

Problem 2 : Bayes Rule

1. What is total probability of phone being defective?

Answer Given, A produces $\rightarrow 20\%$, B $\rightarrow 30\%$ and C $\rightarrow 50\%$.

Then, Probability of phone being produced by factory A is

$$P(A) = 20/100 = 0.2$$

$$\text{Similarly, } P(B) = 30/100 = 0.3$$

$$P(C) = 50/100 = 0.5$$

Also given, probability of defective phone from A is 2%.

$$\Rightarrow P(D/A) = 2/100 = 0.02$$

$$\text{Similarly, } P(D/B) = 1/100 = 0.01$$

$$1) \quad P(D/C) = \frac{0.05}{100} = 0.0005$$

Now, To calculate the probability of a phone being defective, we can use law of total probability & Bayes Rule i.e. $P(A) = \sum_i P(A/B_i) P(B_i)$

$$\Rightarrow P(D) = P(D/A) P(A) + P(D/B) P(B) + P(D/C) P(C)$$

$$= (0.02) * (0.2) + (0.01) * (0.3) + (0.0005) * 0.5$$

$$\boxed{P(D) = 0.00725}$$

2. Probability that this defective phone is from factory A?

Answer To calculate this, we can use Baye's Rule :-

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

Applying this in this problem:-

$$P(A/D) = \frac{P(D/A) P(A)}{P(D)}$$

Date ___ / ___ / ___

$$= \frac{(0.02) * (0.2)}{0.00725} = \frac{0.004}{0.00725}$$

$$\boxed{P(A/D) = 0.5517}$$

3. Probability that this defective phone is from factory B.

Answer: Similarly, using Bayes Rule -

$$P(B/D) = \frac{P(D/B) P(B)}{P(D)} = \frac{(0.01) * (0.3)}{0.00725}$$

$$\Rightarrow \boxed{P(B/D) = 0.4138}$$

4. Probability that this defective phone is from factory C?

Answer: Similarly, again using Bayes Rule -

$$P(C/D) = \frac{P(D/C) P(C)}{P(D)} = \frac{0.0005 * 0.5}{0.00725}$$

$$\Rightarrow \boxed{P(C/D) = 0.03448}$$

Problem-1 : Independence and law of Total probability

1. Compute $P(X=1)$, Given $P(Y=1) = 0.9$, $P(Z=1) = 0.8$,
 $P(X=1 | Y=1, Z=1) = 0.6$, $P(X=1 | Y=1, Z=0) = 0.1$,
and $P(X=1, Y=0) = 0.2$. Y and Z are independent events.

Ans - Using law of total probability - $P(A) = \sum_i P(A|B_i) P(B_i)$

We have 3 cases here when $X=1$ i.e.

$$P(X=1 | Y=1, Z=1) = 0.6$$

$$P(X=1 | Y=1, Z=0) = 0.1$$

$$P(X=1, Y=0) = 0.2$$

Hence, applying law of total probability and summing up these 3 cases, we get \rightarrow

Date _____

$$P(X=1) = P(X=1 \mid Y=1, Z=1) \cdot P(Y=1, Z=1) + \\ P(X=1 \mid Y=1, Z=0) \cdot P(Y=1, Z=0) + \\ P(X=1, Y=0) \cdot P(Y=0)$$

Further, expanding $P(Y=1, Z=1)$ & $P(Y=1, Z=0)$ using conditional probability, we get →

$$P(X=1) = P(X=1 \mid Y=1, Z=1) \cdot P(Y=1) \cdot P(Z=1) + \\ P(X=1 \mid Y=1, Z=0) \cdot P(Y=1) \cdot P(Z=0) + \\ P(X=1 \mid Y=0) \cdot P(Y=0)$$

(as Y and Z are independent events).

$$\Rightarrow P(X=1) = 0.6 * 0.9 * 0.8 + \\ 0.1 * 0.9 * (1 - 0.8) + \\ 0.2 * (1 - 0.9)$$

$$\Rightarrow P(X=1) = 0.47$$

2. Compute Expected value $E[Y]$

Answer: $E[Y] = \sum (y * P(Y=y))$, where summation is for all values of y .

here, we are $y=0$ or $y=1$.

$$\Rightarrow E[Y] = 0 * P(Y=0) + 1 * P(Y=1) \\ = 0 * 0.1 + 1 * 0.9$$

$$E[Y] = 0.9$$

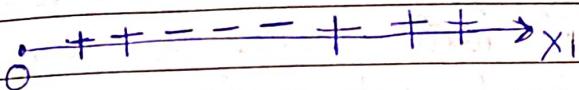
3. Given, instead of 0 and 1, y is taking 115 and 20 as values.
 $P(Y=115) = 0.9$.

$$\Rightarrow P(Y=20) = 1 - 0.9 = 0.1$$

$$E[Y] = 115 * P(Y=115) + 20 * P(Y=20)$$

$$E[Y] = 105.5$$

Problem 3 : Feature Transformation & Kernels .

1.  x_1

Consider 1-D dataset in Fig 1(a). Can you think of a 1-D transformation that will make the points linearly separable

Ans. Yes, we can make the polynomial transformation of degree 2, which transform the points into a parabola. This following equation is for the same which will make these points linearly separable -

$$\phi(x) = x^2 \quad \text{e.g. } \phi(x) = x^2 + c \quad (c \text{ is constant})$$

2. Still consider the same above 1-D dataset (in Fig 1(a)). Can you think of a 2-D transformation that makes points linearly separable.

Ans. Yes, we can have following 2-D transformation →

$$\phi(x, y) = [x, y^2 + c]$$

where x and y are variables & c is constant real value. The 1st dimension x is simply the original data points. The 2nd dimension y^2 , is a quadratic function of the original data point. This means, the transformed data will be able to catch non-linear relationship between original data points.

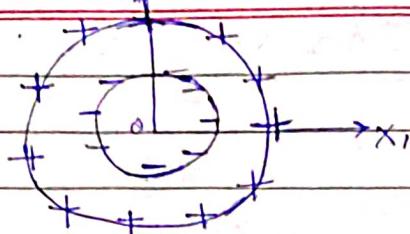
For e.g. → Let's say $x = [1, 2, 3]$

then, by applying $\phi(x, y) = [x, y^2 + 1]$, we get -

$$x_{\text{transformed}} = [[1, 2], [2, 5], [3, 10]]$$

Date ___ / ___ / ___

3.



Consider the 2-D dataset, can you suggest a 1-D transformation that will make the data linearly separable.

Ans yes, we can use the following transformation technique that leverages the radial distance of each data point from the origin. here are the steps-

- a. Calculate the radial distance (r) of each data point from the origin $(0,0)$ in 2D space using the Euclidean distance formulae:

$$r = \text{square root } (x^2 + y^2) = \sqrt{x^2 + y^2}$$

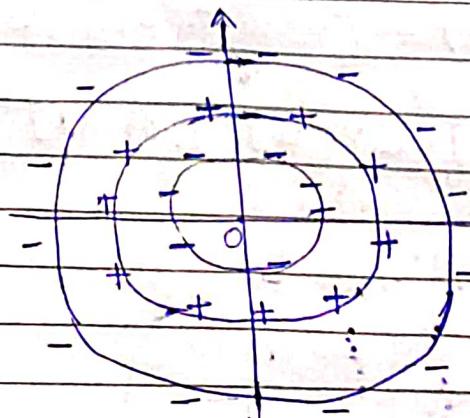
- b. Assign class label (1 and 0) to both classes respectively.

- c. The transformed 1D dataset will consist of pairs of (r, label) , where " r " = radical distance and 'label' is assigned class label (0 or 1)

- d. we can use value $r_{\text{threshold}}$ to classify datapoints as:- If $r \leq r_{\text{threshold}}$, then assign label 1.
& If $r > r_{\text{threshold}}$, assign label 0 .

By choosing appropriate ' $r_{\text{threshold}}$ ' and following these steps, we can achieve linear separation between two classes.

4.



Can you suggest a 2D transformation of the dataset, that makes it linearly separable?

Answer

No, we cannot make datapoints belonging to different class linearly separable here in this case.

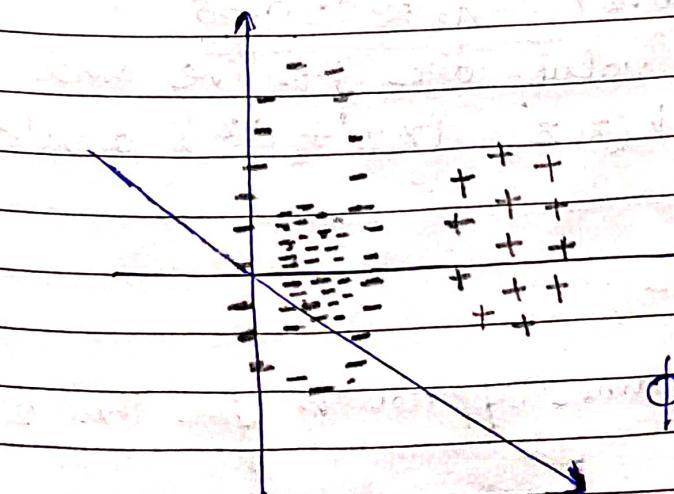
This figure has 3 concentric circles, the innermost and the outermost belongs to the same class, but the middle circle's datapoints belong to a different class.

The linear Kernel which is simple dot product of the input vectors essentially works in the original feature space and would not be suitable here where a transformation of dataset is required.

The RBF Kernel, which implicitly maps data into a higher-dimensional space, can potentially make the data linearly separable by transforming it into a high dimensional space where the classes become separable. This implies RBF will require to transform the dataset which is 2-D in original space into a higher dimensional space.

Similarly, Polynomial Kernel is also incapable to make 2-D transformation to make classes linearly apart. Because, it also maps data into higher dimensional space using suitable polynomial functions.

For this dataset, a transformation from 2-D to 3-D is at least required to make them linearly separable.



This figure depicts the 3-D transformed representation where negative and positive classes are linearly separable

$$\phi: x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2 \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \in \mathbb{R}^3$$

Kernel or Not: For the following functions, prove or disprove that it is a valid Kernel.

$$1. K(x, z) = (xz + 1)^2$$

Ans - To find if the function is a valid Kernel or not, we need to show kernel matrix is symmetric and all of its eigen values are non-negative.

Symmetry :

$$K(x, z) = (xz + 1)^2 = (zx + 1)^2 = K(z, x)$$

$\therefore K$ is symmetric here.

For calculating eigen values, Let's construct a matrix of 3×3 , $X = \{x_1, x_2, x_3\}$, $Z = \{z_1, z_2, z_3\}$

$$\Rightarrow K = \begin{vmatrix} (x_1 z_1 + 1)^2 & (x_1 z_2 + 1)^2 & (x_1 z_3 + 1)^2 \\ (x_2 z_1 + 1)^2 & (x_2 z_2 + 1)^2 & (x_2 z_3 + 1)^2 \\ (x_3 z_1 + 1)^2 & (x_3 z_2 + 1)^2 & (x_3 z_3 + 1)^2 \end{vmatrix}$$

\Rightarrow Let's say $X = \{x_1=1, x_2=2, x_3=3\}$, $Z = \{z_1=1, z_2=2, z_3=3\}$
Substituting the values, we get \rightarrow

$$\Rightarrow K = \begin{vmatrix} 4 & 9 & 16 \\ 9 & 25 & 49 \\ 16 & 49 & 100 \end{vmatrix}$$

\Rightarrow eigen values $\lambda_1, \lambda_2, \lambda_3$ are:-
 $\lambda_1 \approx 0.15, \lambda_2 \approx 8.27, \lambda_3 \approx 120.58$

Since, all the eigen values are positive and K is symmetric, $K(x, z) = (xz + 1)^2$ is a valid Kernel

$$2. K(x, z) = (xz - 1)^3$$

Answer We will use the same approach for this 2nd part two.

Symmetry.

$$K(x, z) = (xz - 1)^3 = (zx - 1)^3 = K(z, x)$$

therefore, K is symmetric

Now, calculating the Kernel matrix, by substituting the same $x_1=1, x_2=2, x_3=3$ & $z_1=1, z_2=2, z_3=3$, we get -

$$K = \begin{vmatrix} (1 \cdot 1 - 1)^3 & (1 \cdot 2 - 1)^3 & (1 \cdot 3 - 1)^3 \\ (2 \cdot 1 - 1)^3 & (2 \cdot 2 - 1)^3 & (2 \cdot 3 - 1)^3 \\ (3 \cdot 1 - 1)^3 & (3 \cdot 2 - 1)^3 & (3 \cdot 3 - 1)^3 \end{vmatrix}$$

$$\Rightarrow K = \begin{vmatrix} 0 & 1 & 8 \\ 1 & 7 & 26 \\ 8 & 26 & 64 \end{vmatrix}$$

Eigen values are $\lambda_1 \approx -3.56, \lambda_2 \approx 22.16, \lambda_3 \approx 45.40$.

In this case, λ_1 is negative, means not all λ_i 's are non negative, which implies

$K(x, z) = (xz - 1)^3$ is not a valid Kernel for set of input points because they do not satisfy the Mercer's condition.

4. Problem 4 : Exponential Family & Geometric Distribution

- Consider the geometric distribution parameterized by ϕ

$$p(y; \phi) = (1-\phi)^{y-1} \phi, y = 1, 2, 3, \dots$$

Show that geometric distribution is the exponential family, and give $b(y), \eta, T(y)$ and $\alpha(n)$.

Proof :- To show that geometric distribution parameterized by ϕ belongs to the exponential family, we need to express its Probability mass Function (PMF) in the canonical form of exponential family. The PMF of geometric distribution is defined as :-

$$p(y; \phi) = (1 - \phi)^{y-1} \phi$$

The canonical form of the distribution is given as -

$$p(y; \theta) = e^{(\eta(y) * T(y) - a(\eta(y)))}$$

where e = exponent,

$\eta(y)$ = The natural parameter

$T(y)$ = The sufficient statistic

$a(\eta)$ = The log partition function.

The value of $\eta(y)$, taking log of PMF, we get -

$$\eta(y) = \log(p(y; \phi)) = (y-1) * \log(1-\phi) + \log(\phi)$$

$$T(y) = y$$

Now, Let's put ϕ in terms of η

$$\rightarrow \phi = e^{(\eta(y) - (y-1) * \log(1-\phi))}$$

And $a(\eta)$ here will be log of sum of probabilities for all values of y and probability distributions sum to 1.

$$\Rightarrow a(\eta) = \log \left(\sum_{y=1}^{\infty} (1-\phi)^{y-1} \phi \right)$$

$$\Rightarrow a(\eta) = \log \left(\phi / (1 - (1-\phi)) \right) \quad (\text{This is geometric series with common ratio of } (1-\phi))$$

$$\Rightarrow a(\eta) = \log \left(\phi / \phi \right)$$

$$= \log 1$$

$$\Rightarrow a(\eta) = 0 \quad \text{or it can converge for infinite series.}$$

$$a(\eta) = \log \left(\frac{1}{1 - e^{-\eta}} \right)$$

Hence, putting values of $\eta(y)$, $T(y)$, $a(\eta)$, we can express the geometric distribution in the canonical form of the exponential family:-

$$P(y; \phi) = e^{(n(y)*T(y) - a(\eta))}$$

$$P(y; \phi) = e^{((y-1)*\log(1-\phi) + \log(\phi)*y - 0)}$$

2. Given a training set $\{(x_n, y_n)\}_{n=1}^N$ and let the log-likelihood of an example be $\log p(y_n | x_n; w)$. rule for learning using a GLM model with Geometric responses y .

Ans. To get the stochastic gradient ascent rule for learning GLM, we need following three assumptions -

1. Exponential Family Distribution - which we have got from 1st part of Problem 4. (Refer that).

2. Link Function : linear predictor $g(\eta) = \log(\eta)$ for the geometric function

3. Linear predictor : $\eta = w^T x$

$$\begin{aligned} \text{Log-Likelihood} &= \log [\exp [-(y-1) \cdot \log(1-\phi) - \log \\ &\quad (1 - \exp [-\log(1-\phi)])]] \\ &= (y-1) \cdot \log(1-\phi) - \log(1 - \exp[-\log(1-\phi)]) \end{aligned}$$

$$= \eta^T (y-1) - \log(1 - \exp(\eta))$$

$$= (w^T x)^T (y-1) - \log(1 - \exp(\eta))$$

$$= (w^T x)^T (y-1) - \log(1 - \exp(w^T x))$$

$$= (x^T w) \cdot (y-1) - \log(1 - \exp(w^T x))$$

On computing gradient ∇w , we get -

$$\nabla_w \log(y/n; w) = (y-1) * \nabla_w \log(1-p) + \nabla_w \log(p)$$

$$\Rightarrow \nabla_w \log(1-p) = -\nabla_w \eta$$

$$\nabla_w \log(p) = \nabla_w \eta$$

Substituting these values, we get -

$$\nabla_w \log p(y/x; w) = (y-1)(-\nabla_w \eta) + \nabla_w \eta$$

$$\Rightarrow [(1-y)*\nabla_w \eta] = (1-y)*x$$

applying stochastic gradient ascent update rule \rightarrow

$$w_{\text{new}} = w_{\text{old}} + \alpha * \nabla_w \log p(y/x; w)$$

$$w_{\text{new}} = w_{\text{old}} + \alpha(1-y)x \quad (\alpha = \text{learning rate})$$

Moreover, we can also write -

$$w + \alpha \left(\frac{\partial \text{log-likelihood}}{\partial w} \right)$$

$$= w + \alpha \left[x^T (y-1) - \frac{\exp(w^T x) \cdot x}{1 - \exp(w^T x)} \right]$$