

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. $\bigcup_i C_i = \{x_1, \dots, x_n\}$

Define $d(x, x') = \|x - x'\| = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}$ Euclidean norm

Cost function:

① J_{avg} Average distance to center

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

$$m_j = \frac{1}{|C_j|} \sum_{x \in C_j} x$$

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

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

↑
cluster center of C_j

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

Sum of squared intra-cluster distance

$$④ J_{max} = \max_{i=1}^k \max_{x \in C_i} d(x, m_i)$$

Max distance to cluster center

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

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

K means algorithm is to minimize one of the 4 cost function

By HW2 we know the procedure can min J_{avg}^2
of iteration $\leq n^{kd}$ degrees of freedom

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

Linear = PCA

Nonlinear:

MDS

Unsupervised Learning

= unlabeled data"

Def: If you only have input x , you want to learn the data structure of it

K means ++

Initialization: m_2 uniformly at random

Choose m_i ($i=3, \dots, k$) with

$$p(m_i = x_j) = \frac{(D_{i-1}(x_j))^2}{\sum_i (D_{i-1}(x_i))^2}$$

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

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

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

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

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

② Update the cluster assignment = Find the closest cluster center for each x_i and get new clusters

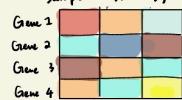
$$r_i \leftarrow \arg \min_{j \in \{1, \dots, k\}} d(x_i, m_j)$$

③ Update the centroids:

$$m_j = \frac{1}{|r_j|} \sum_{i=r_j} x_i$$

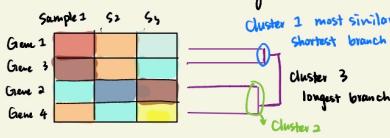
Hierarchical Clustering

Eg. Heatmap - dendograms



Hierarchy

Dendrogram



Similarity ① "Euclidean distance"

② "Manhattan Distance"

$$\begin{array}{|c|c|} \hline & S_1 & S_2 \\ \hline \text{Gene 2} & 1.6 & 0.5 \\ \hline & -0.5 & -1.9 \\ \hline \text{Gene 2} & & \\ \hline \end{array}$$

Difference between Gene 2 and Gene 2
in Sample 1

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

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

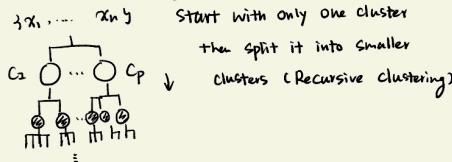
Output: A clustering tree (Dendrogram)

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: $d(C_i, C_j) = \max_{x \in C_i} \min_{x' \in C_j} d(x, x')$

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

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 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 $p(x, z) = p(x|z) \cdot p(z)$

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

Gaussian mixture model

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

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

$$p(x) = \sum_{z=1}^k \pi_z \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{\det(\Sigma_z)}} e^{-\frac{1}{2} (x - \mu_z)^T \Sigma_z^{-1} (x - \mu_z)}$$

Maximum Likelihood

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

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

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

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

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

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

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

Mean = μ

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

$$\hat{\pi}_{zi} = \frac{\pi_{z_i} N(x_i | \hat{\mu}_i, \hat{\Sigma}_i)}{\sum_{j=1}^k \pi_j N(x_i | \hat{\mu}_j, \hat{\Sigma}_j)}$$

Pr(data i th belonging to both Gaussian / cluster) \rightarrow the expected log-likelihood

$$E[L(\theta)] = \sum_{j=1}^k \sum_{i=1}^n \sum_{j=1}^k p\{z_i=j, z_n=j\} \cdots \{z_n=j\} \cdot L(\theta)$$

$$\begin{aligned} &\text{and known } \downarrow \\ &x_i \text{ independent} \\ &= \sum_{j=1}^k \sum_{i=1}^n p\{z_i=j|x_i\} \cdot p\{z_2=j|x_2\} \cdots p\{z_n=j|x_n\} \cdot L(\theta) \end{aligned}$$

$$= \sum_{i=1}^n \sum_{j=1}^k \pi_{z_i} \log \pi_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^k \pi_{z_i} \bar{x}_{ij}^T \Sigma_j^{-1} (\bar{x}_i - \mu_j)$$

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

By using Lagrange method =

$$\hat{\pi}_k = \frac{\sum_{i=1}^n \pi_{z_i}}{n} \quad \begin{aligned} &\text{the new kth weight} \\ &= \text{the average of probabilities} \\ &\text{that a point belongs to kth Gaussian} \end{aligned}$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \hat{\pi}_k \bar{x}_i}{\sum_{i=1}^n \hat{\pi}_k} \quad \begin{aligned} &\text{the new kth mean} \\ &= \text{the weighted average of all points} \end{aligned}$$

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^n \hat{\pi}_k (\bar{x}_i - \hat{\mu}_k)^T (\bar{x}_i - \hat{\mu}_k)}{\sum_{i=1}^n \hat{\pi}_k}$$

Proved in HW 1

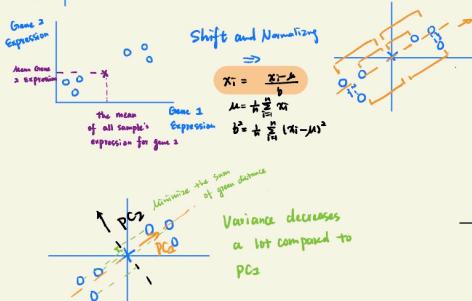
Dimension Reduction

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

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

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

We want $\mathbf{u} = \operatorname{argmax} \langle \mathbf{u}, \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \mathbf{u} \rangle$

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

Advantage:

- ① Find the most relevant subspace for the Data
- ② expensive to Compute
- ③ Components are not Sparse

PCA = 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.

$$\begin{aligned} \text{Gram 2 projection} & \quad \min \left(\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{x}_i^T \mathbf{v}) \mathbf{v}\|^2 \right) \\ & \quad \Leftrightarrow \|\mathbf{x}_i\|^2 = \|\mathbf{x}_i - (\mathbf{x}_i^T \mathbf{v}) \mathbf{v}\|^2 + \|(\mathbf{x}_i^T \mathbf{v}) \mathbf{v}\|^2 \\ \text{Gram 2} & \quad \max \left(\frac{1}{n} \sum_{i=1}^n \|(\mathbf{x}_i^T \mathbf{v}) \mathbf{v}\|^2 \right) \end{aligned}$$

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

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

Continuing this process, we get $\mathbf{u}_2, \dots, \mathbf{u}_d$

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

(Nonlinear)

Multidimensional scaling (MDS)

• Classic MDS:

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

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

s.t.

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

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

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

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

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

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

ii. For $K \in \mathbb{R}^{m \times m}$ 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 X = the first d row of $K^{\frac{1}{2}} B^T$
 $K = X^T X$

propositions:

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

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

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

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

$$\text{③ } \underset{G^*}{\operatorname{argmin}} \| \tilde{G} - G^* \|_{\text{Frob}}^2 = Q \Lambda^* Q^T, \Lambda^* = \text{diag}(\lambda_1, \dots, \lambda_d, 0, \dots, 0)$$

G^* , $\text{rank}(G^*) \leq p$

$$\lambda_1 > \lambda_2 > \dots > \lambda_p$$

Approach:

1. Compute \tilde{G} (Centered Gram Matrix)

2. Compute eigendecomposition $\tilde{G} = Q \Lambda^* Q^T = \tilde{G}$

3. Assuming $\Lambda = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$, $\lambda_1 \geq \dots \geq \lambda_d$

$$\Lambda^* = \text{Diag}(\lambda_1, \dots, \lambda_p, 0, \dots, 0)$$

$$G^* = Q \Lambda^* Q^T$$

$$Y_i = [Q \Lambda^*]_{i,*}^T$$

Locally Linear Embedding (LLE)

Idea = Each point should be approximately reconstructed as a linear combination of its neighbors $x_i = \sum_{j \in \text{km}(i)} w_{i,j} x_j$

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

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

$$y_i \propto \sum_{j \in \text{km}(i)} w_{i,j} y_j$$

Steps for this method:

2. Find the weights for each i :

For each i , we want to minimize

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

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

$$w = (w_j)_{j \in \text{km}(i)} \quad K_{j,j'}^{(i)} = (x_i - x_j)^T (x_i - x_j)$$

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

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

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

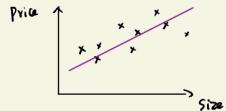
2. Find y_i 's minimize $\sum_{i=1}^n \|y_i - \sum_{j \in \text{km}(i)} w_{i,j} y_j\|^2$

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

Indicator function $i=j, b=1$

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

"Regression problem"

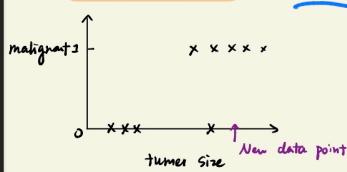


Y is continuous

Supervised Learning

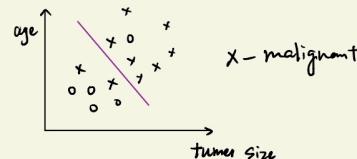
Def: Given input X with labels Y from training sets, we'd like to learn a mapping $h: X \rightarrow Y$ s.t. $h(x)$ is a good predictor for new input x and the corresponding y

"Classification Problem"



Y is discrete

Logistic Regression Algorithm



↓ Projecting to 2D

= Dimension Reduction

Unsupervised

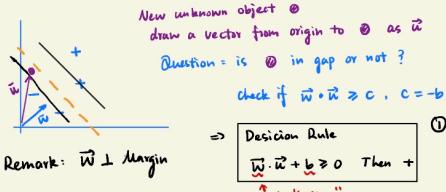
SVM ∞ -dimensional features
= Support Vector Machine

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 =

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

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

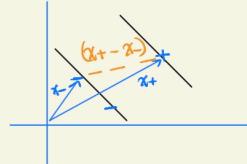
Introduce y_i s.t. $y_i = +1$ For positive samples

$y_i = -1$ For negative samples

$$\left. \begin{array}{l} y_i (\vec{w} \cdot \vec{x}_+ + b) \geq 1 \\ y_i (\vec{w} \cdot \vec{x}_- + b) \geq 1 \end{array} \right\} \Rightarrow y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

Constraint

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



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

$$\begin{aligned} \vec{x}_+ \cdot \vec{w} &= 1 - b \\ -\vec{x}_- \cdot \vec{w} &= 1 + b \end{aligned} \quad y_i := y_i(\vec{x}_i \cdot \vec{w} + b) - 1 = 0 \text{ for } \vec{x}_i \text{ at the margin}$$

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

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

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

mathematically convenient

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

Find the extremum:

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

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum_{i=1}^n y_i \vec{x}_i = 0 \Rightarrow \vec{w} = \sum_{i=1}^n y_i \vec{x}_i$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^n y_i = 0 \Rightarrow \sum_{i=1}^n y_i = 0$$

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

Plug \vec{w} back to L =

$$L = \frac{1}{2} (\sum_{i=1}^n y_i \vec{x}_i) \cdot (\sum_{j=1}^n y_j \vec{x}_j) - \sum_{i=1}^n [y_i (\sum_{j=1}^n y_j \vec{x}_j \cdot \vec{x}_i + b) - 1]$$

$$L = -\frac{1}{2} (\sum_{i=1}^n y_i \vec{x}_i) \cdot (\sum_{j=1}^n y_j \vec{x}_j) - \sum_{i=1}^n y_i \cdot b + \sum_{i=1}^n$$

(4)

$$L = \sum_{i=1}^n y_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \vec{x}_i \cdot \vec{x}_j$$

Optimization only depends on dot product

Decision Rule =

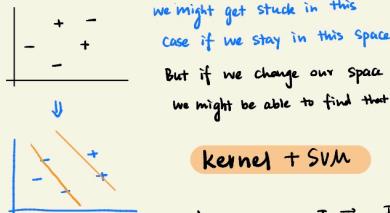
dot product again

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

\vec{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

$$k(x_i, x_j) = \vec{x}_i \cdot \vec{x}_j$$

Inner product of x_i and x_j in another space

Output



Remark: If b is small, shrink the distance between x_i and x_j . We get overfitting

Kernel

positive semidefinite =

A kernel k is positive semidefinite on X =

Given a finite sequence $\{x_1, \dots, x_n\}$, $x_i \in X$ for $i=1, 2, \dots, n$

Kernel Matrix

define $\tilde{k}_{ij} = k(x_i, x_j)$, for any $f \in \mathbb{R}^n$, $f^T k f \geq 0$

$$f^T k f = \sum_{i=1}^n \sum_{j=1}^n f_i k_{ij} f_j = \sum_{i=1}^n \sum_{j=1}^n f_i k(x_i, x_j) \cdot f_j \geq 0$$

↑
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

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

positive semidefinite: $y^T k y = y^T x^T \cdot x y = \|y\|^2 \geq 0$

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. $\theta(x, x') = \text{angle between } x, x' \in \mathbb{R}^n$, $k_{\theta}(x, x') = \cos \theta(x, x')$ is psd
$$k_{\theta}(x, x') = \frac{x \cdot x'}{\|x\| \|x'\|}, \text{ define } \varphi(x) = \frac{x}{\|x\|}$$
$$k_{\theta}(\varphi(x), \varphi(x')) = k_{\theta}(x, x'), \text{ by (i). We get psd.}$$

Gaussian Process

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

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

↑ ↑
mean covariance
function function

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

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}

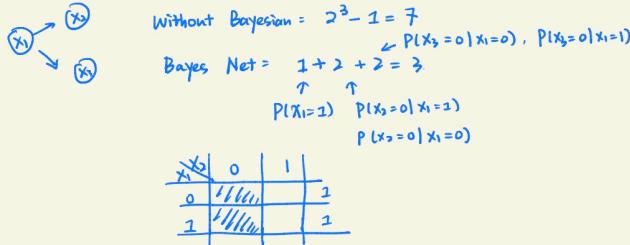
↑ Have a lot of choices for k here

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

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



Markov Chain:

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

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

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-descendants, given its parents
 \downarrow
there's no directed path starting from x to $=$

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"

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

General form of discrete Bayes net:

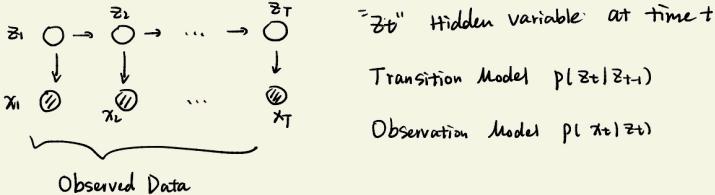
$$P(x) = \prod_{v \in V} P(x_v | \text{pa}(v))$$

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.

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

Questions HMM try to answer:

① What's the probability that a model generated a sequence of observations $O = \{O_1, \dots, O_T\}$ given $\lambda = (A, B, \pi)$

② What sequence of states best explain $O = \{O_1, \dots, O_T\}$?

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

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 were 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

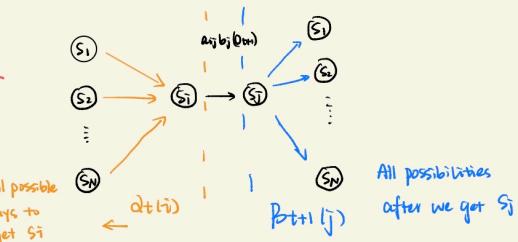
Baum-Welch Algorithm (find local optimal)

= Gradient descent algorithm

Mathematics tool:

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

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



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

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

$$\gamma_t(i) = \sum_{j=1}^N \hat{\gamma}_t(i, j)$$

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

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

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

Data Set = O

$$\text{tools: } \alpha_t(i), \beta_t(i), \gamma_t(i), \hat{\gamma}_t(i, j)$$

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

$$\begin{aligned} \bar{a}_{ij} &= \frac{\text{Expected number of transition from } s_i \text{ to } s_j}{\text{expected number of transition from } s_i} \\ &= \frac{\sum_{t=1}^{T-1} \hat{\gamma}_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \end{aligned}$$

$$\begin{aligned} \bar{b}_j(k) &= \frac{\text{Expected # of times in state } j \text{ observing } v_k}{\text{Expected # of times in state } j} \\ &= \frac{\sum_{t=1}^T \gamma_t(i)}{\sum_{t=1}^T \gamma_t(i)} \end{aligned}$$

$$\text{Given } \alpha_t(i), \beta_t(i), \gamma_t(i), \hat{\gamma}_t(i, j)$$

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

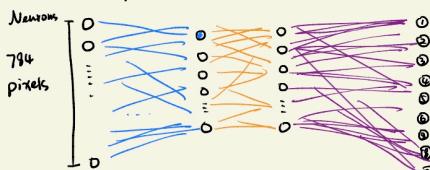
$$p(x) = \frac{1}{Z} \cdot \prod_{C \in \text{Clique}(G)} \phi_C(x_C)$$

- 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:

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

$\textcircled{2}$ $b(l) = \sum_i w_{li} a_i + b_l$
 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) ,

$$E(\hat{y}, y) = \frac{1}{2} \sum_{i=1}^m (\hat{y}_{(i)} - y_{(i)})^2,$$

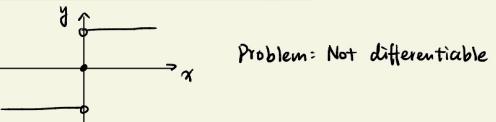
$\hat{y}_{(i)}$ - the output of i th neuron in the output layers

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

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

$$\textcircled{2} \text{ ReLu}(a) = \max(0, a)$$

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



Problem: Not differentiable

$$\textcircled{4} \text{ Linear: } b(t) = t \text{ "identity function"}$$



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 1 layer.

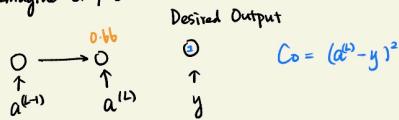
Backpropagation

(Gradient Descent Method)

- Stochastic gradient descent

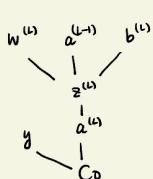
Subdivide data into "mini batches"
and compute each step by the gradient
of the mini batch

Imagine only one neuron for each layer



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

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



$$\frac{\partial C_0}{\partial w^{(L)}} = \text{the sensitivity of cost function regarding small change in weights}$$

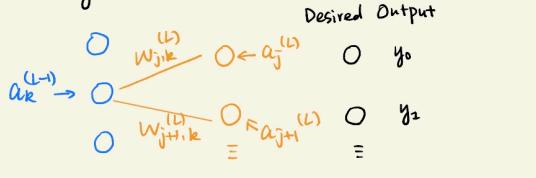
$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial C_0}{\partial a^{(L)}}$$

Chain Rule

$$\frac{\partial C_0}{\partial w^{(L)}} = a^{(L-1)} \cdot b'(z^{(L)}) \cdot 2(a^{(L)} - y)$$

$$\frac{\partial C_0}{\partial a^{(L-1)}} = \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial C_0}{\partial a^{(L)}} = w^{(L)} b'(z^{(L)}) 2(a^{(L)} - y)$$

Imagine multiple neurons from one layer to next:



$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

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

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

$$\frac{\partial C_0}{\partial w_{j,k}^{(L)}} = \frac{\partial z_j^{(L)}}{\partial w_{j,k}^{(L)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial C_0}{\partial a_j^{(L)}}$$

$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial C_0}{\partial a_j^{(L)}}$$

Numerical Differentiation

(Finite difference - Taylor Expansion)

$$\frac{\partial f}{\partial x_i} \approx \frac{f(x+h\hat{e}_i) - f(x)}{h}$$

\hat{e}_i 1st unit vector
h step size

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

* Symbolic differentiation: Sympy

Automatic differentiation

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

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

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 $w_0 = y$

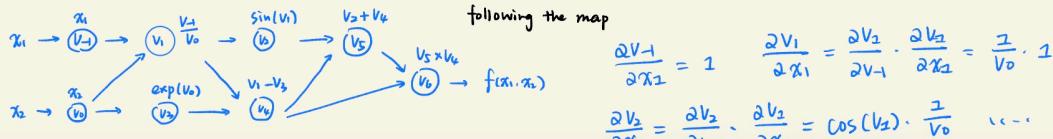
$$\text{Seed: } \frac{\partial x}{\partial x} = 1$$

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 $w_0 = x$

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

2-D case (Forward Accumulation)

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



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

= Jacobian Matrix

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

First-Order Matrix

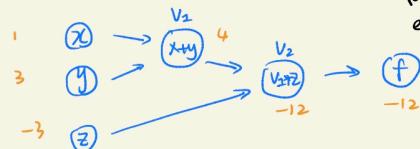
$$= \begin{bmatrix} \nabla^T f_1 \\ \vdots \\ \nabla^T f_m \end{bmatrix}$$

When $n \ll m$
Forward mode is idea

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)

$$\frac{\partial V_1}{\partial x} = 1 \quad \frac{\partial V_2}{\partial x_1} = \frac{\partial V_2}{\partial V_1} \cdot \frac{\partial V_1}{\partial x_1} = z \cdot 1 = z$$

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

We can compute $\frac{\partial f}{\partial x_i}$ by following the map

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

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