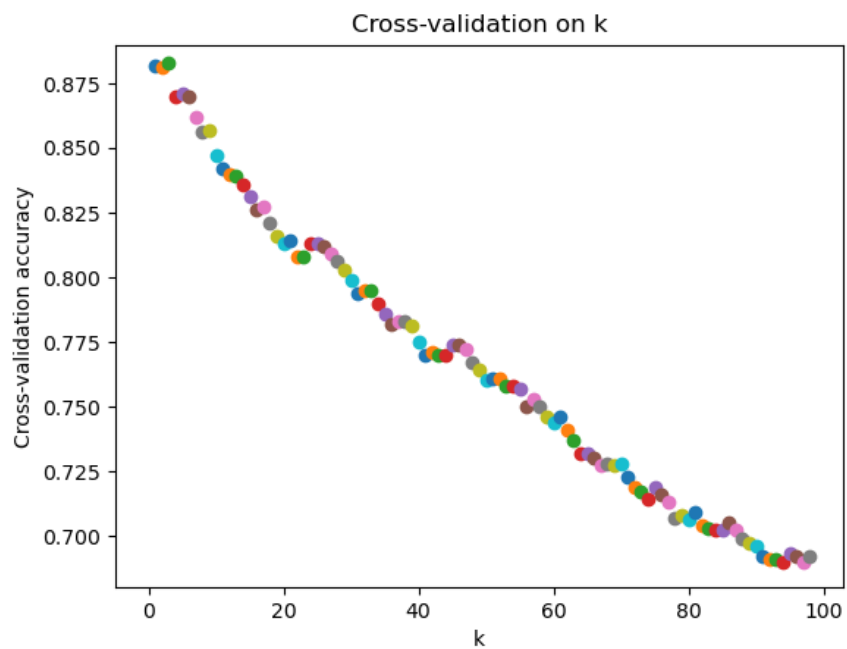
Nearest Neighbor:

1

.The KNN accuracy for $k = 10$. its got 882/1000 correct labeling. and the accuracy rate is 0.882000

2

The best K found is $k = 4$ with 883/1000 correct labeling.and the accuracy rate is 0.883000



It was my first time "LaTeXing" hope it was find (: