

ML TECHNIQ



WEEK 1

→ Broad Paradigms

- ↳ Supervised learning: data, features, and labelling. Classification, regression, ranking
- ↳ Unsupervised learning: only datapoints. Clustering, representational learning
- ↳ Sequential learning: make a decision → get feedback → make decision ... Reinforcement learning

→ Representation Learning

- ↳ Sub-category of Unsupervised learning
- ↳ Goal: given a set of data points, understand something useful about them

- ↳ Problem:

Input: $\{u_1, u_2, \dots, u_n\} \quad u_i \in \mathbb{R}^{d_{\text{features}}}$

Output: some "compressed" representation of the dataset

- ↳ Example:

$$\text{Data: } \left\{ \begin{bmatrix} -7 \\ -14 \end{bmatrix}, \begin{bmatrix} 2.5 \\ 5 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\}$$

$$\text{Representative: } \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \text{Coefficients: } \{-7, 2.5, 0.5, 0\}$$

- ↳ Example.

$$\text{Data: } \left\{ \begin{bmatrix} -5 \\ -5 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 7 \\ 7 \end{bmatrix}, \begin{bmatrix} 10 \\ 10 \end{bmatrix}, \begin{bmatrix} 5 \\ 6 \end{bmatrix} \right\}$$

$$\text{Rep: } \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\} \quad \text{Co-eff: } \{(-5, -5), (1, 1), (7, 7), (10, 10), (5, 6)\}$$

The only way to compress: project $\begin{bmatrix} 5 \\ 6 \end{bmatrix}$ onto the line $y = x$.

- ↳ Goal: find the line that has the least "reconstruction" error.

Dataset: $\{u_1, u_2, \dots, u_n\} \quad u_i \in \mathbb{R}^d$

$$\begin{aligned} \text{ERROR}(\text{line}, \text{dataset}) &= \sum_{i=1}^n \text{error}(\text{line}, u_i) \\ &\stackrel{\text{represented using } w \text{ as well}}{=} \sum_{i=1}^n \text{length}^2(u - (u^T w)w) \end{aligned}$$

$$\begin{aligned} f(w) &= \frac{1}{n} \sum_{i=1}^n \|u - (u^T w)w\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n (u - (u^T w)w)^T (u - (u^T w)w) \\ &= \frac{1}{n} \sum_{i=1}^n u^T u - (u^T w)^2 + (u^T w)^2 \\ &= \frac{1}{n} \sum_{i=1}^n u^T u - (u^T w)^2 \end{aligned}$$

$\rightarrow u^T u$ is a constant. So instead of minimizing $f(w)$, minimize $g(w)$

$$g(w) = \frac{1}{n} \sum_{i=1}^n -(u^T w)^2 \quad \longrightarrow \text{maximize } g(w)$$

$$\begin{aligned} \max_w g(w) &= \frac{1}{n} \sum_w (w^T u_i)(u_i^T w) = \frac{1}{n} \sum_w w^T (u_i u_i^T) w \\ &= w^T \underbrace{\left(\frac{1}{n} \sum_{i=1}^n u_i u_i^T \right)}_C w \end{aligned}$$

$$\text{equivalently: } \max_{\substack{w: \\ \|w\|=1}} w^T C w \quad ; \text{ where } C = \frac{1}{n} \sum_{i=1}^n u_i u_i^T$$

L covariance matrix

↳ The residue data (error from reconstruction) might also have information instead of just being error
 Possible algorithm to deal with this:

Input: $\{u_1, u_2, \dots, u_n\} u_i \in \mathbb{R}^d$
 → De-mean the data: $u_i = u_i - \mu \mathbf{1}_d$
 → Find the best line $w_1 \in \mathbb{R}^d$
 → Replace u_i with $u_i - (u_i^\top w_1) w_1$
 → Repeat to obtain w_2

↳ By continuing this procedure, we get

$$\{w_1, w_2, \dots, w_d\} \text{ s.t. } \underbrace{\|w_k\|^2 = 1}_{\text{orthonormal vectors}} + k \text{ and } w_i^\top w_j = 0 \quad i \neq j$$

↳ Residues after 2 rounds: $\{u_i - (u_i^\top w_1) w_1 - (u_i^\top w_2) w_2, \dots\}$ looks exactly like the gram-schmidt process
 " " d rounds: $\forall i \quad u_i - [(u_i^\top w_1) w_1 + \dots + (u_i^\top w_d) w_d]$

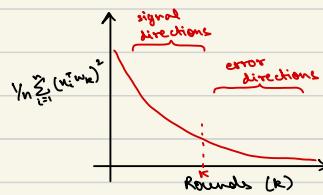
↳ Q. What if data "approximately" in low-dimensional space?

For any $w \in \mathbb{R}^d$ s.t. $\|w\|^2 = 1$

$$\forall i \quad \|u_i\|^2 = \|u_i - (u_i^\top w) w\|^2 + \|(u_i^\top w) w\|^2$$

$$Y_n \sum_{i=1}^n \|u_i\|^2 = Y_n \sum_{i=1}^n \|u_i - (u_i^\top w) w\|^2 + \underbrace{Y_n \sum_{i=1}^n \|(u_i^\top w) w\|^2}$$

The larger the value of this term, the better the fit.



WEEK 2

→ Issues with PCA

① Time complexity of finding the eigenvalues and eigenvectors
 ↳ typically $O(d^3)$

② Only looks at linear relationships between features.

→ Time-complexity issue with PCA

↳ large d [$d \gg n$] , where $d \rightarrow \# \text{features}$ $n \rightarrow \# \text{datapoints}$

$$\mathbf{x} = \begin{bmatrix} 1 & 1 & 1 \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \\ 1 & 1 & 1 \end{bmatrix} \quad \mathbf{x} \in \mathbb{R}^{d \times n} \quad \mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$$

$$\text{Covariance matrix} = \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

↳ let w_k be the k^{th} largest eigenvalue of \mathbf{C} . corresponding to λ_k

$$\mathbf{C} w_k = \lambda_k w_k$$

$$(\frac{1}{n} \sum \mathbf{x}_i \mathbf{x}_i^T) w_k = \lambda_k w_k$$

$$w_k = \frac{1}{n \lambda_k} \sum_{i=1}^n (\mathbf{x}_i^T w_k) \cdot \mathbf{x}_i$$

↳ w_k is a linear combination of data points.

$$w_k = \mathbf{X} \alpha_k \text{ for some } \alpha_k \in \mathbb{R}^n$$

$$\begin{bmatrix} 1 & 1 & 1 \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} \alpha_{k1} \\ \alpha_{k2} \\ \vdots \\ \alpha_{kn} \end{bmatrix}$$

dim mat = $d \times n$

$$(\mathbf{X}^T \mathbf{X}) \alpha_k = (n \lambda_k) \alpha_k \quad \text{--- (1)}$$

$$1 = \alpha_k^T (\mathbf{X}^T \mathbf{X}) \alpha_k \quad \text{--- (2)}$$

Any vector that satisfies (1) & (2) is the α_k

↳ Algorithm

Input: $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \quad \mathbf{x}_i \in \mathbb{R}^d \quad [d \gg n]$

Steps:

① Compute $\mathbf{K} = \mathbf{X}^T \mathbf{X} \quad \mathbf{K} \in \mathbb{R}^{n \times n}$

② Compute eigendecomposition of \mathbf{K}

↳ eigenvectors $\{\beta_1, \dots, \beta_n\}$ corresponding to eigenvalues $\{n\lambda_1, \dots, n\lambda_n\}$ $\longrightarrow O(n^3)$

③ Set $\alpha_k = \frac{\beta_k}{\sqrt{n \lambda_k}} \quad \forall k$

④ $w_k = \mathbf{X} \alpha_k \quad \forall k$

→ Non-linear relationships

↳ Transform features from low dimension \mathbb{R}^d to high dimension \mathbb{R}^D

$$x \in \mathbb{R}^d \xrightarrow{\phi} \phi(x) \in \mathbb{R}^D$$

$\rightarrow \phi(x)^T \phi(x)$ for eigen decomposition

↳ Issue: there might be too many transformations if the power of features is high.

In general: d features, p^{th} power $\rightarrow \sum_i C_i = O(d^p)$

→ Kernel Functions

$$\begin{array}{ccc} \text{↳ } \phi(x_i) & \phi(x_j) & \rightarrow \square \rightarrow \phi(x_i)^T \phi(x_j) = K_{ij} \\ \uparrow & \uparrow & \nearrow ? \\ x_i & x_j & \end{array}$$

↳ Example:

$$x = [f_1, f_2] \quad x' = [g_1, g_2]$$

$$\text{Consider the function: } (x^T x + 1)^2 = (f_1 g_1 + f_2 g_2 + 1)^2$$

$$= [f_1^2 \ f_2^2 \ 1 \ \sqrt{2} f_1 f_2 \ \sqrt{2} f_1 \ \sqrt{2} f_2] \begin{bmatrix} g_1^2 \\ g_2^2 \\ 1 \\ \sqrt{2} g_1 g_2 \\ \sqrt{2} g_1 \\ \sqrt{2} g_2 \end{bmatrix}$$

$$= \phi(x)^T \phi(x')$$

$$\text{where, } \phi\left(\begin{bmatrix} a \\ b \end{bmatrix}\right) = \begin{bmatrix} a^2 \\ b^2 \\ 1 \\ \sqrt{2}ab \\ \sqrt{2}a \\ \sqrt{2}b \end{bmatrix}$$

Note: $(x^T x + 1)^2$ computes the dot-product in a "transformed space".

$$\text{↳ } k(x, x') = (x^T x + 1)^p \quad \text{for some } p \geq 1 \quad x \in \mathbb{R}^d$$

\rightarrow can be shown to be a valid function

$$\text{i.e., } \exists \phi: \mathbb{R}^d \rightarrow \mathbb{R}^D \text{ s.t. } k(x, x') = \phi(x)^T \phi(x')$$

$$\text{↳ } k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \quad \text{for some } \sigma > 0$$

↑ Radial Basis Function

\rightarrow can be shown to be valid

\rightarrow ϕ in this case maps x to an infinite dimensional space.

\rightarrow think of mapping a datapoint to a "function"

and dot-products b/w functions become integrals.

↳ Kernel Functions: Any Function $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ which is a "valid" map

↳ Given a function, how to determine if it's a valid kernel?

Method ①: Exhibit a map ϕ explicitly. Hard sometimes.

Method ②: Mercer's Theorem

↳ $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a valid kernel iff:

① k is symmetric i.e., $k(x, x') = k(x', x)$

② For any dataset $\{x_1, \dots, x_n\}$, the matrix $K \in \mathbb{R}^{n \times n}$ where

$K_{ij} = k(x_i, x_j)$ is Positive Semi-Definite

↳ All eigenvalues ≥ 0

→ Kernel PCA

↳ Algorithm:

Input: $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$, kernel $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Steps:

① Compute $K \in \mathbb{R}^{n \times n}$, where $K_{ij} = K(x_i, x_j) \quad \forall i, j$

② Compute β_1, \dots, β_k eigenvectors and $n\lambda_1 \geq \dots \geq n\lambda_k$ eigenvalues of K
and normalise to get $\alpha_k = \frac{\beta_k}{\sqrt{n\lambda_k}}$

③ $w_k = \phi(x) K_k \rightarrow$ defeats the purpose of kernel because it needs $\phi(x)$

↳ shouldn't compute $w_k \Rightarrow$ we cannot reconstruct eigenvectors of co-variance mat.

↳ but we can compute the "compressed" representation

$$\phi(x_i)^T w_k = \phi(x_i)^T \left(\sum_{j=1}^n \alpha_{kj} K_{ij} \right)$$

$$= \sum_{j=1}^n \alpha_{kj} \phi(x_i)^T \phi(x_j) = \sum_{j=1}^n \alpha_{kj} K_{ij}$$

modified ③ Compute $\sum_{j=1}^n \alpha_{kj} K_{ij} \quad \forall k$

$$x_i \rightarrow \left[\sum_{j=1}^n \alpha_{1j} K_{ij}, \sum_{j=1}^n \alpha_{2j} K_{ij}, \dots, \sum_{j=1}^n \alpha_{kj} K_{ij} \right]$$



WEEK 3

→ Clustering

↳ Goal: $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$

→ Partition data into k different clusters

↳ datapoints: x_1, x_2, \dots, x_n

cluster indicator: z_1, z_2, \dots, z_n where $z_i \in \{1, \dots, k\}$

↳ Performance measurement of clustering:

$$F(z_1, \dots, z_n) = \sum_{i=1}^n \|x_i - \mu_{z_i}\|_2^2, \text{ where } \mu_{z_i} = \text{mean/average of } z_i \text{ cluster}$$

$$\mu_k = \frac{\sum_{i=1}^n x_i \mathbb{1}(z_i=k)}{\sum_{i=1}^n \mathbb{1}(z_i=k)} \quad \text{where } \mathbb{1}(u) = \begin{cases} 1 & \text{if true} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Goal: } \min_{\{z_1, \dots, z_n\}} \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2$$

→ Lloyd's Algorithm / K-means clustering

↳ Initialization

$$z_1^*, z_2^*, \dots, z_n^* \stackrel{\text{iteration } k}{\in} \{1, \dots, k\}$$

→ Algorithm runs until convergence

① Compute means

$$\hookrightarrow \forall k \quad \mu_k^* = \frac{\sum_{i=1}^n x_i \mathbb{1}(z_i^* = k)}{\sum_{i=1}^n \mathbb{1}(z_i^* = k)}$$

② Reassignment

$$\hookrightarrow \forall i \quad z_i^{(k+1)} = \underset{k}{\operatorname{argmin}} \|x_i - \mu_k^*\|^2$$

↳ Convergence

→ Let $x_1, x_2, \dots, x_n \in \mathbb{R}^d$

$$v^* = \underset{v \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^n \|x_i - v\|^2 \quad \rightarrow \quad v^* = \frac{1}{n} \sum_{i=1}^n x_i$$

→ Because every reassignment reduces the aggregate distance from the mean of each cluster and its datapoints, every iteration reduces the objective function.

Because there are a finite number of clusters, the algorithm has to converge.

↳ Nature of clusters

→ Let's say Lloyd's algorithm clusters data points into 2 clusters with mean μ_1 and μ_2

↳ For points in cluster 1 : $\|x - \mu_1\|^2 \leq \|x - \mu_2\|^2$

$$\begin{aligned} &= \|x\|^2 + \|\mu_1\|^2 - 2x^\top \mu_1 \leq \|x\|^2 + \|\mu_2\|^2 - 2x^\top \mu_2 \\ &= x^\top (\mu_2 - \mu_1) \leq \frac{\|\mu_1\|^2 - \|\mu_2\|^2}{2} \end{aligned}$$

→ Cluster regions are intersection of half spaces.

↳ Voronoi Regions

↳ Initialization

→ One method is to pick k-means uniformly at random and then run the algorithm several times.

→ K-means++ algorithm

① choose first mean μ_1 uniformly at random

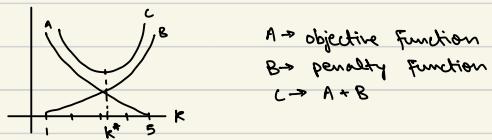
② for $k = 2, \dots, k$: choose μ_k probabilistically proportional to score . $S(x) = \min_{\mu \in \text{means}} \|x - \mu\|^2 + \nu$

↳ Scores are given to each datapoint by looking at how far is the respective datapoint to the closest mean (from all the means that have already been picked).

↳ Choice of K

→ Find K that has the smallest "objective function + penalty (K)"

$$\text{objective function} = \sum_{i=1}^n \|x_i - \mu_k\|^2$$



→ Criterion for picking the penalty function

↳ A.I.C: Akaike Information Criterion $[2K - 2 \log(L(\theta^*))]$

↳ B.I.C: Bayesian Information Criterion $[K \log(n) - 2 \log(L(\theta^*))]$

WEEK 4

Unsupervised learning

- Rep. learning - PCA / kernel
- Clustering - Lloyd's / K-means
- Estimation
 - ↳ Max Likelihood
 - ↳ Bayesian Modeling

→ Maximum Likelihood Estimation

↳ Likelihood function:

$$L(p, \{x_1, x_2, \dots, x_n\}) = P(x_1, x_2, \dots, x_n; p) \quad \begin{matrix} \text{underlying} \\ \text{parameter} \end{matrix}$$

$$= \prod_{i=1}^n P(x_i; p)$$

$$\hat{P}_{ML} = \arg \max_p \prod_{i=1}^n P(x_i; p) = \arg \max_p \log \left(\prod_{i=1}^n P(x_i; p) \right)$$

→ derivative of above function then set it to 0.

↳ assuming Gaussian distribution with μ and σ^2 :

$$L(\mu, \sigma^2, \{x_1, x_2, \dots, x_n\}) = \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \right]$$

$$\hat{\mu}_{ML} = \arg \max_{\mu} \sum_{i=1}^n -(x_i - \mu)^2$$

→ Bayesian Estimation

↳ Think of the parameter to estimate as a random variable.

↳ Munch → codified using a probability distribution over $\theta \rightarrow P(\theta) \leftarrow$ Prior Dist.

$$\text{Update} \rightarrow \cdots \cdots \cdots \cdots \cdots \cdots \rightarrow P(\theta | \text{Data}) \leftarrow \text{Posterior Dist.}$$

↳ Bayes' Theorem

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

A → Parameter θ
B → Data $\{x_1, x_2, \dots, x_m\}$

$$\underbrace{P(\theta | \{x_1, x_2, \dots, x_n\})}_{\text{Posterior}} = \frac{\text{Likelihood}}{\underbrace{P(\{x_1, \dots, x_n\} | \theta)}_{\text{Doesn't depend on } \theta}} \cdot \underbrace{P(\theta)}_{\text{Prior}}$$

↳ Beta prior

$$f(p; \alpha, \beta) = (p)^{\alpha-1} (1-p)^{\beta-1}$$

$$\hookrightarrow P(\theta | \text{Data}) \propto P(\text{Data} | \theta) \cdot P(\theta)$$

→ Gaussian Mixture Models

6 steps to model a generative distribution with multiple gaussians:

① Pick which mixture a data point comes from

② Generate a data point from that mixture

\rightarrow Step 1 : Generate a mixture component among $\{1, \dots, k\}$

$z_i \in \{1, \dots, k\} \rightarrow z_i$ indicates which mixture
with datapoint comes from

$P(Z_i = l) = \pi_l$, where $\sum_{i=1}^k \pi_i = 1$ and $0 \leq \pi_i \leq 1 \forall i$

\rightarrow Step 2 : Generate $x_i \sim N(\mu_{zi}, \sigma_{zi}^2)$

↳ Parameters to be estimated:

$$\rightarrow \pi = [\pi_1, \pi_2, \dots, \pi_k]$$

$$\rightarrow \forall k (\mu_k, \sigma_k^2)$$

$$\text{Total} = 2k + k - 1 = 3k - 1$$

→ Maximum Likelihood of GMM

$$\hookrightarrow L\left(\frac{\mu_1, \dots, \mu_n}{\sigma_1^2, \dots, \sigma_n^2}, \frac{x_1, \dots, x_n}{\pi_1, \dots, \pi_n}\right) = \prod_{i=1}^n f_{\text{mix}}\left(\frac{x_i; \mu_i, \sigma_i^2}{\pi_i}\right)$$

$$= \prod_{i=1}^n \left[\sum_{k=1}^K \pi_k \cdot f(x_i; \mu_k, \sigma_k^2) \right]$$

↪ Not possible to solve $\log L(\theta)$ analytically.

↪ For convex functions: $f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b)$, where $\lambda \in [0,1]$

For concave functions: $f(\lambda a + (1-\lambda)b) \geq \lambda f(a) + (1-\lambda)f(b)$

Jensen's Inequality: $f\left(\sum_{k=1}^K \lambda_k a_k\right) \geq \sum_{k=1}^K \lambda_k f(a_k)$

↪ Log of GMM

$$\rightarrow \log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k \exp\left(\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}\right) \cdot \frac{1}{\sqrt{2\pi}\sigma_k} \right)$$

Log is a concave function

↪ Introduce for every datapoint i , the parameters

$$\{\lambda_1^i, \lambda_2^i, \dots, \lambda_K^i\} \text{ s.t. } \sum_{k=1}^K \lambda_k^i = 1, 0 \leq \lambda_k^i \leq 1 \forall i, k$$

$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \lambda_k^i \left(\frac{\pi_k \exp\left(\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}\right)}{\lambda_k^i} \right) \right)$$

by Jensen's Inequality: $\log L(\theta) \geq \text{modified-log } L(\theta, \lambda)$

$$\underbrace{\text{modified-log } L(\theta, \lambda)}_{\text{gives a lower bound for}} = \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left(\frac{\pi_k \exp\left(\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}\right)}{\lambda_k^i} \right)$$

log of sum is greater than sum of log because log is a concave function.

↪ Fixing λ will allow us to maximise

w.r.t θ and vice versa.

↪ Fix λ and maximise over θ

$$\max_{\theta} \text{modified-log } L(\theta, \lambda)$$

$$= \max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \left[\lambda_k^i \log \pi_k - \lambda_k^i \frac{(x_i - \mu_k)^2}{2\sigma_k^2} - \lambda_k^i \log \sqrt{2\pi} \sigma_k \right]$$

↪ Take derivative w.r.t. μ, σ to get:

$$\hat{\mu}_k^{\text{MML}} = \frac{\sum_{i=1}^n \lambda_k^i x_i}{\sum_{i=1}^n \lambda_k^i}; \hat{\sigma}_k^{\text{MML}} = \frac{\sum_{i=1}^n \lambda_k^i (x_i - \hat{\mu}_k^{\text{MML}})^2}{\sum_{i=1}^n \lambda_k^i}; \hat{\pi}_k^{\text{MML}} = \frac{\sum_{i=1}^n \lambda_k^i}{n}$$

↪ Can think of λ_k^i as the probability i -th datapoint belongs to the k -th "cluster". Then, $\hat{\mu}_k^{\text{MML}}$ and $\hat{\sigma}_k^{\text{MML}}$ are just sample mean and variance for k -th "cluster"

↪ can find π_k & λ using Lagrangian multiplier method.

↪ Fix θ and maximise over λ

$$\hat{\lambda}_k^{\text{MML}} = \frac{\left(\frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}\right) \right) \cdot \pi_k}{\sum_{k=1}^K \left(\frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}\right) \cdot \pi_k \right)}$$

→ EM Algorithm

① Initialise $\theta^0 = \begin{Bmatrix} \mu_1^0, \dots, \mu_n^0 \\ \sigma_1^0, \dots, \sigma_n^0 \\ \pi_1^0, \dots, \pi_n^0 \end{Bmatrix}$

② Until convergence ($\|\theta^{t+1} - \theta^t\| \leq \epsilon$) tolerance parameter

$$\lambda^{t+1} = \arg \max_{\lambda} \text{modified log L}(\theta^t, \lambda) \quad \longrightarrow \text{Expectation step}$$

$$\theta^{t+1} = \arg \max_{\theta} \text{modified log L}(\theta, \lambda^{t+1}) \quad \longrightarrow \text{Maximisation step}$$

→ EM produces "soft clustering", whereas Lloyd's algorithm produces hard clustering.

→ EM takes variances into account, which becomes covariances in higher dimensions.



WEEK 5

→ Supervised learning

- ↳ Input : $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$ ← Features / Attributes
- $\{y_1, \dots, y_n\}$ ← Labels
 - ↳ $\{+1, -1\}, \{0, 1\}$ → Binary classification
 - ↳ $\{0, 1, \dots, k\}$ → Multiclass classification
 - ↳ $\in \mathbb{R}$ → Regression

→ Regression

- ↳ Input / Training data : $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d ; \{y_1, \dots, y_n\} \quad y_i \in \mathbb{R}$
- Goal. learn $h: \mathbb{R}^d \rightarrow \mathbb{R}$
- Measure 'goodness' of a function h : $\text{error}(h) = \sum_{i=1}^n (h(x_i) - y_i)^2$

- ↳ $\text{error}(h)$ can be as small as 0. The h that achieves $\text{error}(h)=0$ is
- st. $h(x_i) = y_i \quad \forall i$

- ↳ Assuming h is linear : $h_w: \mathbb{R}^d \rightarrow \mathbb{R}$ st. $h_w(x) = w^T x$

Goal:

$$\min_w \sum_{i=1}^n (h_w(x_i) - y_i)^2 \rightarrow \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\begin{bmatrix} X^T \\ \vdots \\ -x_1 \\ -x_2 \\ \vdots \\ -x_n \\ \hline \end{bmatrix}_{n \times d} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}_{d \times 1} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad = \min_{w \in \mathbb{R}^d} \|X^T w - y\|_2^2$$

↪ $\min_{w \in \mathbb{R}^d} (x^T w - y)^T (x^T w - y)$ ← unconstrained optimisation

↪ Solution: set gradient to zero

$$f(w) = (x^T w - y)^T (x^T w - y)$$

$$\nabla f(w) = 2(x x^T) w - 2(x y)$$

$$(x x^T) w^* = x y$$

$$w^* = (x x^T)^{-1} (x y)$$

Depends on a covariance-like matrix, like PCA.

→ $(x x^T)^{-1}$ is a bad matrix. Like a covariance matrix.

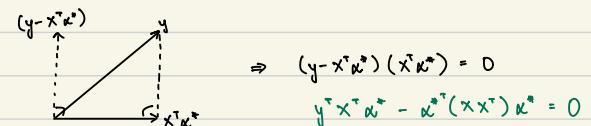
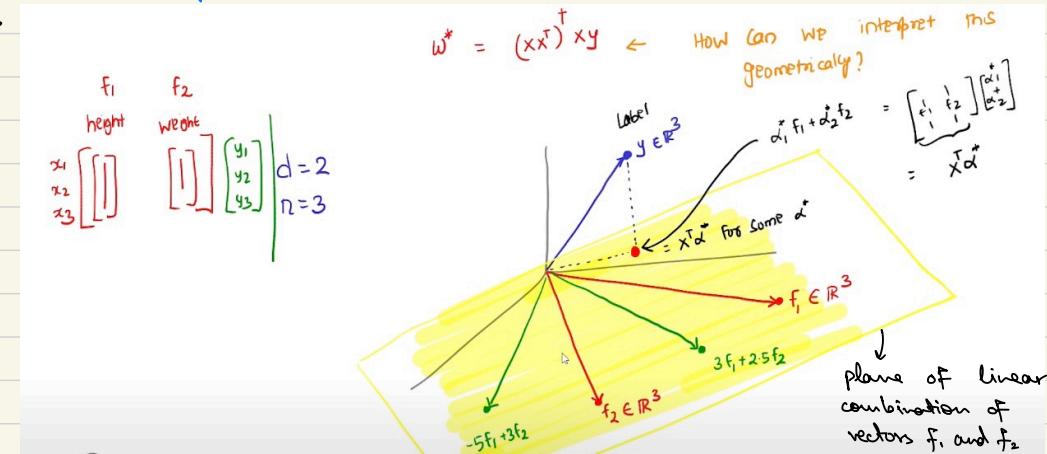
$$(AB)^T = (B^T A^T)$$

$$f(w) = (w^T x - y^T) (x^T w - y)$$

$$= w^T x x^T w - w^T x y - y^T x^T w + y^T y$$

$$= 2(\frac{1}{2} w^T x x^T w - w^T x y)$$

→ Geometric interpretation of Linear Regression



Recall, $w^* = (x x^T)^+ x y$

↪ Substitute $w^* = \alpha^*$: $y^T x^T ((x x^T)^+ x y) - ((x x^T)^+ x y)^T (x x^T) ((x x^T)^+ x y)$

↪ Conclusion: $x^T w^*$ is the projection of the labels onto the subspace spanned by the features.

→ Gradient Descent

↳ $w^* = (x x^T)^{-1} x y$ inverse computation expensive $O(d^3)$

↳ we can apply gradient descent algorithm to obtain w^* .
 ↳ $w^{t+1} = w^t - \eta^t \nabla f(w^t)$

↳ Gradient descent for linear regression:

$$w^{t+1} = w^t - \eta^t [(x x^T) w - x y]$$

↳ If d (dimensions) is large, computing inverse of $(x x^T)^{-1}$ was computationally expensive that's why use gradient descent algorithm.

↳ If n (# of datapoints) is large, computing $x x^T$ is expensive. Use stochastic gradient descent.

↳ Stochastic Gradient Descent

↳ For $t=1, \dots, T$

→ At each step, sample k datapoints uniformly at random from set of all the points.

→ Pretend the sampled datapoints is the entire dataset and take a gradient step w.r.t. the sampled points.

→ After T rounds, we $\underbrace{w_{\text{avg}}^T = \frac{1}{T} \sum_{t=1}^T w^t}_{\text{guaranteed to converge to optima with high probability.}}$ → average of all w

→ Kernel Regression

↳ $w^* = (x x^T)^{-1} x y$

→ w^* must lie in the span of data points.

↳ if it didn't lie on the plane, then it cannot be the optimal solution because the projection of w^* would further minimize the error

→ w^* is some linear combination of the datapoints

$$w^* = x \alpha^* \text{ for some } \alpha^* \in \mathbb{R}^n$$

$$\Rightarrow x \alpha^* = (x x^T)^{-1} x y$$

$$(x x^T) x \alpha^* = x y$$

$$x^T (x x^T) x \alpha^* = x^T x y$$

$$\underbrace{(x^T x)^{-1}}_K x^T x \alpha^* = x^T x y \Rightarrow K^{-1} x^T x \alpha^* = K y \Rightarrow \alpha^* = K^{-1} y \rightarrow \text{ER}'''$$

↳ Prediction:

$$\text{For some } x_{\text{test}} \in \mathbb{R}^d, w^T \phi(x_{\text{test}})$$

$$\text{where, } w^* = \arg \min w \sum_i (\phi(x_i) - y_i)^2$$

$$w^T \phi(x_{\text{test}}) = \left(\sum_{i=1}^n \alpha_i^* \phi(x_i) \right)^T \phi(x_{\text{test}})$$

$$= \sum_{i=1}^n \alpha_i^* K(x_i, x_{\text{test}})$$

how important
is i -th datapoint
towards w^*

→ kernel function is as if it's asking "how similar are these two datapoints?"

→ Probabilistic view of Linear Regression

↳ $x \in \mathbb{R}^d$, $y \in \mathbb{R}$ Dataset: $\{(x_1, y_1), \dots, (x_n, y_n)\}$

→ Assume the label is generated as: $y|x \sim w^T x + e$

↳ can view this as an estimation problem

↳ Maximum Likelihood

$$L(w; x_1, \dots, x_n) = \prod_{i=1}^n \exp\left(\frac{-(w^T x_i - y_i)^2}{2\sigma^2}\right) \cdot \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$\log L(w; x_1, \dots, x_n) = \sum_{i=1}^n \frac{-(w^T x_i - y_i)^2}{2\sigma^2} \cdot \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$\max_w \sum_{i=1}^n -(w^T x_i - y_i)^2 = \min_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\hat{w}_{ML} = w^* = (X^T X)^{-1} X^T y$$

→ Maximum Likelihood estimator assuming zero-mean gaussian noise is same as Linear Regression with squared error!

WEEK 6

→ Goodness of ML estimator for linear regression

↳ $w^* = \hat{w}_{ML}$

$$y|x = w^T x + e, \text{ where } e \sim N(0, \sigma^2)$$

$$\Rightarrow y|x \sim N(w^T x, \sigma^2)$$

↳ $E[\|\hat{w}_{ML} - w\|^2]$

↳ over randomness in y

↳ will converge to the expected value of this random variable

$$E[\|\hat{w}_{ML} - w\|^2] = \sigma^2 \cdot \text{trace}((X^T X)^{-1})$$

average deviation of b/w estimator \hat{w}_{ML} and true w .

variance of the noise that is added to w^* .

→ cannot do much about σ but can reduce $\text{trace}((X^T X)^{-1})$.

$$\Rightarrow \text{trace}((X^T X)^{-1}) \quad \text{tr}(A) = \sum_i \lambda_i$$

Let eigenvalues of $(X^T X)$ be $\{\lambda_1, \dots, \lambda_d\}$

eigenvalues of $(X^T X)^{-1}$ are $\{\lambda_1^{-1}, \dots, \lambda_d^{-1}\}$

→ Mean squared error of \hat{w}_{ML} :

$$E[\|\hat{w}_{ML} - w\|^2] = \sigma^2 \left(\sum_{i=1}^d \lambda_i^{-1} \right), \text{ where } \lambda_i \text{ are the eigenvalues of } (X^T X)$$

→ Cross-validation

↳ Consider the estimator: $\hat{w}_{\text{new}} = \frac{(X^T X + \lambda I)^{-1} X^T Y}{\text{rank}}$

For some matrix A, let eigenvalues: $\{\lambda_1, \dots, \lambda_d\}$

$$AV_i = \lambda_i V_i$$

What are eigenvalues of $(A + \lambda I)$

$$(A + \lambda I)V_i = AV_i + \lambda V_i$$

$$= \lambda_i V_i + \lambda V_i$$

$$= (\lambda_i + \lambda) V_i$$

$$\text{trace}((X^T X + \lambda I)^{-1}) = \sum_{i=1}^d \frac{1}{\lambda_i + \lambda}$$

→ Existence theorem: $\exists \lambda \in \mathbb{R}$ s.t. $\hat{w}_{\text{new}} = (X^T X + \lambda I)^{-1} X^T Y$ has lesser mean squared error than \hat{w}_{ML}

↳ Find λ by cross-validation

Train set	Validation set
80%	20%

↳ Find \hat{w}_{new} on training set.

↳ Train on the training set and check for error on validation set.

↳ Pick λ that gives the least error.

↳ K-Fold cross-validation

F_1	F_2	F_K
-------	-------	-----	-----	-------

↳ Train on folds $\{F_1, \dots, F_{i-1}, F_{i+1}, \dots, F_K\}$

Then validate on F_i

↳ Pick λ that gives the least average error.

→ Bayesian Modeling for Linear Regression

↳ Need a prior on w .

Likelihood: $y|X \sim N(w^T X, 1)$

A choice for prior: $w \sim N(0, \underset{\in \mathbb{R}^d}{\underbrace{\gamma^2 I}})$ covariance matrix $\mathbb{R}^{d \times d}$ $\begin{bmatrix} \gamma^2 & \gamma^2 & \dots \\ \gamma^2 & \gamma^2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$

Posterior:

$$\begin{aligned} P(w | \{(x_i, y_i), \dots, (x_n, y_n)\}) &\propto \underbrace{P(\{(x_i, y_i), \dots, (x_n, y_n)\} | w)}_{\text{Likelihood}} \cdot \underbrace{P(w)}_{\text{Prior}} \\ &\propto \left[\prod_{i=1}^n \exp\left(-\frac{(y_i - w^T x_i)^2}{2}\right) \right] \cdot \underbrace{\left[\prod_{i=1}^n \exp\left(-\frac{(w_i - 0)^2}{2\gamma^2}\right) \right]}_{\exp\left(-\frac{\|w\|^2}{2\gamma^2}\right)} \end{aligned}$$

Max A Priori estimate:

$$\hat{w}_{\text{MAP}} = \arg \max_w \log \left(\prod_{i=1}^n \exp\left(-\frac{(y_i - w^T x_i)^2}{2}\right) \right) = \arg \max_w \sum_{i=1}^n -\frac{(y_i - w^T x_i)^2}{2} - \frac{\|w\|^2}{2\gamma^2}$$

$$\hat{w}_{\text{MAP}} = \arg \min_w \frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \frac{\|w\|^2}{2\gamma^2}$$

→ Take gradient set to 0.

$$\nabla f(w) = (X^T X)w - XY + \frac{w}{\gamma^2}$$

→ Set $\nabla f(w) = 0$

$$\hat{w}_{\text{MAP}} = (X^T X + \frac{1}{\gamma^2} I)^{-1} X^T Y$$

↳ MAP estimation for linear regression with a gaussian prior for w is equivalent to the \hat{w}_{new} we used in cross-validation.

→ Ridge Regression

$$\hookrightarrow \hat{W}_{ML} = \underset{W}{\operatorname{argmin}} \sum_{i=1}^n (W^T x_i - y_i)^2$$

$$\hat{W}_R = \underset{W}{\operatorname{argmin}} \sum_{i=1}^n (W^T x_i - y_i)^2 + \lambda \|W\|^2 \quad \xrightarrow{\text{Ridge Regression}}$$

Loss Regulariser

↳ The goal of the regulariser is to push W towards zero. This helps in picking the W which has the least exposures to redundant features.

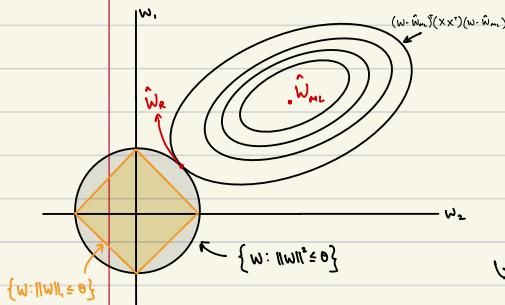
↳ " λ " in the regulariser is $(\frac{1}{\gamma^2})$ where γ^2 is the variance. In the scenario where most features are redundant, a lot of W 's will be zero, meaning low variance. If γ^2 is small, then λ is going to be large, which means it's going to pay a large penalty for increasing the length of W .

$$\hookrightarrow \min_{W \in \mathbb{R}^d} \sum_{i=1}^n (W^T x_i - y_i)^2 + \lambda \|W\|^2 \quad \text{--- (A)}$$

is equivalent to $\min_{W \in \mathbb{R}^d} \sum_{i=1}^n (W^T x_i - y_i)^2 \text{ s.t. } \|W\|^2 \leq \theta$ (B)
depends on θ

→ For every choice of $\lambda > 0$, $\exists \theta$ s.t. the optimal solutions of problems (A) and (B) coincide

Assuming $W = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$
 $\|W\|^2 \leq \theta \Rightarrow w_1^2 + w_2^2 \leq \theta$



→ Loss of linear regression at \hat{W}_{ML} : $f(W_{ML}) = \sum_{i=1}^n (\hat{W}_{ML}^T x_i - y_i)^2$

Consider the set of all W st. $f(W) = f(\hat{W}_{ML}) + c$; where $c > 0$

$$\hookrightarrow S_c = \{W : f(W) = f(\hat{W}_{ML}) + c\}$$

$$\text{every } W \in S_c \text{ satisfies } \|X^T W - Y\|^2 = \|X^T \hat{W}_{ML} - Y\|^2 + c$$

IF $X^T X = I$,
 $\|W - \hat{W}_{ML}\|^2 = c \leftarrow \{ \Rightarrow (W - \hat{W}_{ML})^T (X^T)(W - \hat{W}_{ML}) = c \leftarrow$

some constant that depends on c , (X^T) , \hat{W}_{ML} and not on W .

↳ Conclusion:

→ Ridge regression pushes weight values towards zero. But does not necessarily make it zero.

→ Lasso Regression

↳ Alternate way of regularization: use $\|\cdot\|_1$ norm instead of $\|\cdot\|_2$ norm.

$$\|W\|_1 = \sum_{i=1}^n |w_i|$$

$$\hookrightarrow \min_{W \in \mathbb{R}^d} \sum_{i=1}^n (W^T x_i - y_i)^2 + \lambda \|W\|_1 \quad = \quad \min_{W \in \mathbb{R}^d} \sum_{i=1}^n (W^T x_i - y_i)^2 \text{ s.t. } \|W\|_1 \leq \theta$$

↳ increases the chances of picking a W which has weights as zeros instead of really small values.

↳ Points to note:

→ Lasso does not have a closed-form solution.

→ Sub-gradient methods usually used to solve Lasso.

↳ Subgradient

→ A vector $g \in \mathbb{R}^d$ is a sub-gradient of $f: \mathbb{R}^d \rightarrow \mathbb{R}$ at a point $x \in \mathbb{R}^d$ if:

$$\forall z \quad f(z) \geq f(x) + g^T(z-x)$$

→ If the function f to minimize is a convex function, then sub-gradient descent algorithm converges



WEEK 7

→ Introduction to Binary Classification

↳ Input: $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \quad y_i \in \{0, 1\} \cup \{-1, +1\}$

Goal: $h: \mathbb{R}^d \rightarrow \{0, 1\}$

$$\text{Loss: } \text{Loss}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i) \quad \mathbb{1}(z) = \begin{cases} 1 & \text{if true} \\ 0 & \text{otherwise} \end{cases}$$

$$\min_{w \in H_{\text{linear}}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h_w(x_i) \neq y_i) \quad \text{where } H_{\text{linear}} = \{h_w: h_w(x) = \text{sign}(w^T x)\}$$

↳ NP-hard problem in general

↳ Can we use linear regression to solve classification problem?

$$\{(x_1, y_1), \dots, (x_n, y_n)\} \rightarrow \boxed{\text{Lin. reg.}} \rightarrow w \in \mathbb{R}^d \rightarrow h_w: \mathbb{R}^d \rightarrow \{0, 1\}$$

→ Regression is sensitive to the datapoints and not just the "side" on which data lies w.r.t. separator.

→ K-nearest Neighbour

↳ Given a test point $x_{\text{test}} \in \mathbb{R}^d$. Find the closest point x^* to x_{test} in the training set. Predict $y_{\text{test}} = y^*$

↳ Issue: can get affected by outliers

Fix: look at the distance from k different neighbours.

↳ K-NN algorithm

→ Given x_{test} , find the k-closest points in training set - $(x_1^*, x_2^*, \dots, x_k^*)$

→ Predict $y_{\text{test}} = \text{majority}(y_1^*, y_2^*, \dots, y_k^*)$

→ Too few k will result in outliers impacting the algorithm and too high k will make the algorithm very general (high loss)

→ Choosing k: cross-validate → try diff. k and see performance on validation set

→ Issues with K-nn:

① choosing a distance function

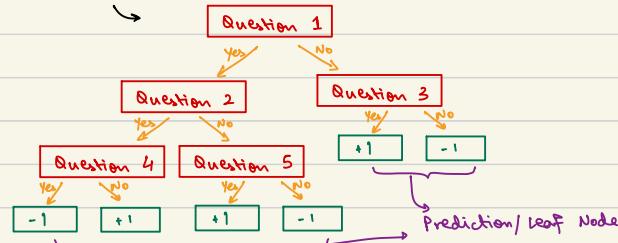
② Prediction is computationally expensive because of sorting

③ No model is learned - cannot throw data after "learning"

→ Introduction to Decision Trees

↳ Input: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ s.t. $x_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$

Output: Decision Tree

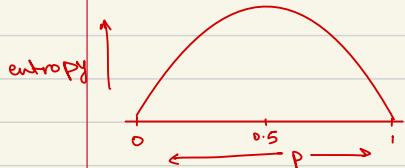


Prediction: Traverse through the tree to reach a leaf node.

Question: (feature, value) pair. Example: height ≤ 180 cm?

↳ How to measure "goodness" of a question?

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

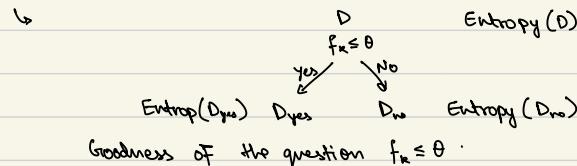


$$\begin{array}{c} f_k \leq \theta \\ \text{yes} \quad \text{no} \end{array}$$

$$D_{\text{yes}} = \{(x_1, y_1), (x_2, y_2), \dots\} \quad D_{\text{no}} = \{(x_3, y_3), (x_4, y_4), \dots\}$$

→ measure of "impurity" for a set of labels $\{y_1, \dots, y_n\}$

$$\hookrightarrow \text{Entropy } (y_1, \dots, y_n) = \text{Entropy } (p) = - (p \log p + (1-p) \log (1-p))$$



Goodness of the question $f_k \leq \theta$.

$$\text{Information Gain}(\text{Feature}, \text{value}) = \text{Entropy}(D) - [\gamma \text{Entropy}(D_{\text{yes}}) + (1-\gamma) \text{Entropy}(D_{\text{no}})]$$

$$\text{where } \gamma = \frac{|D_{\text{yes}}|}{|D|} = \frac{\text{cardinality of } D_{\text{yes}}}{\text{cardinality of } D}$$

↳ measures how much of the impurity is being removed by asking this question.

→ Decision Tree Algorithm

↳ Discretize each feature in [min, max] range

↳ Pick the question that has the highest Information Gain.

↳ Repeat the procedure for D_{yes} , D_{no} .

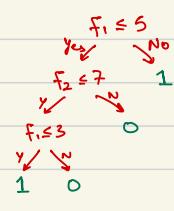
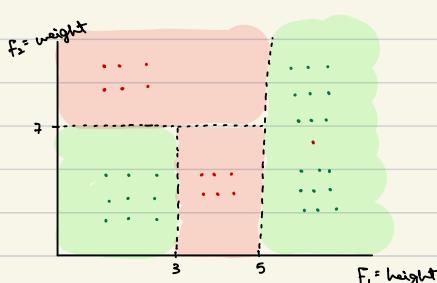
↳ Points:

→ can stop growing a tree if a node becomes "sufficiently" pure.

→ Depth of the tree (# of questions) is a hyperparameter.

→ Alternate measures for "goodness" of a question. E.g. Gini Index

Decision Boundary



→ Generative and Discriminative Models

↳ Types of Modelling

→ Assumption: there is some distribution D over $X \times Y_m$ labels

$$\{x_1, \dots, x_n\}$$

$$\{y_1, \dots, y_n\}$$

$$(x_i, y_i) \sim D$$

↙ features

→ implies there is some probabilistic mechanism that generates the data.

→ The link b/w Train and Test

↳ Classification

↳ ① Generative Model

② Discriminative Model

Generative

→ $P(x, y)$. There is some probability/density for every feature-value pair

example: There is a sentence and the task is to tell whether the sentence is in Bengali or not

→ Generative model is to model how Bengali sentences are generated. Modelling x themselves
Joint Modelling.

Discriminative

→ $P(y|x)$. Probability of a label given the feature.

→ Your goal is only to differentiate if the sentence is of the Bengali class or not.
Don't necessarily need to know how Bengali sentences are generated.

→ If you have reasons to believe how the exact features are generated, then it is a good idea to do generative modelling. If, on the other hand, you have no idea how the features are generated, you're better off with Discriminative.



WEEK 8

→ Generative model-based algorithm

$$\hookrightarrow \text{Data} = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad x_i \in \{0,1\}^d \quad y_i \in \{0,1\}$$

example: spam classification of e-mails

$x_i \in \{0,1\}^d$, where $d = \# \text{ of words in a dictionary}$

e.g. email = "Hello, how are you?"

$$P(x, y) = P(\text{"Hello, how are you"}, \text{spam}) = ?$$

$$P(x, y) = P(x) \cdot P(y|x) = \underbrace{P(y) \cdot P(x|y)}$$

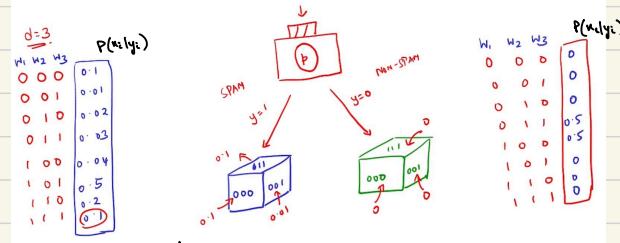
better way of modelling

↪ Generative Story

Steps:

① Decide the label by tossing a coin. $P(y_i=1) = p$

② Decide the features using the label is step 1 by $P(x_i|y_i)$



↪ # of parameters in the model

$$= 1 + [(2^d - 1)] + [(2^d - 1)] \rightarrow P(x_i|y_i)$$

label, whether an email is spam or not. For the blue die (spam world) For the green die (spam world)

$$= 15 \text{ for the example above}$$

→ If we know all the 15 parameters, then we have the ability to generate a new (x, y) pair.

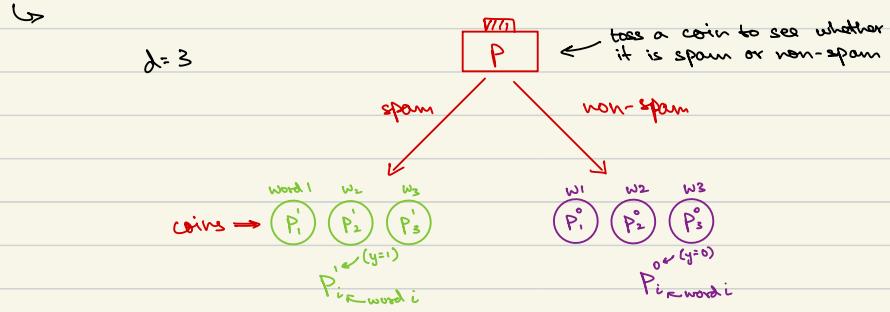
$$= 1 + (2^d - 1) + (2^d - 1)$$

$$= 1 + 2(2^d - 1)$$

$$= 2^{d+1} - 1 \text{ parameters}$$

↪ Issue: Too many parameters!

→ Alternate generative model-based algorithm



$$\begin{aligned} \# \text{ parameters} &= 1 + d + d \\ &= 2d + 1 \leftarrow \text{much better than } (2^{d+1}-1) \end{aligned}$$

→ The difference is that the features in this algorithm are independent to each other. The features are orthogonal.
→ Class Conditional Independence Assumption!

↳ Steps of the algorithm:

$$\textcircled{1} \quad P(y=1) = P$$

$$\textcircled{2} \quad P(n=[f_1, f_2, \dots, f_d] | y) = \prod_{i=1}^d (P_i^y)^{f_i} (1 - P_i^y)^{1-f_i}$$

→ Naive Bayes Algorithm

↳ Parameters to estimate: $P, \{P_i^1, \dots, P_i^d\}, \{P_i^0, \dots, P_i^d\}$

Maximum Likelihood estimates:

$$\hat{P} = \frac{1}{n} \sum_{i=1}^n y_i \rightarrow \text{Fraction of spam emails in the dataset}$$

$$\hat{P}_j^y = \frac{\sum_{i=1}^n \mathbb{1}(f_j^i = 1, y_i = y)}{\sum_{i=1}^n \mathbb{1}(y_i = y)} \quad \left. \begin{array}{l} \text{Fraction of } y\text{-labelled emails} \\ \text{that contain the } j\text{th word} \end{array} \right\} \quad \text{# of emails with label } y$$

↳ Prediction

Given $X^{\text{test}} \in \{0,1\}^d$, what is \hat{y}^{test} ?

$$P(y^{\text{test}}=1 | X^{\text{test}}) > P(y^{\text{test}}=0 | X^{\text{test}}) \Rightarrow \hat{y}^{\text{test}} = 1, 0 \text{ otherwise}$$

↳ How to obtain $P(y|n)$ from $P(y)$ and $P(n|y)$

$$\rightarrow P(y^{\text{test}}=1 | n^{\text{test}}) = \frac{P(X^{\text{test}} | y^{\text{test}}=1) \cdot P(y^{\text{test}}=1)}{P(X^{\text{test}})}$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$P(X^{\text{test}} | y^{\text{test}}=1) \cdot P(y^{\text{test}}=1) = \left[\prod_{i=1}^d (\hat{P}_i^1)^{f_i} (1 - \hat{P}_i^1)^{1-f_i} \right] \cdot \hat{P}$$

compare with

$$P(X^{\text{test}} | y^{\text{test}}=0) \cdot P(y^{\text{test}}=0) = \left[\prod_{i=1}^d (\hat{P}_i^0)^{f_i} (1 - \hat{P}_i^0)^{1-f_i} \right] \cdot (1 - \hat{P})$$

↳ Pitfalls of Naive Bayes

→ If a word does not appear in the train set

but appears in a test datapoint, $\hat{p}_i^i = 0$ and $\hat{p}_i^o = 0$

↳ this will result in y_{test} not being classified.

Occurrence of one word is throwing away all the information from other words.

↳ Possible Fix: Laplace Smoothing

↳ Add two "pseudo" emails with all words

present - one email has label 0 and another 1.

$$x_{\text{pseudo}}^i = [1, 1, 1, \dots, 1] \quad y_{\text{pseudo}}^i = 1$$

$$x_{\text{pseudo}}^o = [1, 1, 1, \dots, 1] \quad y_{\text{pseudo}}^o = 0$$

→ This way every word will have occurred at least once in the test dataset.

↳ Decision Function:

$$\text{Given } x_{\text{test}}, y_{\text{test}} = 1 \text{ if } \frac{P(y_{\text{test}}=1|x_{\text{test}})}{P(y_{\text{test}}=0|x_{\text{test}})} \geq 1$$

$$\Rightarrow \log \left[\frac{P(x_{\text{test}}|y_{\text{test}}=1) \cdot P(y_{\text{test}}=1)}{P(x_{\text{test}}|y_{\text{test}}=0) \cdot P(y_{\text{test}}=0)} \right] \geq 0$$

$$x_{\text{test}} = [f_1, f_2, \dots, f_d]$$

$$\Rightarrow \log \left[\prod_{i=1}^d \frac{(\hat{p}_i^i)^{f_i} (1-\hat{p}_i^i)^{1-f_i} \cdot \hat{p}}{(\hat{p}_i^o)^{f_i} (1-\hat{p}_i^o)^{1-f_i} \cdot (1-\hat{p})} \right] \geq 0$$

$$\Rightarrow \left[\sum_{i=1}^d f_i \log \left(\frac{\hat{p}_i^i}{\hat{p}_i^o} \right) + (1-f_i) \log \left(\frac{1-\hat{p}_i^i}{1-\hat{p}_i^o} \right) + \log \left(\frac{\hat{p}}{1-\hat{p}} \right) \right] \geq 0$$

$$\Rightarrow \underbrace{\sum_{i=1}^d f_i}_{\text{features}} \left(\log \left(\frac{\hat{p}_i^i}{\hat{p}_i^o} \right) \right) + \underbrace{\log \left(\frac{1-\hat{p}_i^i}{1-\hat{p}_i^o} \right)}_{\text{constants}} + \underbrace{\log \left(\frac{\hat{p}}{1-\hat{p}} \right)}_{\text{constants}} \geq 0$$

Decision function is of the form:

$$\text{Predict } y_{\text{test}} = 1 \text{ if } w^T x_{\text{test}} + b \geq 0$$

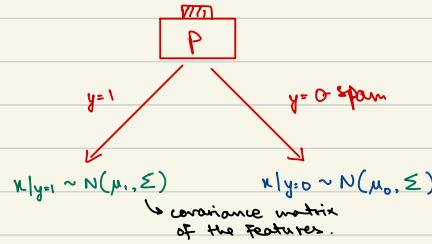
$$\text{where } w_i = \log \left(\frac{\hat{p}_i^i (1-\hat{p}_i^o)}{\hat{p}_i^o (1-\hat{p}_i^i)} \right) \text{ and } b = \log \left(\frac{(1-\hat{p}_i^i)}{(1-\hat{p}_i^o)} \right) + \log \left(\frac{\hat{p}}{1-\hat{p}} \right)$$

Conclusion: decision function of naive bayes is linear!

Gaussian Naive Bayes

→ Data: $\{(x_i, y_i), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \{0, 1\}$

A generative story:



parameters: p, μ_1, μ_0, Σ

Max Likelihood estimate for parameters:

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\hat{\mu}_1 = \frac{\sum_{i=1}^n \mathbb{1}(y_i=1) \cdot x_i}{\sum_{i=1}^n \mathbb{1}(y_i=1)} \quad \leftarrow \begin{array}{l} \text{sample mean of} \\ \text{data points labelled 1} \end{array}$$

$$\hat{\mu}_0 = \frac{\sum_{i=1}^n \mathbb{1}(y_i=0) \cdot x_i}{\sum_{i=1}^n \mathbb{1}(y_i=0)}$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^T$$

→ Prediction

$$P(y_{\text{test}} | x_{\text{test}}) \propto \underbrace{P(x_{\text{test}} | y_{\text{test}})}_{f(x_{\text{test}}; \hat{\mu}_y, \hat{\Sigma})} \cdot \underbrace{P(y_{\text{test}})}_{\hat{p}}$$

$$y_{\text{test}} = 1 \quad f(x_{\text{test}}; \hat{\mu}_1, \hat{\Sigma}) \cdot \hat{p} \geq f(x_{\text{test}}; \hat{\mu}_0, \hat{\Sigma}) \cdot (1 - \hat{p})$$

$$\exp(-(x_{\text{test}} - \hat{\mu}_1)^T \hat{\Sigma}^{-1} (x_{\text{test}} - \hat{\mu}_1)) \cdot \hat{p} \geq \exp(-(x_{\text{test}} - \hat{\mu}_0)^T \hat{\Sigma}^{-1} (x_{\text{test}} - \hat{\mu}_0)) \cdot (1 - \hat{p})$$

↓ simplification by taking a log !

Predict $y_{\text{test}} = 1$ if:

$$\underbrace{(\hat{\mu}_1 - \hat{\mu}_0)^T \hat{\Sigma}^{-1} x_{\text{test}}}_{w^T} + \underbrace{\hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1 - \hat{\mu}_0^T \hat{\Sigma}^{-1} \hat{\mu}_0}_{b} + \log\left(\frac{1 - \hat{p}}{\hat{p}}\right) \geq 0$$

$$w^T x_{\text{test}} + b \geq 0$$

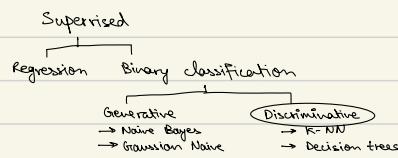
⇒ Decision function is again linear!

↪ assumption of Σ is same in both classes is imp.

→ If the covariance matrix is identity matrix, the contours are circular, not elliptical.



WEEK 9



→ Discriminative Models for classification

↳ How to model $P(y=1|u)$

Simplest assumption:

$$P(y=1|u) = \begin{cases} 1 & \text{if } W^T u \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

↳ Linear Separability Assumption:

$$\exists W \in \mathbb{R}^d \text{ s.t. } \text{sign}(W^T u_i) = y_i \quad \forall i$$

→ Perception Learning Algorithm

↳ Input: $\{(u_1, y_1), \dots, (u_n, y_n)\}$ $u_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$

$$\rightarrow W^0 = 0 \in \mathbb{R}^d \quad [0 \dots 0]$$

Until convergence:

→ Pick (u_i, y_i) pair from dataset

→ If $\text{sign}(W^T u_i) = y_i$, do nothing

Else $W^{t+1} = W^t + \eta_i y_i$ ← Update rule

↳ Update Rule:

Mistake type 1

→ Predicted = 1 ← $W^T u_i \geq 0$

Actual = -1 ← $y_i = -1$

Mistake type 2

→ Predicted = -1 ← $W^T u_i < 0$

Actual = 1 ← $y_i = +1$

$$W^{t+1} = W^t + \eta_i y_i$$

$$(W^{t+1})^T u_i = (W^t + \eta_i y_i)^T u_i$$

$$= W^T u_i + \underbrace{\eta_i y_i \|u_i\|^2}_{\geq 0}$$

Negative

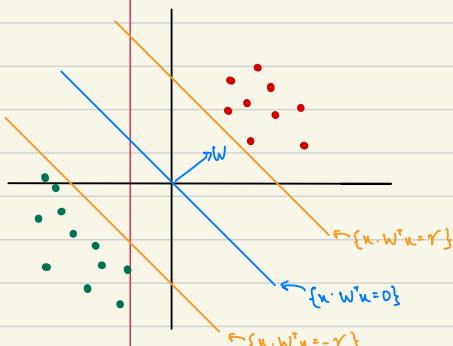
$$(W^{t+1})^T u_i = \underbrace{W^T u_i}_{\leq 0} + \underbrace{\eta_i y_i \|u_i\|^2}_{> 0}$$

- Update rule pushes W in the right direction for u_i
- Fixing it locally might break something else for some other datapoint.

↳ Assumption: Linear separability with γ -margin

→ Dataset $\{(u_1, y_1), \dots, (u_n, y_n)\}$ is linearly separable with γ margin if

$$\exists W^* \in \mathbb{R}^d \text{ s.t. } ((W^*)^T u_i) y_i \geq \gamma \quad \forall i ; \gamma > 0$$



↳ Assumption: Radius Assumption

$$\rightarrow \|x\|_2 \leq R \text{ for some } R > 0$$

↳ all the points fall within a ball of radius R

↳ Assumption: Without loss of generality, assume $\|W^*\| = 1$

→ if you rescale W^* by $\|W^*\|$, γ also gets rescaled.

↳ Proof of Perceptron

→ Analysis of "mistakes" made by perceptron

Update happens only when there is a mistake

Say W^l is the current guess and a mistake happens w.r.t (x, y)

$$\begin{aligned} W^{l+1} &= W^l + x \cdot y \\ \|W^{l+1}\|^2 &= \|W^l + x \cdot y\|^2 \\ &= (W^l + x \cdot y)^T (W^l + x \cdot y) \\ &= \|W^l\|^2 + 2(W^l \cdot x)y + \|x\|^2 \cdot y^2 \\ &\stackrel{\leq 0}{\substack{\text{because} \\ \text{mistake}}} \quad \stackrel{\leq R^2}{\substack{\text{radius} \\ \text{assumption}}} \end{aligned}$$

$$\begin{aligned} \|W^{l+1}\|^2 &\leq \|W^l\|^2 + R^2 \\ \|W^{l+1}\|^2 &\leq (\|W^l\|^2 + R^2) + R^2 \end{aligned}$$

$$\|W^{l+1}\|^2 \leq \|W^l\|^2 + lR^2 \quad \text{--- (1)} \quad \stackrel{\text{# of mistakes}}{\curvearrowleft}$$

$$\begin{aligned} \rightarrow (W^{l+1})^T W^* &= (W^l + x \cdot y)^T W^* \\ &= W^l W^* + (W^l \cdot x) y \end{aligned}$$

$$(W^{l+1})^T W^* \geq W^l W^* + \gamma$$

$$(W^{l+1})^T W^* \geq W^l W^* + l\gamma \quad \text{--- (2)}$$

$$\rightarrow l\gamma \leq (W^{l+1})^T W^*$$

$$\Rightarrow (l\gamma)^2 \leq ((W^{l+1})^T W^*)^2 \leq \|W^{l+1}\|^2 \cdot \|W^*\|^2 \quad \text{--- (3)} \quad \text{From Cauchy-Schwarz}$$

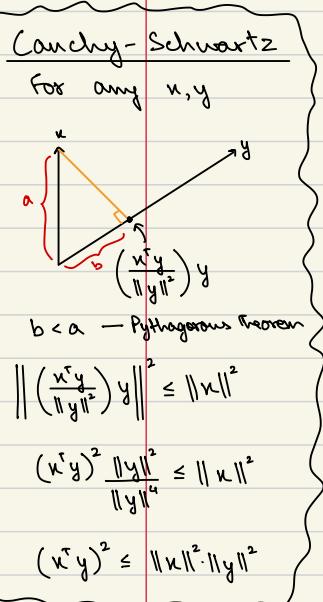
$$\Rightarrow \|W^{l+1}\|^2 \geq l^2 \gamma^2 \quad \text{--- (3)}$$

$$\Rightarrow l^2 \gamma^2 \leq \|W^{l+1}\|^2 \leq lR^2$$

$$\Rightarrow l^2 \gamma^2 \leq lR^2$$

$$\Rightarrow l \leq \frac{R^2}{\gamma^2} \quad \begin{matrix} \leftarrow \text{Radius} \\ \leftarrow \text{Margin} \end{matrix} \quad \text{--- Radius Margin Bound}$$

→ # of mistakes is bounded because $\gamma > 0$



→ Modelling Class Probabilities

↳ For perceptron: $P(y=1|x) = 1 \text{ if } W^T x \geq 0, 0 \text{ otherwise}$

↳ simple linear model

$$z = W^T x \rightarrow \text{score}$$

→ We want such that larger the value of z , score, higher the probability of being classified as +1.

$$\rightarrow g: \mathbb{R} \rightarrow [0, 1] \quad g(z) = 0.5 \text{ if } z=0$$

$$g(z) \rightarrow 1 \text{ as } z \rightarrow \infty$$

$$g(z) \rightarrow 0 \text{ as } z \rightarrow -\infty$$

Sigmoid / Logistic Function $\cdot \quad g(z) = \frac{1}{1 + e^{-z}}$

→ Logistic Regression

↳ Model: $P(y=1|x) = \frac{1}{1 + e^{-W^T x}} = g(W^T x)$ Data: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ $x_i \in \mathbb{R}^d$
 $y_i \in \{0, 1\}$

Find W using Max. Likelihood

$$L(W; \text{Data}) = \prod_{i=1}^n (g(W^T x_i))^{y_i} \cdot (1 - g(W^T x_i))^{(1-y_i)}$$

$$\log L(W; \text{Data}) = \sum_{i=1}^n y_i \log(g(W^T x_i)) + (1-y_i) \log(1 - g(W^T x_i))$$

$$= \sum_{i=1}^n y_i \log\left(\frac{1}{1 + e^{-W^T x_i}}\right) + (1-y_i) \log\left(\frac{e^{-W^T x_i}}{1 + e^{-W^T x_i}}\right)$$

$$= y_i \log(1) - y_i \log(1 + e^{-W^T x_i}) + \log(e^{-W^T x_i}) - \log(1 + e^{-W^T x_i}) - y_i \log(e^{-W^T x_i}) + y_i \log(1 + e^{-W^T x_i})$$

$$= -W^T x_i + y_i (W^T x_i) - \log(1 + e^{-W^T x_i})$$

$$= \sum_{i=1}^n (1-y_i)(-W^T x_i) - \log(1 + e^{-W^T x_i})$$

Goal: $\max_w \log L(W; \text{Data}) \rightarrow$ no closed-form solution

→ Use gradient ascent

$$\nabla \log L(W) = \sum_{i=1}^n (1-y_i)(-x_i) - \left(\frac{e^{-W^T x_i}}{1 + e^{-W^T x_i}}\right)(-x_i)$$

$$= \sum_{i=1}^n y_i x_i - x_i \left(\frac{1}{1 + e^{-W^T x_i}}\right)$$

$$= \sum_{i=1}^n x_i \left(y_i - \frac{1}{1 + e^{-W^T x_i}}\right)$$

$$W_{t+1} = W_t + \eta_x \nabla \log(L)$$

$$= W_t + \eta_x \left[\sum_{i=1}^n x_i \left(y_i - \frac{1}{1 + e^{-W^T x_i}} \right) \right]$$

$$\underbrace{\left(y_i - \frac{1}{1 + e^{-W^T x_i}} \right)}_{\{0, 1\}} \rightarrow g(W^T x_i)$$

$\nabla \log L$ is some linear combination of the dataset.

↳ Prediction: $\hat{y}_{\text{test}} = \text{sign}(W^T x_{\text{test}})$



WEEK 10

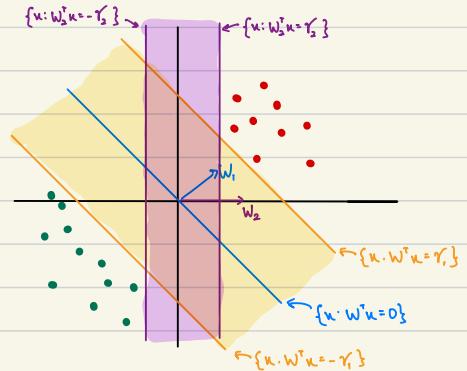
→ Perceptions and Margin

↳ Quality of Final Solution

→ W and W_2 are both acceptable separators.

→ # of mistakes depend on the best possible W^* margin

→ The W that perception outputs need not be the one with the largest margin γ



↳ Maximise Margin

Goal: formulate margin maximisation

$$\max_{W, \gamma} \gamma \text{ s.t. } (W^T x_i) y_i \geq \gamma + i$$

→ This is not useful because W can grow arbitrarily

$$\max_{W, \gamma} \gamma \text{ s.t. } (W^T x_i) y_i \geq \gamma + i \text{ and } \|W\|^2 = 1$$

→ If you fix $\|W\|^2$, γ varies and if you fix γ , $\|W\|^2$ can be scaled.

$$\max_{W} \text{width}(W) \text{ s.t. } (W^T x_i) y_i \geq 1 + i$$

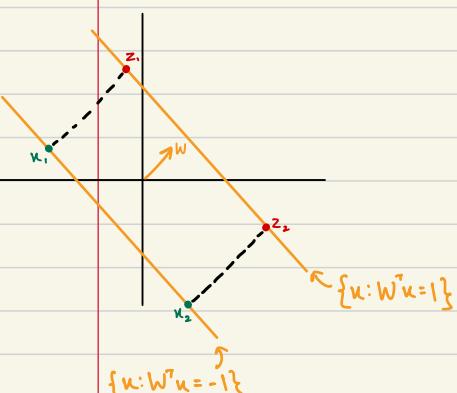
→ Define width(W).

$$\min_z \frac{1}{2} \|x - z\|^2 \text{ s.t. } W^T z = 1, \text{ where } W^T x = -1$$

→ solution to the optimisation: $\text{width}(W) = \frac{2}{\|W\|^2}$

$$\rightarrow \text{Final formulation: } \max_{W} \frac{2}{\|W\|^2} \text{ s.t. } (W^T x_i) y_i \geq 1 + i$$

$$= \min_{W} \frac{\|W\|^2}{2} \text{ s.t. } (W^T x_i) y_i \geq 1 + i$$



→ Constrained Optimisation

$$\hookrightarrow \min_w f(w) \quad \text{s.t. } g(w) \leq 0$$

$$L(w, \alpha) = f(w) + \alpha g(w)$$

Fix any w . Consider $\max_{\alpha \geq 0} L(w, \alpha) = \max_{\alpha \geq 0} f(w) + \alpha g(w)$

$$= \begin{cases} \infty & \text{if } g(w) > 0 \\ f(w) & \text{if } g(w) \leq 0 \end{cases}$$

$$\rightarrow \min_w \left[\max_{\alpha \geq 0} f(w) + \alpha g(w) \right]$$

If $f(w)$ and $g(w)$ are convex functions,

then

$$\min_w \left[\max_{\alpha \geq 0} f(w) + \alpha g(w) \right] = \max_{\alpha \geq 0} \left[\min_w f(w) + \alpha g(w) \right]$$

↳ for multiple constraints

$$\min_w f(w) \quad \text{s.t. } g_i(w) \leq 0 \quad \forall i=1, \dots, k$$

$$= \min_w \left[\max_{\alpha_i \geq 0, i} f(w) + \sum_{i=1}^k \alpha_i g_i(w) \right]$$

→ Perceptrons and Margins

$$\hookrightarrow \min_w \frac{\|w\|^2}{2} \quad \text{s.t. } (w^T x_i) y_i \geq 1 \quad \forall i=1, \dots, n$$

$\underbrace{\quad}_{\text{quadratic in } w} \quad \underbrace{\quad}_{\text{linear in } w}$

$$f(w) = \frac{\|w\|^2}{2} \quad g_i(w) = 1 - (w^T x_i) y_i$$

$$L(w, \alpha) = \frac{\|w\|^2}{2} + \sum_{i=1}^n \alpha_i (1 - (w^T x_i) y_i)$$

$$\min_w \max_{\alpha \geq 0} L(w, \alpha) \equiv \max_{\alpha \geq 0} \min_w L(w, \alpha)$$

① ②

① will fix the w and maximise over $\alpha \geq 0$.

② will fix the $\alpha \geq 0$ and minimize over w ; makes the minimisation of w an unconstrained minimisation (solved by setting the gradient to 0).

Example using ②:

$$\text{Fix } \alpha = \begin{bmatrix} 2 \\ 3 \end{bmatrix}; \quad \min_w \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - (w^T x_i) y_i)$$

$$X = \begin{bmatrix} 1 & | & x_1 & | & x_2 & | & \dots & | & x_n \end{bmatrix} \quad Y = \begin{bmatrix} y_1 & 0 \\ 0 & y_2 \\ \vdots & \vdots \end{bmatrix} \quad \alpha = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad W^* = X^T Y \alpha$$

$$0 = w + \sum_{i=1}^n \alpha_i (-x_i y_i)$$

$$W^* = \sum_{i=1}^n \alpha_i (x_i y_i)$$

≥ 0 α_i $\{ \pm 1 \}$

Substitute w into objective:

$$x^T 1 - \frac{1}{2} (x^T \alpha)^2 \quad 1 = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\hookrightarrow \text{Dual Problem: } \max_{\alpha \geq 0} \alpha^T 1 - \frac{1}{2} \alpha^T y^T X^T X y \alpha$$

\hookrightarrow Advantages:

\rightarrow Dual variable dimension \mathbb{R}^n , while primal problem dimension in \mathbb{R}^d .

\rightarrow Dual constraints are "easier"

\rightarrow The objective of dual depends on $X^T X$ and so can be "kernelized"

\rightarrow Support Vector Machine

\hookrightarrow Revisiting Lagrangian:

$$\min_w \left[\max_{\alpha \geq 0} f(w) + \alpha g(w) \right] \stackrel{\text{Primal}}{=} \max_{\alpha \geq 0} \left[\min_w f(w) + \alpha g(w) \right] \stackrel{\text{Dual}}{=}$$

\rightarrow say w^* is the primal solution and α^* is the dual solution

$$\underbrace{\max_{\alpha \geq 0} f(w^*) + \alpha g(w^*)}_{f(w^*)} = \min_w f(w) + \alpha^* g(w)$$

$$\Rightarrow f(w^*) \leq f(w^*) + \alpha^* g(w^*)$$

$$\Rightarrow \alpha^* g(w^*) \geq 0 \quad \text{--- (1)}$$

\rightarrow but we know $\alpha^* \geq 0$ and $g(w^*) \leq 0$

$$\Rightarrow \alpha^* g(w^*) \leq 0 \quad \text{--- (2)}$$

\rightarrow Only way (1) & (2) are true: $\underbrace{\alpha^* g(w^*)}_{\text{Complementary Slackness}} = 0$

\hookrightarrow In our problem:

$$\alpha_i^* (1 - (w^T u_i) y_i) = 0 \quad \forall i$$

$$\Rightarrow \text{If } \alpha_i^* > 0 \Rightarrow 1 - (w^T u_i) y_i = 0 \Rightarrow \underbrace{(w^T u_i) y_i = 1}_{\text{Supporting hyperplane}}$$

Only the points on the "Supporting" hyperplane can contribute to w^*

w^* is a sparse linear combination of the data points

\hookrightarrow Prediction:

$$\begin{aligned} \text{Given } u_{\text{test}}, \quad w^T u_{\text{test}} &= \left(\sum_{i=1}^n \alpha_i^* u_i \cdot y_i \right)^T u_{\text{test}} \\ &= \sum_{i=1}^n \alpha_i^* y_i (u_i^T u_{\text{test}}) \end{aligned}$$

In the kernel,

$$w^T \phi(u_{\text{test}}) = \sum_{i=1}^n \alpha_i^* y_i k(u_i, u_{\text{test}})$$

\uparrow
 > 0 only for support vector

→ Soft-Margin SVM

↳ How to deal with outliers?

Insight: Make every W feasible. Earlier, only the W that correctly classified the datapoints were feasible.

↳ Fix any W . W classifies some points correctly and misclassifies some other points

For the incorrectly classified points "pay penalty" to go to the "correct" side.

↳ Modified formulation:

$$\min_{W, \epsilon} \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \epsilon_i \quad \stackrel{\geq 0}{\text{Hyperparameter}}$$

$$\text{s.t. } (W^T \mathbf{x}_i) y_i + \epsilon_i \geq 1 \quad \forall i \\ \epsilon_i \geq 0$$

→ If $C = 0$, $W = 0 \in \mathbb{R}^d$ is the solution

If $C = \infty$, W is the same as strict linear separator



WEEK 11

→ Dual Formulation For Soft-Margin SVM

↳ Primal Problem

$$\min_{W, \epsilon} \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \epsilon_i \quad \text{s.t. } (W^T \mathbf{x}_i) y_i + \epsilon_i \geq 1 \quad \forall i \quad \epsilon_i \geq 0 \quad \forall i$$

Lagrangian:

$$L(W, \epsilon, \alpha, \beta) = \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \epsilon_i + \sum_{i=1}^n \alpha_i [1 - (W^T \mathbf{x}_i) y_i - \epsilon_i] + \sum_{i=1}^n \beta_i (-\epsilon_i)$$

$$\min_{W, \epsilon} \left[\max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} L(W, \epsilon, \alpha, \beta) \right] \equiv \max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} \left[\min_{W, \epsilon} L(W, \epsilon, \alpha, \beta) \right]$$

$$\frac{\partial L}{\partial W} = 0 \Rightarrow W_{\alpha, \beta}^* = \sum_{i=1}^n \alpha_i \mathbf{x}_i y_i \quad \text{--- ①}$$

$$\frac{\partial L}{\partial \epsilon_i} = 0 \Rightarrow C + \alpha_i (-1) + \beta_i (1) = 0 \Rightarrow \alpha_i + \beta_i = C \quad \text{--- ②}$$

Substitute $W_{\alpha, \beta}^*$ into the Lagrangian:

$$\text{Dual Problem: } \max_{\substack{\alpha \geq 0 \\ \beta \geq 0 \\ \alpha + \beta = C}} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T \mathbf{y}^T \mathbf{X}^T \mathbf{X} \mathbf{y} \alpha$$

↪ w.p. β is effectively restricting α to C .

↳ Dual Problem: $\max_{\alpha \in \mathbb{R}^n} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T \mathbf{y}^T \mathbf{X}^T \mathbf{X} \mathbf{y} \alpha$

\rightarrow If $C=0$, $\alpha^* = 0 \in \mathbb{R}^n \Rightarrow W^* = \sum \alpha_i y_i x_i = 0$

\rightarrow If $C=\infty$, \Rightarrow Hard-margin SVM

\rightarrow Kernelisable because of $\mathbf{X}^T \mathbf{X}$

↳ Complementary Slackness

let (W^*, ε^*) be the primal optimal solutions

let (α^*, β^*) " " dual "

$$\textcircled{1} \quad \forall i \quad \alpha_i^* (1 - W^* x_i y_i - \varepsilon_i^*) = 0$$

$$\textcircled{2} \quad \forall i \quad \beta_i^* (\varepsilon_i^*) = 0$$

Cases possible: (dual problem)

$$\textcircled{1} \quad \alpha_i^* = 0 \Rightarrow \beta_i^* = C \Rightarrow \varepsilon_i^* = 0 \Rightarrow W^* x_i y_i \geq 1$$

$\Leftrightarrow W^*$ classifies (x_i, y_i) correctly

$$\textcircled{2} \quad \alpha_i^* \in (0, C) \Rightarrow 0 < \beta_i^* < C \Rightarrow \varepsilon_i^* = 0 \Rightarrow W^* x_i y_i = 1$$

$\Leftrightarrow x_i$ lies on the supporting hyperplane

$$\textcircled{3} \quad \alpha_i^* = C \Rightarrow \beta_i^* = 0 \Rightarrow \varepsilon_i^* \geq 0 \Rightarrow (1 - W^* x_i y_i - \varepsilon_i^*) = 0$$

$\Leftrightarrow W^* x_i y_i \leq 1 \Rightarrow x_i$ is being misclassified by W^* or isn't classified with enough margin

Cases possible: (primal problem)

$$\textcircled{1} \quad W^* x_i y_i < 1 \Rightarrow \varepsilon_i^* > 0 \Rightarrow \beta_i^* = 0 \Rightarrow \alpha_i^* = C$$

$\Leftrightarrow x_i$ is misclassified

$$\textcircled{2} \quad W^* x_i y_i = 1 \Rightarrow \varepsilon_i^* \geq 0 \Rightarrow \alpha_i^* \in [0, C]$$

$\Leftrightarrow x_i$ is on the supporting hyperplane

$$\textcircled{3} \quad W^* x_i y_i > 1 \Rightarrow \alpha_i^* = 0$$

$\Leftrightarrow x_i$ which is classified with margin > 1 does not contribute to W^*

↳ Summary

Dual:

$$\textcircled{1} \quad \alpha_i^* = 0 \Rightarrow W^* x_i y_i \geq 1$$

$$\textcircled{2} \quad \alpha_i^* \in (0, C) \Rightarrow W^* x_i y_i = 1$$

$$\textcircled{3} \quad \alpha_i^* = C \Rightarrow W^* x_i y_i \leq 1$$

Primal:

$$\textcircled{1} \quad W^* x_i y_i < 1 \Rightarrow \alpha_i^* = C$$

$$\textcircled{2} \quad W^* x_i y_i = 1 \Rightarrow \alpha_i^* \in [0, C]$$

$$\textcircled{3} \quad W^* x_i y_i > 1 \Rightarrow \alpha_i^* = 0$$

→ Algorithms so far

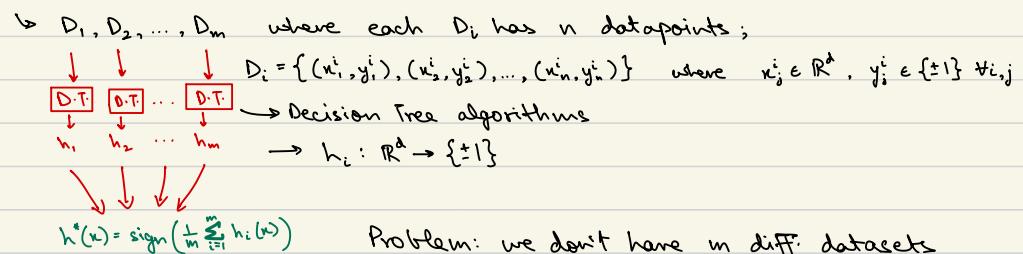
- ↳ Generative: Naive Bayes, Gaussian Naive Bayes
- Discriminative: K-NN, Decision trees, Logistic Regression, Perceptron, SVM

→ Meta classifiers / Ensemble classifiers

- ↳ Weak learners → Strong learners
(better than random)
- ↳ Underfitting → high bias: even if the data changes a little, the model will vary within a small range.
- Overfitting → high variance: if the data changes a little bit, the prediction is going to change by a lot.
- ↳ Bias · tells you how good is your assumed structure actually representing the truth.
- ↳ error = bias + variance

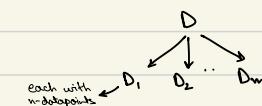
→ Bagging: Bootstrap Aggregation

- ↳ $x_1, x_2, \dots, x_n \sim N(\mu, 1)$
- $\hat{\mu}_1 = x_1, \hat{\mu}_2 = x_2, \dots, \hat{\mu}_n = x_n \quad \left\{ \begin{array}{l} \text{these are all} \\ \text{unbiased estimators} \end{array} \right.$
- Unbiased estimators: $E[\hat{\theta}] = \theta$
- $E[\hat{\mu}_i] = E[x_i] = \mu$; similarly, $E[\hat{\mu}_{ml}] = \mu$
- Using $\hat{\mu}_{ml}$ instead of $\hat{\mu}_1, \dots, \hat{\mu}_n$ reduces "variance"



- ↳ Input: $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ $x_i \in \mathbb{R}^d$ $y_i \in \{\pm 1\}$ &

Bootstrapping: sampling with replacement



create m-different datasets using bootstrapping and then train different models on each dataset

Probability that a point doesn't get picked: $1 - 1/n$

" " " " " " in any bag: $(1 - 1/n)^n$

Probability that a point appears in any bag: $1 - (1 - 1/n)^n \approx 67\%$. For $\gg n$

→ Bagging reduces variance

↳ Random Forest

→ Bag of Decision Trees (typically overfit trees)

→ Feature Bagging

↳ when building a decision tree, pick a subset of features using the concept of impurity.

→ Boosting

↳ also weak learner → strong learner

Decision Stumps : $f_k < \theta$: 1-level or 2-level

↳ AdaBoost Algorithm

Input : $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ $x_i \in \mathbb{R}^d$ $y_i \in \{\pm 1\}$

→ every datapoint gets assigned a weight

Initialize : $D_0(i) = \frac{1}{n}$
iteration \uparrow
with datapoint

for $t=1, \dots, T$:

→ h_t = Input (S, D_t) to a weak learner

→ $\tilde{D}_{t+1}(i) = \begin{cases} D_t(i) \cdot e^{h_t(x_i)} & \text{if } h_t(x_i) \neq y_i \\ D_t(i) \cdot e^{-h_t(x_i)} & \text{if } h_t(x_i) = y_i \end{cases}$

→ $D_{t+1}(i) = \frac{\tilde{D}_{t+1}(i)}{\sum_j \tilde{D}_{t+1}(j)}$

Output : h_1, h_2, \dots, h_T . $h^*(x) = \text{sign}(\sum_t \alpha_t h_t(x))$

$$\alpha_t = \ln \sqrt{\frac{1 - \text{err}(h_t)}{\text{err}(h_t)}}$$



WEEK 12

Classification Loss Functions

→ Data: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ $x_i \in \mathbb{R}^d$, $y_i \in \{\pm 1\}$

Goal: $h: \mathbb{R}^d \rightarrow \{\pm 1\}$

Performance measure: $\sum_{i=1}^n \mathbf{1}(h(x_i) \neq y_i)$

If $h \in H_{\text{linear}}$,

$$\min \sum_{i=1}^n \mathbf{1}(h(x_i) = y_i) \equiv \min_{\substack{w \in \mathbb{R}^d \\ b \in \mathbb{R}}} \sum_{i=1}^n \mathbf{1}(\text{sign}(w^T x_i + b) \neq y_i)$$

↳ NP-HARD

$$x \in \mathbb{R}^d, x' = [x \mid 1] \in \mathbb{R}^{d+1}$$

$$w \in \mathbb{R}^d, w' = [\underbrace{w_1, \dots, w_d}_w \mid \underbrace{w_{d+1}}_b] \quad w'^T x' = w^T x + b$$

→ loss-function view

→ every point suffers a loss of 1 or 0

$$\mathbf{1}(h(x) \neq y) = \mathbf{1}(h(x) \cdot y < 0)$$

$$(x, y) \in \mathbb{R}^d \times \{\pm 1\} \quad h: \mathbb{R}^d \rightarrow \{\pm 1\}$$

$$h_w(x) = \text{sign}(w^T x)$$

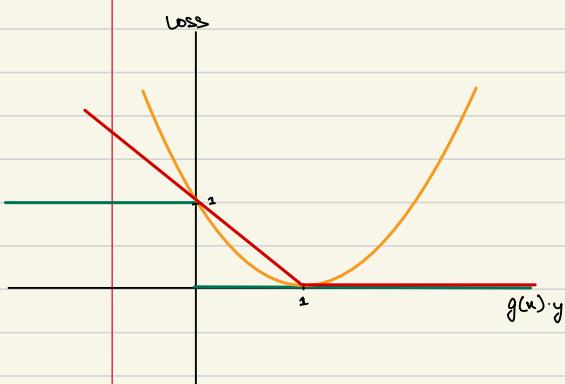
→ Algorithm

→ Using Regression for classification

$$\begin{aligned} \text{Loss}(g, (x, y)) &= (g(x) - y)^2 \quad \text{--- ①} \quad \text{where } g(x) = w^T x + b \\ &= (g(x) - 1)^2 \quad \text{--- ②} \\ &= g(x)^2 + 1 - 2g(x) \cdot y \end{aligned}$$

→ ① and ② are not graphically same although they're algebraically equivalent.

→ Therefore, regression is not a great algorithm for binary classification



→ Support Vector Machine

$$\begin{aligned} \min_{w, \epsilon} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \epsilon_i \quad &\text{st. } (w^T x_i) \cdot y_i + \epsilon_i \geq 1 \quad \text{and } \epsilon_i \geq 0 \\ &\text{Regularizer} \\ &= \min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - (w^T x_i) \cdot y_i) \\ &\text{margin, Model Dependent} \quad \text{Data-dependent; Loss} \end{aligned}$$

→ Convex function which is also a good approximation for the 0-1 loss function.

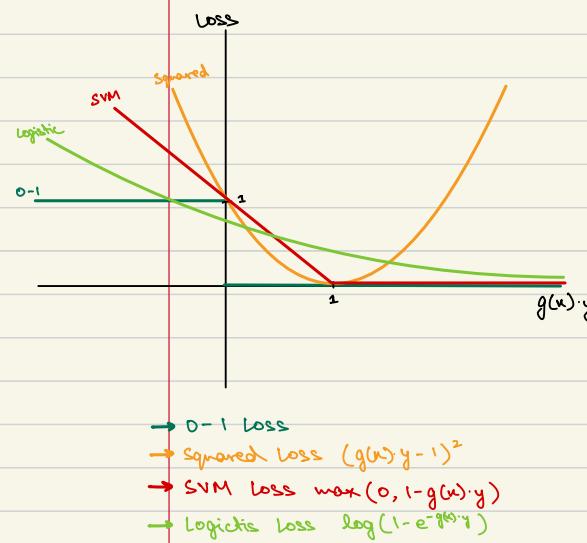
→ A.R.O. Hinge Loss

→ SVM: Regularization + Hinge Loss

→ 0-1 Loss

→ Squared Loss $(g(x) \cdot y - 1)^2$

→ SVM Loss $\max(0, 1 - g(x) \cdot y)$



→ Conclusion

→ 0-1 loss is NP-hard.

→ Diff. algorithms use different "surrogate" loss.

→ Surrogates are convex.

→ Perceptron

$$W_{t+1} = W_t + \nu_t y_t$$

$$\text{modified hinge loss}(w, (x, y)) = \max(0, -(w^T x)y)$$

$$\nabla_w \text{hinge} = \begin{cases} -\nu_t y & (w^T x)y < 0 \\ 0 & (w^T x)y > 0 \\ [-1, 0] \nu_t y & (w^T x)y = 0 \end{cases}$$



$$\text{gradient descent: } W_{t+1} = W_t - \eta \nabla_w \text{loss}(W_t)$$

$$\text{Perceptron: } W_{t+1} = W_t - \underbrace{\nu_t y_t}_{\eta = 1}$$

→ Perceptron can be interpreted as stochastic gradient descent with modified hinge loss and step size = 1

→ Boosting

$$\text{loss}(h, (x, y)) = e^{-y h(x)}$$

→ exponential loss

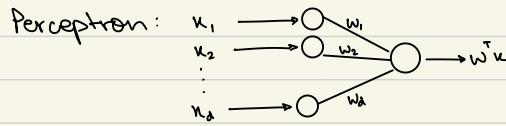
→ Neural Networks

↳ Inspired by Perceptrons

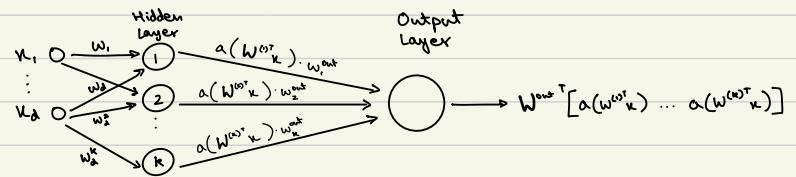
Doesn't use convex loss functions

$$\hookrightarrow \mathbf{x} \in \mathbb{R}^d \rightarrow \text{sign}(\mathbf{w}^\top \mathbf{x})$$

\downarrow
 $[x_1, x_2, \dots, x_d]$



Neural Network:



Parameters: $\{w^{(1)}, \dots, w^{(L)}, w^{out}\}$ $w^{(l)} \in \mathbb{R}^d$, $w^{out} \in \mathbb{R}^k$

Prediction:

$$\hat{y} = \sum_{i=1}^k w_i^{out} \alpha(w^{(i)^T} \mathbf{x})$$

Activation function

↳ examples of activation function

$$① \alpha(z) = \frac{1}{1 + e^{-z}} \quad \text{Sigmoid}$$

$$② \alpha(z) = \max(0, z) \quad \text{Rectified Linear Unit}$$

↳ Parameter estimation:

$$\begin{aligned} \text{Regression} : L(\text{NN}(\mathbf{x}_i, \theta), y_i) & \xrightarrow{\{w^{(1)}, \dots, w^{(L)}, w^{out}\}} \\ & = \sum_n (\text{NN}(\mathbf{x}_i, \theta) - y_i)^2 \end{aligned}$$

→ learn θ using Gradient descent

→ gradient computed using Back-Propagation algorithm, which uses the chain rule.

↳ converges to local minima

Classification: Cross entropy loss

↳ compares $\hat{y}_i = P(y| \mathbf{x}_i)$ with y_i

↳ Conclusion:

→ NN learns local minima of non-convex functions

→ Typically works very well in practice, especially for unstructured data.



Week 4 GA

$$\textcircled{7} \quad \frac{(0.054)(0.7)}{(0.242)(0.3) + (0.054)(0.7)} = 0.342$$

$$\textcircled{8} \quad \begin{bmatrix} \lambda_1^1 & \lambda_2^1 \\ \vdots & \vdots \\ \lambda_1^4 & \lambda_2^4 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 4 \\ 4 & 1 \end{bmatrix} \begin{bmatrix} 0.66 & 0.34 \\ 0.41 & 0.59 \\ 0.21 & 0.79 \\ 0.09 & 0.91 \end{bmatrix} \quad z = \begin{bmatrix} 1 \\ 2 \\ 2 \\ 2 \end{bmatrix}$$

$$\textcircled{9} \quad \hat{\mu}_1 = \frac{\sum \lambda_i^1 \cdot x_i}{\sum \lambda_i^1} = 1.8029$$

Week 7 GA

$$\textcircled{3} \quad \begin{aligned} d(x_1) &= \sqrt{1+4} = \sqrt{5} & 1 \\ d(x_2) &= \sqrt{16+16+36+81} = \sqrt{149} & 0 \\ d(x_3) &= \sqrt{4+1+9} = \sqrt{14} & 1 \\ d(x_4) &= \sqrt{1+1} = \sqrt{2} & 0 \\ d(x_5) &= \sqrt{81+36+16+1} = \sqrt{134} & 0 \end{aligned}$$

Week 8 GA

$$\textcircled{10} \quad P(y=1|x) \propto P(x|y) \cdot P(y)$$

$$\begin{aligned} &\underline{P(y=1|x)} \\ &= P(\xi_1=0|y=1) \cdot P(\xi_2=1|y=1) \cdot P(y=1) \\ &= (0.9)(0.02)(0.3) \\ &= 0.0054 \end{aligned}$$

$$\begin{aligned} &\underline{P(y=0|x)} \\ &= P(\xi_1=0|y=0) \cdot P(\xi_2=1|y=0) \cdot P(y=0) \\ &= (0.8)(0.3)(0.7) \\ &= 0.168 \end{aligned}$$

