

Spectral Clustering (Graph based) - cosine distance
 From similarity matrix find Adjacency matrix (W) with nearest neigh.

$$Obj = cut(A_1, A_2, \dots, A_k) = \frac{1}{2} \sum_{k=1}^K W(A_k, A_k) \Rightarrow W(A, B) = \sum_{i \in A, j \in B} W_{ij}$$

 This is NP hard problem (No optimal solution)
 Graph Laplacian $\rightarrow L \triangleq D - W$ (D : Diagonal matrix [sum of each row in W])
 D represents degree of corresponding Node. \rightarrow constant
Properties (L): ① Each row sums to zero ② 1 Eigen vector with eigen value 0
 ③ Symmetric and positive semi-definite, ④ Has N nonnegative eigen value vectors ⑤ K comp then L has K eigen vector with eigen values 0
K-Means: Revise cluster centers using centroid or medoid (more time needed)
 using Euclidean provide circular clusters. Use Trail and error to set K

$$Obj \rightarrow J(C, R) = \sum_{i=1}^n \sum_{k=1}^K R_{ik} \|x_i - c_k\|^2 \quad (\because R_{ik} = \begin{cases} 1 & \text{if } x_i \in C_k \\ 0 & \text{otherwise} \end{cases})$$

 very sensitive to initialization. choose random, farthest, second farthest
Strengths: Simple, can be extended to other type of data, Easy to parallelise
Weakness: Circular clusters, Choose K, Not guaranteed to be optimal, hard clustering, Assign every data point to exactly one cluster
 for HC probability will be 1 for one cluster or for others.

PCA - Dimensionality Reduction (loss of information (minimize this))
Why DR: Reduce data (storage, compute, transport)
 Visualization: other algo do DR by preserving nearest neighbor property (notice)
 ① Distance between points in dataset
 improves data (important features, address correlation, remove)

PCA tries to preserve the variance of the data.
 Project data in maximal variance direction. I want a unit vector u such that when I project along unit vector my variance should be maximised. For PCA we need data to be mean centered. For this you have to subtract mean for each row to make it mean centered.

$$variance = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (\text{variance of reduced data along } u) \quad S = \frac{1}{N} \sum_{i=1}^N (x_i - u)(x_i - u)^T$$

$$S = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$$

 S-covariance matrix, use Lagrangian multiplier (λ) (Find eigen value, vector of S)
 $Su = \lambda u$ (u eigen vector of S)
 Obj. fun $\rightarrow u^T S u$ or $u^T S u$ (A value) \rightarrow The best embedding will be the eigen vector with largest eigen value.
First PC: Eigen vector of the (S) with largest eigen value. $1 \leq \lambda \leq 0$
Second PC: variance of each PC is given by λ_i | Variance captured by PC $\sum_{i=1}^N x_i^2 \times 100$
 orthogonal to first one.

Reinforcement Learning: An agent in an environment, The environment has a state and the environment will give reward on the action taken by the agent to change the state. Agent will always tries to improve the reward.
 Ex: Playing chess, Go, Robot learning to walk, Every day activity, Controller adjustment
 Parameters of an engineered system (engine, oil refinery) real time.
Exploration: learning from prev experiences and using that to the advantage to maximize the reward to work properly in this case.

Exploration: First testing the unknown and then improving the performance to that new environments so that you can actually leverage that and exploit that in future (trade off between exploration and exploitation)
 Agent operating in an env, State of the agent at time t - $S_t \in S$, Action taken by agent at time t - $A_t \in A(S_t)$, Reward at time t - $R_t \in R$, Policy - π (decision making rule)
 $\pi(s)$: $S \rightarrow A$ (Action at a given state)
Adv of RL: Learn optimal policy π that maximizes the reward (will maximize reward)
 ① Policy ② Reward signal - used by the env to inform the agent of the reward at a given time step (current state and action) ③ Value function: Expected total reward starting at a given state (hard to determine, long term desirability of a state), (A state might have a small reward but a high value - might be followed by a sequence of states with high reward)
 ④ Model - Allows us to model the env (predict next steps and next rewards) | learning parameters
 - Using the value of current state to update the value of prev state: $V(S_t) \leftarrow V(S_t) + \alpha [V(S_{t+1}) - V(S_t)]$ α step size parameter (function like α values are NN are combined with RL to handle high dimension state action policy and generalize across states, enabling RL to handle complex real world problems)

Difference from evolutionary: operates at a policy level, Evaluates policy over many games and then chooses the next policy, for each game, only the final outcome is considered. **Value function/learning**: Evaluates individual states during the course of play.
MDP: The central assumption is the decision that we take at the current step is only influenced by the action we have taken in the prev. **MDP**: The mathematical Idealized formulation of the RL problem. (Two types of learning ① Episodic ② Continuing)
 P defines the dynamic fun of the MDP. A four arguments: $S \times R \times S \times A \rightarrow [0, 1]$
Episodic reward: save the reward for a seq. of actions that machine took in response to some of states. Expected return with discount factor, γ (short sighted)
 $G_t = R_t + \gamma V_{t+1} + \gamma V_{t+2} + \dots + \gamma V_T$ ($\gamma = 0$ maximizing immediate reward (short sighted) $\gamma = 1$ more weight to future)

if y is linear, then it is Linear
 it is **Robust Regression**. Linear Regression is susceptible to outliers
Logistic Regression: If $\frac{1}{1+e^{-w^T x}} > 0.5$ then $y=1$ (or) $y=0$
 $\frac{1}{1+e^{-w^T x}} > \frac{1}{2} \Rightarrow 2 > 1+e^{-w^T x} \Rightarrow 1 > e^{-w^T x} \Rightarrow 0 > -w^T x$
 $w^T x > 0$ **Perceptron** if $w^T x > 0$, then $y=1$ else $y=0$

LR, P, SVM have same decision function, different learning strategy. **Decision function in SVM** $y = w^T x + b$
Best Boundary: Boundary that separates the data the most
 $w^T x + b = 0 \rightarrow$ line of the Margin $= \frac{2}{\|w\|}$
 $w^T x + b = -1 \rightarrow$ we want to learn w, b such that:
 for all points x_i for which $y_i = 1, w^T x_i + b > 0$
 for all points x_i for which $y_i = -1, w^T x_i + b < 0$
 $x_1 \Rightarrow \frac{2}{\|w\|}$ as large as possible

The distance from the origin to the Margin $= \frac{b}{\|w\|}$
Misclassification: To create a threshold not sensitive to outliers we use misclassification
Soft margin: The threshold which is determined after considering misclassification. Soft margin is also called SVC.
 SVC can handle outliers. And can handle overlapping data.
Kernel function: This helps to choose (or) draw SVC at a higher dimension.
Kernel Trick: compare relationships b/w observations as they would appear in higher dimension without talking them to higher dimension.

Polynomial Kernel function: This contains a parameter degree which will affect the complexity of the non-linear decision boundary.
 Minimize $f(x, y) = 2 - x^2 - 2y^2$, s.t $h(x, y) = x + y - 1 = 0$. Here we apply Lagrangian multiplier (β) \rightarrow minimize $L(x, y, \beta) = f(x, y) + \beta h(x, y)$
 for multiple constraint minimize $L(x, y, \beta) = f(x, y) + \sum \beta_i h_i(x, y)$

$$\frac{\partial L}{\partial x} = 0 \Rightarrow -2x + \beta = 0, \frac{\partial L}{\partial y} = 0 \Rightarrow -4y + \beta = 0, \frac{\partial L}{\partial \beta} = 0 \Rightarrow x + y - 1 = 0$$

$$x = \frac{\beta}{2}, y = \frac{\beta}{4}$$

$$x = \frac{4}{6}, y = \frac{4}{12}$$

 Inequality Constraints = minimize $f(x, y) = x^2 + y^3$
 s.t $g(x) = x^2 - 1 \leq 0$

Handling both constraints: minimize $f(w)$, s.t $g_i(w) \leq 0, h_i(w) = 0$
 $L(w, \alpha, \beta) = f(w) + \sum_{i=1}^I \alpha_i g_i(w) + \sum_{j=1}^J \beta_j h_j(w)$ s.t $\alpha_i \geq 0$

Primal optimization: $Op(w) = \max_{\alpha, \beta} L(w, \alpha, \beta)$
 $p^* = \min_w Op(w) = \min_w \max_{\alpha, \beta} L(w, \alpha, \beta)$ p^* solving is equivalent to solving the constrained optimization problem.
Dual Optimization: $D_d(w, \beta) = \min_w L(w, \alpha, \beta) \Rightarrow d^* = \max_{\alpha, \beta} D_d(w, \beta)$
 $d^* = \max_{\alpha, \beta} \min_w L(w, \alpha, \beta)$ $d^* \leq p^*$ weak duality theorem

Adv of Dual Optimization: 1) Easier to solve than other optimization problem 2) Determines some interesting facts about SVMs
 $d^* = p^*$ \Rightarrow strong duality, when 1) data is convex 2) constraints are affined \rightarrow if $d^* - p^*$ is very small then dual can be used as a proxy for primal.
KKT conditions: for $d^* = p^* = L(w^*, \alpha^*, \beta^*) \Rightarrow \frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial w} = 0$
 $\frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial \alpha_i} = 0, \alpha_i^* g_i(w^*) = 0, g_i(w^*) \leq 0, \alpha_i^* \geq 0$
 $\frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial \beta_j} = 0$
 min $f(w) = \frac{\|w\|^2}{2}$ s.t $y_1(w^T x_1 + b) - 1 \geq 0, x_1 = 1, y_1 = 1$
 $y_2(w^T x_2 + b) - 1 \geq 0, x_2 = 2, y_2 = 1$
 $\Rightarrow -(w^T x_1 + b) - 1 \geq 0, (w^T x_2 + b) - 1 \geq 0 \Rightarrow -w_1 - w_2 - b - 1 \geq 0$
 $2w_1 + 2w_2 + b - 1 \geq 0 \Rightarrow L = \frac{w_1^2 + w_2^2}{2} + \alpha_1(-w_1 - w_2 - b - 1) + \alpha_2(2w_1 + 2w_2 + b - 1)$

$$\frac{\partial L}{\partial w_1} = w_1 - \alpha_1 + 2\alpha_2 = 0, \frac{\partial L}{\partial w_2} = -w_1 - w_2 - b - 1 = 0$$

$$\frac{\partial L}{\partial \alpha_1} = -w_1 - w_2 - b - 1 = 0, \frac{\partial L}{\partial \alpha_2} = 2w_1 + 2w_2 + b - 1 = 0$$

$$w_1 = w_2 = 1, b = -3$$

→ Posterior \propto Likelihood \times Prior | Best hypothesis = maximum likelihood method
 → $2 \times (d-1)$ parameters for two classes
 → All features are independent and each variable can be assumed to be a Bernoulli random variable.

→ we reduce to 2D parameter
 → $P(Y=1|X=x) = P(X=1|Y=1)P(Y=1)$
 $P(X=1|Y=1) = \frac{P(X=1, Y=1)}{P(Y=1)}$
 $P(X=1|Y=0) = \frac{P(X=1, Y=0)}{P(Y=0)}$
 $P(Y=1) = P(X=1|Y=1)P(Y=1) + P(X=0|Y=1)P(Y=1)$
 $P(Y=0) = P(X=1|Y=0)P(Y=0) + P(X=0|Y=0)P(Y=0)$
 $\theta_1 = P(X=1|Y=1)$
 $\theta_0 = 1 - \theta_1$

review: $\log \text{likelihood}$
 $L(w) = \sum_{i=1}^N y_i \log \theta_i + (1-y_i) \log (1-\theta_i)$
 $\frac{dL(w)}{dw} = \sum_{i=1}^N (y_i - \theta_i) x_i$
 $w_{k+1} = w_k + \eta \frac{dL(w)}{dw}$

1) role of a prior: a prior represents the initial belief about the probabilities of classes before considering any evidence from the data. It helps incorporate prior knowledge and serves as a baseline for updating probabilities using Bayes's theorem. It makes sure that current probabilities will not have full control on final outcome.
 2) limitation: It assumes feature independence, which rarely happens in real world. 1) performs poorly on small datasets due to unreliable probability estimates 2) sensitive to noise data and irrelevant features

3) Bernoulli: because it simplifies classification into binary outcomes, aligning with the algorithmic assumption of independent feature contribution. This makes it computationally efficient and well suited for binary classification tasks.

4) conditionally independent → $q(d-1)$ totally → $2(d)$
 for multiclass, use multinomial

5) why only 2D: we assume features independent, which means avoid the need to account for feature interaction. with d features and 2 classes, we only estimate d parameters per class, resulting in 2D
 Hessian is the second order derivative of object function
 $H = - \sum_{i=1}^N \theta_i (1-\theta_i) x_i x_i^T$

→ More N → MORE D

linear regression (iid → independent, identical)

Probabilistic Interpretation: y is assumed to be normally distributed. σ^2 controls the impact of margin and margin error

$y \sim N(w^T x, \sigma^2)$
 random mean variance
 $P(Y|X) = N(w^T x, \sigma^2)$
 Log likelihood or MLE (maximum)
 $w_2 = (X^T X)^{-1} X^T Y$
 $\sigma^2 = \frac{1}{N} (Y - X w)^T (Y - X w)$

Geometric:
 $y = w^T x = w_1 x_1 + w_2 x_2 + \dots + w_d x_d$
 $w = w_1 x_1 + w_2 x_2 + \dots + w_d x_d$
 $x \rightarrow [1 \ x]$
 $y \rightarrow N \times 1$
 $X \rightarrow (N \times d)$
 $Y \rightarrow N \times 1$

minimize squared loss:
 Least Squares estimate
 $J(w) = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2$
 Gradient descent
 $\frac{\partial J(w)}{\partial w_j} = \sum_{i=1}^N (w^T x_i - y_i) x_{ij}$
 $w_j = w_j - \alpha \frac{\partial J(w)}{\partial w_j}$
 Newton's method
 setting η is sometimes tricky, too large, wrong results, too small, slow convergence

issues: 1) susceptible to outlier 2) too simplistic - underfitting 3) unstable in presence of correlated input attributes

Bayesian linear regression we assume a prior on w

→ multivariate gaussian & an appropriate prior

$P(w|D) = \frac{P(D|w)P(w)}{\text{total}}$ prior
 diff ways to predict y^* given x^* :
 1) MLE: $P(y^*|x^*) = N(w_{MLE}^T x^*, \sigma_{MLE}^2)$
 2) for map: $P(y^*|x^*) = N(w_{map}^T x^*, \sigma_{MLE}^2)$
 3) full bayesian treatment
 $P(y^*|x^*) = \int N(w^T x^*, \sigma_{MLE}^2) P(w|D) dw$
 → impact of prior on w :
 1) Penalizes large values of w 2) translates to a zero mean gaussian prior

for non linear curves: $\phi(x) \rightarrow x^2$
 basis function → $\phi(x) = [1, x, x^2, \dots, x^d]$

Ridge regression: controls overfitting
 $\hat{w} = \arg \min J(w) + \alpha \|w\|^2$
 Optim for Non separable case
 $J(w) = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2 + \frac{\lambda}{2} \|w\|^2$
 $w = (\lambda I + X^T X)^{-1} X^T Y$
 how to control overfitting:
 1) Use simpler models (might lead to underfitting)
 2) Use regularized complex models
 $\hat{\theta} = \arg \min J(\theta) + \lambda R(\theta)$ penalty paid for complexity of model
 → Ridge helps in reducing impact of correlated inputs

lasso: $\hat{w} = \arg \min J(w) + \lambda \|w\|_1$
 1) Helps in feature selection inducing sparsity
 2) Some of the features be zero and some non zero from which we can ignore zero, that is feature selection
 3) prior is Laplace distribution
 $P(w) = \frac{1}{2b} \exp \left| -\frac{|w - \mu|}{b} \right|$
 VC: analysis of how well a line fits various datasets

review:
 1) why $y = w^T x$: because this form assume linear relation b/w inputs & outputs and also forces the line to pass through origin. many models include an intercept to allow flexibility and also involve η
 2) diff of class and regre:
 class: output type is discrete categories or labels, objective is assign labels to data points.
 regre: output is continuous numerical values, objective is predict values by fitting a function

robust regression → handles outlier better
 $J(w) = \sum_{i=1}^N |y_i - w^T x_i|$
 why robust better for outlier:
 This doesn't have square in the loss function, hence impact of error can be minimized. The errors caused by outlier do not have impact

Generative model
 Discriminative model
 better as these need lesser training examples for PAC learning