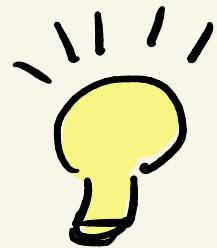


Machine Learning

Mindmap



K means clustering = Flat clustering

Input = n data points $\{x_1, \dots, x_n\}$ and # of cluster = k

Output = k disjoint sets C_1, \dots, C_k s.t. $\bigcup_{i=1}^k C_i = \{x_1, \dots, x_n\}$

Define $d(x, x') = \|x - x'\| = \left[\sum_{i=1}^d (x_i - x'_i)^2 \right]^{\frac{1}{2}}$ Euclidean norm

Cost function:

① J_{avg} Average distance to center

$$J_{avg} = \frac{1}{n} \sum_{j=1}^k \sum_{x \in C_j} d(x, m_j)$$

② J_{avg}^2 Average square distance to cluster

$$J_{avg}^2 = \frac{1}{n} \sum_{j=1}^k \sum_{x \in C_j} d(x, m_j)^2$$

↑
Cluster center of C_j

$$\textcircled{3} J_{IC} = \frac{1}{k} \sum_{j=1}^k \frac{1}{|C_j|} \sum_{x \in C_j} \sum_{x' \in C_j} d(x, x')^2$$

Sum of squared intra-cluster distance

$$\textcircled{4} J_{max} = \max_{j=1}^k \max_{x \in C_j} d(x, m_j)$$

Max distance to cluster center

$$\text{HW1: } 2J_{avg}^2 = J_{IC}$$

$$\text{Hints: } d(x, x')^2 = (x - x')^T (x - x')$$

$$d(x, m_j)^2 = (x - \frac{1}{|C_j|} \sum_{x' \in C_j} x')^T (x - \frac{1}{|C_j|} \sum_{x' \in C_j} x')$$

K means algorithm is to minimize

one of the 4 cost function

Clustering

Agglomerative clustering

Hierarchical clustering

Mixture of Gaussians

Model based clustering

Dimension Reduction

Input = $x_1, x_2, \dots, x_n \in \mathbb{R}^d$

Output = $y_1, y_2, \dots, y_n \in \mathbb{R}^p$, $p \ll d$

Linear = PCA

Nonlinear =

MDS

Unsupervised Learning

= unlabeled data"

Def = If you only have input X , you want
to learn the data structure of it

Specific procedure = r_i denote as the cluster x_i is assigned to

$$\min J_{avg}^2 = \min \sum_{j=1}^k \sum_{x \in C_j} d(x, m_j)^2 = \min \sum_{i=1}^n d(x_i, m_{r_i})^2$$

① Assign k cluster center randomly m_1, \dots, m_k (First step for the whole procedure)

② Update the cluster assignment =

$$r_i \leftarrow \arg \min_{j \in \{1, 2, \dots, k\}} d(x_i, m_j)$$

Find the closest cluster center for each x_i and get new clusters

By HW1
we know
the procedure
can min J_{avg}^2

of iteration

$\leq n^{kd}$ degree of freedom

③ Update the centroids:

$$m_j = \frac{1}{|C_j|} \sum_{i=r_i=j} x_i$$

K means ++

* Initialization: m_2 uniformly at random

Choose m_j ($j=2 \dots k$) with

$$P(m_j = x_i) = \frac{(D_{i-1}(x_j))^2}{\sum_j (D_{i-1}(x_j))^2}$$

$$D_{i-1}(x) = \min_{p \in \{1, \dots, k\}} \|x - m_p\|^2$$

The points more far away from previous cluster centroids have higher prob to be new centroid

Hierarchical Clustering

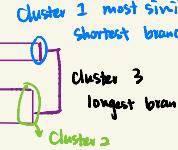
E.g. Heatmap - dendograms

	Sample 1	S ₂	S ₃
Gene 1	Red	Orange	Cyan
Gene 2	Cyan	Blue	Dark Red
Gene 3	Red	Orange	Cyan
Gene 4	Orange	Cyan	Yellow

Hierarchy

	Sample 1	S ₂	S ₃
Gene 1	Red	Orange	Cyan
Gene 3	Red	Orange	Cyan
Gene 2	Cyan	Blue	Dark Red
Gene 4	Orange	Cyan	Yellow

Dendrogram



Similarity ① "Euclidean distance"

② "Manhattan Distance"

$$\sqrt{(1.6+0.5)^2 + (0.5+1.9)^2}$$

Difference between Gene 1 and Gene 2
in Sample 1

	S ₂	S ₂
Gene 2	1.6	0.5
Gene 2	-0.5	-1.9

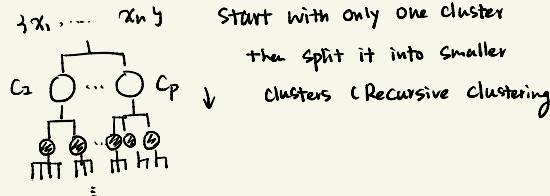
Input: $x_1, \dots, x_n \in \mathbb{R}^d$

Output: A clustering tree (Dendrogram)

Advantage: No need to specify # of clusters

You can cut the tree with k = ?
or the height in your code

Divisive Algorithm



Agglomerative Algorithm

① Start with n data points as n clusters

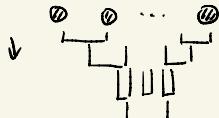
② merge clusters pairwise until they get $\{x_1, \dots, x_n\}$

Merge Method =

i. Complete linkage = $d(C_i, C_j) = \max_{x \in C_i, x' \in C_j} d(x, x')$

ii. Single linkage = $d(C_i, C_j) = \min_{x \in C_i, x' \in C_j} d(x, x')$

iii. Average linkage = $d(C_i, C_j) = \frac{1}{|C_i|} \cdot \frac{1}{|C_j|} \cdot \sum_{x \in C_i} \sum_{x' \in C_j} d(x, x')$



Mixture of Gaussian Model

"Answering How likely is it for a datapoint to belong to a certain cluster"

Input = n data points $x_1, \dots, x_n \in \mathbb{R}^d$ and # of cluster = k

Output = k disjoint sets C_1, \dots, C_k s.t. $\bigcup_{i=1}^k C_i = \{x_1, \dots, x_n\}$

(x_i, z_i)

- Observed i-th datapoint $x_i \in \mathbb{R}^d$

- $z_i \in \{1, 2, \dots, k\}$ the cluster assignment of the i-th data point (Hidden)

★ Assumption = "Generative modeling"

$$(x_i, z_i) \sim p_\theta$$

(x_i, z_i) is drawn independently from some probability distribution p_θ

We model the joint distribution as $p(x, z) = p(x|z) \cdot p(z)$

$$p(x) = \sum_{z=1}^k p(x|z) = \sum_{j=1}^k p(x|z_j) \cdot p(z_j)$$

Gaussian mixture model

Given a data point $x \in \mathbb{R}^d$, $z = 1, 2, \dots, k$

$$p(x, z) = \pi_z \cdot N(x|z; \mu_z, \Sigma_z)$$

$$p(x) = \sum_{z=1}^k \pi_z N(x|z; \mu_z, \Sigma_z), \text{ where } \sum_{z=1}^k \pi_z = 1$$

Maximum Likelihood

Likelihood for $\{(x_1, z_1), \dots, (x_n, z_n)\}$

$$L(\theta) = \log(p(x_1, z_1) \cdot p(x_2, z_2) \dots p(x_n, z_n))$$

$$= \sum_{i=1}^n \log p(x_i, z_i)$$

$$= \sum_{i=1}^n \log [\pi_{z_i} \cdot N(x_i|z_i; \mu_{z_i}, \Sigma_{z_i})]$$

$$= \text{constants} + \sum_{i=1}^n \log \pi_{z_i} - \frac{1}{2} \sum_{i=1}^n \log |\Sigma_{z_i}| - \sum_{i=1}^n (x_i - \mu_{z_i})^\top \Sigma_{z_i}^{-1} (x_i - \mu_{z_i})$$

Multivariate Normal Distribution $N(x; \mu, \Sigma)$

$$\frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}$$

Variance = Σ Mean = μ

$$\Sigma = \frac{1}{(2\pi)^{\frac{d}{2}} (2\pi)^{\frac{1}{2}}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}^{\frac{d}{2}}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Estimation Step =

$$\hat{\pi}_{ik} = \frac{\pi_k \cdot N(x_i|\hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{j=1}^k \pi_j N(x_i|\hat{\mu}_j, \hat{\Sigma}_j)}$$

Pr(data i-th belongs to j-th Gaussian/cluster)

the expected log-likelihood

$$E(L(\theta)) = \sum_{j=1}^k \sum_{i=1}^n \sum_{j=1}^k P(z_i=j | x_i, z_i=j_1, \dots, z_n=j_n) \cdot L(\theta)$$

z_i independent

$$\text{and known} = \sum_{j=1}^k \sum_{i=1}^n P(z_i=j_1 | x_i) \cdot P(z_2=j_2 | x_i) \dots P(z_n=j_n | x_i) \cdot L(\theta)$$

$$= \sum_{i=1}^n \sum_{j=1}^k \pi_{ij} \log \pi_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^k \sum_{j=1}^k \pi_{ij} (x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j)$$

Maximize $E(\theta)$ so we can update parameters $\hat{\pi}_k, \hat{\mu}_k, \hat{\Sigma}_k$

By using Lagrange method =

$$\hat{\pi}_k = \sum_{i=1}^n \frac{\pi_{ik}}{n} = \text{the new k-th weight}$$

= the average of probabilities
that a point belongs to k-th Gaussian

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \hat{\pi}_{ik} x_i}{\sum_{i=1}^n \hat{\pi}_{ik}} = \text{the new k-th mean}$$

= the weighted average of all points

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^n \hat{\pi}_{ik} (x_i - \hat{\mu}_k)^\top (x_i - \hat{\mu}_k)}{\sum_{i=1}^n \hat{\pi}_{ik}}$$

PROVED in HW 1

Dimension Reduction

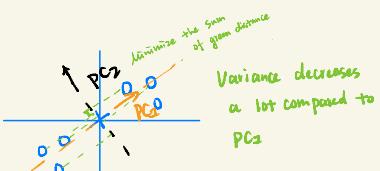
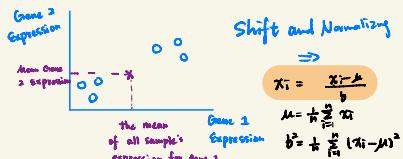
PCA (Linear)

= Find the most relevant subspace for the data" — dimension reduction

After normalizing the data. We want to keep their variance as much as possible in its subspace

= Variance " — Variance matrix $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$

Example in 2D



How to find principal components

$$\text{Maximize Variance} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{u})^2 = \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i^T \mathbf{u}, \mathbf{x}_i^T \mathbf{u} \rangle \\ = \frac{1}{n} \sum_{i=1}^n \mathbf{u}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{u} = \mathbf{u}^T \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{u}$$

\Rightarrow As $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$ is symmetric

We want $\mathbf{u} = \arg \max \langle \mathbf{u}, \hat{\Sigma} \mathbf{u} \rangle$

$\Rightarrow \mathbf{u}_1 = \mathbf{p}_d$

Advantage=

① Find the most relevant subspace for the Data

- Disadvantage=
- ① expensive to Compute
 - ② Components are not Sparse

PC₂ = "Major axis of variation"

the direction on which data's variation is maximized. In other words, find a unit vector \mathbf{u} s.t. the data is projected onto this direction with largest variance.

$$\text{Min} \left(\frac{1}{n} \sum_{i=1}^n \| \mathbf{x}_i - (\mathbf{x}_i^T \mathbf{v}) \mathbf{v} \|^2 \right) \\ \Leftrightarrow \| \mathbf{v} \|^2 = \| \mathbf{x}_i - (\mathbf{x}_i^T \mathbf{v}) \mathbf{v} \|^2 + \| (\mathbf{x}_i^T \mathbf{v}) \mathbf{v} \|^2 \\ \text{Max} \left(\frac{1}{n} \sum_{i=1}^n \| (\mathbf{x}_i^T \mathbf{v}) \mathbf{v} \|^2 \right)$$

With \mathbf{u}_1 , we can project all data points

onto subspace $\text{Span}\{\mathbf{p}_1, \dots, \mathbf{p}_{d-1}\}$

Continuous this process. We get $\mathbf{u}_2, \dots, \mathbf{u}_d$

And $\hat{\Sigma} = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^T$

(Nonlinear)

Multidimensional scaling (MDS)

• Classic MDS:

Input: n data points $\{ \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n \} \in \mathbb{R}^d$

Output: n lower dimensional points $\{ \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n \} \in \mathbb{R}^d$
 $p \leq d-1$

$$\min E_{\text{MDS}} = \| D - D^* \|_{\text{Frob}}^2 = \sum_{i,j} (D_{i,j} - D_{i,j}^*)^2$$

$$D_{i,j} = \| \mathbf{x}_i - \mathbf{x}_j \|^2 \text{ and } D_{i,j}^* = \| \mathbf{y}_i - \mathbf{y}_j \|^2$$

Gram matrix of $\{ \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \}$ $\mathbf{x}_i \in \mathbb{R}^d$

$$G_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle = \mathbf{x}_i^T \mathbf{x}_j$$

$$G = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \dots & \mathbf{x}_1^T \mathbf{x}_n \\ \vdots & \ddots & \vdots \\ \mathbf{x}_n^T \mathbf{x}_1 & \dots & \mathbf{x}_n^T \mathbf{x}_n \end{pmatrix}$$

i. Rank(G) ≤ d $\quad \text{Rank}(AB) \leq \min\{\text{Rank}(A), \text{Rank}(B)\}$

ii. For $K \in \mathbb{R}^{n \times n}$ symmetric positive semi-definite matrix of rank r, $r \leq d$, we can find n points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ s.t. $G = K$.

(HW2) Define \mathbf{X} = the first d rows of $A^T \mathbf{D}^*$
 $\mathbf{K} = \mathbf{X}^T \mathbf{X}$

propositions:

$$\text{① } E_{\text{MDS}} = \| \tilde{G} - G^* \|_{\text{Frob}}^2 = \sum (G_{ij} - G_{ij}^*)^2$$

\tilde{G} = centered Gram matrix, $G_{ij} = (\mathbf{x}_i - \mu)^T (\mathbf{x}_j - \mu)$

G^* = Gram matrix of $\{ \mathbf{y}_1, \dots, \mathbf{y}_n \}$

$$\text{② } \tilde{G} = -\frac{1}{n} P D P^T, \quad P = I - \frac{1}{n} \mathbf{1} \mathbf{1}^T, \quad D_{ij} = \| \mathbf{x}_i - \mathbf{x}_j \|^2$$

$$\text{③ } \underset{G^*}{\text{argmin}} \| \tilde{G} - G^* \|_{\text{Frob}}^2 = \mathbf{Q} \Lambda^* \mathbf{Q}^T, \quad \Lambda^* = \text{diag}(\lambda_1, \dots, \lambda_p, 0, \dots, 0) \\ \text{rank}(G^*) \leq p \quad \lambda_1 > \lambda_2 \geq \dots \geq \lambda_p$$

Approach:

1. Compute \tilde{G} (Centered Gram Matrix)
2. Compute eigendecomposition $\mathbf{Q} \Lambda^* \mathbf{Q}^T = \tilde{G}$
3. Assuming $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$, $\lambda_1 \geq \dots \geq \lambda_d$
 $\Lambda^* = \text{diag}(\lambda_1, \dots, \lambda_p, 0, \dots, 0)$, $G^* = \mathbf{Q} \Lambda^* \mathbf{Q}^T$
4. $\mathbf{y}_i = [\mathbf{Q} \Lambda^{1/2}]_{i,:}^T$

Locally Linear Embedding (LLE)

Idea = Each point should be approximately reconstructed as a linear combination

$$\text{of its neighbors } x_i = \sum_{j \in \text{Knn}(i)} w_{i,j} x_j$$

$$\sum_{j \in \text{Knn}(i)} w_{i,j} = 1.$$

Now we want to find $y_1, \dots, y_n \in \mathbb{R}^p$ s.t

$$y_i \approx \sum_{j \in \text{Knn}(i)} w_{i,j} y_j$$

Steps for this method:

1. Find the weights for each i .

For each i , we want to minimize

$$\Phi = \|x_i - \sum_{j \in \text{Knn}(i)} w_{i,j} x_j\|^2 \text{ s.t. } \sum_{j \in \text{Knn}(i)} w_{i,j} = 1$$

$$\Rightarrow \Phi = \left\| \sum_{j \in \text{Knn}(i)} w_{i,j} (x_i - x_j) \right\|^2 = w^T K^{(i)} w$$

$$w = (w_j)_{j \in \text{Knn}(i)} \quad k_{j,j}^{(i)} = (x_i - x_j)^T (x_i - x_j)$$

$$L(\lambda) = \Phi - \lambda \left(\sum_{j \in \text{Knn}(i)} w_{i,j} - 1 \right) = w^T K^{(i)} w - \lambda (w^T \mathbf{1} - 1)$$

$$\frac{\partial L}{\partial w_j} = [2 K^{(i)} w - \lambda \mathbf{1}]_j = 0 \Rightarrow w = \lambda (K^{(i)})^{-1} \mathbf{1}$$

$$\Rightarrow w = \frac{(\lambda K^{(i)})^{-1} \mathbf{1}}{\|(\lambda K^{(i)})^{-1} \mathbf{1}\|_2}$$

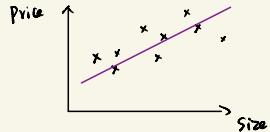
2. Find y_i 's minimize $\sum_{i=1}^n \|y_i - \sum_{j \in \text{Knn}(i)} w_{i,j} y_j\|^2$

$$\text{s.t. } \sum_{i=1}^n y_i = 0 \text{ and } \frac{1}{n} \sum_{i=1}^n y_i y_i^T = I \quad (\text{Make Problem well-posed})$$

Indicator function $i=j, b=1$

$$\Psi = \sum_{i=1}^n \|y_i - \sum_{j \in \text{Knn}(i)} w_{i,j} y_j\|^2 = \sum_{i,j} w_{i,j} G_{i,j}, \quad w_{i,j} = b_{i,j} - w_{i,j} - w_{i,i} + \sum_k w_{k,i} w_{k,j}$$

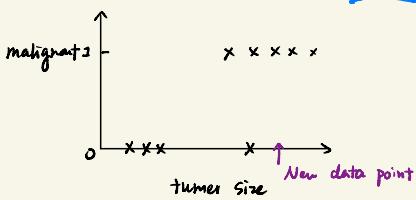
"Regression problem"



Y is continuous

Supervised Learning

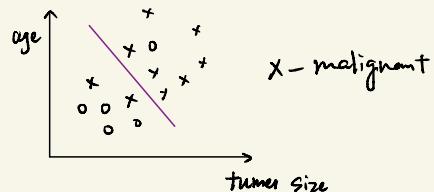
"Classification Problem"



Y is discrete

Def: Given input x with labels Y from training sets, we'd like to learn a mapping $h: x \rightarrow Y$ s.t $h(x)$ is a good predictor for new input x and the corresponding y

Logistic Regression Algorithm



SVM ∞ -dimensional features
Support vector Machine

↓ Projecting to 2D

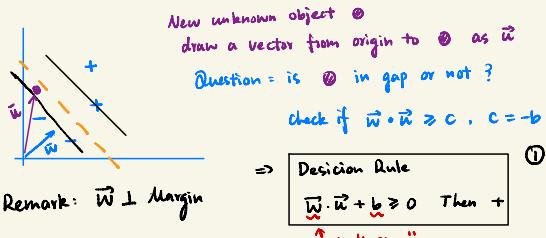
= Dimension Reduction Unsupervised

As input x could be high-dimension
in real life

SVM Support Vector Machine =

Idea: "separating data with a large gap"

Binary Classification:



For all the known samples =

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

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

Introduce y_i s.t. $y_i = +1$ For positive samples

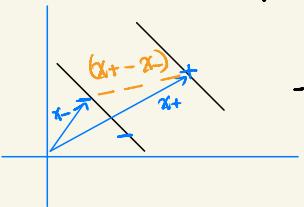
$y_i = -1$ For negative samples

$$\begin{aligned} y_i (\vec{w} \cdot \vec{x}_i + b) &\geq 1 \\ y_i (\vec{w} \cdot \vec{x}_i + b) &\geq 1 \end{aligned} \quad \Rightarrow \quad y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

Constraint

②

$y_i (\vec{w} \cdot \vec{x}_i + b) - 1 = 0$, for x_i in the cutter



$$\text{width} = (x_+ - x_-) \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

$$x_+ \cdot \vec{w} = 1 - b \quad \text{y} := y_i (\vec{w} \cdot \vec{x}_i + b) - 1 = 0 \text{ for } x_i \text{ at the margin}$$

$$-x_- \cdot \vec{w} = 1 + b$$

$$\text{width} = (1 - b + 1 + b) \cdot \frac{1}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}$$

Maximize width \Rightarrow Maximize $\frac{1}{\|\vec{w}\|} \Rightarrow \min \|\vec{w}\|^2$

$$\min \frac{1}{2} \|\vec{w}\|^2 \quad \text{③}$$

mathematically convenient

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

Find the extremum:

$$\frac{\partial L}{\partial \vec{w}} = 0$$

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum \alpha_i y_i \vec{x}_i = 0 \Rightarrow \boxed{\vec{w} = \sum \alpha_i y_i \vec{x}_i}$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i y_i = 0 \Rightarrow \boxed{\sum \alpha_i y_i = 0}$$

\vec{w} is a linear combination of \vec{x}_i samples

Plug \vec{w} back to L :

$$L = \frac{1}{2} (\sum \alpha_i y_i \vec{x}_i) \cdot (\sum \alpha_j y_j \vec{x}_j) - \sum \alpha_i [y_i (\sum \alpha_j y_j \vec{x}_j \cdot \vec{x}_i + b) - 1]$$

$$L = -\frac{1}{2} (\sum \alpha_i y_i \vec{x}_i) \cdot (\sum \alpha_j y_j \vec{x}_j) - \sum \alpha_i y_i \cdot b + \sum \alpha_i$$

④

$$L = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

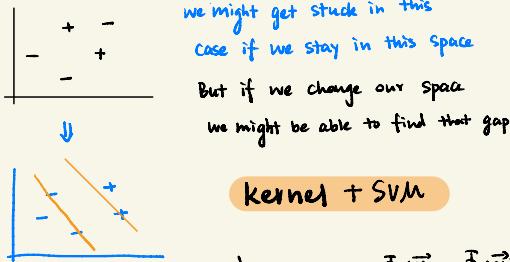
Optimization only depends on dot product

Decision Rule = dot product again

$$\sum_i \alpha_i y_i \vec{x}_i \cdot \vec{w} + b \geq 0$$

\vec{x}_i : x_i is at the margin

Remark: this question is in convex space
so we won't get stuck in local minimum.
we get global minimum

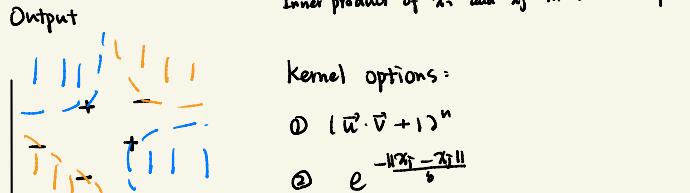


kernel + SVM

← Inspiration = Dot product

$$K(\vec{x}_i, \vec{x}_j) = \vec{\Phi}(\vec{x}_i) \cdot \vec{\Phi}(\vec{x}_j)$$

Inner product of \vec{x}_i and \vec{x}_j in another space



Remark: If b is small, shrink the distance between \vec{x}_i and \vec{x}_j . We get overfitting

Kernel

positive semidefinite =

A kernel k is positive semidefinite on X =

Given a finite sequence $\vec{x}_1, \dots, \vec{x}_n$, $\vec{x}_i \in X$ for $i=1, 2, \dots, n$

kernel Matrix define $K_{ij} = k(\vec{x}_i, \vec{x}_j)$, for any $f \in \mathbb{R}^n$, $f^T K f \geq 0$

$$f^T K f = \sum_j \sum_i f_i K_{ij} f_j = \sum_i \sum_j f_i k(\vec{x}_i, \vec{x}_j) \cdot f_j \geq 0$$

$\overset{\text{↑}}{\text{Vector in } \mathbb{R}^n}$

Symmetric = $k(x, x') = k(x', x)$, for all $x, x' \in X$

Positive = $k(x, x') \geq 0$ for all $x, x' \in X$

Example: Positive semidefinite + Symmetric but not positive

$$k(x, x') = x \cdot x', \quad X = \mathbb{R}, \quad k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$$

$$\text{positive semidefinite: } y^T K y = y^T x^T x y = \|x y\|^2 \geq 0$$

positive + symmetric but not positive Semidefinite:

$$k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, \quad k(x, x') = 1 - \cos(x, x'), \quad X = \{0, \pi\}$$

$$K = \begin{bmatrix} 1 - \cos(0), & 1 - \cos(\pi) \\ 1 - \cos(\pi), & 1 - \cos(0) \end{bmatrix} = \begin{bmatrix} 0 & 2 \\ 2 & 0 \end{bmatrix}, \quad \lambda = 2, -2$$

Proposition for PSD Kernel

- i. If k' is a psd kernel on \mathcal{X}' , given any function $\varphi: \mathcal{X} \rightarrow \mathcal{X}'$
 $k(x, x') = k'(\varphi(x), \varphi(x'))$, then k is psd kernel. (HW3)
- ii. If k_1 and k_2 are psd kernels on \mathcal{X} , $k(x, x') = k_1(x, x') + k_2(x, x')$
 is also psd kernel
- iii. $k_{\oplus}((x_1, x_2), (x'_1, x'_2)) = k_1(x_1, x_2) + k_2(x'_1, x'_2)$ is psd
- iv. $\alpha(x, x') = \text{angle between } x, x' \in \mathbb{R}^n$, $k_{\alpha}(x, x') = \cos \alpha(x, x')$ is psd

$$k_{\alpha}(x, x') = \frac{x \cdot x'}{\|x\| \|x'\|}, \text{ define } \varphi(x) = \frac{x}{\|x\|}$$

$$k_{\alpha}(\varphi(x), \varphi(x')) = k_{\alpha}(x, x'), \text{ by (i). we get psd.}$$

Gaussian Process

Def: A finite collection of random variables have a joint Gaussian distribution.

$f(x) \sim GP(m(x), k(x, x'))$ iff for all $N \in \mathbb{N}$, $(f(x_1), \dots, f(x_N)) \sim N(f_m, k)$ with

\uparrow	\uparrow
mean	covariance
function	function

$$f_{xi} = \mu(x_i) \text{ and } k_{ij} = \text{cov}(f(x_i), f(x_j)) = k(x_i, x_j)$$

Rmk: μ can be any function: $\mathbb{R} \rightarrow \mathbb{R}$

$k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is any positive semi-definite kernel on \mathbb{R}

↑ Have a lot of choices for k here

$$y = f(x) + \epsilon \quad \text{"Noise function"} \quad \epsilon \sim N(0, b^2)$$

Directed Graphical Models (Bayes nets)

Pros: Using Bayes network can significantly reduce the # of probabilities needed when describing joint distribution of (x_1, \dots, x_N)

Example: x_1, x_2, x_3 are all binary R.V



$$\text{Without Bayes net} = 2^3 - 1 = 7$$

$$\text{Bayes Net} = 1 + 2 + 2 = 3$$

$\uparrow \quad \uparrow$
 $P(x_1=1) \quad P(x_2=0|x_1=1)$
 $P(x_3=0|x_1=0)$

x_1	0	1
0	111111	1
1	111111	1

Markov Chain:

$$\text{Def: } P(x_t | x_1, \dots, x_{t-1}, x_{t+1}, \dots) = P(x_t | x_1, \dots, x_{t-1})$$

x_t should only depend on what happened in the past

Def: = k 'th order Markov Chain"

$$P(x_t | x_1, x_2, \dots, x_{t-k}, x_{t+1}, \dots) = P(x_t | x_{t-k}, \dots, x_{t-1})$$

A first order Markov Chain is said to be "stationary" if $P(x_t | x_{t-1})$ is independent from t , $P(x_t | x_{t-1}) = \mu_{x_t, x_{t-1}}$

= D-Separation"

Determine which variables are independent in a Bayes net?

Bayes Net assumption: \exists Each variable is conditionally independent of its non-descents, given its parents
↓
there's no directed path starting from x to y

Step: ① Draw ancestral graph
(all nodes in X , all ancestor of any node in X and all edges between those nodes)

② "Moralize" the parents
(Each pair of variables with a common children draw line between them)

③ "Disorient" the graph (replacing \Rightarrow with $=$)
④ Delete given and their edges

Multiplication Rule:

$$P(A_1 \cap \dots \cap A_n) = P(A_1) \cdot P(A_2 | A_1) \cdot P(A_3 | A_1 \cap A_2) \dots$$

= V structure "

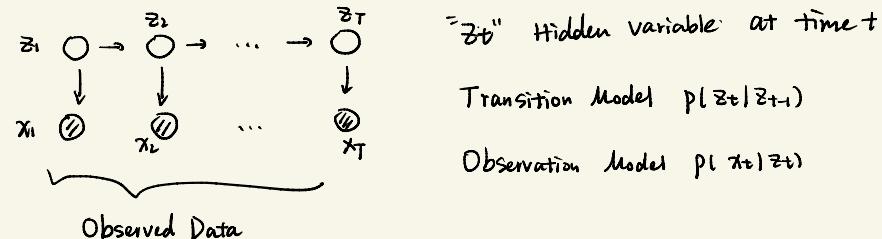
$$\begin{array}{c} \textcircled{A} \quad \textcircled{B} \\ \swarrow \quad \searrow \\ \textcircled{C} \end{array}$$
$$\begin{aligned} & P(X_A, X_B, X_C) \\ & = P(X_A, X_C) \cdot P(X_B, X_C) \\ & = P(X_A) \cdot P(X_C | X_A) \cdot P(X_B) \cdot P(X_C | X_B) \end{aligned}$$

= Inverse V structure"

$$\begin{array}{c} \textcircled{C} \\ \swarrow \quad \searrow \\ \textcircled{A} \quad \textcircled{B} \end{array}$$
$$\begin{aligned} & P(X_A, X_B, X_C) \\ & = P(X_C) \cdot P(X_A | X_C) \cdot P(X_B | X_C) \end{aligned}$$

HMM Hidden Markov Model =

A First-order discrete HMM:



Discrete HMM (General)

$$\text{HMM} = \{N, M, A, B, \pi\}$$

A: transition matrix = Hidden

Store transition probability from state i to state j
every time step, a robot must move from one state to another (it's okay to stay in place)

B: A vector of seeing the different observations
given you're in a particular state

π = probability of starting in a particular state
and then moving through the whole states
and making observations from time 1 to time T

General form of discrete Bayes net:

$$P(\pi) = \prod_{v \in V} P(x_v | \pi_{\text{pa}(v)})$$

Simpson's Paradox =

1. A graphical model can't capture all variables that might possibly relevant.
2. We can't know the Causal relationship based on observational study. We need randomized controlled trials.

Questions HMM try to answer:

- What's the probability that a model generated a sequence of observations $O = \{O_1, \dots, O_T\}$ given $\lambda = (A, B, \pi)$

- What sequence of states best explain $O = \{O_1, \dots, O_T\}$?

Find $\hat{\lambda} = \{q_1, q_2, \dots, q_T\}$?

- ③ Given $O = \{O_1, \dots, O_T\}$ Baum-Welch algorithm

How do we learn parameters that would generate them?
 $\lambda = (A, B, \pi)$?

Forward Parameter What came in the past

$$\alpha_t(i) = P(O_1, \dots, O_t, q_t = s_i | \lambda)$$

All observations capture from time 1 to t and we're at s_i state, given λ

Backward Parameter: what'll come in the future

$$\beta_t(i) = P(O_{t+1}, O_{t+2}, \dots, O_T | q_t = s_i, \lambda)$$

All observations after time t given current state is s_i and λ

$$\gamma_t(i) = P(q_t = s_i | O, \lambda)$$

given we know what comes before and after t

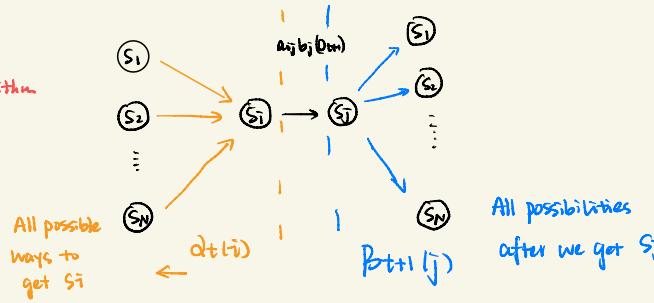
Baum-Welch Algorithm (find local optimal)

= Gradient descent algorithm

Mathematics tool =

$$\xi_t(i, j) = P(q_t = s_i, q_{t+1} = s_j | O, \lambda)$$

Prob of transition from s_i to s_j at time t given all observations and λ



$$\xi_t(i, j) = \frac{\alpha_t(i) \cdot a_{ij} \cdot b_j(O_{t+1}) \cdot \beta_{t+1}(j)}{P(O | \lambda)}$$

$$= \frac{\alpha_t(i) \cdot a_{ij} \cdot b_j(O_{t+1}) \cdot \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) \cdot a_{ij} \cdot b_j(O_{t+1}) \cdot \beta_{t+1}(j)}$$

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

$$E(\# \text{ of times } s_i \text{ is visited}) = \sum_{t=1}^{T-1} \gamma_t(i) \cdot 1$$

$$E(\# \text{ of transition from } s_i \text{ to } s_j) = \sum_{t=1}^{T-1} \xi_t(i, j)$$

$$\lambda = (A, B, \pi)$$

$$\text{Data Set} = O$$

$$\text{tools} = \alpha_t(i), \beta_t(i), \gamma_t(i), \xi_t(i, j)$$

$$\bar{\lambda} = ?$$

$$\bar{a}_{ij} = \frac{\text{Expected number of transition from } s_i \text{ to } s_j}{\text{expected number of transition from } s_i}$$

$$= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$\bar{b}_j(k) = \frac{\text{Expected # of times in state } j \text{ observing } v_k}{\text{Expected # of times in state } j}$$

$$= \frac{\sum_{t=1}^T \gamma_t(i)}{\sum_{t=1}^T \sum_{v_k \in O} \gamma_t(i)}$$

Given $\alpha_t(i), \beta_t(i), \gamma_t(i), \xi_t(i, j)$

We can produce $\bar{\lambda} = (\bar{A}, \bar{B}, \bar{\pi})$

\Rightarrow iterative steps

Undirected Graphical Models (Markov Random Fields)

Def: the general form of joint distribution

$$P(x) = \frac{1}{Z} \cdot \prod_{C \in \text{Clique}(G)} \phi_C(x_C)$$

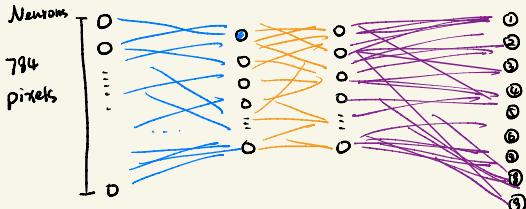
- ϕ_C = clique potential (A positive function)

- Z = normalizing factor s.t. $\sum_x \frac{1}{Z} \prod_{C \in \text{Clique}(G)} \phi_C(x_C) = 1$

Neuro Network

Example: Digit Recognition

First layer Hidden layer Output layer



Neuron - Function that takes in input to a number
 ↑
 = Activation"

Activation functions =

① Sigmoid function: $s(x) = \frac{1}{1+e^{-x}} \in [0,1]$

$$o = b + \sum_i w_{i,j} a_i + b_j$$

A point
on the
next layer
points
from current
layers

b : only activate meaningfully when weighted sum > b

$$\text{We have } 784 \times 16 + 16 \times 16 + 16 \times 10 + 16 \times 2 + 10 = 13002$$

weights
Bias

numbers for weights and bias.

Learning means to find right weights and bias

Cost function / Loss function - Machine learning minimizes Cost function

i. Feed-forward NN for regression

$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, Output y is continuous

for one example (x, y) ,

$$E(\hat{y}, y) = \frac{1}{2} \sum_{i=1}^m (a_i(x) - y_i)^2,$$

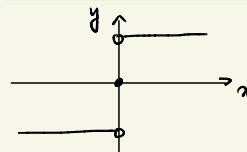
$a_i(x)$ - the output of i th neuron in the output layers

ii. Feed-forward NN for classification (y is discrete)

$$f: \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}, E(\hat{y}, y) = -\log \left(\frac{e^{a_i(x)}}{\sum_{j=1}^k e^{a_j(x)}} \right)$$

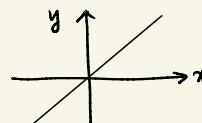
② ReLu(a) = $\max(0, a)$

③ Hard Shreshold $b(t) = \text{sgn}(t)$



Problem: Not differentiable

④ Linear: $b(t) = t$ = "identity function"



Problem:

1. it's not possible to use back propagation as derivation is 1.

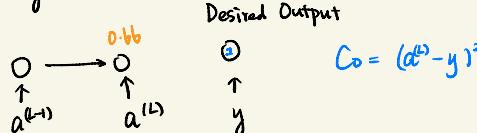
2. No matter the number of layers in the neural network, a linear activation function turns it into 2 layers

Backpropagation

(Gradient Descent Method)

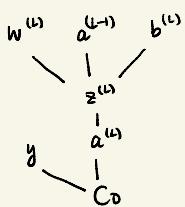
- Stochastic gradient descent
Subdivide data into "mini batches" and compute each step by the gradient of the mini batch

Imagine only one neuron for each layer



$$a^{(L)} = b(w^{(L)}a^{(L-1)} + b^{(L)})$$

$$\text{let } z^{(L)} = w^{(L)}a^{(L-1)} + b^{(L)}, \text{ then } a^{(L)} = b(z^{(L)})$$



$\frac{\partial Co}{\partial w^{(L)}}$ = the sensitivity of cost function regarding small change in weights

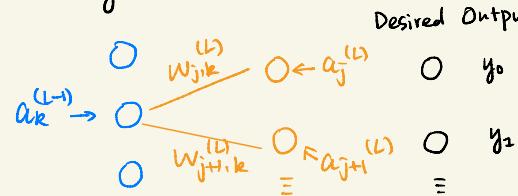
$$\frac{\partial Co}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial Co}{\partial a^{(L)}}$$

Chain Rule

$$\frac{\partial Co}{\partial w^{(L)}} = a^{(L-1)} \cdot b'(z^{(L)}) \cdot 2(a^{(L)} - y)$$

$$\frac{\partial Co}{\partial a^{(L-1)}} = \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial Co}{\partial a^{(L)}} = w^{(L)} b'(z^{(L)}) 2(a^{(L)} - y)$$

Imagine multiple neurons from one layer to next:



$$Co = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

$$z_j^{(L)} = \sum_k^n w_{j,k}^{(L)} a_k^{(L-1)}$$

$$a_j^{(L)} = b(z_j^{(L)})$$

$$\frac{\partial Co}{\partial w_{j,k}^{(L)}} = \frac{\partial z_j^{(L)}}{\partial w_{j,k}^{(L)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial Co}{\partial a_j^{(L)}}$$

$$\frac{\partial Co}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial Co}{\partial a_j^{(L)}}$$

Numerical Differentiation

(Finite difference - Taylor Expansion)

$$\frac{\partial f}{\partial x_i} \approx \frac{f(x+he_i) - f(x)}{h}$$

$h \approx \text{step size}$

Requires $O(n)$ evaluations = x_1, \dots, x_n

* Symbolic differentiation: Sympy

Automatic differentiation =

Chain Rule: $y = f(g(h(x)))$

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial w_0} \cdot \frac{\partial w_0}{\partial w_1} \cdot \frac{\partial w_1}{\partial w_2} = \frac{\partial f(w_3)}{\partial w_3} \cdot \frac{\partial g(w_2)}{\partial w_2} \cdot \frac{\partial h(x)}{\partial x}$$

Forward Accumulation computes $\frac{\partial w_i}{\partial x} = \frac{\partial w_i}{\partial w_{i-1}} \cdot \frac{\partial w_{i-1}}{\partial x}$ with $w_0 = y$

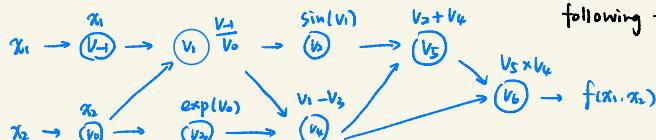
Seed = $\frac{\partial y}{\partial x} = 1$

Reverse Accumulation computes $\frac{\partial y}{\partial w_i} = \frac{\partial y}{\partial w_{i-1}} \cdot \frac{\partial w_{i-1}}{\partial w_i}$ with $w_0 = x$

Seed = $\frac{\partial y}{\partial y} = 1$

2-D case (Forward Accumulation)

$$\text{eg. } f(x_1, x_2) = [\sin(\frac{x_1}{x_2}) + \frac{x_1}{x_2} - e^{x_2}] \times [x_1 - e^{x_2}]$$



$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$

= Jacobian Matrix

$$J_f = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \dots & \frac{\partial f}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

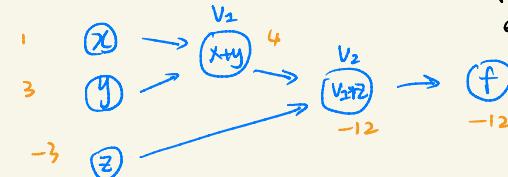
First-Order Matrix

$$= \begin{bmatrix} \nabla^T f_1 \\ \vdots \\ \nabla^T f_m \end{bmatrix}$$

Remark: Second-Order Matrix = Hessian Matrix

Simple example = $(x+y)*z = f$

Forward Pass (Orange Path) = find the value for mathematical expression.



Backward Pass (To compute the gradient)

$$\frac{\partial V_1}{\partial x} = 1 \quad \frac{\partial V_2}{\partial V_1} = \frac{\partial V_2}{\partial V_1} \cdot \frac{\partial V_1}{\partial x} = z \cdot 1 = z$$

$$\frac{\partial V_2}{\partial y} = 1 \quad \frac{\partial V_2}{\partial y} = z \quad \frac{\partial V_2}{\partial z} = V_1 = x+y$$

$$\frac{\partial V_1}{\partial x_2} = 1 \quad \frac{\partial V_1}{\partial x_1} = \frac{\partial V_1}{\partial V_2} \cdot \frac{\partial V_2}{\partial x_2} = \frac{1}{V_0} \cdot 1$$

$$\frac{\partial V_2}{\partial x_1} = \frac{\partial V_2}{\partial V_1} \cdot \frac{\partial V_1}{\partial x_1} = \cos(V_2) \cdot \frac{1}{V_0} \quad \dots$$

When $n \ll m$

Forward mode is idea