

Q-learning + Epsilon Greedy:

Q-Learning algorithm: Find for all $s \in S$ and $a \in A$ a function $Q(s, a)$, which gives a good approximation of $Q^*(s, a)$

1. Start with random values for $Q(s, a)$ (eg all zero)

2. Choose a starting state $s_0 \in S$

3. Look up the current best action in that state. i.e.

$$a_0 = \underset{a \in A}{\operatorname{argmax}} Q(s, a)$$

OR choose a random action $a_0 \in A$ with probability $\epsilon \in [0, 1]$ (Epsilon greedy algo)

4. Apply this action and get a new state s_1 and reward

$$r_0 = R(s_0, a_0, s_1)$$

5. Update the value $Q(s_0, a_0)$ as follows (Bellman Equation)

$$Q(s_0, a_0) = (1 - \alpha) Q(s_0, a_0) + \alpha (r_0 + \gamma \max_{a \in S} Q(s_1, a))$$

Hence, $\alpha \in [0, 1]$ is the 'learning rate'

6. If s_1 is not the terminal state repeat with step 3.



AdaBoost Algo

1. For $i=1$ to N Initialize that data weight $\omega_i^{(i)} = \frac{1}{N}$

2. For $t=1$ to T

- a) Find a classifier $h_t(x)$ by minimizing the weighted error function:

$$e_t = \sum_{i=1}^N \omega_t^{(i)} \times I(y^{(i)} \neq h_t(x^{(i)}))$$

miss classify or weight ~~weight~~

- b) find the weighted error of $h_t(x)$:

$$\epsilon_t = \frac{\sum_{i=1}^N \omega_t^{(i)} \times I(y^{(i)} \neq h_t(x^{(i)}))}{\sum_{i=1}^N \omega_t^{(i)}}$$

error

and a new component is assigned votes based on its

$$\alpha_t = \ln((1 - \epsilon_t) / \epsilon_t)$$

we will assign some points / value to the classifier. Because "a good classifier, that misclassify" and "bad classifier that misclassify", The first one is more significant or bad so we assign big α value to it to maximize its priority than others. And voting time, that classifier will be chosen whose α value is highest.

For weak classifier, $\epsilon_t = 0.5$

$$\ln\left(\frac{1-0.5}{0.5}\right) = 0$$

$\therefore \alpha = 0$ [voting power is reduced]

Strong classifier, $\epsilon_t = 0.4$

$$\ln\left(\frac{1-0.4}{0.4}\right) = \underline{\underline{0.405}}$$

③ The normalized weights are updated:

$$\omega_{t+1}^{(i)} = \omega_t^{(i)} \underbrace{e^{\alpha_t}}_{\text{if correct}} \mathbb{I}(y_t^{(i)} \neq h_t(x_t^{(i)}))$$

[and normalize] + weight added if wrong

e was taken to reduce error exponentially

misclassify $\Rightarrow \alpha=0 \quad e^\alpha=1$.

↑ α \rightarrow less weight

strong

3. Combined classifier $\hat{y} = \text{sign}(H_T(x))$ where $H_T(x) = \sum_{t=1}^m \alpha_t h_t(x)$

→ linear combination of multiple functions

→ weighted sum of additive functions

sigmoid activation not to use with non-linear

SVM

Proof: $\bar{u} \rightarrow$ new value (unknown point)

$\vec{w} \cdot \bar{u} + b \geq 0$ then + [decision Rule]

→ Put additional constraints to calculate

w and b ,

$$\vec{w} \cdot \bar{x}_+ + b \geq 1 \quad \dots \text{+ve}$$

$$\vec{w} \cdot \bar{x}_- + b \leq -1 \quad \dots \text{-ve}$$

→ Carrying two equation is pain.

→ introduce new variable to make math convenient.

y_i such that $y_i = +$ for positive sample

$y_i = -$ for negative "

$$\rightarrow y_i (\vec{w} \cdot \bar{x}_i + b) \geq 1$$

$$\rightarrow y_i (\vec{w} \cdot \bar{x}_i + b) - 1 \geq 0$$

$$y_i (\vec{w} \cdot \bar{x}_i + b) - 1 = 0 \quad [\text{for } x_i \text{ in the gutter/margin}]$$

$$\text{Hence, } w x^+ + b = +1$$

$$\Rightarrow w x^+ + b - 1 = 0 \quad \text{--- (I)}$$

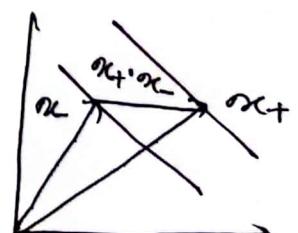
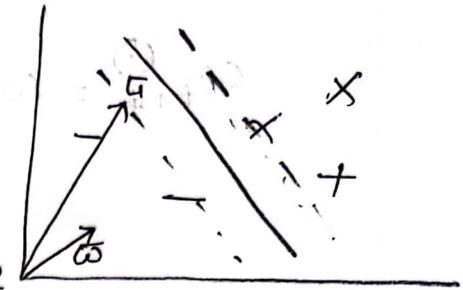
$$w x^- + b + 1 = 0 \quad \text{--- (II)}$$

$$w x^+ + b - 1 = w x^- + b + 1$$

$$\Rightarrow w (x^+ - x^-) = 2 \quad \text{--- (III)}$$

$$\Rightarrow (x^+ - x^-) = \frac{2}{\|w\|}$$

$$\Rightarrow m = \dots$$



From figure $m = (x^+ - x^-)$ is margin
we need a direction so multiply by vector.

$$m = (x^+ - x^-) \frac{w}{\|w\|}$$

$$m = \frac{2}{\|w\|} \quad [\text{from equ (11)}]$$

Maximize the margin:

$$\max \frac{2}{\|w\|} \sim \max \frac{1}{\|w\|} \sim \min \|w\| \sim \min \frac{\|w\|^2}{2}$$

minimize $\frac{\|w\|^2}{2}$ with respect to $y_i(\bar{w} \cdot \bar{x}_i + b) \geq 1$

মানুষ ক্ষেত্রে ক্ষেত্রে তার চাপ্পা আচরণ অর্থে G.D use না করে
Lagrange Multiplier use করা। Optimization problem -
তার চাপ্পা আর্গেন্ট Lm use করা।

$$L = \frac{1}{2} \|w^2\| - \sum \alpha_i [y_i (\bar{w} \cdot \bar{x}_i + b) - 1] \quad (1)$$

$$\begin{aligned} \frac{\partial L}{\partial \bar{w}} &= \bar{w} - \sum \alpha_i y_i x_i = 0 & \frac{\partial L}{\partial b} &= -\sum \alpha_i y_i = 0 \\ \Rightarrow \bar{w} &= \sum \alpha_i y_i x_i & \Rightarrow \sum \alpha_i y_i &= 0. \end{aligned}$$

Put the value of \bar{w} in L eq(1)

$$\begin{aligned} L &= \frac{1}{2} (\sum \alpha_i y_i x_i) (\sum \alpha_j y_j x_j) - \sum \alpha_i y_i x_i (\sum \alpha_j y_j x_j) - \\ &\quad \sum \alpha_i y_i b + \sum \alpha_i \end{aligned}$$

$$= \frac{1}{2} \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i x_j$$

→ parameters are multiple to each other in L so
it's a quadratic optimization problem.

- Now use qp solver to find α
 - all of the math depends on dot product ($x_i \cdot x_j$)
- Decision Rule: \rightarrow dual optimization परिच्छ सोल्व
বাস্তুতি

$$y_i(\bar{w}x_i + b) - \dots$$

$$\sum \alpha_i y_i x_i \bar{w} + b > 0$$

minimizing soft margin

$$2 \cdot \text{loss} + \frac{1}{2} \|\bar{w}\|^2 \text{ where } \text{loss} = \sum_i \alpha_i y_i (\bar{w} \cdot x_i) - b - \frac{1}{2}$$

Bayesian Decision theory

$$\text{Conditional Probability, } f(w_i | x) = \frac{P(x|w_i) P(w_i)}{P(x)}$$
$$= \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

Q8 In Bayes parameter estimation for $\theta = \{\mu\}$, the estimate is calculated from the following equation:

$$\hat{\mu}_n = \frac{n \sigma^2}{n \sigma_0^2 + \sigma^2} * \bar{x}_n + \frac{\sigma^2}{n \sigma_0^2 + \sigma^2} * \mu_0$$

Show that the Bayes estimate of θ is same as that of MLE method when large number of samples is taken.

Ans: Given, $\hat{\mu}_n$ is the MAP estimate of μ

$n \rightarrow$ number of samples

$\sigma^2 \rightarrow$ variance of the Sample mean \bar{x}

$\sigma_0^2 \rightarrow$ Variance Associated with prior distribution of μ

$\bar{x} \rightarrow$ Sample mean.

$\mu_0 \rightarrow$ mean associated with the prior distribution of μ

MLE for μ , assuming a Gaussian distribution, is simply the sample mean \bar{x} , as MLE does not account for prior distribution and maximizes the likelihood function given the observe data.

Now For MAP's equation,

① The term $\frac{n\sigma^2}{n\sigma^2 + \delta^2} \bar{x}$ will dominate as n becomes very large, because $n\sigma^2$ will grow faster than δ^2 , making the fraction approach 1.

② The term $\frac{\delta^2}{n\sigma^2 + \delta^2} \mu_0$ will diminish as n grows, because the denominator will grow much larger than the numerator, making the fraction approach to 0.

Therefore, For $n \rightarrow \infty$, the map estimate $\hat{\mu}_n$, converges to the sample mean \bar{x} , which is the MLE estimate for μ . So it concludes with large n , the prior becomes less influential and MAP estimate effectively becomes the MLE estimate.

MAP estimate is influenced by both data and prior information. But the prior information is not taken into account when the prior is not informative. Thus the prior information does not influence the MAP estimate.

Define log-likelihood function and use Maximum Likelihood Estimation method for estimating the unknown parameters $\Theta = (\mu, \sigma^2)$ of a univariate Gaussian distribution/ also $\Theta = (\mu)$

\Rightarrow For $\Theta = (\mu)$ only mean is unknown

$$P(x_k | \Theta) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x_k - \Theta)^T \Sigma^{-1} (x_k - \Theta) \right]$$

We know, it's easier to maximum if we add \ln ;

$$\begin{aligned} \ln P(x_k | \Theta) &= \ln \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x_k - \Theta)^T \Sigma^{-1} (x_k - \Theta) \right] \\ &= -\frac{1}{2} (x_k - \Theta)^T \Sigma^{-1} (x_k - \Theta) - \frac{1}{2} [\ln(2\pi)^d |\Sigma|] \end{aligned}$$

$$\text{Now, } \nabla_{\Theta} (\ln P(x_k | \Theta)) = 0$$

$$\Rightarrow \frac{\partial}{\partial \Theta} \left\{ -\frac{1}{2} [\ln(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \Theta)^T \Sigma^{-1} (x_k - \Theta) \right\} = 0$$

$$\Rightarrow -\frac{1}{2} (x_k - \Theta) \Sigma^{-1} = 0$$

$$\Rightarrow \Sigma^{-1} (x_k - \Theta) = 0$$

$$\Rightarrow \sum_{k=1}^n \Sigma^{-1} (x_k - \hat{\Theta}) = 0$$

$$\Rightarrow \sum_{k=1}^n (x_k - \hat{\Theta}) = 0$$

$$\Rightarrow \sum_{k=1}^n x_k - \sum_{k=1}^n \hat{\Theta} = 0$$

$$\Rightarrow n \hat{\Theta} = \sum_{k=1}^n x_k$$

$$\hat{\Theta} = \frac{1}{n} \sum_{k=1}^n x_k$$

For unknown mean and variance

Let $\theta = [\mu, \sigma^2]$ $\theta_1 = \mu$ $\theta_2 = \sigma^2$

$$\ln(P(x_k|\theta)) = -\frac{1}{2} \ln[\theta_2] - \frac{1}{2} (x_k - \theta_1)^t \theta_2^{-1} (x_k - \theta_1)$$

$$\nabla_{\theta} (\ln P(x_k|\theta)):$$

$$\text{with respect to } \theta_1 = \frac{1}{\theta_2} (x_k - \theta_1) \quad (1)$$

$$\text{u " " " } \theta_2 = \frac{-1}{2\theta_2} + \frac{(x - \theta_1)^2}{2\theta_2^2}$$

From ①

$$\sum_{k=1}^n \frac{1}{\theta_2} (x_k - \hat{\theta}_1) = 0$$

$$\hat{\theta}_2 = \hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2$$

From (a)

$$\sum_{k=1}^n -\frac{1}{2\hat{\theta}_2} + \frac{(x - \hat{\theta}_1)^2}{2\hat{\theta}_2^2} = 0$$

$$\Rightarrow \sum_{k=1}^n (x - \hat{\theta}_1)^2 = \sum_{k=1}^n \frac{2\hat{\theta}_2^2}{2\hat{\theta}_2}$$

$$\Rightarrow \sum_{k=1}^n (x - \hat{\theta}_1)^2 = \sum_{k=1}^n \hat{\theta}_2 \quad Q_0$$

$$\hat{\theta}_1 = \hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

Derive a decision rule that minimizes the overall risks for a two-class problem.

Lets,

$\alpha_1 \rightarrow$ action \rightarrow model predict feature to go to class w_1

$\lambda_{ij} = \lambda(\alpha_i | w_j)$ predicted i class but actual j class.

\uparrow
Loss function.

Suppose we observe X and take action α_i^* . The Conditional Risk (or expected loss) with taking action α_i^* is,

$$R(\alpha_i | x) = \sum_{j=1}^c \lambda(\alpha_i | \omega_j) P(\omega_j | x)$$

Now for two class,

$$R(\alpha_1|x) = \lambda_{11} P(\omega_1|x) + \lambda_{12} P(\omega_2|x)$$

$$R(\alpha_2|x) = \lambda_{22} P(\omega_2|x) + \lambda_{21} P(\omega_1|x)$$

If $R(\alpha_1|x) < R(\alpha_2|x)$ than $c\alpha_1$ else $c\alpha_2$.
 Right class loss: α_1 is assigned with loss 0.

As it is known, Bayes rule minimizes R by

- ① Computing $R(\alpha_i | s)$ for every α_i given on s .
 - ② Choosing the action α_i^* with minimum $R(\alpha_i | s)$

$$R^* = \min R$$

Discriminant function: $g_i(x) = \ln P(x|\omega_i) + \ln P(\omega_i)$

use $\ln \omega_i$ monotonically increasing, does not change the classification results.

Discriminant function for multivariate Gaussian Density

$$P(x|\omega_i) N(\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$

$$g_i(x) = \ln P(x|\omega_i) + \ln P(\omega_i)$$

$$\text{now, } g_i(x) = -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} \ln 2\pi - \frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) + \ln P(\omega_i)$$

Case 1: (Isotropic Gaussian) $\Sigma = \sigma^2 I$

$$\Sigma = \sigma^2 I = \sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

Description: In an Isotropic Gaussian, the covariance matrix is scaled identity matrix. This means that the variance of different variable are 'equal' and there is no correlation between variables. Features are statistically independent.

We are doing this derivation to find out on which class does a feature (x) belongs to, by finding the ratio of the \ln of the posterior of two classes for x .

$$P(x|y=i) = N(\mu_i, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_i)^T \Sigma^{-1} (x - \mu_i)\right)$$

For class = 1, taking \ln on both sides (μ_1)

$$\begin{aligned} \ln P(x|y=1) &= -\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) - \underbrace{\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma|}_{\text{constant} \rightarrow \text{ignore}} \\ &= -\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \end{aligned}$$

$$= \frac{1}{2} \left(\sum x^T x - 2x^T \mu_1 \Sigma^{-1} + \mu_1^T \mu_1 \Sigma^{-1} \right)$$

$$\begin{aligned} &= (x - \mu_1)^T (x - \mu_1) \\ &= (x - \mu_1)^T (x - \mu_1) \\ &= (x^2 - 2x^T \mu_1 + \mu_1^2) \\ &= (x^T x - 2x^T \mu_1 + \mu_1^T \mu_1) \Sigma^{-1} \end{aligned}$$

Similar for class = 2 just put μ_2

$$\ln P(x|y=2) = -\frac{1}{2} \left(\sum x^T x - 2x^T \mu_2 \Sigma^{-1} + \mu_2^T \mu_2 \Sigma^{-1} \right)$$

To find the ratio of posterior probability,

$$\ln \frac{P(y=1|x)}{P(y=2|x)} = \ln \frac{P(x|y=1) P(y=1)}{P(x|y=2) P(y=2)} \quad \begin{array}{l} \text{evidences are equal} \\ \& \text{prior are } 1 \end{array}$$

$$= \ln P(x|y=1) - \ln P(x|y=2)$$

$$= -\frac{1}{2} \left[\sum x^T x - 2x^T \mu_1 \Sigma^{-1} + \mu_1^T \mu_1 \Sigma^{-1} \right] - \left[-\frac{1}{2} \left[\sum x^T x - 2x^T \mu_2 \Sigma^{-1} + \mu_2^T \mu_2 \Sigma^{-1} \right] \right]$$

$$= -\frac{1}{2} \sum x^T x + x^T \mu_1 \Sigma^{-1} - \frac{1}{2} \mu_1^T \mu_1 \Sigma^{-1}$$

$$+ \frac{1}{2} \sum x^T x - x^T \mu_2 \Sigma^{-1} + \frac{1}{2} \mu_2^T \mu_2 \Sigma^{-1}$$

$$= \alpha (\mu_1^T \Sigma^{-1} - \mu_2^T \Sigma^{-1}) + \left(-\frac{1}{2} \mu_1^T \mu_1 \Sigma^{-1} + \frac{1}{2} \mu_2^T \mu_2 \Sigma^{-1} \right)$$

$$= \alpha w^T + w_0$$

$$= w^T \alpha + w_0$$

Thus it converts into linear equation.

$$\frac{d \mathcal{B}}{d \alpha} = g_i^{\circ}(\alpha) - g_j^{\circ}(\alpha)$$

$$\begin{aligned} w^T \alpha + w_0 &= w^T \alpha_0 + w_0 \\ \Rightarrow w^T (\alpha - \alpha_0) &= 0 \end{aligned}$$

Case 2: Diagonal Gaussian

$$\Sigma_i^{\circ} = \Sigma$$

$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

proof: same as case-1

Case 3: (Full Covariance Gaussian)

Σ_i are arbitrary.

④ Variance is not equal & has the prior

$$\begin{aligned}
 \ln \frac{P(Y=1|x)}{P(Y=2|x)} &= \ln \frac{P(x|Y=1) P(Y=1)}{P(x|Y=2) P(Y=2)} \\
 &= -\frac{1}{2} \left[\cancel{x^T \Sigma_1^{-1} x} - \cancel{2x^T \Sigma_1^{-1} \mu_1} + \cancel{\mu_1^T \Sigma_1^{-1} \mu_1} \right] \\
 &\quad - \frac{1}{2} \left[x^T \Sigma_1^{-1} x + x^T \Sigma_2^{-1} x - 2x^T \Sigma_1^{-1} \mu_1 + 2x^T \Sigma_2^{-1} \mu_2 \right. \\
 &\quad \left. + \mu_1^T \Sigma_1^{-1} \mu_1 - (\mu_2^T \Sigma_2^{-1} \mu_2) \right] - \ln |\Sigma_1| - \ln |\Sigma_2| \\
 &\quad + \ln P(Y=1) - \ln P(Y=2) \\
 &= -\frac{1}{2} \left((\Sigma_1^{-1} - \Sigma_2^{-1}) x^T x \right) - x^T (\mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1}) \\
 &\quad + \omega_0 \\
 &= \omega_{12} x^T x - \omega_{12} x + \omega_0 + c_0 \\
 &= \omega_{12} x^2 - \omega_{12} x + \omega_0 \\
 &= \text{quadratic equation}
 \end{aligned}$$

so from bayesian we can also derive quadratic equation

Recursive - 40

$$\Sigma = \begin{bmatrix} 1.1 & 0.3 \\ 0.3 & 1.9 \end{bmatrix} \quad m_1 = [0 \ 0]^T \quad m_2 = [3 \ 3]^T$$

Classify the vector: $[1.0 \quad 2.2]^\top$

According to Mahalanobis distance classifier;

As they share same covariance matrix

$$g_i(x) = -\frac{1}{2} (x - u_i)^T \sum_{j=1}^n (x - u_j)$$

Avoiding the negative and constant value

$$d_m^2(x, \mu) = (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)$$

$$= \begin{bmatrix} 1.0 & 2.2 \end{bmatrix} \begin{bmatrix} 1.1 & 0.3 \\ 0.3 & 1.9 \end{bmatrix}^{-1} \begin{bmatrix} 1.0 \\ 2.2 \end{bmatrix}$$

$$= \begin{bmatrix} 1.0 & 2.2 \\ 2.2 & 2.2 \end{bmatrix} \begin{bmatrix} 0.95 & -0.15 \\ -0.15 & 0.55 \end{bmatrix} \begin{bmatrix} 1.0 \\ 2.2 \end{bmatrix}$$

$$= \begin{bmatrix} (1 \times 0.95) + (2.2 \times (-0.15)) & (-0.15 \times 1) + (2.2 \times 0.55) \\ \end{bmatrix}$$

$$= [0.62 \quad 1.06] \begin{bmatrix} 1.0 \\ 12.2 \end{bmatrix}$$

$$= 2.952$$

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$\Sigma^{-1} = \frac{1}{1.1 \times 1.9 - (0.3)^2} \begin{bmatrix} 1.9 & -0.3 \\ -0.3 & 1.1 \end{bmatrix}$$

$$\Sigma^{-1} = \begin{bmatrix} 0.95 & -0.15 \\ -0.15 & 0.55 \end{bmatrix}$$

$$\textcircled{5} \quad \left[\begin{array}{c} 1 \cdot 0 \\ 2 \cdot 2 \end{array} \right] \quad \frac{1 \times 2}{2 \times 1}$$

$$d\mathbf{m}^2(\mu_2, \mathbf{x}) = (\mathbf{x} - \mu_2)^T \Sigma^{-1} (\mathbf{x} - \mu_2)$$

$$\text{协方差矩阵 } \Sigma = \begin{bmatrix} 2.0 & -0.8 \\ -0.8 & 2.2 \end{bmatrix}, \quad \Sigma^{-1} = \begin{bmatrix} 0.05 & -0.15 \\ -0.15 & 0.55 \end{bmatrix}$$

i don't know where it comes from

$$= 3.672$$

$$(\mu_2, \mathbf{x}) = (\mathbf{x} - \mu_2)^T \Sigma^{-1} = [1.0 \ 2.2] - [3 \ 3] = [-2 \ -0.8]$$

Zero-One Loss function

$$\lambda(x_i | \omega_j) = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases}$$

কেবলমাত্র কোণের ক্ষেত্রে 1 assign করা যাবে
otherwise 0

$$R(x_i | \mathbf{x}) = \sum_{j=1}^C \lambda(x_i | \omega_j) P(\omega_j | \mathbf{x}) = \sum_{i \neq j} P(\omega_j | \mathbf{x}) = 1 - P(\omega_i | \mathbf{x})$$

So Decision Rule:

Decide ω_1 if $[1 - P(\omega_1 | \mathbf{x})] < [1 - P(\omega_2 | \mathbf{x})]$; otherwise ω_2

Decide ω_1 if $R(x_1 | \mathbf{x}) < R(x_2 | \mathbf{x})$; " " ω_2 .

Decide ω_1 if $P(\omega_1 | \mathbf{x}) > P(\omega_2 | \mathbf{x})$; " " ω_2

General Loss Vs zero-one Loss

Decide ω_1 if,

$$\frac{P(\mathbf{x} | \omega_1)}{P(\mathbf{x} | \omega_2)} > \frac{(\lambda_{12} - \lambda_{21})P(\omega_2)}{(\lambda_{21} - \lambda_{12})P(\omega_1)} \quad \text{vs} \quad \frac{P(\mathbf{x} | \omega_1)}{P(\mathbf{x} | \omega_2)} > \frac{P(\omega_2)}{P(\omega_1)}$$

Case 1: (Isotropic Gaussian)

$$\Sigma = \sigma^2 I = \sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

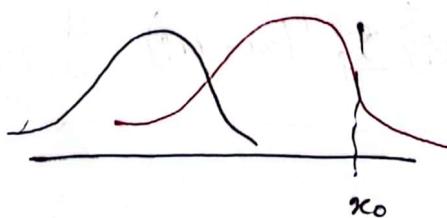
Description: In an Isotropic Gaussian, the covariance matrix is scaled identity matrix. This means that the variance of different variable are 'equal' and there is no correlation between variables. **features are statically independent.

Example: Imagine a circular (sphere) distribution in 2D where the spread of the data is the same along both dimensions.



Properties of Decision boundary:

- It passes through x_0
- It is orthogonal to the line linking the means
 $\mu_1 - \mu_2$
- If $P(\omega_1) \neq P(\omega_2)$ then x_0 shifts away from the most likely category.
- If σ is very small, the position of the boundary is insensitive to $P(\omega_1)$ and $P(\omega_2)$.



Q: When Bayesian distance become minimum distance classifier

→ Case-1
When $P(\omega_i)$ are equal then discriminant becomes,

$$g_i(x) = -\frac{\|x - \mu_i\|^2}{2\sigma^2} + \ln P(\omega_i)$$

$$g_i(x) = -\|x - \mu_i\|^2$$

This is the Euclidean distance.

[Since all feature has same variance & same prior probability it becomes Euclidean D.]

Case 2: (Diagonal Gaussian)

$$\Sigma_i = \Sigma$$

Description: The covariance matrix is diagonal, indicating that there is no correlation between variables but their variances can be different. This allows for modeling different levels of variability along each dimension.

Example: Think of an Ellipse in 2D where the major and minor axis align with the co-ordinate axes, representing different variances along the two dimensions.



$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

Properties of Decision Boundary:

- Passes through x_0
- not orthogonal to the line linking the means.
- If $P(\omega_i) \neq P(\omega_j)$ then x_0 shifts away from most likely category.

Mahalanobis distance classifier

→ When $P(\omega_i)$ are equal, then discriminant becomes

$$g_i(x) = -\frac{1}{2} (x - \mu_i)^T \Sigma^{-1} (x - \mu_i) + \ln P(\omega_i)$$

$$= -\frac{(x - \mu_i)^T (x - \mu_i)}{2 \Sigma} + \ln P(\omega_i)$$

$$= -\frac{\|x - \mu_i\|^2}{2 \sigma^2} + \ln P(\omega_i)$$

$$g_i(x) = -\frac{1}{2} \|x - \mu_i\|^2 \quad [\Sigma = \sigma^2]$$

$$= -\frac{\|x - \mu_i\|^2}{2 \sigma^2} \quad [\text{since all } P(\omega_i) \text{ are equal so } 1]$$

Hence the prior probability does not affect the decision and classification is based solely on the Mahalanobis distance term.

Case 3: (Full Covariance Gaussian)

Σ_i = arbitrary.

Description: In Full C.G., there is no restrictions on the covariance matrix. It can model different variances and arbitrary correlations between variables. The covariance matrix can have off-diagonal elements, capturing complex relations between variables.

$$\begin{bmatrix} 1 & 0.9 \\ 0.9 & 4 \end{bmatrix}$$

→ Clusters have different shapes and sizes (centered at)

Example: Envision an ellipse in 2D that is tilted and stretched, indicating correlation or dependence between the two variables.

$$\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$



D.B is determined by hyperquadrics; setting $g(x) = g_0(x)$ (hyperplane, pair of hyperplanes, hyperspheres, hyperellipsoids...).

Differences:

Case 1 Vs Case 2: The key difference lies in how they handle correlations. Case 1 assumes equal variance and no correlation, while case 2 allows for different variance but still assumes no correlation.

Case 2 Vs Case 3: Case-2 allows different variance but assumes no cross-variable correlations while, ~~case~~ case 3 allows for arbitrary correlations between variables.

Case 1 Vs Case 3: Case-1 assumes equal variances and no correlation, while ~~for~~ case-3 can model different variances and arbitrary correlations, making it more flexible but also requiring more parameters.

Hard margin vs soft margin

Hard margin:

1. Strict boundary with no tolerance of misclassification.
2. Works well when data is well-separable but may struggle with noisy or overlapping data.
3. Finds ω and b such that,

$$\Phi(\omega) = \frac{1}{2} \omega^T \omega \text{ is minimized and for all } \{(x_i, y_i)\}$$

$$y_i (\omega^T x_i + b) \geq 1$$

Soft margin:

1. Allows some misclassification.
2. Creates more flexible boundary, can handle noisy or overlapping data.
3. More robust but might sacrifice a bit of accuracy for improved generalization [misclassified]

$$\Phi(\omega) = \frac{1}{2} \omega^T \omega + C \sum \xi_i \text{ is minimized and for all } \{(x_i, y_i)\}$$

$$y_i (\omega^T x_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for all } i$$

* Parameter C controls overfitting

A small C encourages a wider margin at the expense of allowing some training points to be misclassified, promoting better generalization.

But a large C allows the model closely fit the training data which can lead to overfitting and poor generalization to new data. tries to minimize training error so narrow margin and creates overfitting.

Linear SVM vs Non-Linear SVM

* Finds a_1, \dots, a_N such that

$$Q(a) = \sum a_i - \frac{1}{2} \sum \sum a_i a_j y_i y_j [x_i \cdot x_j]$$

is minimized and

$$\textcircled{1} \quad \sum a_i y_i = 0$$

$$\textcircled{2} \quad 0 \leq a_i \leq C \text{ for all } a_i$$

$$f(x) = \sum a_i y_i x_i^T x + b$$

* Finds a_1, \dots, a_N such that

$$Q(a) = \sum a_i - \frac{1}{2} \sum \sum a_i a_j y_i y_j K(x_i, x_j)$$

is minimized and

$$\textcircled{1} \quad \sum a_i y_i = 0$$

$$\textcircled{2} \quad a_i > 0 \text{ for all } a_i$$

$$f(x) = \sum a_i y_i K(x_i, x) + b$$

~~* *~~ Higher function द्वारा आवृत्ति 1D और जास्तीय dot Product के लिए। dimension

Proofs A kernel function is some function that corresponds to an inner product in some expanded feature space.

2-dimension vector: $x = [x_1 \ x_2]$ let $K(x_i, x_j) = (1 + x_i^T x_j)^2$

Need to prove that $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$; [Linear]

$$K(x_i, x_j) = (1 + x_i^T x_j)^2$$

$$= 1 + x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} \cdot x_{i2} x_{j2} + 2 x_{i1} x_{j1} +$$

$$2 x_{i2} x_{j2}$$

$$= \begin{bmatrix} 1 & x_{i1}^2 & \sqrt{2} x_{i1} x_{i2} & x_{i2}^2 & \sqrt{2} x_{i1} & \sqrt{2} x_{i2} \end{bmatrix}^T$$

$$\begin{bmatrix} 1 & x_{j2}^2 & \sqrt{2} x_{j2} x_{i1} & x_{j1}^2 & \sqrt{2} x_{j1} & \sqrt{2} x_{j2} \end{bmatrix}$$

$$= \phi(x_i)^T \phi(x_j)$$

SVM

Proof for new unknown point \vec{u} are correct.

$$\vec{w} \cdot \vec{u} + b > 0 \quad [\text{From decision Rule}]$$

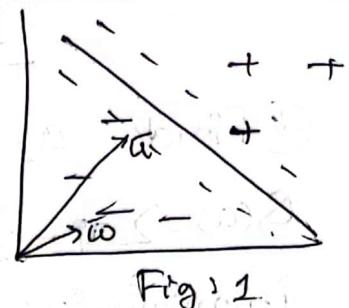


Fig: 1

Now Putting addition constraints to calculate c_0 & b .

$$\vec{w} \cdot \vec{x}_i + b \geq 1 \rightarrow +ve$$

$$\vec{w} \cdot \vec{x}_i + b \leq -1 \rightarrow -ve.$$

To combine both condition, new equation.

$$y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 \quad | \text{cohere, } y_i = + \text{ for +ve}$$

$$y_i = - \text{ for -ve}$$

For gutter value x_i^* ,

$$y_i(\vec{w} \cdot \vec{x}_i^* + b) - 1 = 0$$

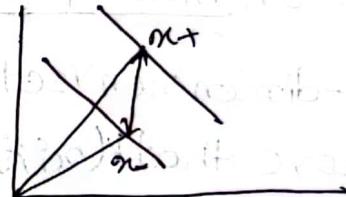


Fig: 2

Again, from Fig:2

$$w x^+ + b = +1 \Rightarrow w x^+ + b - 1 = 0 \quad (i)$$

$$w x^- + b = -1 \Rightarrow w x^- + b + 1 = 0 \quad (ii)$$

From (i) & (ii)

$$w x^+ + b - 1 = w x^- + b + 1$$

$$\Rightarrow w(x^+ - x^-) = 2 \quad (iii)$$

From Fig:2 Margin, $m = x^+ - x^-$ and we need a direct

so multiply by a vector, now m becomes.

$$m = \frac{(x^+ - x^-) \cdot w}{\|w\|} \rightarrow 2$$

$$\Rightarrow m = \frac{2}{\|w\|} \quad [\text{from equ (iii)}]$$

maximize the margins

$$\max \frac{2}{\|w\|} \sim \max \frac{1}{\|w\|} \sim \min \|w\| \sim \min \frac{\|w\|^2}{2}$$

Now minimize margin $\frac{\|w\|^2}{2}$ with respect to $y_i(\bar{w}\bar{x}_i + b) - 1$ with Lagrange multiplier.

$$L = \frac{1}{2} \|w\|^2 - \lambda \left[y_i (\bar{w}\bar{x}_i + b) - 1 \right]$$
$$= \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i (\bar{w}\bar{x}_i + b) - 1] - w$$

$$\frac{\partial L}{\partial \bar{w}} = \bar{w} - \sum \alpha_i y_i x_i = 0 \Rightarrow \bar{w} = \sum \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = - \sum \alpha_i y_i = 0 \Rightarrow \sum \alpha_i y_i = 0 \leftarrow (1)$$

Put the values in L equation (1)

$$L = \frac{1}{2} \left(\sum \alpha_i x_i y_i \right) \left(\sum \alpha_j x_j y_j \right) - \sum \alpha_i y_i x_i \left(\sum \alpha_j y_j x_j \right) - \sum \alpha_i y_i b + \sum \alpha_i$$
$$= -\frac{1}{2} \sum \alpha_i y_i x_i + \sum \alpha_i$$
$$= \sum \alpha_i - \sum \sum_j \alpha_i \alpha_j y_i y_j x_i x_j$$

parameters are multiple of each other so it's a quadratic optimization problem.