

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Mathematical Foundations



Table of contents

1. Linear Algebra
2. Basic Probability
3. Bayes Theorem

Basic Linear Algebra

- Given a vector \mathbf{x} of m dimensions, the transpose \mathbf{x}^t is -

$$\mathbf{x} = [x_1 \ x_2 \ x_3 \ \dots \ \dots \ \dots \ \dots \ x_m] \quad \mathbf{x}^t = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ \vdots \\ \vdots \\ x_m \end{bmatrix}$$

Basic Linear Algebra - Determinant

- Given a 2x2 matrix \mathbf{A} , the determinant $|\mathbf{A}|$ is defined as -

$$\mathbf{A} = \begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix} = a_{00}a_{11} - a_{01}a_{10}$$

Basic Linear Algebra - Determinant

- Given a 2x2 matrix \mathbf{A} , the determinant $|\mathbf{A}|$ is defined as -

$$\mathbf{A} = \begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix} = a_{00}a_{11} - a_{01}a_{10}$$

- Given a 3x3 matrix \mathbf{A} , the determinant $|\mathbf{A}|$ is defined as -

$$\mathbf{A} = \begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \\ a_{20} & a_{21} & a_{22} \end{bmatrix} = a_{00}(a_{11}a_{22} - a_{12}a_{21}) - a_{01}(a_{10}a_{22} - a_{20}a_{12}) + a_{02}(a_{10}a_{21} - a_{11}a_{20})$$

Determinant - Examples

- What is the determinant of $A = \begin{bmatrix} 1 & 7 \\ 2 & -4 \end{bmatrix}$

Determinant - Examples

- What is the determinant of $A = \begin{bmatrix} 1 & 7 \\ 2 & -4 \end{bmatrix}$

$$|A| = 1 \times -4 - 7 \times 2 = -18$$

Determinant - Examples

- What is the determinant of $A = \begin{bmatrix} 1 & 7 \\ 2 & -4 \end{bmatrix}$

$$|A| = 1 \times -4 - 7 \times 2 = -18$$

- What is the determinant of $A = \begin{bmatrix} 2 & 3 & 7 \\ -3 & 4 & 0 \\ 1 & -1 & 6 \end{bmatrix}$

Determinant - Examples

- What is the determinant of $A = \begin{bmatrix} 1 & 7 \\ 2 & -4 \end{bmatrix}$

$$|A| = 1 \times -4 - 7 \times 2 = -18$$

- What is the determinant of $A = \begin{bmatrix} 2 & 3 & 7 \\ -3 & 4 & 0 \\ 1 & -1 & 6 \end{bmatrix}$

$$|A| = 2(4 \times 6 - 0 \times -1) - 3(-3 \times 6 - 1 \times 0) + 7(-3 \times -1 - 4 \times 1) = 95$$

Basic Linear algebra - Inverse of a matrix

- For a matrix A, inverse is denoted by A^{-1}

Basic Linear algebra - Inverse of a matrix

- For a matrix A, inverse is denoted by A^{-1}
- $AA^{-1} = A^{-1}A = I$

Basic Linear algebra - Inverse of a matrix

- For a matrix A, inverse is denoted by A^{-1}
- $AA^{-1} = A^{-1}A = I$
- Why is an inverse of a matrix needed?
 - Because matrices cannot be divided!

Inverse of a matrix - Example 1

Q. Calculate the inverse of $P = \begin{bmatrix} 1 & 2 \\ 3 & -5 \end{bmatrix}$

Inverse of a matrix - Example 1

Q. Calculate the inverse of $P = \begin{bmatrix} 1 & 2 \\ 3 & -5 \end{bmatrix}$

Step 1: Calculate the determinant. Here, $|P| = -11$

Inverse of a matrix - Example 1

Q. Calculate the inverse of $P = \begin{bmatrix} 1 & 2 \\ 3 & -5 \end{bmatrix}$

Step 1: Calculate the determinant. Here, $|P| = -11$

Step 2: Swap diagonal elements and add negative sign to off-diagonal elements.

$$\begin{bmatrix} -5 & -2 \\ -3 & 1 \end{bmatrix}$$

Inverse of a matrix - Example 1

Q. Calculate the inverse of $P = \begin{bmatrix} 1 & 2 \\ 3 & -5 \end{bmatrix}$

Step 1: Calculate the determinant. Here, $|P| = -11$

Step 2: Swap diagonal elements and add negative sign to off-diagonal elements.

$$\begin{bmatrix} -5 & -2 \\ -3 & 1 \end{bmatrix}$$

Step 3: Divide the matrix by the determinant.

$$\begin{bmatrix} \frac{5}{11} & \frac{2}{11} \\ \frac{-3}{11} & \frac{-1}{11} \end{bmatrix}$$

Inverse of a matrix - Example 2

Q. Calculate the inverse of $\begin{bmatrix} 3 & -9 \\ 2 & -6 \end{bmatrix}$

Inverse of a matrix - Example 2

Q. Calculate the inverse of $\begin{bmatrix} 3 & -9 \\ 2 & -6 \end{bmatrix}$

Solution. Here, $|A| = 0$. Therefore, the inverse does not exist!

Such a matrix is called **singular matrix**!

Probability

- A probability space is a triplet (Ω, \mathcal{B}, P) that is used to model a process or an experiment with random outcomes.

Probability

- A probability space is a triplet (Ω, \mathcal{B}, P) that is used to model a process or an experiment with random outcomes.
 - Ω - Sample space

Probability

- A probability space is a triplet (Ω, \mathcal{B}, P) that is used to model a process or an experiment with random outcomes.
 - Ω - Sample space
 - \mathcal{B} - Collection of subsets of Ω

Probability

- A probability space is a triplet (Ω, \mathcal{B}, P) that is used to model a process or an experiment with random outcomes.
 - Ω - Sample space
 - \mathcal{B} - Collection of subsets of Ω
 - P - probability

Conditional Probability

- Let (Ω, \mathcal{B}, P) be a probability space and let $H \in \mathcal{B}$ with $P(H) > 0$. For any $B \in \mathcal{B}$, $P(B|H)$ is defined as-

$$P(B|H) = P(BH) / P(H)$$

and call $P(B|H)$ the conditional probability of B given H .

Total Probability Rule

- Let (Ω, \mathcal{B}, P) be a probability space, and let $\{H_j\}$ be pairwise disjoint events in \mathcal{B} (i.e. $H_j H_k = \emptyset, \forall j \neq k$) and $\bigcup_{j=1, \dots, \infty} H_j = \Omega$.

Suppose $P(H_j) > 0, \forall j$, then,

$$P(B) = \sum_{j=1, \dots, \infty} P(H_j)P(B|H_j)$$

Bayes Theorem

- Let (Ω, \mathcal{B}, P) be a probability space, and let $\{H_j\}$ be pairwise disjoint events in \mathcal{B} (i.e. $H_j H_k = \emptyset, \forall j \neq k$) and $\bigcup_{j=1, \dots, \infty} H_j = \Omega$. and $P(H_j) > 0, \forall j$. We have, $\forall B \in \mathcal{B}$ and $P(B) > 0$,

Bayes Theorem

- Let (Ω, \mathcal{B}, P) be a probability space, and let $\{H_j\}$ be pairwise disjoint events in \mathcal{B} (i.e. $H_j H_k = \emptyset, \forall j \neq k$) and $\bigcup_{j=1, \dots, \infty} H_j = \Omega$. and $P(H_j) > 0, \forall j$. We have, $\forall B \in \mathcal{B}$ and $P(B) > 0$,

$$P(H_j|B) = \frac{P(H_j) P(B|H_j)}{\sum_{i=1, \dots, \infty} P(H_i) P(B|H_i)}, \quad \forall j$$

Bayes Theorem

- Let (Ω, \mathcal{B}, P) be a probability space, and let $\{H_j\}$ be pairwise disjoint events in \mathcal{B} (i.e. $H_j H_k = \emptyset, \forall j \neq k$) and $\bigcup_{j=1, \dots, \infty} H_j = \Omega$. and $P(H_j) > 0, \forall j$. We have, $\forall B \in \mathcal{B}$ and $P(B) > 0$,

$$\text{Posterior } \Rightarrow P(H_j|B) = \frac{\text{Prior } \rightarrow P(H_j) \text{ Likelihood } \rightarrow P(B|H_j)}{\sum_{i=1, \dots, \infty} P(H_i)P(B|H_i)}, \quad \forall j$$

↑
Evidence

Bayes Theorem

Consider two events A and B, then the joint probability is-

$$P(AB) = P(B|A)P(A)$$

$$\Rightarrow P(AB) = P(A|B)P(B)$$

$$\Rightarrow P(B|A)P(A) = P(A|B)P(B)$$

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

Bayes Theorem - Example 1

Q. A test is developed to detect a disease that 0.1% of the population have. The test is 99% effective in detecting an infected person. However, it gives a false positive result for 0.5% of cases. Find the probability that a person actually has the disease if the person tests positive?

Sol.: Let X be the event that a person has the disease & Y be the event that the test result is true. $P(X) = 0.001$, $P(Y|X) = 0.99$, $P(Y|\sim X) = 0.005$, $P(\sim X)=1-P(X)=0.991$. We need to find $P(X|Y)$. Using Bayes theorem and total probability rule, we have

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y|X)P(X)+P(Y|\sim X)P(\sim X)}$$

Substituting the above values, we get $P(X|Y) = 0.165$

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Supervised Learning



Table of contents

1. Project Part - 1
2. Density Estimation using MLE
3. Generative v/s Discriminative Models

Project Part 1 - Density Estimation and Supervised Learning

- Due on: 17 Feb, 11:59 PM MST.
- Binary classification using Fashion-MNIST dataset.
- Two classes - T-shirt (label 0) and Trouser (label 1)
- 28x28 grayscale images
- Training set: "Tshirt": 6000; "Trouser": 6000.
- Testing set: "Tshirt": 6000; "Trouser": 1000.



Project Part 1 - Tasks

- Calculate two features of each image - average and standard deviation.
- Estimate parameters for 2D Gaussian distribution for each class.
- Implement Naive Bayes to perform classification
- Train Logistic regression using gradient ascent to perform classification

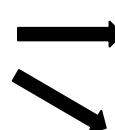
Project Part 1 - Deliverables

- Code:
 - Acceptable file types are .py/.m or .zip.
 - Well-commented code. Be sure to read through the directions carefully to ensure you have included all necessary parts in your code.
- Report:
 - Acceptable File types: .pdf
 - Length: 2-5 A4 pages
 - Content: Refer Canvas

Project Part 1 - Languages/Software

- What can you use?
 - Python/MATLAB to code.
 - Libraries or in-built functions to manipulate the data i.e. NumPy, SciPy libraries etc.
- What cannot be used?
 - Packages like scikit-learn which have ready-to-use algorithms.

What is supervised learning?

- Data - <sample, label> pairs
- Objective - Learn from the given data such that label can be predicted for a new data sample.
- Based on labels -  Regression
Classification

Density Estimation

- Estimating underlying probability density function, based on the training data
- Parametric: each class of images (the feature vectors) may be modeled by a density function $p_\theta(x)$ with parameter θ .
- Non-parametric: makes no assumption about the distribution of data or modeling the data without any parameters

Maximum Likelihood Estimation

- Given some training data and assuming a parametric model $p(x|\theta)$; what specific θ will fit/explain the data best?
- To consider all the samples denoted by $D=\{x_1, x_2, \dots, x_n\}$, assume that all the samples are i.i.d - independent and identically distributed.
- So, data likelihood represented by $L(\theta)$ is -

$$L(\theta) = P(D|\theta) = \prod_i P(x_i|\theta)$$

Maximum Likelihood Estimation

- ① Likelihood function
- ② take derivative
- ③ Set to 0.

- Maximum Likelihood Estimation (MLE): Finding the parameter that maximizes the likelihood function

$$\hat{\theta} = \operatorname{argmax}_{\theta} p(x|\theta) \rightarrow L1$$

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} \prod_{i=1}^n p(x_i|\theta) \\ &\stackrel{?}{=} \operatorname{argmax}_{\theta} \log \prod_{i=1}^n p(x_i|\theta) = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log p(x_i|\theta)\end{aligned}$$

MLE - Example 1

$$\theta = \{\mu, \sigma^2\}$$

Given n i.i.d. samples $\{x_i\}$ from the 1-D normal distribution $N(\mu, \sigma^2)$, find the MLE for μ and σ^2 .

① Likelihood func

$$P(D|\theta) = \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$P(D|\theta) = \prod_{i=1}^n P(x_i|\theta)$$

$$\begin{aligned} \log P(D|\theta) &= \underbrace{\log \prod_{i=1}^n}_{\text{log}} P(x_i|\theta) = \underbrace{\log \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n}_{\text{log}} \sum_{i=1}^n \underbrace{\frac{(x_i-\mu)^2}{2\sigma^2}}_{-\frac{1}{2\sigma^2} \log \sigma^2} \end{aligned}$$

$$(\mu, \sigma^2)$$

MLE - Example 1

② Take derivative w.r.t μ

$$\begin{aligned}\frac{\partial}{\partial \mu} \log P(D|\theta) &= 0 - \sum_{i=1}^n \frac{\partial}{\partial \mu} \frac{f(x_i - \mu)}{\sigma^2} \\ &= \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma^2}\end{aligned}$$

③ Equate derivative to 0.

$$\begin{aligned}\sum_{i=1}^n (x_i - \mu) &= 0 \\ \sum_{i=1}^n x_i - n\mu &= 0 \\ n\bar{x} - n\mu &= 0 \\ \bar{x} - \mu &= 0\end{aligned}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \Rightarrow \text{Sample mean}$$

$$\sum_{i=1}^n x_i - \frac{\sum_{i=1}^n x_i}{n} = 0$$

Generative vs Discriminative Models

- Generative models -
 - Learn $P(y)$ and $P(x|y)$.
 - Ex: Bayesian classifier, Naive Bayes.
- Discriminative models -
 - Directly learn $P(y|x)$
 - Ex: Logistic Regression

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Supervised Learning



Table of contents

- 1. MLE solution for variance**
- 2. Naive Bayes**
- 3. Logistic Regression**

MLE - Example 1 (Continued)

$$\sigma \rightarrow \sum \begin{cases} \sigma^2 \\ \sigma^2 \end{cases}$$

Given n i.i.d. samples $\{x_i\}$ from the 1-D normal distribution $N(\mu, \sigma^2)$, find the MLE for μ and σ^2 .

$$\log P(D|\mu, \sigma) = \underbrace{-n \log(\sigma \sqrt{2\pi})}_{\text{log}} - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}$$
$$\frac{\partial}{\partial \sigma} n P(D|\mu, \sigma) = -\frac{n}{\sigma} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^3}$$

Equate it to 0,

$$-\frac{n}{\sigma} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^3} = 0$$

MLE - Example 1 (Continued)

$$\frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2} = \frac{n}{\sigma^2}$$

$$\left\{ \hat{\sigma}_{\text{MLE}}^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \right\} \text{Biased}$$

$$\text{Bias} - E(\hat{\sigma}_{\text{MLE}}^2) - \sigma^2 = 0 \quad \text{Unbiased}$$

$$P(X|Y)$$

Naive Bayes

$$P(X_1, X_2, X_3 | Y)$$

2^3 2^3

$$Y=0$$

$$Y=1$$

- The "naive" conditional independence assumption: each feature is (conditionally) independent of every other feature, given the label, i.e.,
 $p(x_i | \{x_j \text{ for any } j \neq i\}, y) = p(x_i | y)$
- How does this simplify?
 - Consider the previous example again: d-dimensional binary features, and y is also binary.

Bayesian \rightarrow No. of probabilities $P(X_1, X_2, X_3 | Y) = 2^{3+1}$

\nexists d-dimensional $X \rightarrow 2^{d+1} \doteq 4d$

$$P(X_1, X_2, X_3 | Y) = P(X_1 | Y) \cdot P(X_2 | Y) \cdot P(X_3 | Y)$$
$$Y=0 \rightarrow P(X_1=0 | Y=0) \cdot P(X_2=0 | Y=0) \cdot P(X_3=0 | Y=0) = 6 \cdot 6 = 36$$
$$Y=1 \rightarrow P(X_1=1 | Y=1) \cdot P(X_2=1 | Y=1) \cdot P(X_3=1 | Y=1) = 6 \cdot 6 = 36$$

Naive Bayes

- The predicted label is given by -

$$\hat{y} = \operatorname{argmax}_y P(y) \prod_{i=1}^d p(x_i|y)$$

$P(y|x) = P(x|y) \cdot P(y)$

$P(x)$

- Model parameters- x_1, x_2, x_3, y

$$\frac{P(x_1=0|y=0)}{P(x_1=0|y=1)}$$

$$\frac{P(x_2=0|y=0)}{P(x_2=1|y=1)}$$

$$\frac{P(x_3=0|y=0)}{P(x_3=1|y=1)}$$

$$P(x=0)$$

$y = \text{binary}$
 $x_i = \text{two values}$
 $(2d+1)$

Naive Bayes - Example

X

Y

Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

$$P(Y = \text{Yes}) = \frac{3}{4}$$
$$P(x_1 = \text{Sunny} | Y = \text{Yes}) = \frac{3}{3}$$

$(2d+1) = 13$ independent parameters

$$P(X|Y)$$

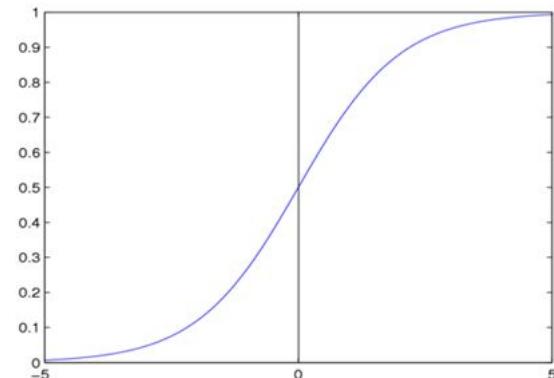
Logistic Regression

$$\underline{\underline{w^T x \geq 0}}$$

- Training set: n labelled samples $\langle \underline{\underline{x(i)}}, \underline{\underline{y(i)}} \rangle$
- Use the logistic function for modeling $\overbrace{P(y|x)}$, considering only the case of $y \in \{0,1\}$

$$\underline{\underline{P(y=0|x) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}}} = 1 - \sigma(\underline{\underline{w^T x}})$$

$$\underline{\underline{P(y=1|x) = \frac{\exp(w_0 + \sum_{i=1}^d w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}}} = \sigma(\underline{\underline{w^T x}})$$



$$\text{Sigmoid} \quad \underline{\underline{\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t}}}$$

Logistic Regression

$w \rightarrow$ model parameter

- Model parameters -

$$P(Y|X) = \left[1 - \sigma(w^T x) \right]^{1-y} \left[\sigma(w^T x) \right]^y$$
$$P(y^{(1)}, y^{(2)}, \dots, y^{(n)} | x^{(1)}, x^{(2)}, \dots, x^{(n)}) = \prod_{i=1}^n \left[1 - \sigma(w^T x^{(i)}) \right]^{1-y^{(i)}} \left[\sigma(w^T x^{(i)}) \right]^{y^{(i)}}$$

Argmax

$$\hat{L}(w) = \sum_{i=1}^m [1 - y^{(i)}] \log [1 - \sigma(w^T x^{(i)})] + y^{(i)} \log [\sigma(w^T x^{(i)})]$$

Gradient Ascent

Update Equation -



$$w_{t+1} = w_t + \alpha \nabla_w L(w)$$

Updated weights Current weight Learning rate Gradient of the loss function

The equation $w_{t+1} = w_t + \alpha \nabla_w L(w)$ is displayed. Above the equation, arrows point from the terms to their definitions: 'Updated weights' points to w_{t+1} , 'Current weight' points to w_t , 'Learning rate' points to α , and 'Gradient of the loss function' points to $\nabla_w L(w)$.

Descent -

Gradient Ascent

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Supervised Learning

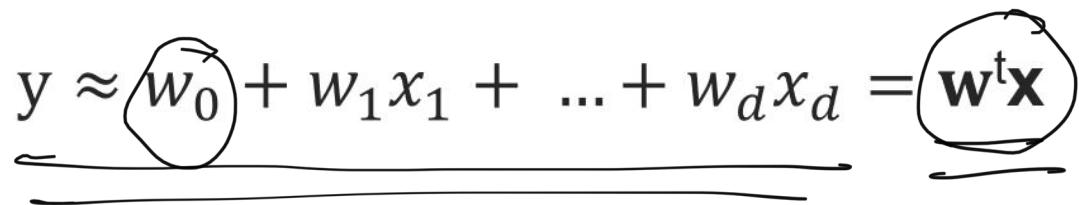


Table of contents

1. Linear Regression
2. Q/A Session

Linear Regression

- Regression - A training set of n samples $\langle x^{(i)}, y^{(i)} \rangle$ where $y^{(i)}$ is a continuous “label” (or target value) for $x^{(i)}$
- Linear regression - modeling the relation between y and x via a linear function

$$y \approx w_0 + w_1 x_1 + \dots + w_d x_d = \mathbf{w}^t \mathbf{x}$$


Linear Regression

- The error is given as -

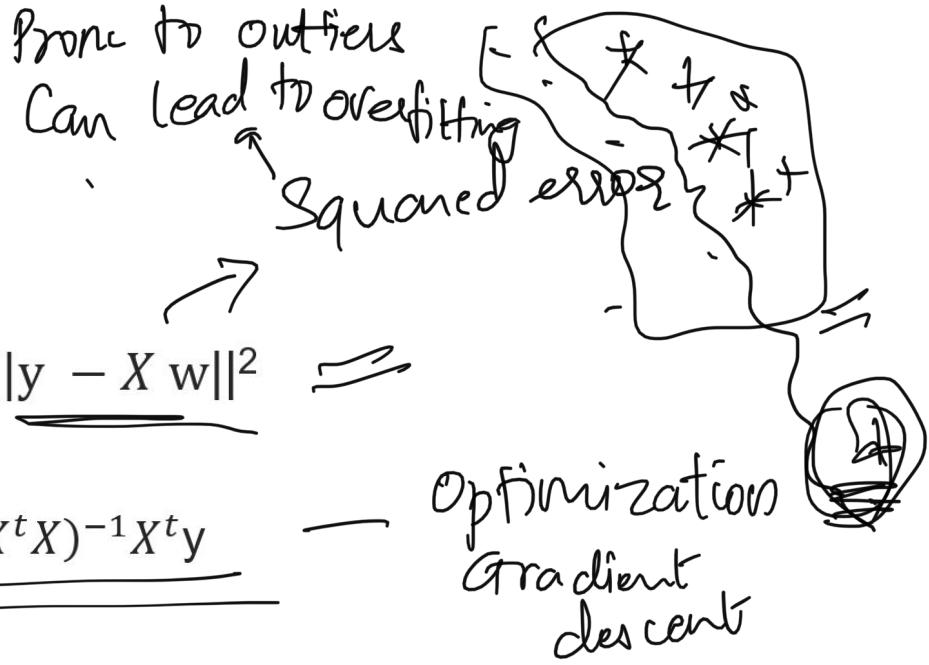
$$\text{argmin} \|\mathbf{e}\|^2 = \|\mathbf{y} - \mathbf{Xw}\|^2$$

- Weights can be found using -

$$\hat{\mathbf{w}} = \underline{(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}}$$

- How to generalize?

$$\underline{\mathbf{y}} = \underline{w_0 + w_1 \phi_1(\mathbf{x}) + \dots + w_{M-1} \phi_{M-1}(\mathbf{x})}$$



Regularization in Linear Regression

Squared error + $\lambda \|w\|^1$

$$E_D(w) + \lambda E_W(w)$$

Squared error

hyperparameter

$\|w\|^p$

$\begin{cases} 1, & \text{Ridge} \\ 2, & \text{Lasso} \end{cases}$

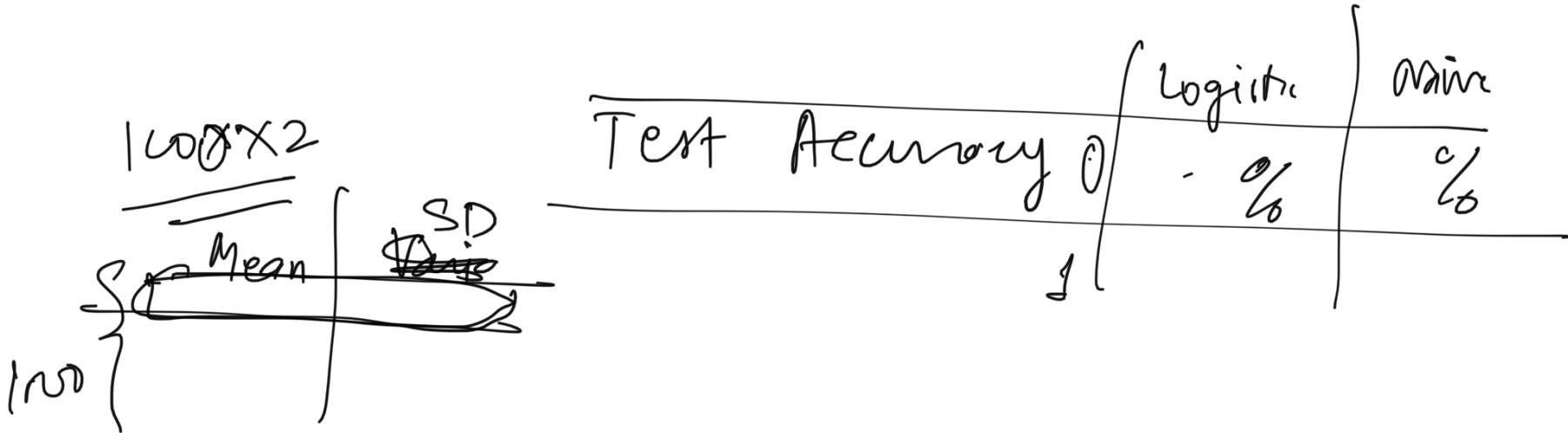
Regularization term

Why regularize and how it works?

- Used to avoid overfitting of the network
- Regularization terms shrinks the w estimates.



Covariance M
[]



Questions?

Project Part - 1

After feature extraction - 6000×2
 Training size = 12000×2 784
 Test data size = 2000×2

① Estimation parameters - np. cov

One class $\mu = [\mu_1, \mu_2]$ $\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$?

$\Sigma = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{bmatrix}$ Class 1 - μ_1, Σ_1
Class 0 - μ_0, Σ_0

Hint: The features are independent.

Class 0 - Avg, standard deviation $\begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{bmatrix}$

$\mu, \sigma \leftarrow 1\text{-D Gaussian}$ 1-D Gaussian

Class 1 - Avg Std
1-D Gaussian

784

P-1

② Naive Bayes -

$$P(Y|X) \sim P(X|Y) \cdot P(Y)$$

From the data

$$\frac{P(X_1|Y) \cdot P(X_2|Y)}{P(Y)}$$

$$\underline{P(Y=0|X)} > \underline{P(Y=1|X)} \rightarrow \text{Class 0}$$

otherwise Class = 1

③ Logistic Regression - Weights w $P(Y|X)$

Step 1 : Initialize w randomly.

Step 2 : $Z = \sigma(w^T x)$ $\sigma(t) = \frac{1}{1+e^{-t}}$

Step 3 : Calculate the cost function.

Log-likelihood funⁿ -

$$L(w) = \sum_{i=1}^m y^{(i)} \log Z^{(i)} + (1-y^{(i)}) \log(1-Z^{(i)})$$

$\in \mathbb{R}^n$

Step 4 : Gradient - ascent

$$w_t^{j+1} = w_t^j + \eta \frac{\partial L(w)}{\partial w^j} \quad \text{np. randn}$$

$$\frac{\partial L(w)}{\partial w^j} = \frac{\partial L(w)}{\partial z} \cdot \frac{\partial z}{\partial p} \cdot \frac{\partial p}{\partial w^j} \quad \Rightarrow \text{Using chain rule.}$$

$$\frac{\partial L(w)}{\partial z} = \frac{y}{z} - \frac{(1-y)}{1-z} \quad - \textcircled{1}$$

$$\frac{\partial z}{\partial p} = z(1-z) \quad - \textcircled{2}$$

$$\frac{\partial p}{\partial w^j} = x^j \quad - \textcircled{3}$$

$$\frac{\partial L(w)}{\partial w^j} = \frac{y(1-z) - z(1-y)}{z(1-z)} \cdot z(1-z) x^j$$

$$= (y - z) x^j$$

Prediction $z = \sigma(w^T x)$

$z > 0.5 \rightarrow \text{Class 1}$

otherwise $\rightarrow \text{Class 0}$

Covariance matrix - Class 0

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{bmatrix}$$

$12000 \begin{bmatrix} x_1 & x_2 \end{bmatrix}^2$ - Training data

$6000 \begin{bmatrix} x_1 & x_2 \end{bmatrix}^2$ Class 0

$$6000 \begin{bmatrix} x_1 & x_2 \end{bmatrix}^2$$

$$N_1, \mu_1, \Sigma_1$$

Class 0

$$\mu_1, \mu_2, \Sigma$$

$$N_2, \mu_2, \Sigma$$

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Linear Machines & SVM



Table of contents

1. Quick Recap
2. Linear Discriminant Function
3. Support Vector Machines (SVM)

Quick Recap - Logistic Regression

- Data - X (features) and Y (labels)
- Learn $\underline{P(Y|X)}$ assuming a logistic function —Discriminative model.
- $g(x) = \underline{\underline{w^T x}}$ is called the linear discriminant function.

$$\underbrace{\underline{\underline{w^T x}}}_w$$

Linear Discriminant functions

- Two types of notations-

$$g(x) = w^T x \quad \text{or} \quad g(x) = w^T x + w_0$$

\downarrow weights $=$
 $w \quad w_0$ bias / threshold

- For 2-class problem, learn w such that -

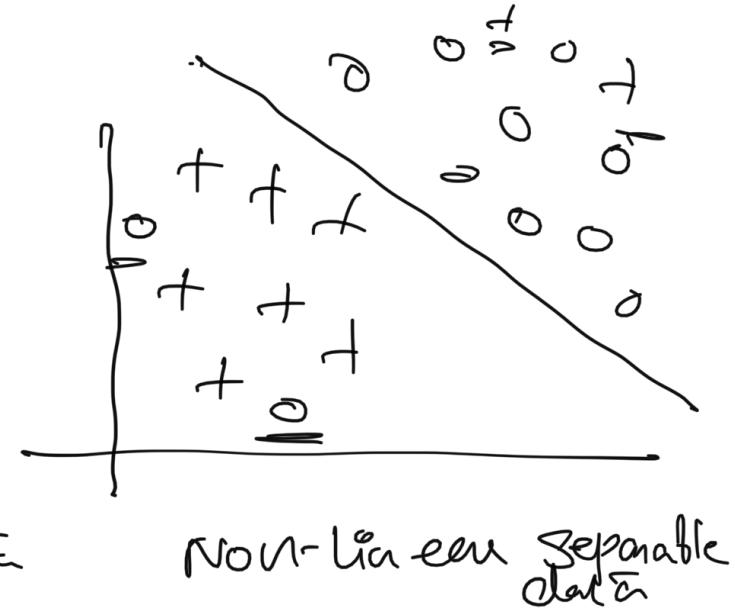
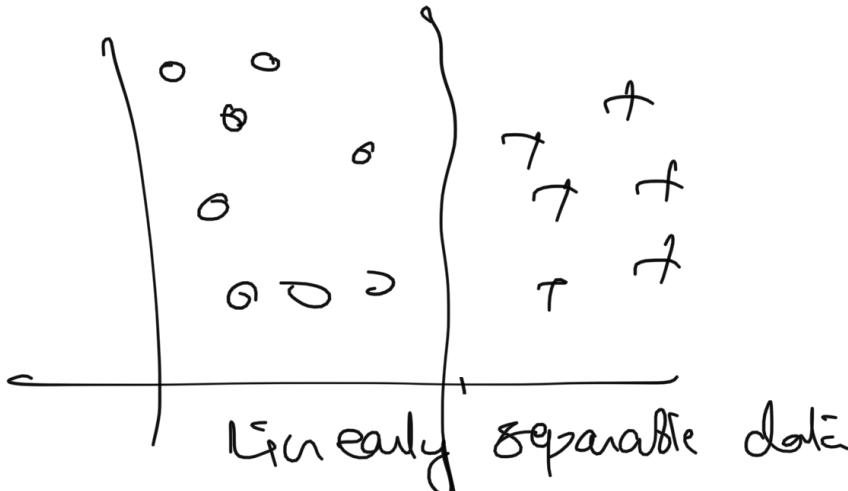
$$w^T x \geq 0 \rightarrow \text{Class 1}$$

$$w^T x \leq 0 \rightarrow \text{Class 0}.$$

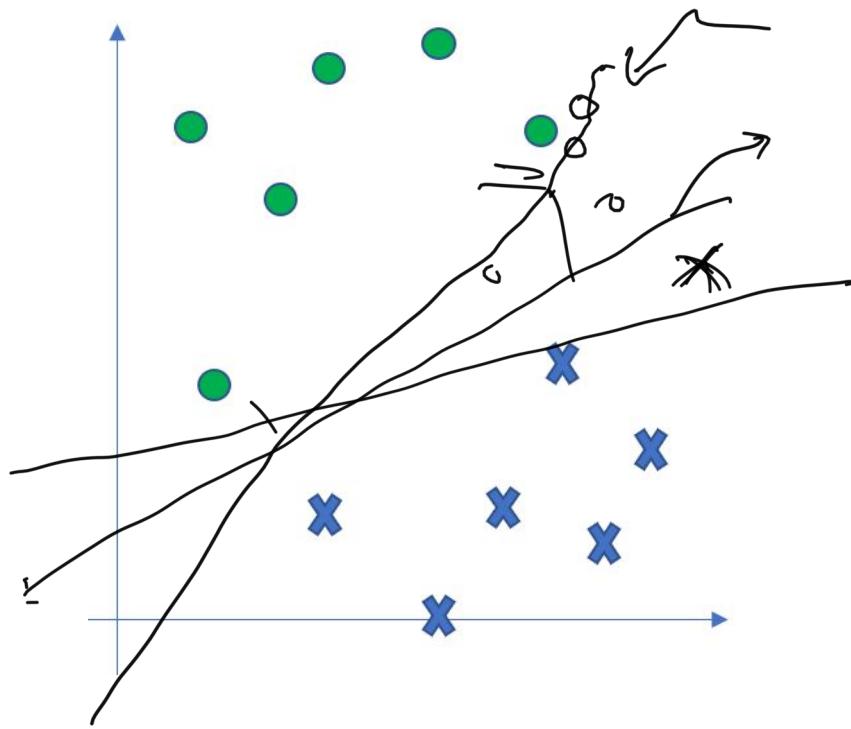
$$\omega^T b$$

What is Linear Separability?

- If there is atleast one solution of w such that $g(x)$ classifies all the samples in the training data, then the data is said to be linearly separable.



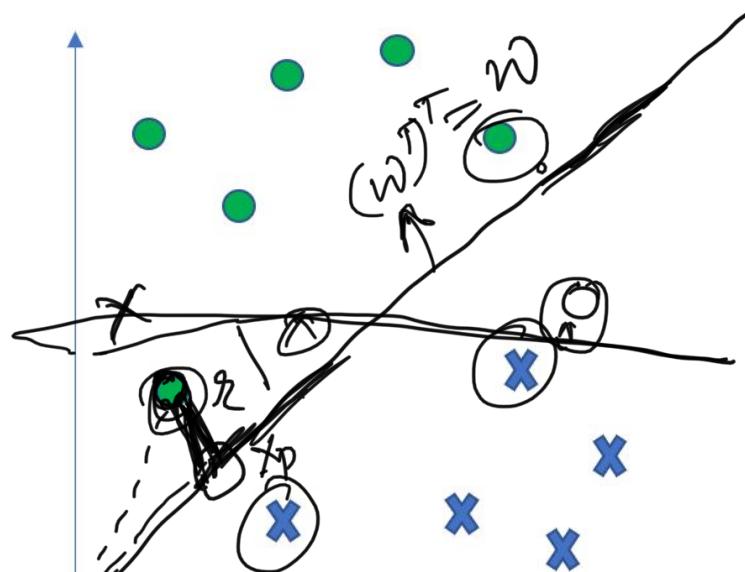
Is the solution unique?



Concept of Margin

$$g(x) = w^T x + w_0 = 0$$

$$\begin{aligned} z &= g(x) \\ \|w\| \end{aligned}$$



$$x = x_p + \gamma \frac{w}{\|w\|}$$

$$x_p = x - \gamma \frac{w}{\|w\|}$$

$$\underline{g(x_p) = 0}$$

$$w^T x - \gamma \frac{w^T w}{\|w\|} + w_0 = 0$$

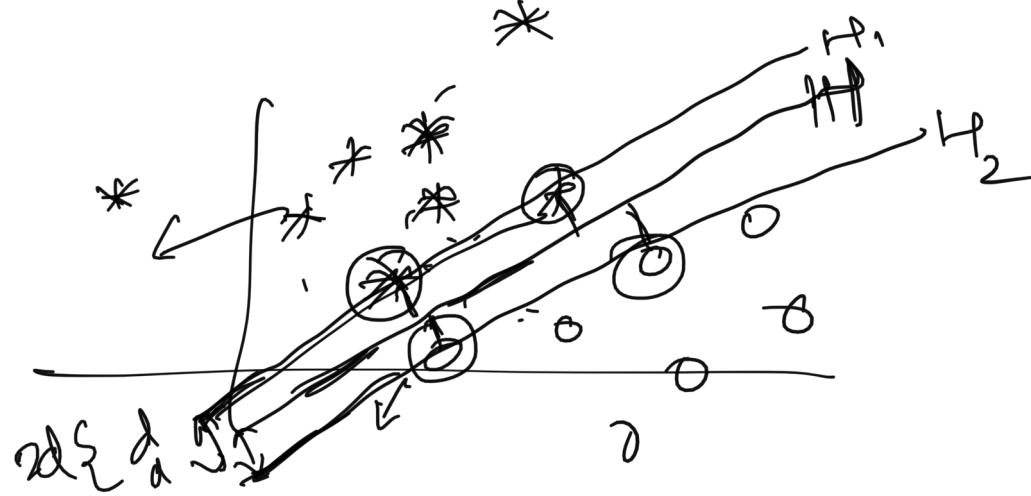
$$g(x) = \frac{w^T}{\|w\|} \left(x - \gamma \frac{w}{\|w\|} \right) + w_0 = 0$$

$$w^T x_p + w_0 = 0$$

$$\left(x - \gamma \frac{w}{\|w\|} \right) + w_0 = 0$$

Concept of Margin

- Margin of a sample x (w.r.t. the decision plane) is defined as the distance from x to the plane.
- For a classifier, the margin should be as large as possible for better performance.

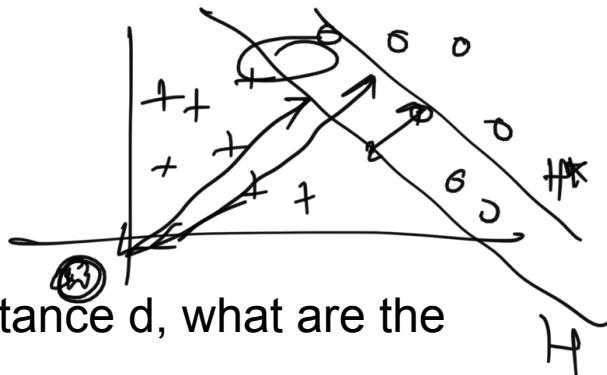


Support Vector Machines (SVM)

⇒ Maximum Margin classifier.

- Key idea - To find the decision boundary such that the margin is maximized.
- Data - $\langle x^{(i)}, y^{(i)} \rangle$, $y^{(i)} \in \{-1, 1\}$, $x^{(i)} \in R^d$, for all $i=1, \dots, n$
- Assume that the data is linearly separable

SVM - Problem Formulation



- Given separating plane $H: \underbrace{w^T x + b}_w = 0$ and distance d , what are the equations for H_1 and H_2 ?
- Consider H^* plane with equation - $H^* = \underbrace{w^T x + b}_{\|w\|d} = \|w\|d$

$$d(\text{origin}, H) = \frac{g(0)}{\|w\|} = \frac{b}{\|w\|}$$

$$d(\text{origin}, H^*) = \frac{g(0)}{\|w\|} = \frac{b - \|w\|d}{\|w\|}$$

$$d(H, H^*) = d$$

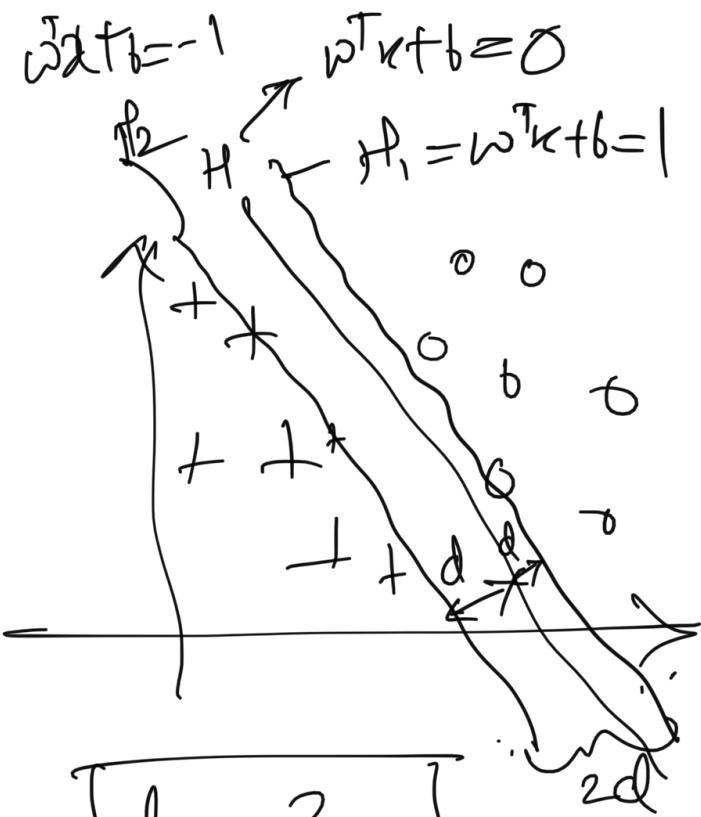
SVM - Problem Formulation

$$H_1 = \underline{\underline{w^T x + b}} = \underline{\underline{\|w\| d}}$$

$$H_2 = \underline{\underline{w^T x + b}} = \underline{\underline{-\|w\| d}}$$

$$\underline{\underline{H = w^T x + b = 0}}$$

$$\left\{ \begin{array}{l} H_1 = \underline{\underline{w^T x + b = 1}} \\ H_2 = \underline{\underline{w^T x + b = -1}} \\ H = \underline{\underline{w^T x + b = 0}} \end{array} \right\}$$



SVM - Problem Formulation

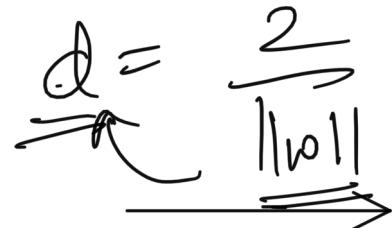
- Plane equations-

$$H_1$$

$$H_2$$

$$H$$

- Margin -

$$d = \frac{2}{\|w\|}$$


$$d = \frac{2}{\|w\|}$$

SVM - Problem Formulation

$$\{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \|\mathbf{w}\| \text{ or } \{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2$$

Subject to

$$\begin{aligned} & \mathbf{w}^T \mathbf{x}^{(i)} + b \geq 1 && \text{for } y^{(i)} = +1 \\ & \mathbf{w}^T \mathbf{x}^{(i)} + b \leq -1 && \text{for } y^{(i)} = -1 \end{aligned} \quad \left. \right\}$$

The constraints can be combined into:

$$\underline{y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1} \geq 0 \quad \forall i$$

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Linear Machines & SVM



Table of contents

- 1. SVM Dual Lagrangian Formulation**
- 2. Kernel Trick**
- 3. Soft-Margin case**

$$d = \frac{2}{\|\omega\|}$$

SVM - Problem Formulation

$$\{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \|\mathbf{w}\| \text{ or } \{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2$$

\Leftarrow

Subject to

$$\mathbf{w}^t \mathbf{x}^{(i)} + b \geq 1 \quad \text{for } y^{(i)} = +1 \quad \underline{\underline{—}}$$

$$\mathbf{w}^t \mathbf{x}^{(i)} + b \leq -1 \quad \text{for } y^{(i)} = -1 \quad \underline{\underline{—}}$$

The constraints can be combined into:

$$\underline{\underline{y^{(i)}(\mathbf{w}^t \mathbf{x}^{(i)} + b) - 1 \geq 0 \quad \forall i}} \quad \underline{\underline{—}}$$

SVM Lagrangian Dual Formulation

$$L(w, b, \alpha) = \frac{1}{2} \underbrace{\|w\|^2}_{w^T w} - \sum_{i=1}^n \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1] \quad \textcircled{1}$$

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \Rightarrow w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \quad \text{Lagrangian multipliers}$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^n \alpha_i y_i \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{--- } \textcircled{3}$$

Plugging $\textcircled{2}$ in $\textcircled{1}$,

$$L(w, b, \alpha) = \frac{1}{2} \sum_{i=1}^n \alpha_i y^{(i)} (x^{(i)})^T \sum_{j=1}^n \alpha_j y^{(j)} (x^{(j)})$$

SVM Lagrangian Dual Formulation

$$-\sum_{i=1}^n \alpha_i y^{(i)} (\underbrace{x^{(i)} \sum_{j=1}^n \alpha_j y^{(j)} x^{(j)}}_{+}) - \sum_{i=1}^n \alpha_i y^{(i)} b =$$

$$+\sum_{i=1}^n \alpha_i.$$

$$= -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (\underbrace{x^{(i)} \top x^{(j)}}) + \sum_{i=1}^n \alpha_i.$$

$$L(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{\langle x^{(i)}, x^{(j)} \rangle}_{< x^{(i)}, x^{(j)} >}$$

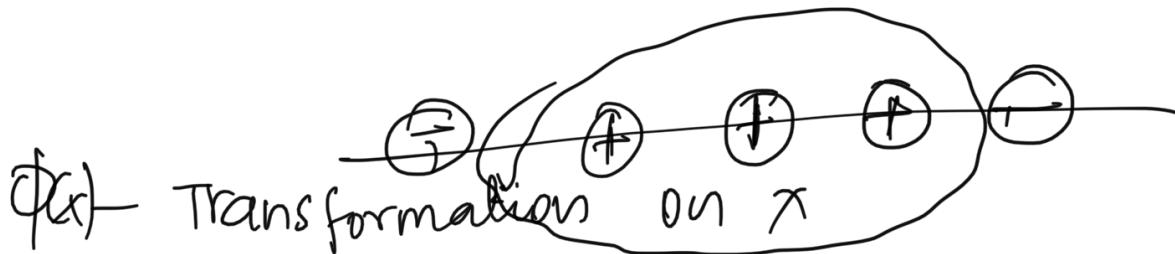
$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = x_1 y_1 + x_2 y_2 + x_3 y_3 \\ = \bar{x} \cdot \bar{y} \cdot \cos \theta$$

What if data is not linearly separable?

$$2 \rightarrow 4$$

$$4 \rightarrow 16$$

$$8 \rightarrow 64$$



$$L(\omega, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

$$K(x^{(i)}, x^{(j)})$$

similarity between
data points

What if data is not linearly separable?

- Kernel trick!

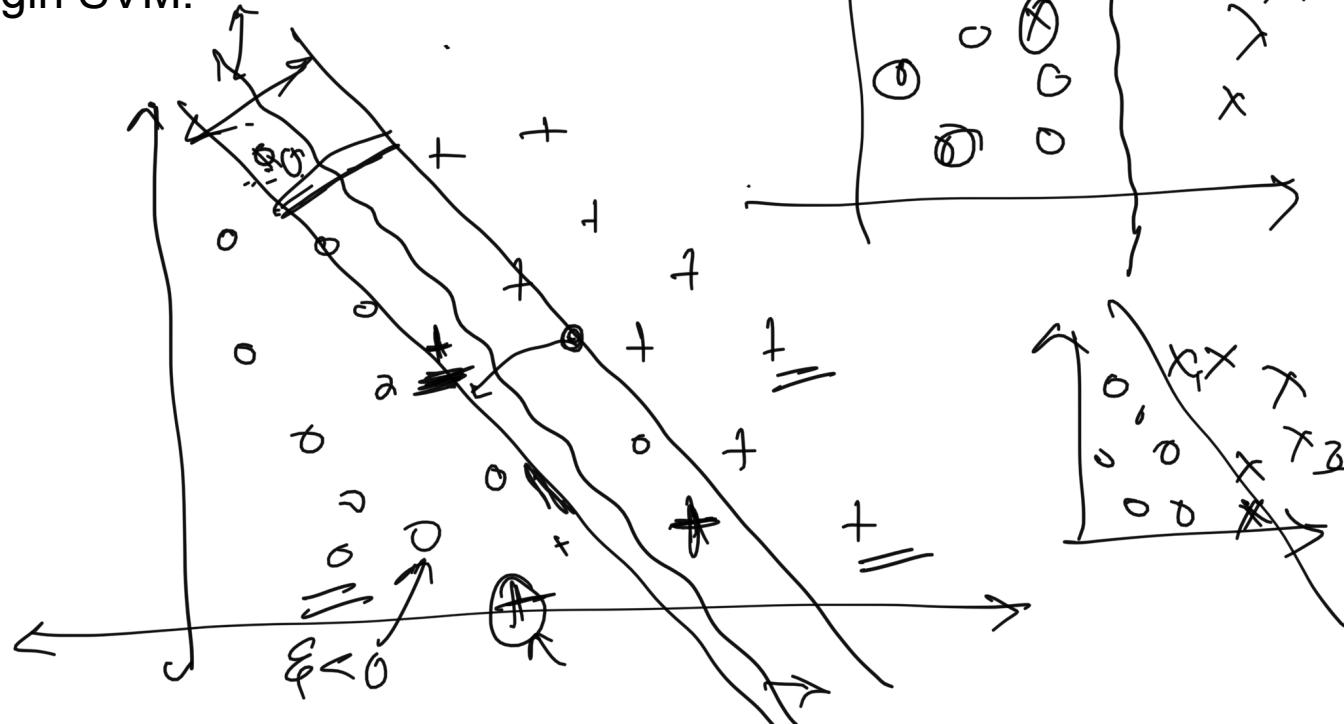
$$L(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} K(\underline{\underline{x}}^{(i)}, \underline{\underline{x}}^{(j)})$$

↳ Reduces computation

↳ Don't have to find the transformation of

What if the data is still not separable?

- Soft-margin SVM!



$$\xi = 0.8$$

Soft-Margin SVM Formulation

→ Hard-margin SVM

$$\boxed{\{w^*, b^*\} = \operatorname{argmin}_{w,b} \frac{1}{2} \|w\|^2 + C(\sum_i \xi_i)}$$

~~$\xi_i = 0$~~

$\xi_i = 0$
- Correctly
classified
samples

subject to

$$0 < \xi_i < 1$$

$$w^t x^{(i)} + b \geq 1 - \xi_i \text{ for } y^{(i)} = +1$$

$$w^t x^{(i)} + b \leq -1 + \xi_i \text{ for } y^{(i)} = -1$$

$$\xi_i \geq 0, \forall i$$

~~$\xi_i \geq 0$~~

- On the correct
side of decision
bound but
within the margin

$\xi_i = 0$ (Support
vectors)

$\xi_i > 1$ completely
misclassified
samples.

Effect of C parameter

$C \Rightarrow \infty$ \Rightarrow Hard-margin SVM

$C = 0$ \Rightarrow ~~Decision boundary with~~ very high misclassification.

Will SVM overfit?

↳ Robust to overfitting.

How to overcome overfitting?

- Tune C parameter.
- Tune the parameters used in kernels.
- Change the kernels

$$y = \frac{1}{1 + e^{w^T x}} \quad 0-1$$

Difference between SVM and Logistic Regression

Support Vector Machines	Logistic Regression
<ul style="list-style-type: none"> ① Margin, no probabilities ② Sparse solution ③ Robust to overfitting ④ Hinge loss 	<ul style="list-style-type: none"> ① Probabilities ② Not sparse ③ Can overfit ④ Log likelihood loss

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

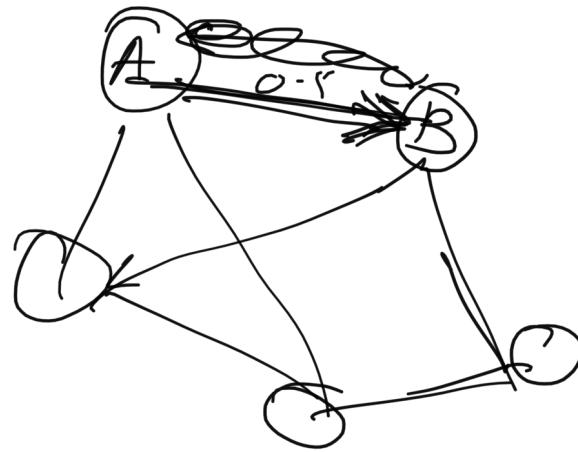
Instructor: Nupur Thakur

Graphical Models



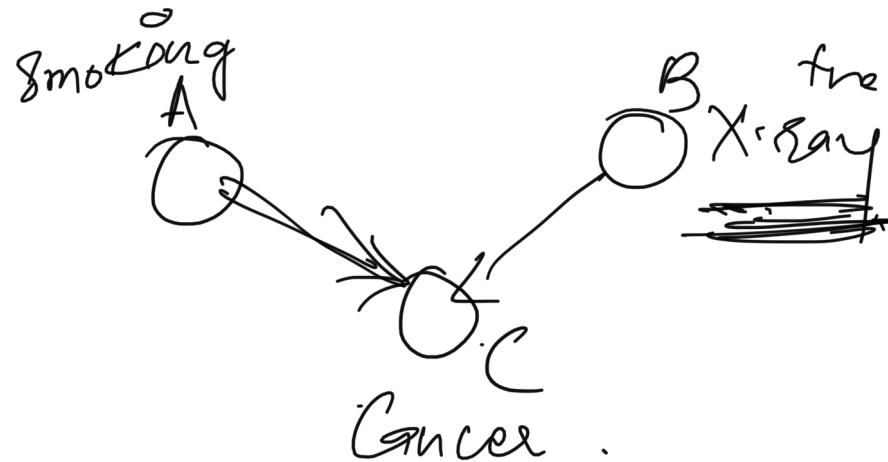
Table of contents

1. Bayesian Networks
2. Hidden Markov Models

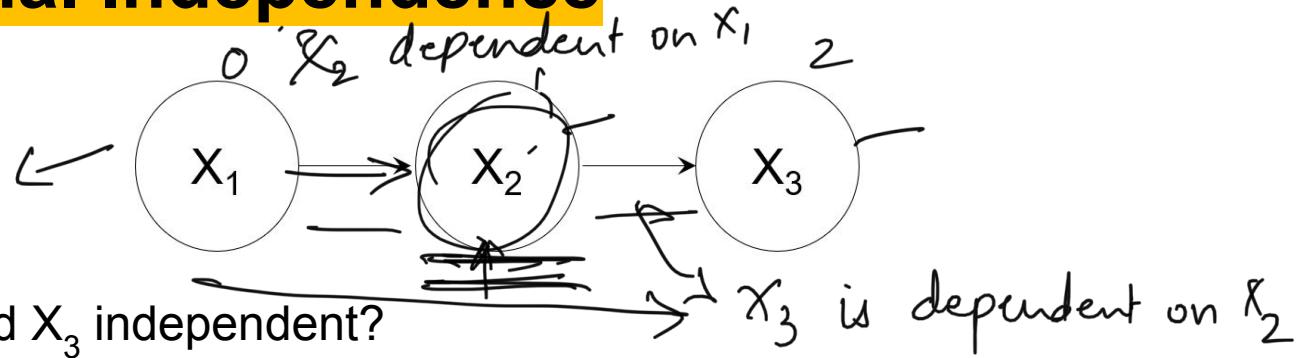


Bayesian Networks

- A probabilistic graphical model representing a set of variables and their conditional dependencies via a directed acyclic graph.

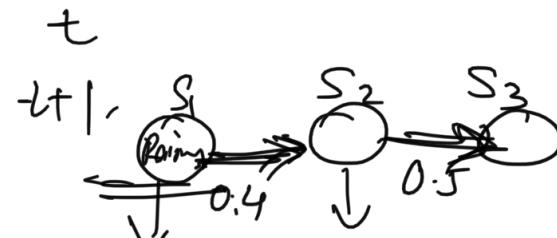


Conditional Independence



- Are X_1 and X_3 independent?
 - No.
- Are they conditionally independent given X_2 ?
 - Yes.

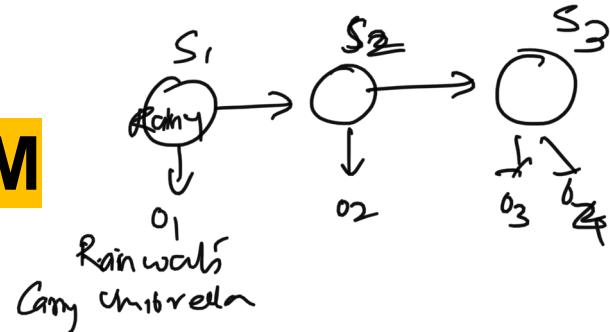
Hidden Markov Model (HMM)



- Dynamic Bayesian network - modeling the process indexed by time.
 - [Rainy days]
- Two assumptions -
 - First-order Markov chain - $P(s^t = S_j | s^{t-1} = S_i) —$
 - $a_{ij} = P(s^t = S_j | s^{t-1} = S_i)$, $1 \leq i, j \leq N$, for any t — Stationary.

$$P(S_3 | S_2) = 0.5$$

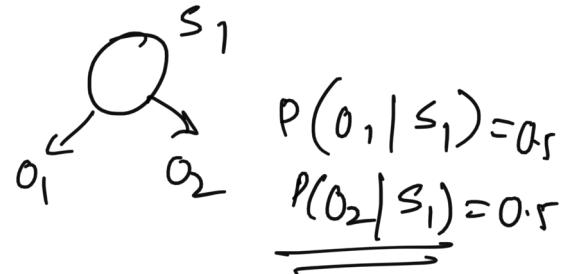
Specifying HMM



- Θ - set of hidden states $= \{s_1, s_2, s_3\}$
- Ω - set of output observations $= \{o_1, o_2, o_3, o_4\}$
- π - initial state distribution

$$\circ \quad \pi = \{\pi_i\}, \pi_i = P(s^1 = s_i) \quad = \underline{s_2}$$
$$\underline{\{0, 1, 0\}} \quad \underline{\underline{P(s'_1 = s_2) = 1}}$$

Specifying HMM



- State transition matrix - A — Shape of A = $n \times n$
 - $A = \{a_{ij}\}$ where $a_{ij} = P(s^t = S_j | s^{t-1} = S_i)$ for $1 \leq i, j \leq N$
- Observation/emission probability matrix B
 - $B = \{b_{jk}\} = P(o^t = v_k \text{ at } t | s^t = S_j)$ where v_k is the k^{th} symbol in Ω

$n \rightarrow$ states $k \rightarrow$ observations

$$S_1 \begin{bmatrix} o_1 & o_2 \\ 0.5 & 0.5 \end{bmatrix} K$$

Problems in HMM

- For a given HMM $\Lambda = \{\Theta, \Omega, A, B, \pi\}$

state transition matrix
emission prob. matrix

- Estimation of model parameters — A, B
- Given an observation sequence $O = \{o^1, o^2, \dots, o^k\}$, what is the most likely state sequence $S = \{s^1, s^2, \dots, s^k\}$ that has produced O ? — Decoding.
- How likely is an observation O ? — likelihood.

HMM Parameter Estimation

$$S_1(O_1) \rightarrow S_2(O_2)$$

- Given labeled data - state and observation

$$\overbrace{S_1(O_1) \rightarrow S_2(O_2) \rightarrow S_3(O_1)}^{\text{state sequence}} \quad -$$

$a^{ij}(S_i|S_j) = \frac{\text{number of } (s^t = S_i, s^{t-1} = S_j)}{\text{number of } S_j}$

$\text{transition probability}$

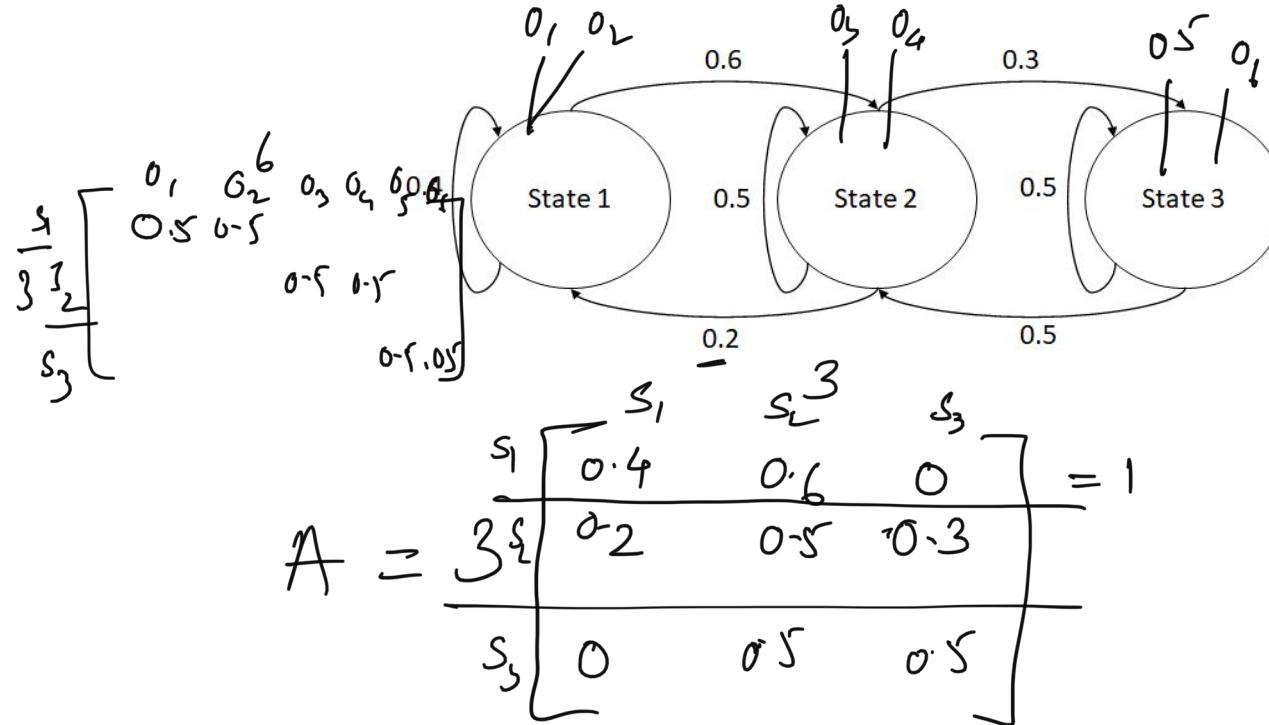
$$= b^{or}(o_r|S_j) = \frac{\text{number of } (o^t = o_r, s^t = S_j)}{\text{number of } S_j}$$

$\text{emission probability}$

Ques. Given data - sequence of observations.

Forward-Backward Algorithm.

HMM Parameter Estimation Example



Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Graphical Models



Table of contents

- 1. Problems in HMM - Parameter Estimation**
- 2. Problems in HMM - State Estimation**
- 3. Problems in HMM - Evaluate P(O)**

Problems in HMM

- For a given HMM $\Lambda = \underline{\{\Theta, \Omega, A, B, \pi\}}$
 - Estimation of model parameters
 - Given an observation sequence $O = \{o^1, o^2, \dots, o^k\}$, what is the most likely state sequence $S = \{s^1, s^2, \dots, s^k\}$ that has produced O ?
 -
 - How likely is an observation O ?

HMM Parameter Estimation

- Given labeled data - state and observation

A

$$t(S_i|S_j) = \frac{\text{number of } (s^t = S_i, s^{t-1} = S_j)}{\text{number of } S_j}$$

==

B

$$e(o_r|S_j) = \frac{\text{number of } (o^t = o_r, s^t = S_j)}{\text{number of } S_j}$$

==

- Given observation sequences -

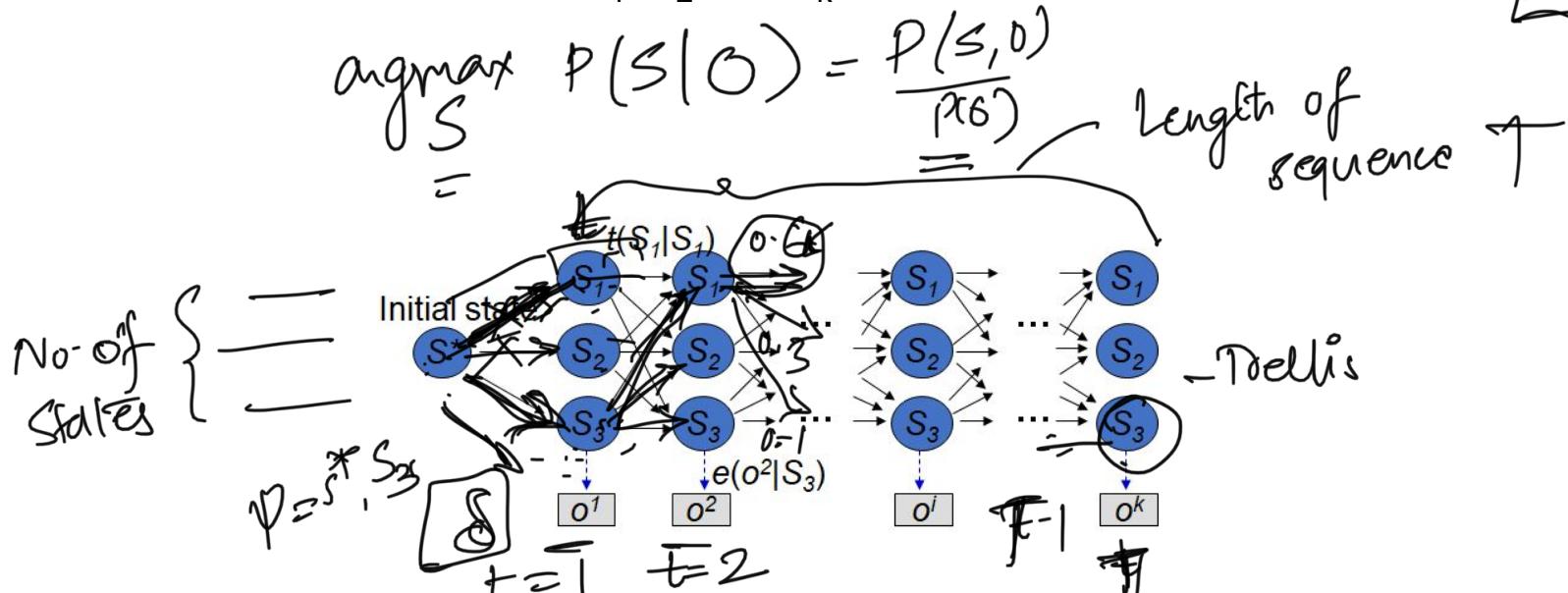
Forward-Backward Algorithm - Expectation-Maximization Approach.

HMM State Estimation

Decoding

$$\begin{aligned} P(S^* | O) &= \frac{P(O_1 | S_1) P(S_1 | S^*)}{P(O_1)} \\ &= \frac{P(O_1 | S_1) P(S_1 | S^*)}{P(O_1 | S^*)} \\ &= P(S^* | O) \end{aligned}$$

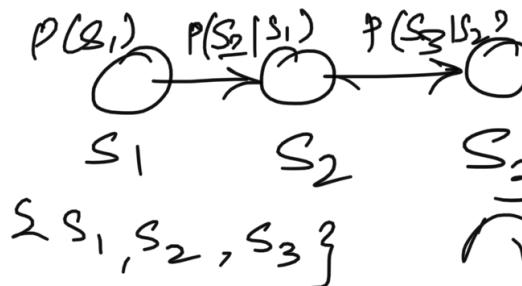
- Given an observation (sequence) $O = \{o_1, o_2, \dots, o_k\}$, what is the most likely state sequence $S = \{s_1, s_2, \dots, s_k\}$ that has produced O ?



HMM State Estimation

$$P(S, O) = P(O|S) \cdot P(S)$$

$$\begin{aligned}
 &= P(O^1 \dots O^K | S_1 \dots S_K) \\
 &= \prod_{i=1}^K P(O^i | O^1 \dots O^{i-1}, S^1 \dots S^{i-1}) \prod_{j=1}^K P(S^j | S^1 \dots S^{j-1}) \\
 &= \prod_{i=1}^K P(O^i | S^i) \quad \text{Emission probability} \quad \prod_{j=1}^K P(S^j | S^{j-1}), \\
 P(S, O) &= \prod_{i=1}^K P(O^i | S^i) \prod_{j=1}^K P(S^j | S^{j-1}) \quad \text{Transition probability}
 \end{aligned}$$



HMM State Estimation - Viterbi Algorithm

Initialization

$$\delta_{S_i}(1) = \underbrace{t(S_i|s^*)}_{\leftarrow} e(o^1|S_i), \quad \forall S_i \in \Theta \quad \leftarrow$$

Induction:

for $2 \leq t \leq k$, do

$$\begin{aligned}\delta_{S_i}(t) &= \max_{S_j} \underbrace{t(S_i|S_j)}_{\leftarrow} e(o^t|S_i) \underbrace{\delta_{S_j}(t-1)}_{\leftarrow} \\ \psi_{S_i}(t) &= \operatorname{argmax}_{S_j} \underbrace{t(S_i|S_j)}_{\leftarrow} e(o^t|S_i) \underbrace{\delta_{S_j}(t-1)}_{\leftarrow}\end{aligned}$$

Termination:

– The probability of the best state sequence: $\max_{S_j} \underbrace{\delta_{S_j}(k)}_{\leftarrow}$

– The best last state: $\hat{s}^k = \operatorname{argmax}_{S_j} \underbrace{\delta_{S_j}(k)}_{\leftarrow}$

– Back trace to get other states:

$$\hat{s}^t = \underbrace{\psi_{\hat{s}^{t+1}}(t)}_{\leftarrow}, \text{ for } t = k-1, \dots, 1.$$

HMM - Evaluate P(O)

- How likely is an observation O?

$$O = \{o_1, \dots, o_k\}$$

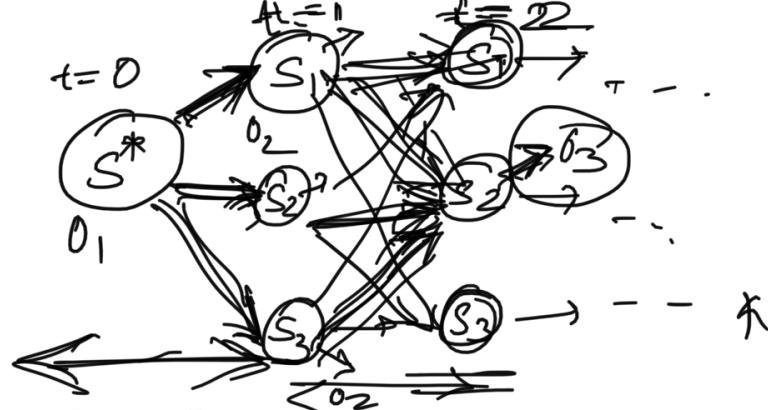
$$P(O) = \sum_{S} \overline{P(S, O)}$$

- forward Algorithm — $\alpha^{(t-1)}$ at t
- Backward Algorithm — $\beta^{(t+1)}$ to T
at t^{th} step

HMM - Evaluate P(O)

- Forward Algorithm

Initialization: $\alpha_{S_i}(1) = t(S_i|s^*)e(o^1|S_i), \forall S_i \in \Theta$

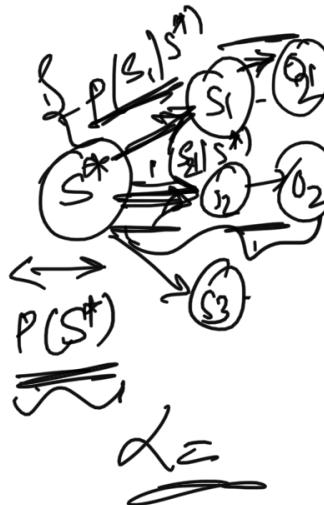


Induction:
for $2 \leq t \leq k$, do

$$\alpha_{S_i}(t) = \sum_{S_j} t(S_i|S_j)e(o^t|S_i)\alpha_{S_j}(t-1)$$

Termination:

$$P(O) = \sum_{S_j} \alpha_{S_j}(k)$$



Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Midterm Review

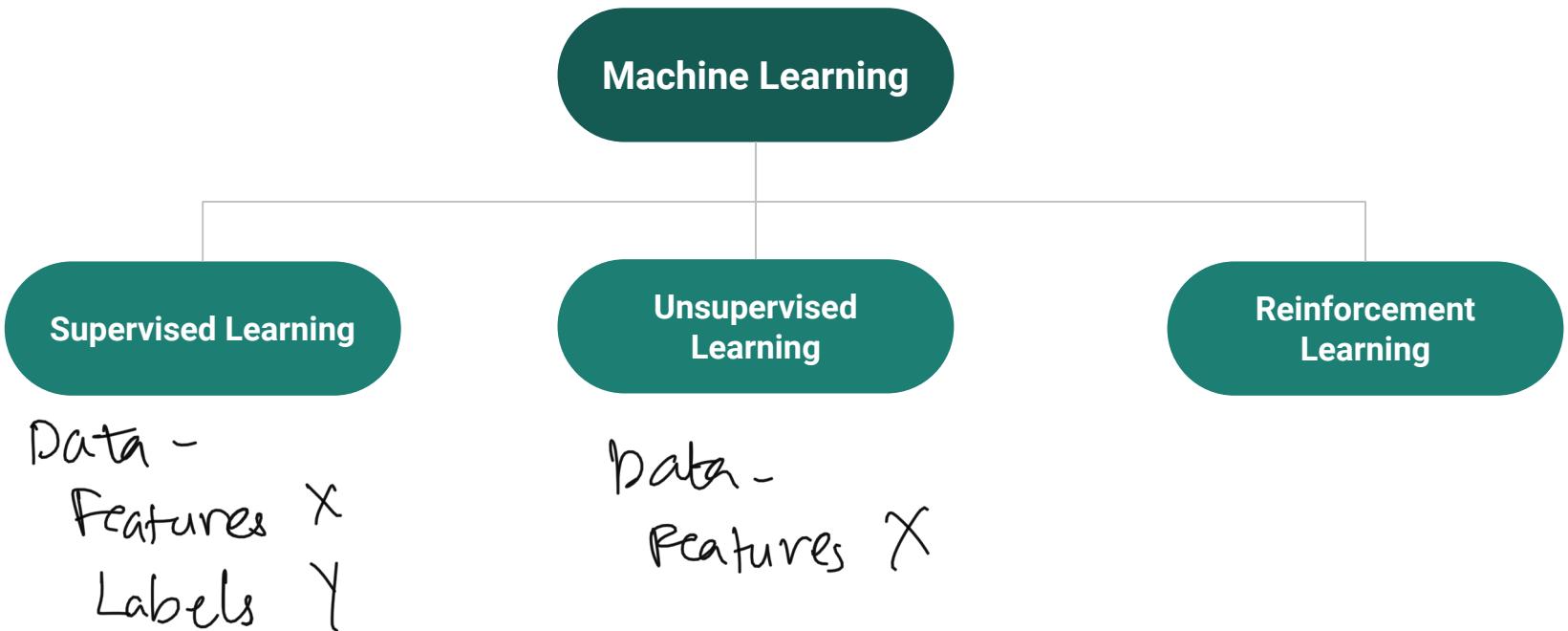


Table of contents

1. Concepts Review

2. Sample Questions

Types of Machine Learning



$$\text{Bias} - \underbrace{E[\hat{\theta}]}_{\text{Estimated Parameter}} - \theta = 0 \quad E[\hat{\theta}] = \theta$$

Maximum Likelihood Estimation

$$\underline{\underline{Ex}} - \underline{\underline{E[\hat{\mu}]}} = \mu \Rightarrow \text{Unbiased}$$

$$\underline{\underline{\hat{\mu}}} \neq \mu \Rightarrow \text{biased.}$$

$E[\mu] = \int x p(x) dx$ unbiased

- Given some training data and assuming a parametric model $\underline{\underline{p(x|\theta)}}$; what specific θ will fit/explain the data best?
- To consider all the samples denoted by $D = \{x_1, x_2, \dots, x_n\}$, assume that all the samples are i.i.d - independent and identically distributed.
- So, data likelihood represented by $L(\theta)$ is -

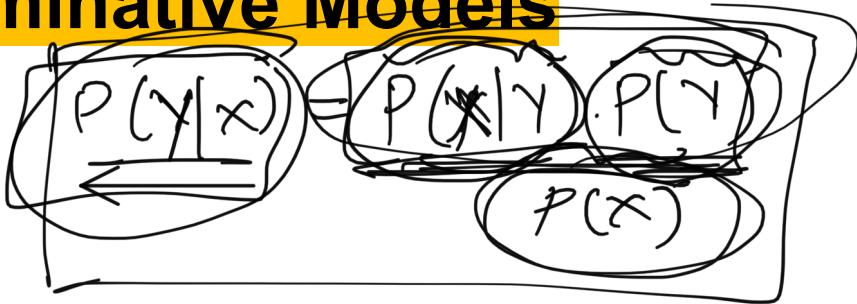
$$L(\theta) = P(D|\theta) = \prod_i P(x_i|\theta)$$

$$\boxed{\hat{\theta} = \operatorname{argmax}_{\theta} p(x|\theta)}$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

Generative vs Discriminative Models

- Generative models -
 - Learn $P(y)$ and $P(x|y)$.
 - Ex: Bayesian classifier, Naive Bayes.
- Discriminative models -
 - Directly learn $P(y|x)$
 - Ex: Logistic Regression



$$P(y|x)$$

Naive Bayes

Supervised

$$\mathbf{x} = \{x_1, \dots, x_d\}$$
$$y = 1$$

- The "naive" conditional independence assumption: each feature is (conditionally) independent of every other feature, given the label, i.e.,
 $p(x_i | \{x_j \text{ for any } j \neq i\}, y) = p(x_i | y)$
- The predicted label is given by -

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(y) \prod_{i=1}^d p(x_i | y)$$

Predicted label.

Linear Regression

- Regression - A training set of n samples $\langle \underline{x}^{(i)}, \underline{y}^{(i)} \rangle$ where $\underline{y}^{(i)}$ is a continuous “label” (or target value) for $\underline{x}^{(i)}$
- Linear regression - modeling the relation between y and x via a linear function

$$y \approx w_0 + w_1 x_1 + \dots + w_d x_d = \mathbf{w}^t \mathbf{x}$$

- The error is given as - $\|e\|^2 = \|y - \underline{\mathbf{X}^t \mathbf{w}}\|^2$

Logistic Regression

$$P(y|x) = \text{logistic function} \\ w^T x \\ \equiv \\ w^T x \geq 0$$

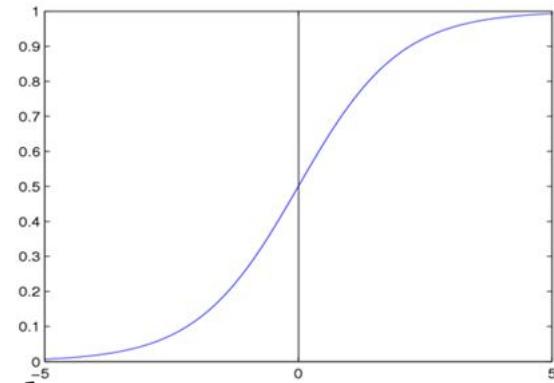
- Training set: n labelled samples $\langle \underline{x}(i), \underline{y}(i) \rangle$
- Use the logistic function for modeling $P(y|x)$, considering only the case of $y \in \{0,1\}$

$$P(y=0|x) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}$$

$$P(y=1|x) = \frac{\exp(w_0 + \sum_{i=1}^d w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}$$

→ Gradient ascent

$$w_{t+1} = w_t + \alpha \frac{\partial L(\omega)}{\partial \omega}$$



$$\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t}$$

Support Vector Machines (SVM)



- Key idea - To find the decision boundary such that the margin is maximized.
- Data - $\langle x^{(i)}, y^{(i)} \rangle$, $y^{(i)} \in \{-1, 1\}$, $x^{(i)} \in R^d$, for all $i=1, \dots, n$
- Plane equations-

$$\begin{cases} w^T x + b = 1 \\ w^T x + b = -1 \end{cases} \quad \text{Margin hyperplane.}$$

- Margin - $w^T x + b = 0$ — Decision boundary.

$$d = \frac{2}{\|w\|}$$

SVM - Problem Formulation

$$\{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \|\mathbf{w}\| \text{ or } \{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2$$

Subject to

$$\begin{cases} \mathbf{w}^t \mathbf{x}^{(i)} + b \geq 1 & \text{for } y^{(i)} = +1 \\ \mathbf{w}^t \mathbf{x}^{(i)} + b \leq -1 & \text{for } y^{(i)} = -1 \end{cases}$$

The constraints can be combined into:

$$y^{(i)}(\mathbf{w}^t \mathbf{x}^{(i)} + b) - 1 \geq 0 \quad \forall i$$

Soft-Margin SVM Formulation

subject to

$$\{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2 + C(\sum \xi_i)$$

Parameter C controls the penalty.

slack variable

Very high values \rightarrow Hard margin SVM

$$\mathbf{w}^t \mathbf{x}^{(i)} + b \geq 1 - \xi_i \text{ for } y^{(i)} = +1$$

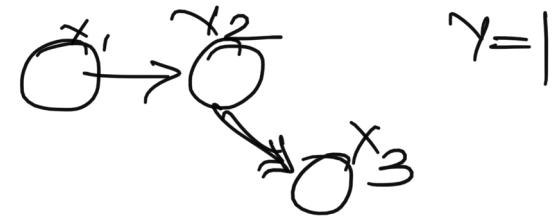
$$\mathbf{w}^t \mathbf{x}^{(i)} + b \leq -1 + \xi_i \text{ for } y^{(i)} = -1$$

$$\xi_i \geq 0, \forall i$$

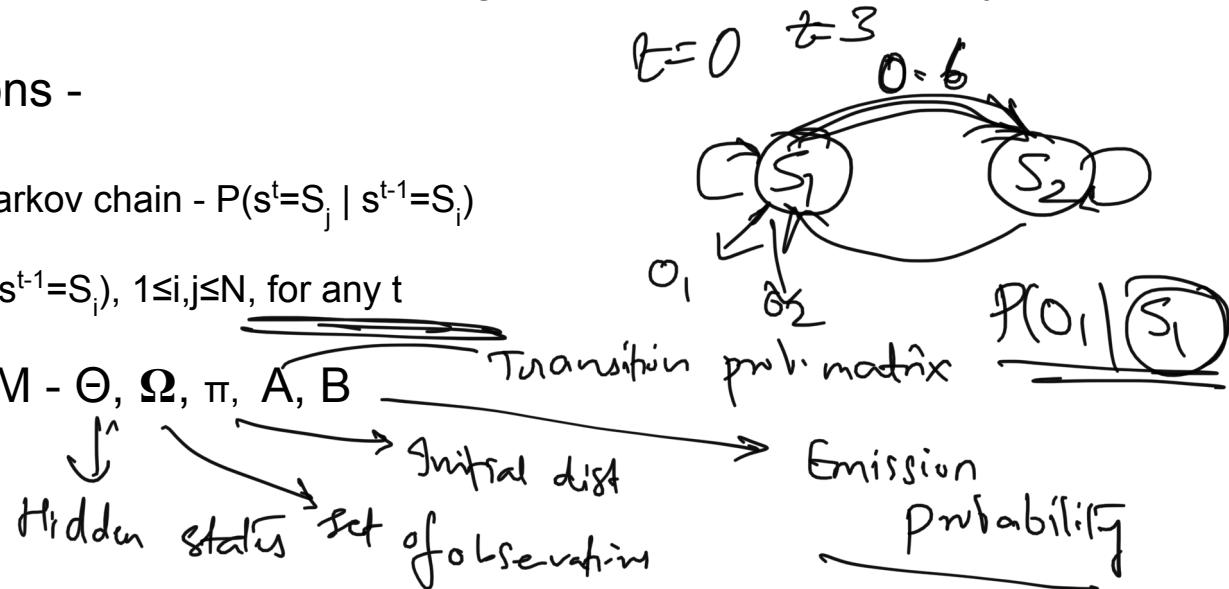
Very low values \rightarrow Larger misclassifications

$$\phi(x_i - x_j)$$

Hidden Markov Model (HMM)



- Dynamic Bayesian network - modeling the process indexed by time.
- Two assumptions -
 - First-order Markov chain - $P(s^t=S_j | s^{t-1}=S_i)$
 - $a_{ij} = P(s^t=S_j | s^{t-1}=S_i)$, $1 \leq i, j \leq N$, for any t
- Specifying HMM - $\Theta, \Omega, \pi, A, B$



Problems in HMM

- For a given HMM $\Lambda = \{\Theta, \Omega, A, B, \pi\}$
 - Estimation of model parameters
 - Given an observation sequence $O = \{o^1, o^2, \dots, o^k\}$, what is the most likely state sequence $S = \{s^1, s^2, \dots, s^k\}$ that has produced O ? $\xrightarrow{\text{Decoding}}$
 - How likely is an observation O ? $\xrightarrow{P(O)}$

Trained
HMM

HMM Parameter Estimation

- Given labeled data - state and observation

$$\text{A} \rightarrow t(S_i|S_j) = \frac{\text{number of } (s^t = S_i, s^{t-1} = S_j)}{\text{number of } S_j}$$

$$\text{B} \rightarrow e(o_r|S_j) = \frac{\text{number of } (o^t = o_r, s^t = S_j)}{\text{number of } S_j}$$

- Observation sequence — Forward-Backward Algo.

HMM State Estimation

- Given an observation (sequence) $O = \{o_1, o_2, \dots, o_k\}$, what is the most likely state sequence $S = \{s_1, s_2, \dots, s_k\}$ that has produced O ?

Viterbi

Initialization

$$\delta_{S_i}(1) = t(S_i|s^*)e(o^1|S_i), \quad \forall S_i \in \Theta$$

Induction:

for $2 \leq t \leq k$, do

$$\delta_{S_i}(t) = \max_{S_j} t(S_i|S_j)e(o^t|S_i)\delta_{S_j}(t-1)$$

$$\psi_{S_i}(t) = \operatorname{argmax}_{S_j} t(S_i|S_j)e(o^t|S_i)\delta_{S_j}(t-1)$$

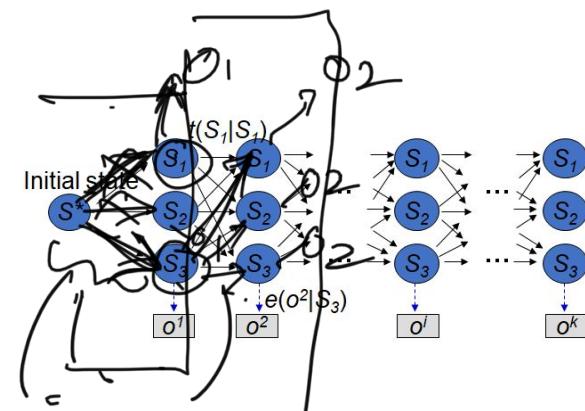
Termination:

- The probability of the best state sequence: $\max_{S_j} \delta_{S_j}(k)$

- The best last state: $\hat{s}^k = \operatorname{argmax}_{S_j} \delta_{S_j}(k)$

- Back trace to get other states:

$$\hat{s}^t = \psi_{\hat{s}^{t+1}}(t), \text{ for } t = k-1, \dots, 1.$$



s^*, s_3, s_1

HMM - Evaluate P(O)

- How likely is an observation O?

$$P(O) = \sum_s P(S, O)$$

s

- Forward algorithm-



Initialization: $\alpha_{S_i}(1) = t(S_i|s^*)e(o^1|S_i), \quad \forall S_i \in \Theta$

Induction:

for $2 \leq t \leq k$, do

$$\alpha_{S_i}(t) = \sum_{S_j} t(S_i|S_j)e(o^t|S_i)\alpha_{S_j}(t-1)$$

Termination:

$$P(O) = \sum_{S_j} \alpha_{S_j}(k)$$

Sample Questions

Problem 1:

$$x_1 = 0 \quad y = 1$$

$$x_1 = 1 \quad y = 0$$

$$\{ P(x_1=1|y=1) \} \\ P(x_1=1|y=0)$$

$$1 - [P(x_1=0|y=1)] \quad P(x_1=0|y=0)$$

The following data is used for training Naive Bayes binary classifier. The last column is the binary class label; Each of the first 4 columns is a binary feature, and each row is a training example.

X1	X2	X3	X4	Y
1	0	0	0	1
0	1	1	0	1
1	0	0	1	0
0	1	1	1	0
1	1	1	0	1

$$P(x_2=0|y=1) = \frac{2}{5}$$
$$P(y=0) = 2/5$$
$$P(x_1=1|y=1) = 1/5$$
$$P(x_1=0|y=1) = 2/5$$
$$P(x_1=0|y=0) = 4/5$$

Problem 2:

State true/false along with justification:

1. Logistic Regression may give us a non-linear classifier, depending on how the training examples are distributed in the feature space.

Ans - False because logistic regression is a linear classifier, producing a linear decision boundary.

Problem 3:

Multiple Choice Questions:

Which of the following statements is true about HMM?

- a. Given some initial state and transition matrix, it is possible that there exists a state that you can never achieve.
- b. The elements of a column of the state transition matrix always sum to 1.
- c. The sum of a row of the transition matrix may not be 1.

Ans ~ (9)

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Unsupervised Learning



Table of contents

- 1. K-means clustering**
- 2. Project Part- 2**

What is unsupervised learning?

- Given a training set of n unlabeled samples - $x(i)$

- What can we learn from the samples?

EM algorithm

- Estimate the overall distribution of the data without knowing their label.



- Figure out the groupings of the samples (if any).

— Clustering (K-means)

- Identify some features that may be more important than others.

Feature selection

Finding clusters

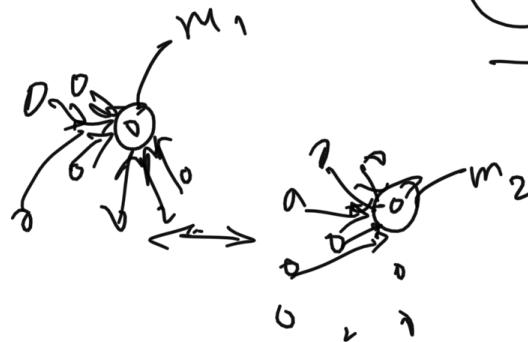
- How to represent the clusters?
 - Centroid (mean of all samples)
- Which cluster a sample should be assigned to (e.g., membership)?
 - Some similarity measure like Euclidean dist.
- What similarity measure to use?
 - Different measures
 - Euclidean dist.
 - Cosine similarity.

Clustering Objective function

- The sum-of-squared-error criterion/cost-

$$J_e = \sum_{i=1}^C \sum_{x \in D_i} \|x - m_i\|^2$$

(centroid of
ith cluster.)


$$m_i = \frac{1}{n_i} \sum_{x \in D_i} x$$

Reduce the intra-cluster distance
increase the inter-cluster distance.

K-Means Clustering

Given: n samples, a number k .

No. of clusters .

Begin

initialize $\underline{\mu}_1, \underline{\mu}_2, \dots, \underline{\mu}_k$ (randomly selected)

do classify n samples according to
nearest μ_i \Rightarrow Membership

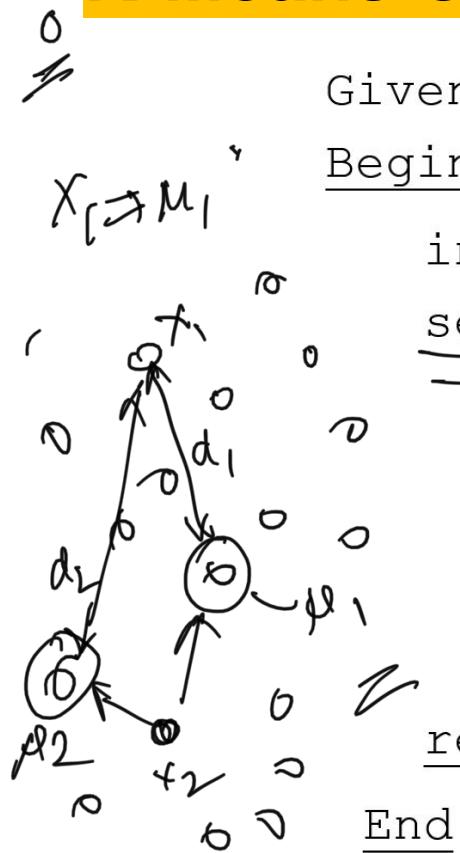
recompute μ_i

Assignment

until no change in μ_i

return $\mu_1, \mu_2, \dots, \mu_k$

End



K-Means Clustering - Some scenarios

- What happens if the distance of a point is same from more than one centroid?

Randomly assigned one of the clusters.

- Hand assignment

- What happens in case of outliers?

Sensitive to outliers.

K-Means Clustering - Some scenarios

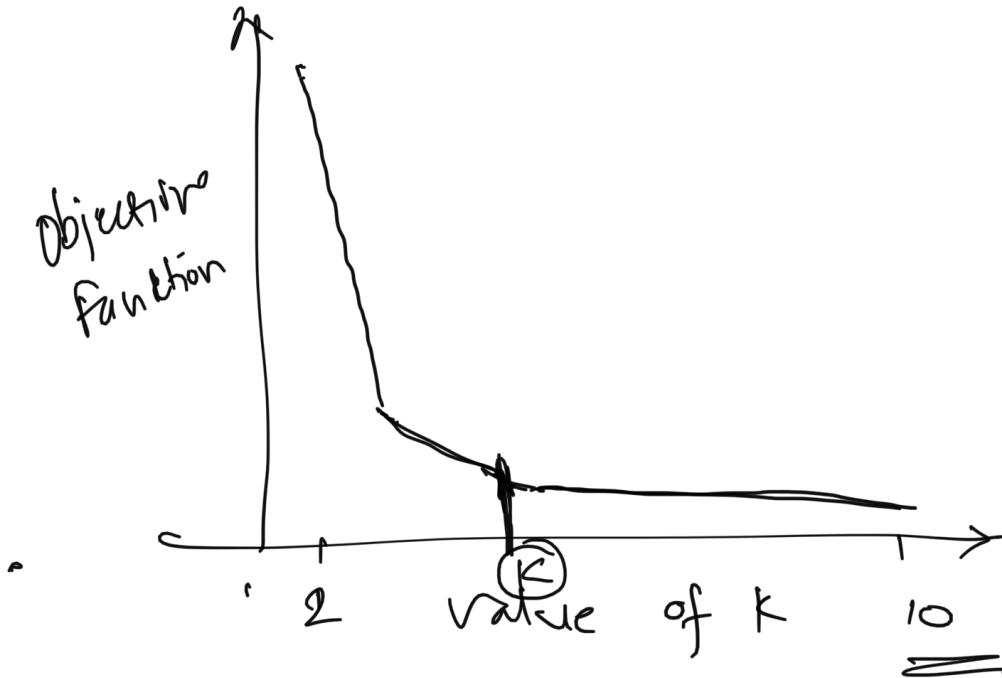
- Different initializations?
 - Might lead to sub-optimal results.
- Different number of clusters than k?

How to choose the k?

- Elbow method.

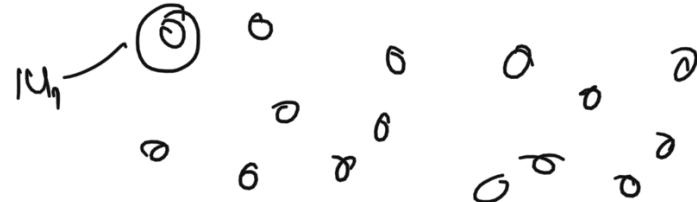
$$\begin{array}{l} 10,000 \\ \underline{\underline{k=1}} \\ 1-9999 \\ K=n \end{array}$$

$$1 \leq K < n$$



How to deal with sensitivity to initialization problem?

- ① Run the algo for few epochs
= & take the average result
- ② first centroid randomly. → Sensitive to outliers.
- ③ Kmeans++ \Rightarrow Smart Initialization
=



K-Means Clustering - Pros and Cons

Pros

- ① Easy to implement
- ② fast.

Cons

- ① Sensitive to initialization
- ② " " outliers.
- ③ Hard Assignment of data points -
- ④ Choose no. of clusters manually.

Project Part 2 - Unsupervised Learning (K-means)

1st time - $k=2 \Rightarrow$ Obj.
 $\overbrace{\quad\quad\quad}^{k=3}$
 \Leftrightarrow

$k=2, 3, 4, 5, 6, 7, 8, 9, 10$ $k=10$

- Apply k-means algorithm the given dataset of 2-D points.

- Initialization strategies:

\nearrow 2 times - $\overbrace{\quad\quad\quad}^{k=2-10}$ Obj. ↑
Value

- Strategy 1: randomly pick the initial centers from the given samples.
- Strategy 2: pick the first center randomly; for the i-th center ($i > 1$), choose a sample
(among all possible samples) such that the average distance of this chosen one to all previous ($i-1$) centers is maximal.

\nearrow 2 times. $\rightarrow k=2-10$

Total - 4 plots

Project Part 2 - Unsupervised Learning (K-means)

- Objective function-

$$\sum_{i=1}^k \sum_{x \in D_i} \|x - \mu_i\|^2 \rightarrow SSE$$

- Number of clusters (k) - from 2-10
- Plot the objective function value vs. the number of clusters k.

Project Part 2 - Deliverables

— Due on 26 March
(Arizona time) 11:59 PM

- Well-commented code
 - Python / MATLAB / Jupyter notebook (.ipynt)
- A report that summarizes the results and includes all the plots.
 - .PDF format .(typed, not handwritten)
- Please submit the code and the report as separate files on Canvas. **Do**

not zip them.

· py
· py
· pdf

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Unsupervised Learning



Table of contents

1. Gaussian Mixture Models

2. Expectation-Maximization Algorithm

Project - strategy 2 $K=3$ $\sum \|x - \tilde{y}_i\|^2$

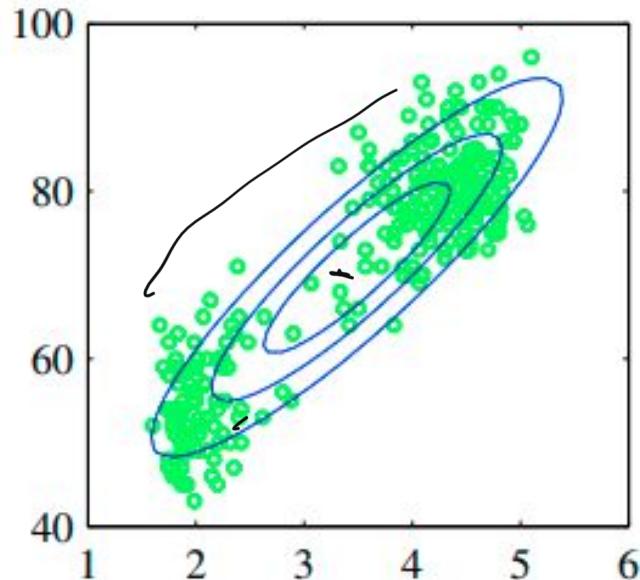
1st center - randomly initialize

2nd center - max. dist from 1st center.

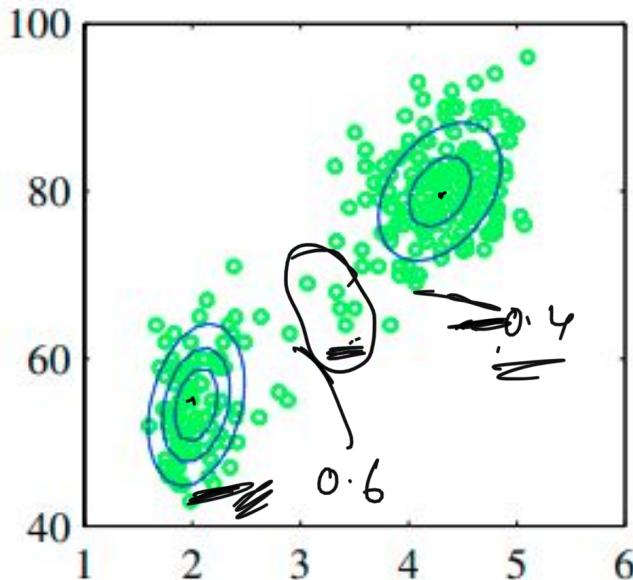
3rd center - max. arg. dist from 1st & 2nd center

Which distribution better represents the given data?

{ Gaussian

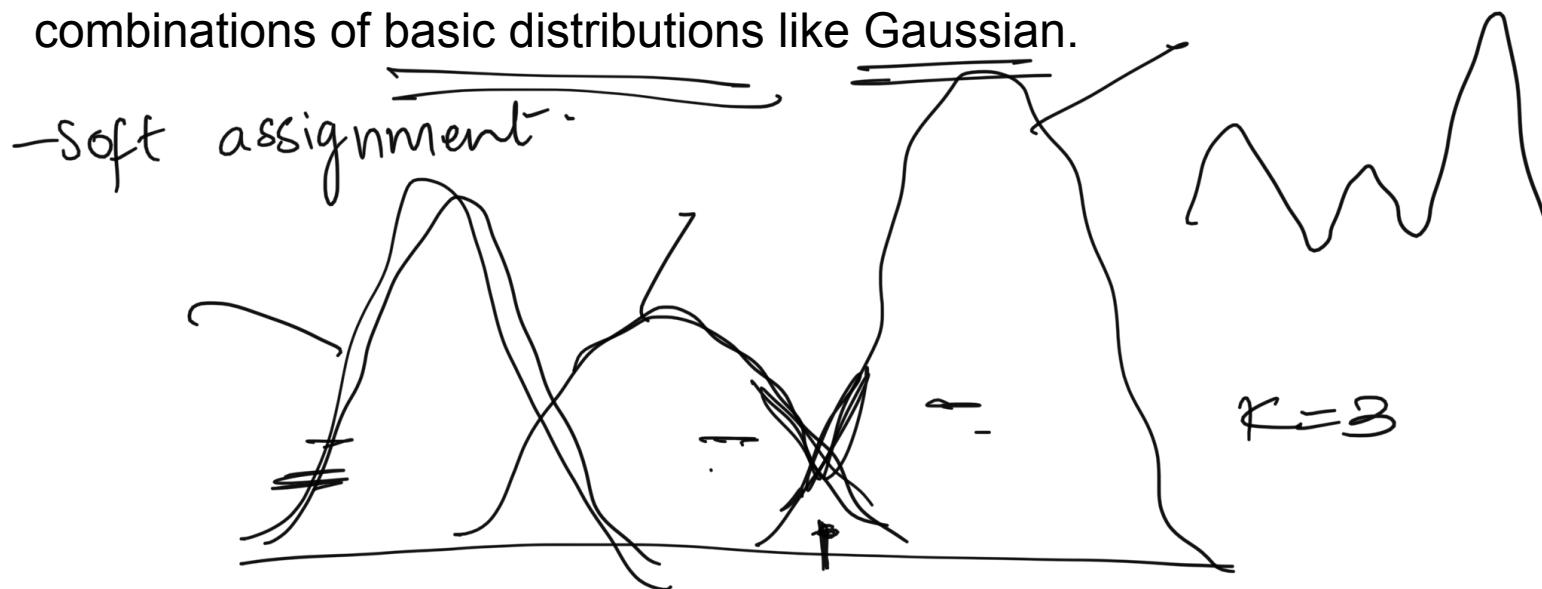


Better - 2 Gaussian



Mixture Models

- Mixture distributions are probabilistic models formed by linear combinations of basic distributions like Gaussian.



Gaussian Mixture Models

- Superposition of K Gaussian densities -

$$p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

$p(x) \geq 0$ $N(x | \mu_k, \Sigma_k) \geq 0 \Rightarrow \pi_k \geq 0$

$\sum_{k=1}^K \pi_k = 1$

$0 \leq \pi_k \leq 1$

Component of mixture model.

following conditions of probability

PDF of each Gaussian

$$\text{log-likelihood} = \ln P(x|N, \Sigma, \pi) = \sum_{n=1}^N \ln \left(\sum_k \pi_k N \left(\underline{\underline{\underline{\underline{x}}}} \mid \underline{\underline{\underline{\underline{\mu_k}}}}, \underline{\underline{\underline{\Sigma_k}}} \right) \right)$$

Gaussian Mixture Models

$$p(x) = \sum_k p(k) \cdot p(x|k) \iff$$

$p(k|x) \cdot \frac{\pi_k}{\sum_k \pi_k}$ Prior $N(x|\mu_k, \Sigma_k)$
 conditional probability

$$\gamma_k(x) = \frac{p(k) \cdot p(x|k)}{\sum_k p(k) \cdot p(x|k)} \quad \} \text{ - From Bayes theorem.}$$

responsibilities

$$= \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_k \pi_k N(x|\mu_k, \Sigma_k)}$$

$\pi_k \Rightarrow \text{No. of components in the mixture model}$

EM

$k=4$

EM algorithm for GMM

- \curvearrowleft k-means initialization
 \curvearrowright k-means centroids
Sample covariance
 $\frac{\text{No. of samples}}{\text{in each cluster}}$
 $\frac{1}{\text{Total samples}}$
1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
 2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \quad (9.23)$$

EM algorithm for GMM

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\underline{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \underline{\gamma(z_{nk})} \underline{\mathbf{x}_n} \quad (9.24)$$

$$\underline{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \underline{\gamma(z_{nk})} (\underline{\mathbf{x}_n} - \underline{\mu}_k^{\text{new}}) (\underline{\mathbf{x}_n} - \underline{\mu}_k^{\text{new}})^T \quad (9.25)$$

$$\underline{\pi}_k^{\text{new}} = \frac{N_k}{N} \quad (9.26)$$

where

$$N_k = \sum_{n=1}^N \underline{\gamma(z_{nk})}. \quad (9.27)$$

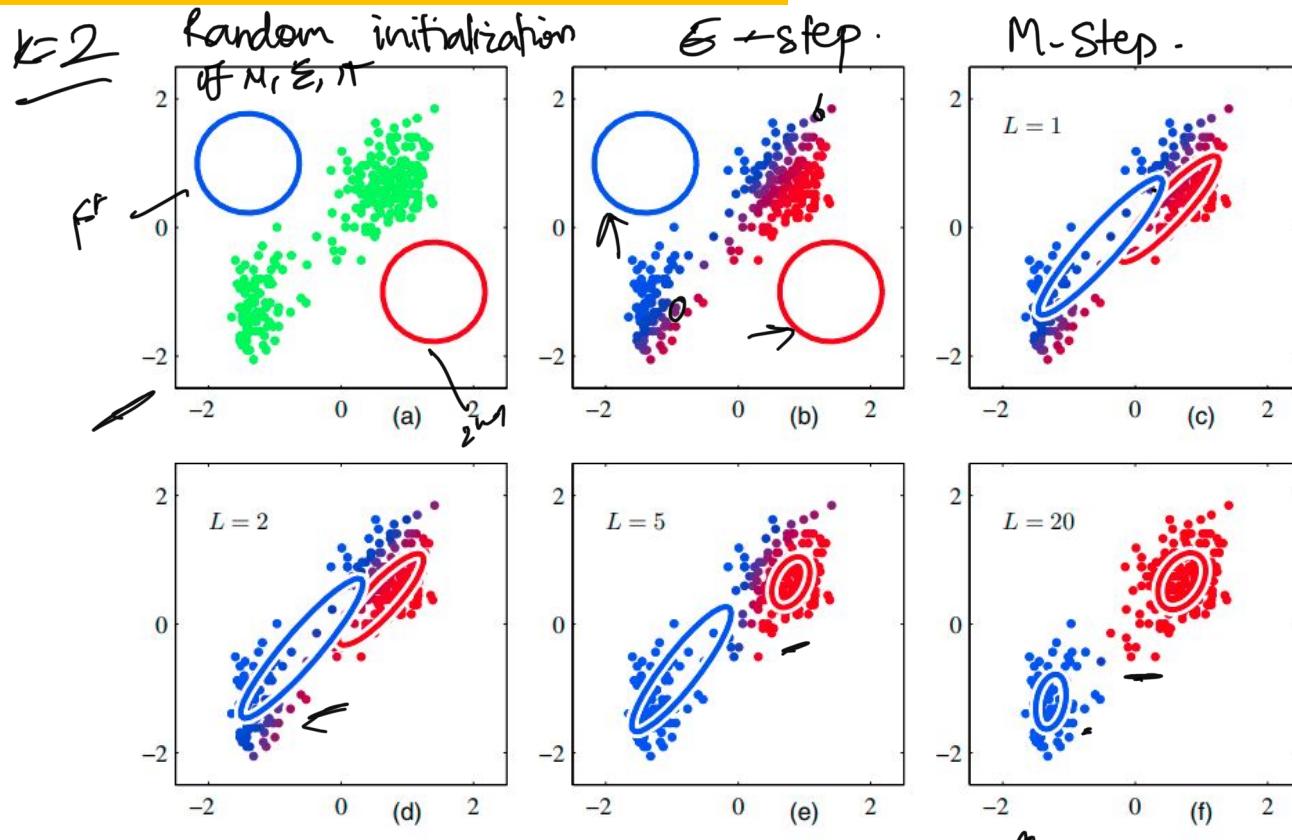
4. Evaluate the log likelihood

MLE

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (9.28)$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

Illustration of EM for GMM



Questions?

CSE 575: Statistical Machine Learning (Spring 2021)
Instructor: Nupur Thakur

Spectral Clustering

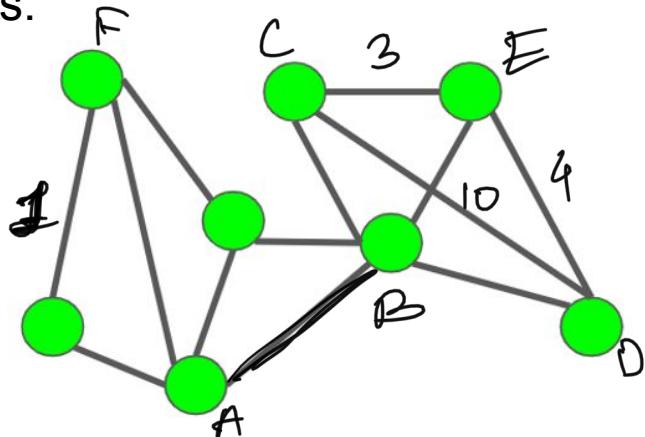


Table of contents

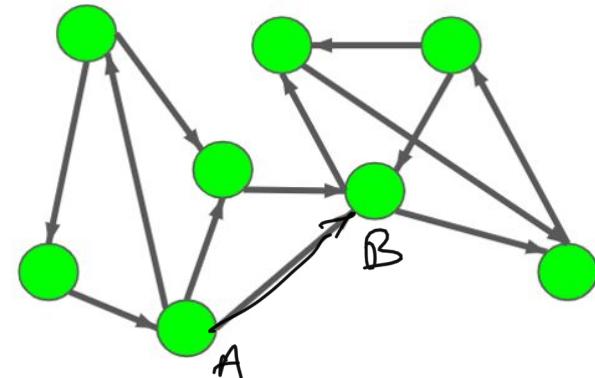
- 1. Graph Representation**
- 2. Clustering as a Graph Partition Problem**
- 3. Graph Cuts**

What is a graph?

- A graph $G = (V, E)$ is defined by V , a set of N vertices, and E , a set of edges.



Undirected graph



Directed graph

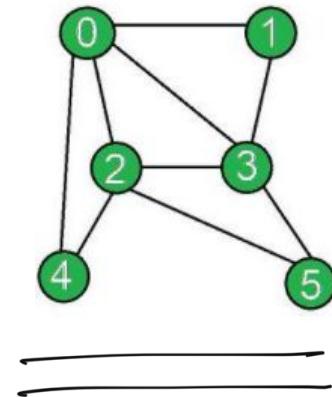
Graph Representation - Adjacency Matrix W

- For an undirected graph - N vertices

$N \times N$ []

Adjacency matrix - 6×6

0	1	2	3	4	5
0	0	1	1	1	0
1	1	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0

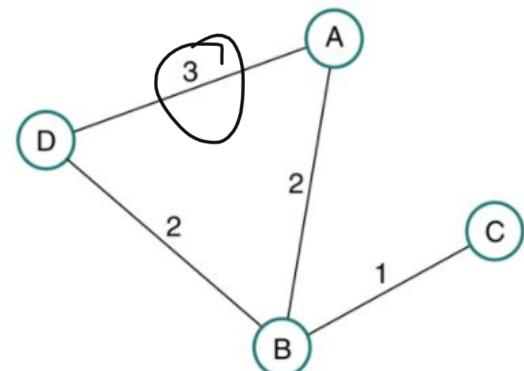


Graph Representation - Weighted Adjacency Matrix

$N \times N$ matrix

4×4 matrix

$$\begin{array}{ccccc} & A & B & C & D \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \left[\begin{matrix} 0 & 2 & 0 & 3 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \right] \end{array}$$

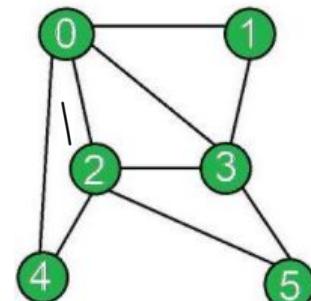


Graph Representation - Degree Matrix D

- Degree of a node is the number of edges incident to the node.
 - $N \times N$ matrix for N nodes.
 - Diagonal matrix

$$\begin{matrix} & 0 & 1 & 2 & 3 & 4 & 5 \\ 0 & 4 & & & & & \\ 1 & & 0 & & & & \\ 2 & & & 2 & & & \\ 3 & & & & 0 & & \\ 4 & & & & & 4 & \\ 5 & & & & & & 4 \end{matrix}$$

Undirected graph.



Graph Representation - Graph Laplacian L

- Degree matrix*
- Unnormalized*
- Adjacency matrix*
- $L = D - W$
 - Properties of L -
 - L is symmetric and positive semi-definite
 - L has n non-negative, real-valued eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$
 - The smallest eigenvalue is 0, the corresponding eigenvector is the 1-vector (all elements being 1)

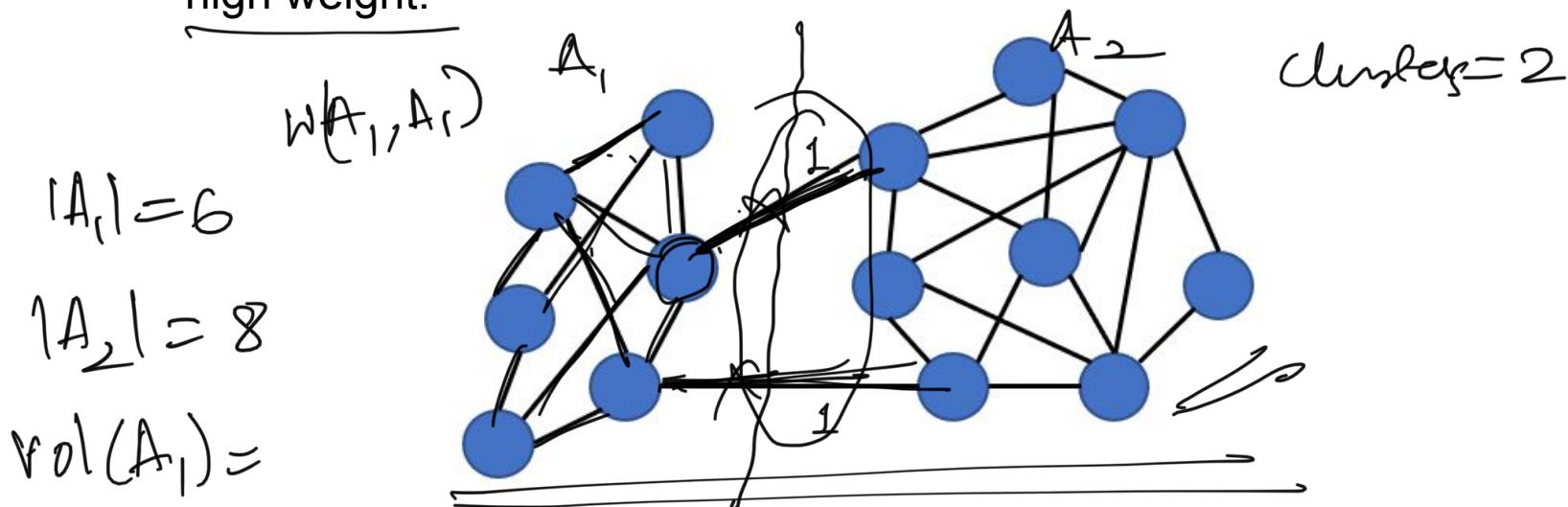
Different ways to build similarity graphs

- ϵ -neighborhood graph $\text{dist}(A, B) < \epsilon \Rightarrow$ Edge betw A & B
 - unweighted graph -
- k -nearest neighbor graph -
 - mutually k -nearest neighbor- property is satisfied both ways
- fully connected graph

Clustering – Increase inter-cluster distance
Decrease intra-cluster distance.

Clustering as a Graph Partition

- Find a partition of a graph such that the edges between different groups have a very low weight while the edges within a group have high weight.



MinCut Problem

- Given a similarity graph with adjacency matrix W , the simplest and most direct way to construct a partition of the graph is to solve the mincut problem.

~~question~~ ✓

$$\underset{k=2}{\text{minimize}} \quad \text{cut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k W(A_i, \overline{A}_i)$$

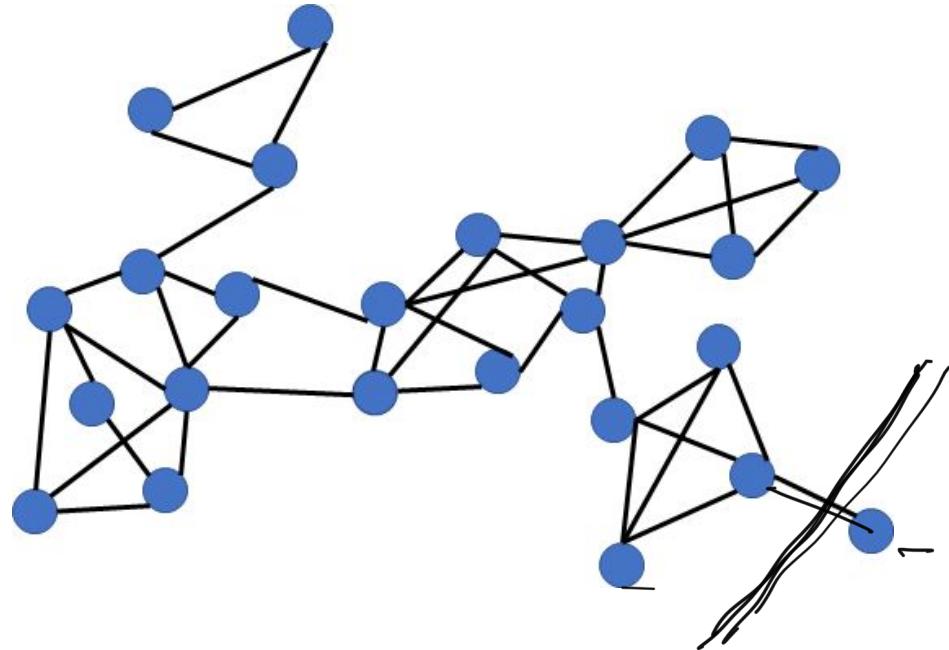
Complement of A_i

$k=2$

$$\text{cut}(A_1, A_2) = \frac{1}{2} \sum_{i=1}^2 W(A_i, \overline{A}_i)$$

MinCut Problem - Major Drawback

- It does not balance the size of the partitions



$|A_i| \Rightarrow$ No. of vertices/nodes in
a cluster A_i .

Other types of cuts

- To overcome the drawback of MinCut, one obvious solution is to specify that the partitions are reasonably large.

$$\underline{\text{RatioCut}}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \overline{A}_i)}{|A_i|} = \sum_{i=1}^k \frac{\text{cut}(A_i, \overline{A}_i)}{|A_i|}$$

Normalized
cut.

$$\underline{\text{Ncut}}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \overline{A}_i)}{\text{vol}(A_i)} = \sum_{i=1}^k \frac{\text{cut}(A_i, \overline{A}_i)}{\text{vol}(A_i)}$$

$$\underline{\text{MinMaxCut}}(A_1, \dots, A_k) := \sum_{i=1}^k \frac{\text{cut}(A_i, \overline{A}_i)}{W(A_i, A_i)}$$

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Spectral Clustering



Table of contents

1. Properties of Graph Laplacian
2. Formulation of MinCut

Other types of cuts

- To overcome the drawback of MinCut, one obvious solution is to specify that the partitions are reasonably large.

$$\Leftarrow \text{RatioCut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{|A_i|} = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{|A_i|}$$

$$\text{Ncut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{\text{vol}(A_i)} = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)}.$$

No. of
vertices in
a cluster A_i

$$\text{MinMaxCut}(A_1, \dots, A_k) := \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{W(A_i, A_i)}$$

Page No. 25
of the doc in Recommended Resources

Properties of Graph Laplacian

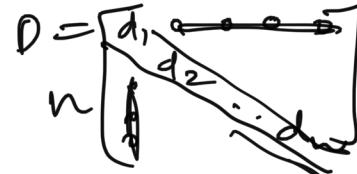
D *degree matrix*

- $L = D - W \rightarrow$ Adjacency Matrix

- Properties of L -

- For every vector $f \in \mathbb{R}^n$ we have

$$f' L f \geq 0$$



$$f' L f = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2$$

Proof -

By the definition of d_i ,

$$\begin{aligned} f' L f &= f' D f - f' W f = \sum_{i=1}^n d_i f_i^2 - 2 \sum_{i,j=1}^n f_i f_j w_{ij} \\ &= \underbrace{\frac{1}{2} \left(\sum_{i=1}^n d_i f_i^2 - 2 \sum_{i,j=1}^n f_i f_j w_{ij} + \sum_{j=1}^n d_j f_j^2 \right)}_{\text{Element of } D} - \underbrace{\frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2}_{\text{Element in } W} \\ &\stackrel{\text{by } (a-b)^2 = a^2 - 2ab + b^2}{=} 0 \end{aligned}$$

Properties of Graph Laplacian

- Properties of L -
 - L is symmetric and positive semi-definite

D-Symmetric
W-Symmetric

From the first property, $f' L f \geq 0$

Positive semi-definite matrix.

Properties of Graph Laplacian

- Properties of L -

- The smallest eigenvalue is 0, the corresponding eigenvector is the 1-vector (all elements being 1)

Proof -

$$L = D - W = \underbrace{\begin{pmatrix} d_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & d_n \end{pmatrix}}_{D} - \underbrace{\begin{pmatrix} w_{11} & \cdots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{n1} & \cdots & w_{nn} \end{pmatrix}}_{W} = \begin{pmatrix} d_1 - w_{11} & -w_{12} & \cdots & -w_{1n} \\ \vdots & \ddots & & \vdots \\ -w_{n1} & \cdots & \sum_{j=1}^n w_{nj} - w_{nn} & \end{pmatrix}$$

$$= \begin{pmatrix} \sum_{j=1}^n w_{1j} - w_{11} & \cdots & -w_{1n} \\ \vdots & \ddots & \vdots \\ -w_{n1} & \cdots & \sum_{j=1}^n w_{nj} - w_{nn} \end{pmatrix}.$$

$$\Rightarrow \begin{pmatrix} \sum_{j=1}^n w_{1j} - w_{11} & \cdots & -w_{1n} \\ \vdots & \ddots & \vdots \\ -w_{n1} & \cdots & \sum_{j=1}^n w_{nj} - w_{nn} \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

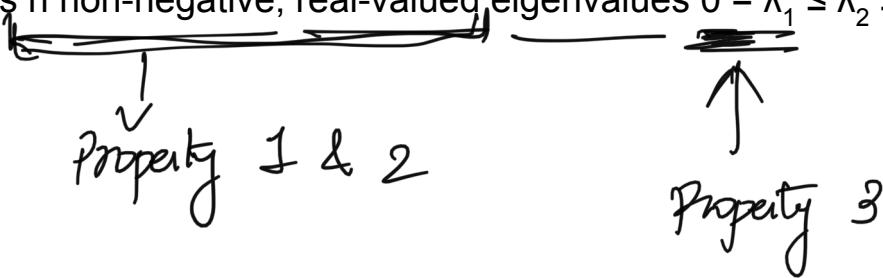
$$n \times 1 \begin{pmatrix} \sum_{j=1}^n w_{1j} - w_{11} & \cdots & -w_{1n} \\ \vdots & \ddots & \vdots \\ -w_{n1} & \cdots & \sum_{j=1}^n w_{nj} - w_{nn} \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$L\mathbf{x} = \lambda \mathbf{x}$

Properties of Graph Laplacian

- Properties of L -

- L has n non-negative, real-valued eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$



MinCut Problem

- Given a similarity graph with adjacency matrix W , the simplest and most direct way to construct a partition of the graph is to solve the mincut problem.

Minimize $\text{cut}(A_1, \dots, A_k) := \underbrace{\frac{1}{2} \sum_{i=1}^k W(A_i, \bar{A}_i)}$

MinCut for k=2

- Given \mathbf{W} and a cluster membership vector \mathbf{q} ,

$$q_i = \begin{cases} 1 & i \in \text{Cluster A} \\ -1 & i \in \text{Cluster B} \end{cases}$$

$$\boxed{\mathbf{q} = \underset{\mathbf{q} \in [-1, 1]^n}{\operatorname{argmin}} \frac{1}{4} \sum_{i,j} (q_i - q_j)^2 w_{i,j}}$$

$n \rightarrow$ no. of samples
 $2^n \rightarrow$ exponential complexity.

Cut size

Relaxation Approach

$$\left\{ \begin{array}{l} \mathbf{q} = \underset{\mathbf{q}}{\operatorname{argmin}} J = \underset{\mathbf{q}}{\operatorname{argmin}} \mathbf{q}^T (\mathbf{D} - \mathbf{W}) \mathbf{q} \\ \text{subject to } \sum_{i=1}^n q_i^2 = n \end{array} \right.$$

eigenvector
corresponding
to second
smallest
eigenvalue
of \mathbf{L} is the sol

$q_i \Rightarrow$ any real value betw $-1, 1$

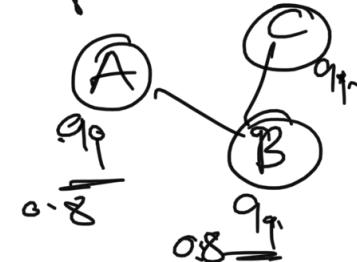
$$\mathbf{q} = [0.8 \quad -0.2 \quad 0.1]$$

$$\downarrow$$

$$k=2$$

$$\{q > 0 \Rightarrow 1 \quad q \leq 0 \Rightarrow -1\}$$

$$q \leq 0 \Rightarrow -1$$



Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Dimensionality Reduction



Table of contents

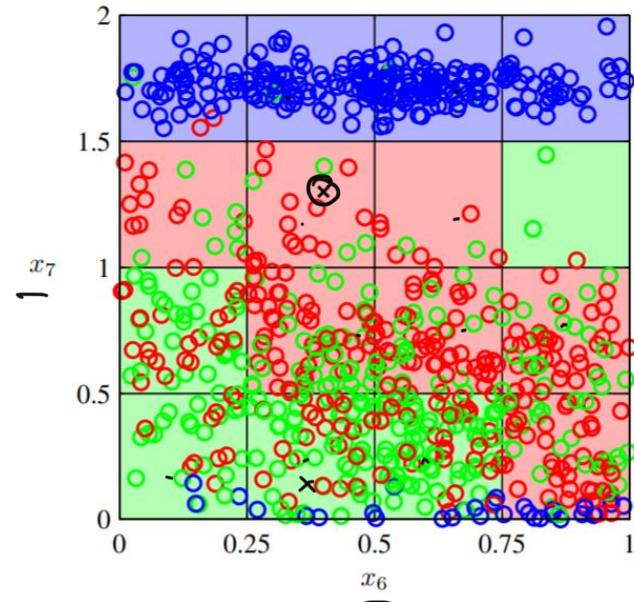
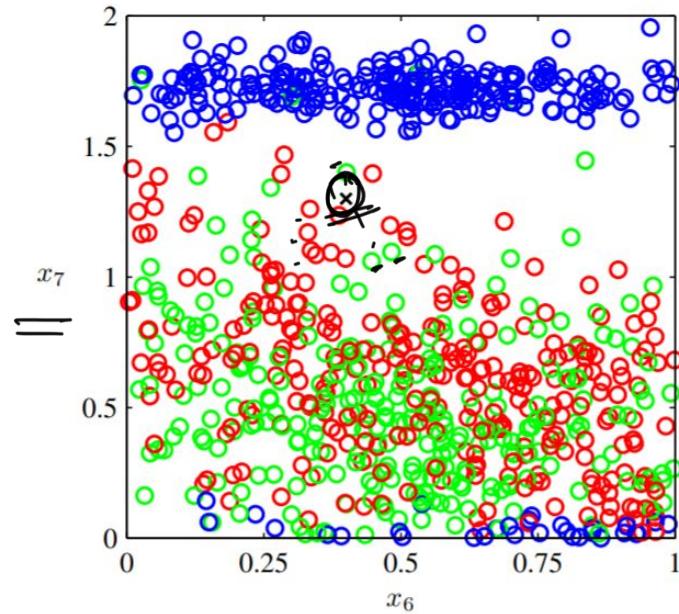
1. Dimensionality Reduction
2. Curse of Dimensionality
3. Principal Component Analysis

Dimensionality reduction

$$\begin{array}{c} \overline{784} \Rightarrow 28 \times 28 \\ \downarrow \\ \boxed{2} \quad 3 \end{array}$$

- Given N data points in a high-dimensional space (in the order of tens of thousands of dimensions), project them into some low-dimensional space.
- Why project them into some low-dimensional space?
 - Curse of dimensionality.

Curse of dimensionality



Principal Component Analysis (PCA)

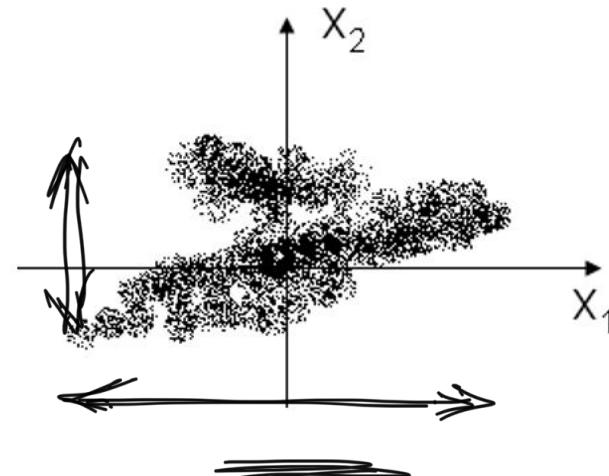


- Basic idea -

2D to 1D

Different ways -

- (1) Random
- (2) Discard the less descriptive
- (3) Project onto a space such that the max-variance is obtained.



PCA Problem Formulation

d dimensions
Problem

Given n samples $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ in d -dimensional space,
find a direction \mathbf{e}_1 , such that the projection of D onto \mathbf{e}_1 gives
the largest variance (compared with any other direction).

$\mathbf{e}_1 \Rightarrow$ unit vector, $\|\mathbf{e}_1\| = 1$

Projections - $y_i = \mathbf{x}_i \cdot \mathbf{e}_1$ \rightarrow inner product.

Mean of projections - $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{e}_1 = \bar{\mathbf{x}} \cdot \mathbf{e}_1$

PCA Problem Formulation

Variance of projections -

$$\begin{aligned}\sigma^2(e_i) &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^n [(x_i - \bar{x}) e_i]^2 \\ &= \sum_{j=1}^d \sum_{k=1}^d \left[\frac{1}{n} \sum_{i=1}^n (x_{i,j} - \bar{x}_{i,j})(x_{i,k} - \bar{x}_{i,k}) \right] e_j e_k \\ &= \underbrace{\sum_{j=1}^d \sum_{k=1}^d e_j^T e_k}_{\text{Covariance}} G_{jk} \quad \text{where } G \rightarrow \text{Covariance matrix.}\end{aligned}$$

$$\underbrace{\sigma^2(e_i)}_{=} = \underbrace{e^T C e}_{=} \Rightarrow \text{Matrix form.}$$

PCA Problem Formulation

$$e = \underset{e}{\operatorname{argmax}} \sigma^2(e) \quad \text{subject to } \|e\|=1$$

$$\underline{F}(e) = \underset{e}{\operatorname{argmax}} \underbrace{\sigma^2(e)}_{\downarrow} - \lambda \underbrace{(e^T e - 1)}_{\text{Lagrange's multipliers}}$$

$$\frac{\partial F(e)}{\partial e} = \underline{\underline{e^T C e}} - \underline{\underline{2\lambda e}} = 0 \Rightarrow \underline{\underline{C e = \lambda e}}$$

$e \rightarrow$ eigenvector of covariance matrix C

$\lambda \rightarrow$ eigenvalue corresponding to eigenvector e .

PCA - Principal Components

- What are principal components?
 - e_i vectors.
 - First principal component - eigenvector with the largest eigenvalue.
- How many principal components to keep?

Total variance = sum of variance in all the projections.

$$d' \Rightarrow \text{new no. of dimensions} \quad \text{where } d' \ll d$$
$$= \sum_{j=1}^d \lambda_j$$

$\approx \frac{\sum_{j=1}^{d'} \lambda_j}{\sum_{j=1}^d \lambda_j}$

$$D = (n \times d) \quad n \times d'$$

PCA Algorithm

$d \rightarrow$ total no. of dimensions of data.

1. Compute the $\underline{\underline{d \times d}}$ sample covariance matrix $\underline{\underline{C}}$
2. Find the eigenvalues and corresponding eigenvectors of $\underline{\underline{C}}$
3. Project the original data onto the space spanned by the eigenvectors
 - The projection may be done onto a d' -dimensional subspace spanned by the first d' eigenvectors (ordered by the eigenvalue in descending order)
 - $\underline{d'}$ is determined by the desired accuracy

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Neural Networks & Deep Learning

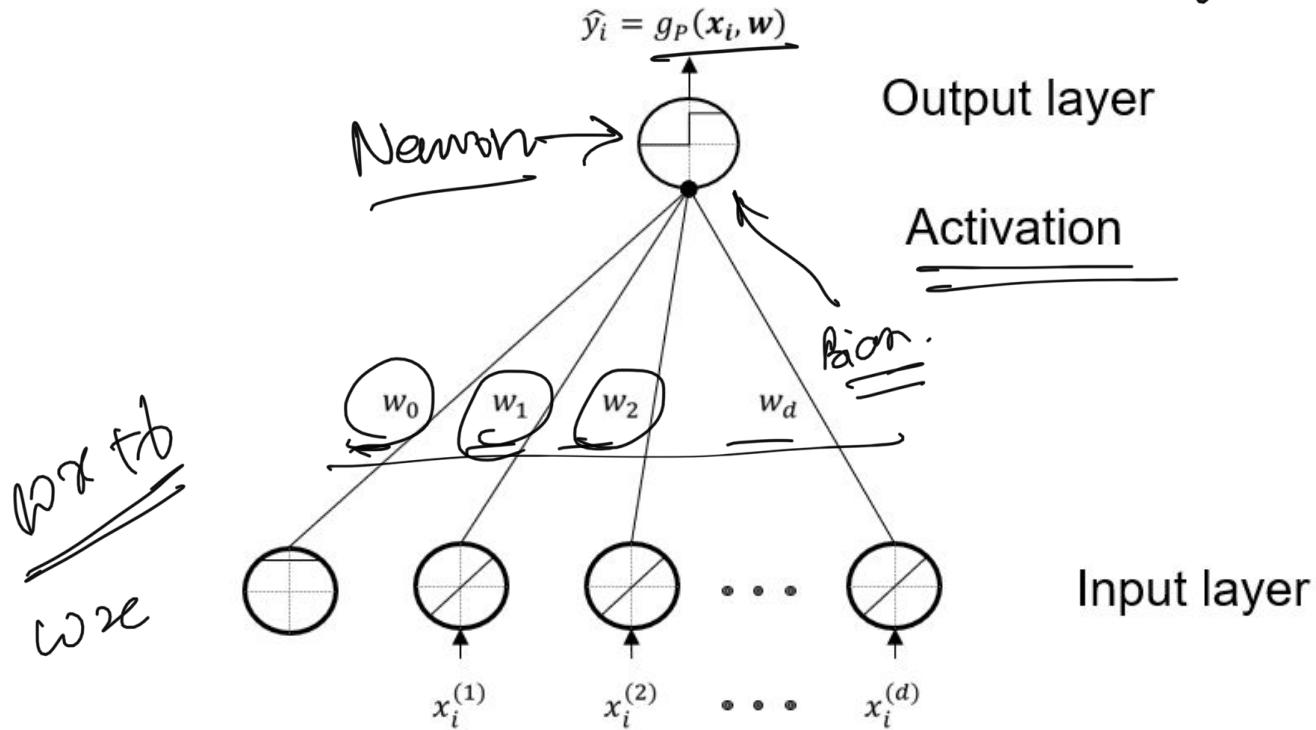


Table of contents

1. Perceptron Learning
2. XOR Problem using MLP
3. Back-propagation

What is a Perceptron?

Linear binary classifier.



Learning in Perceptron

Iterate for t until a stop criterion is met

{

for each sample x_i with label y_i :

{

compute the output \hat{y}_i of the network

estimate the error of the network $e(w(t)) = y_i - \hat{y}_i$

update the weight $w(t+1) = w(t) + e(w(t))x_i$

}

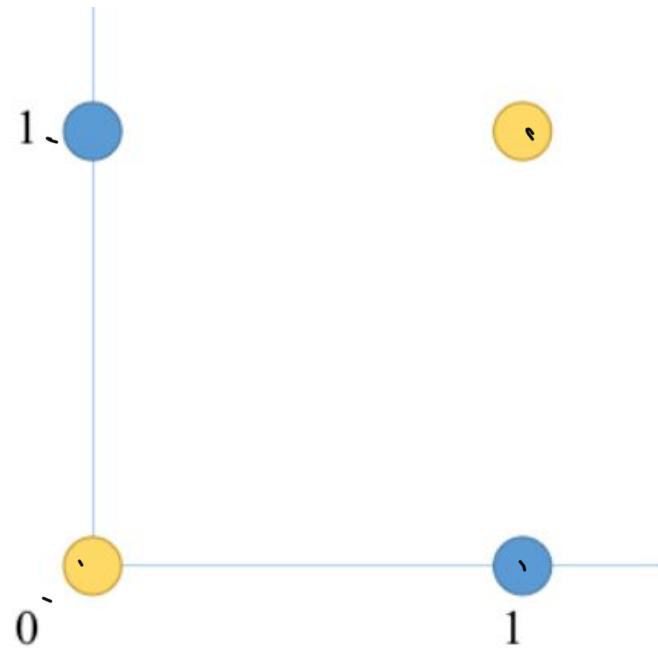
$t++$

}

Max iterations
→ does not change
 $\|w_{old} - w_{new}\| < \epsilon$
ground truth (labels).



XOR problem



XOR problem - Using MLP (Multi-layer perceptron) .

$$w_{t+1} = w_t - \eta \nabla_w J$$

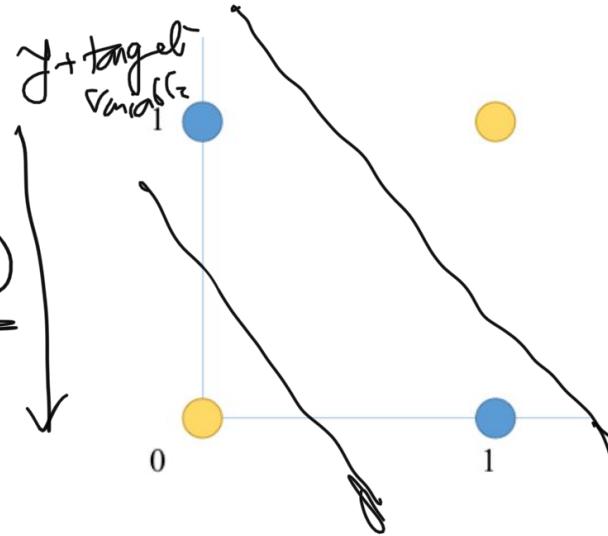
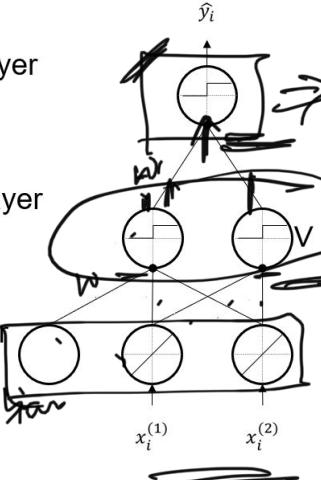
Output layer

Hidden layer

$$z_1 = w_1 x + b_1$$

Input layer

$$z_2 = w_2 x + b_2$$



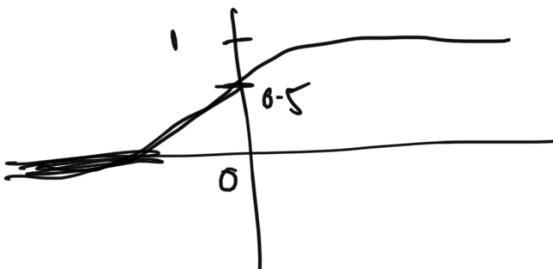
- How do we learn weights? Can perceptron learning be applied?

~ Gradient descent

~ Differentiable.

~ Backpropagation .

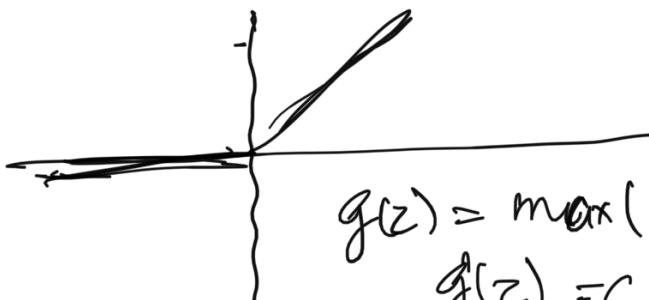
Activation Functions



Sigmoid

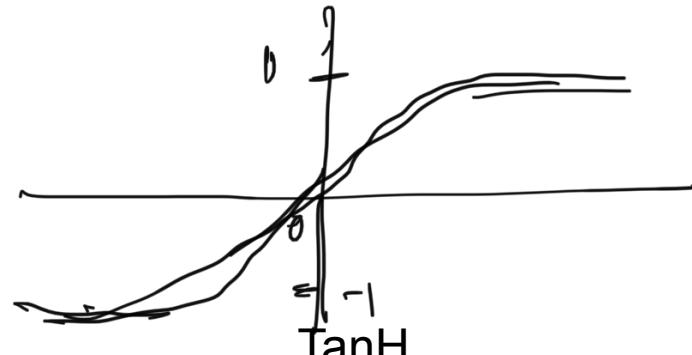
$$g(z) = \frac{1}{1+e^{-z}}$$

$$g'(z) = g(z)(1-g(z))$$



ReLU
(Rectified Linear Unit)

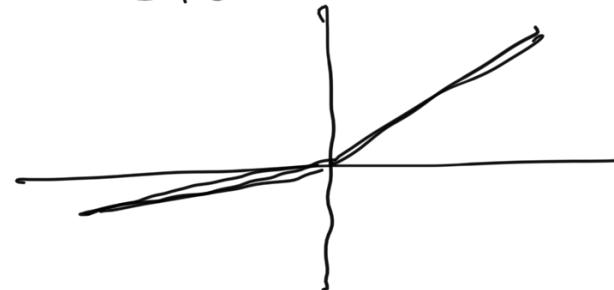
$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$



TanH

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g(z) = 1 - (g(z))^2$$



Leaky ReLU

$$g(z) = \max(\epsilon z, z)$$

Why non-linear activation functions are important?



$$\underbrace{w_1x + b_1}_{w_2(\sigma(w_1x + b_1))} \rightarrow \underbrace{w_2(w_1x + b_1) + b_2}_{w'(w_1x + b_1)}$$

- Bring non-linearity to the network.

Handling Multiple Classes

$$\left[\begin{array}{c} \underline{0.05} \\ \underline{0.05} \\ \textcircled{0.9} \end{array} \right]$$

- No. of output neurons - No. of classes.
- Activation function - Softmax activation -
 - Bring the values b/w 0 & 1.
 - argmax of ^{softmax} probabilities
- Error Function-

$C \rightarrow$ no. of classes.

$$\frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}$$

Normalisation Constant

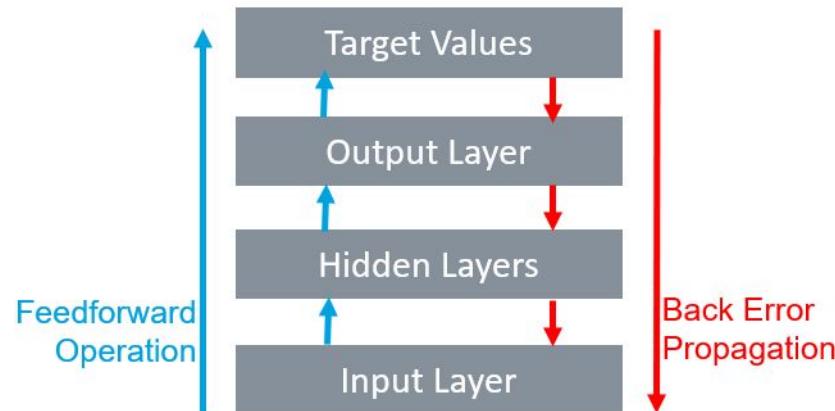
Cross-entropy (Binary cross entropy)

$$L = - \sum_{i=1}^n \sum_{j=1}^C (\gamma_i) \log p(y_i=j | x_i) \rightarrow \begin{cases} 1 & \gamma_i=j \\ 0 & \text{otherwise} \end{cases}$$

softmax probability.

MLP Learning - Back Propagation

- Key idea - Properly distribute error computed from output layer back to earlier layers to allow their weights to be updated in a way that reduce the error.

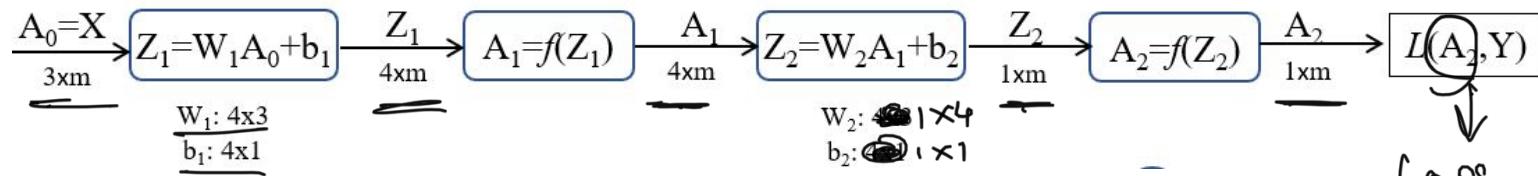
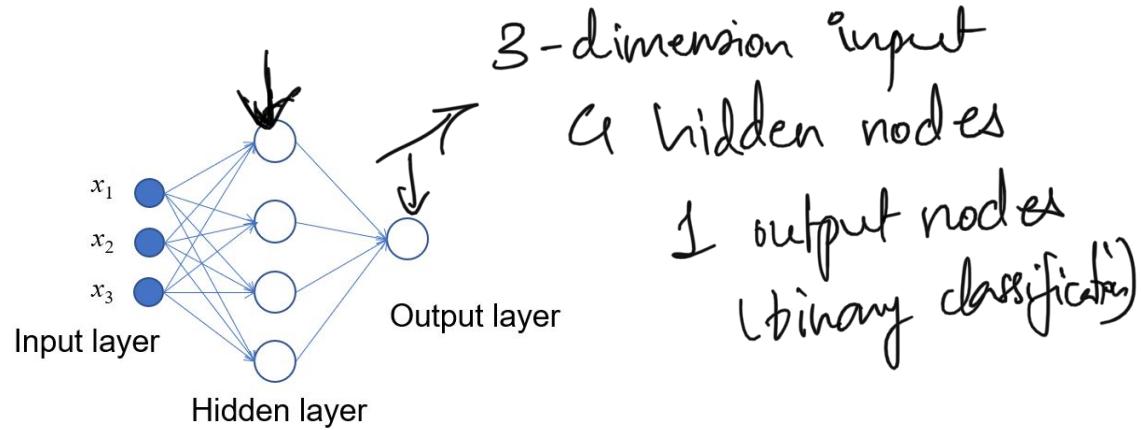


MLP Learning - Back Propagation

4×3 (4×1) 4×3
 $\xrightarrow{=}$ $\xrightarrow{\text{Broadcasting}}$

Feed-Forward Phase-

$m \rightarrow$ total no. of samples -
activation funⁿ \rightarrow sigmoid.



$$Z_1 = W_1 A_0 + b_1$$

\downarrow
 $4 \times m$

\downarrow
 2×3

\downarrow
 $3 \times m$

\downarrow
 4×1

$$Z_2 = f_2 A_1 + b_2$$

\downarrow
 $1 \times m$

\downarrow
 1×4

\downarrow
 $4 \times m$

\downarrow
 1×1

Loss
function

MLP Learning - Back Propagation

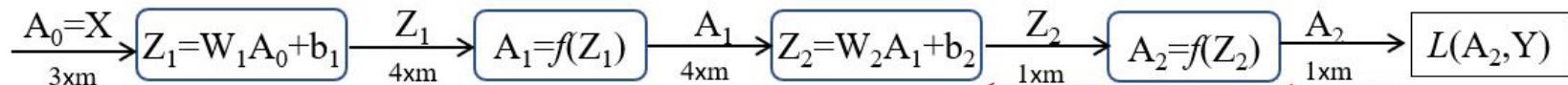
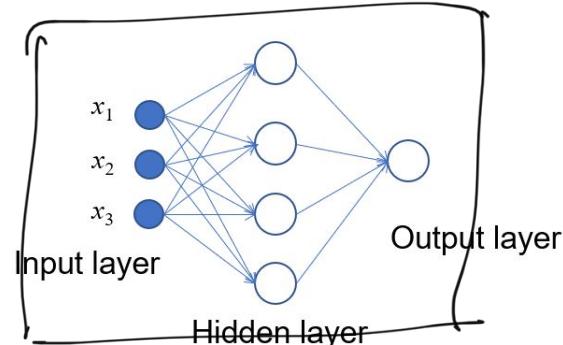
Back propagation phase-

$$w_{t+1} = w_t - \eta \nabla_w L$$

Parameters -

$$W_1, b_1, W_2, b_2$$

L ← Loss function



$$W_1: 4 \times 3$$

$$b_1: 4 \times 1$$

$$W_2: 1 \times 4$$

$$b_2: 1 \times 1$$

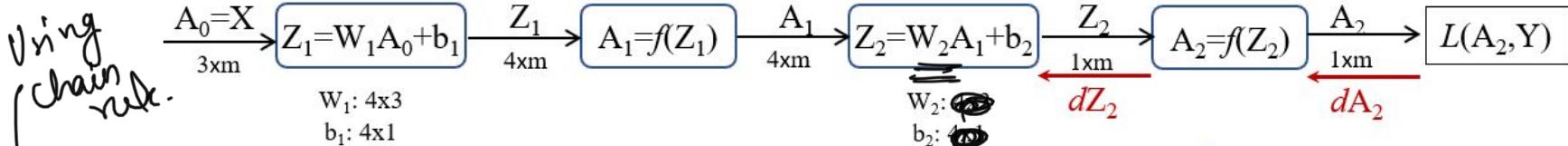
$$\frac{\partial L}{\partial W_1}, \quad \frac{\partial L}{\partial b_1}, \quad \frac{\partial L}{\partial W_2}, \quad \frac{\partial L}{\partial b_2}$$

$\frac{\partial L}{\partial Z_2} = \frac{\partial L}{\partial A_2}$

\Rightarrow Chain rule.

Activation \rightarrow Sigmoid, Output + Binary loss \rightarrow BCE

MLP Learning - Back Propagation



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}, \quad \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2}, \quad \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

$$\frac{\partial z_2}{\partial A_1} = W_2^T, \quad \frac{\partial A_1}{\partial z_1} = A_1(1-A_1)$$

$$L = -(y \log A_2 + (1-y) \log (1-A_2))$$

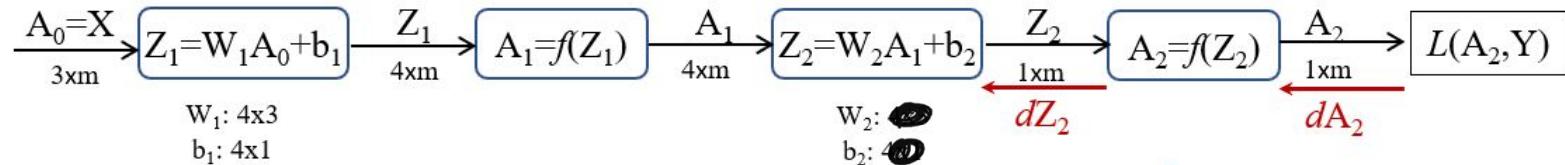
$$\frac{\partial L}{\partial A_2} = \frac{A_2 - y}{A_2(1-A_2)}, \quad \frac{\partial z_2}{\partial w_2} = A_1^T$$

$$\frac{\partial z_1}{\partial w_1} = A_0^T, \quad \frac{\partial z_1}{\partial b_1} = [1, \dots, 1]^T$$

$$\frac{\partial A_2}{\partial z_2} = A_2(1-A_2), \quad \frac{\partial z_2}{\partial b_2} = [1, 1, \dots, 1]^T$$

1xm 4xm 4xm 4xm

MLP Learning - Back Propagation



$$\frac{\partial L}{\partial w_2} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 4} \cdot \underbrace{A_1^T}_{m \times 4}, \quad \frac{\partial L}{\partial b_2} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 1} \cdot \underbrace{[1 \dots 1]^T}_{m \times 1}$$

$$\frac{\partial L}{\partial w_1} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 1} \cdot \underbrace{(w_2)^T}_{4 \times 1} \cdot \underbrace{A_1}_{1 \times 4} \cdot \underbrace{A_0^T}_{m \times 3}$$

$$\frac{\partial L}{\partial w_1} = \underbrace{(w_2)^T}_{4 \times 1} \underbrace{(A_2 - Y)}_{1 \times m} \underbrace{A_1}_{4 \times m} \underbrace{(1-A_1)}_{4 \times m} \underbrace{A_0^T}_{m \times 3}$$

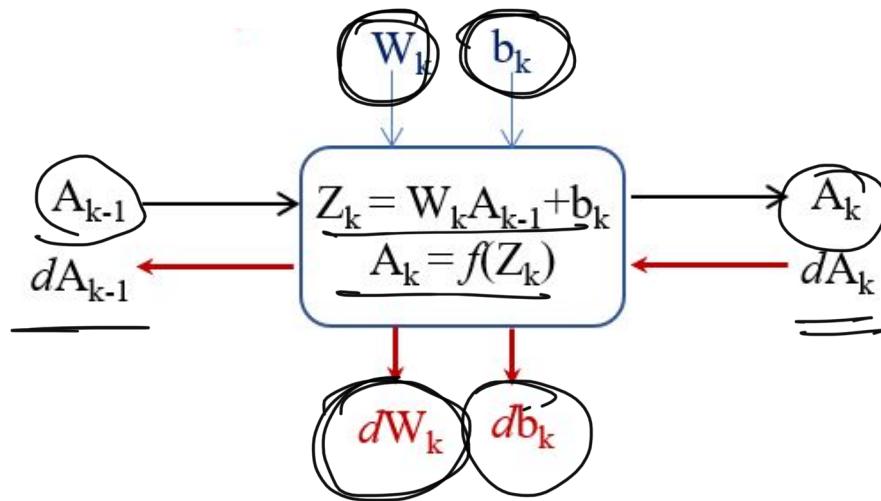
elementwise multiplication

$$\frac{\partial L}{\partial b_1} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 1} \cdot \underbrace{w_2^T}_{4 \times 1} \underbrace{A_1}_{1 \times 4} \underbrace{(1-A_1)}_{1 \times m} \underbrace{[1 \dots 1]^T}_{m \times 1}$$

$$\frac{\partial L}{\partial b_1} = \underbrace{w_2^T}_{4 \times 1} \underbrace{(A_2 - Y)}_{1 \times m} \underbrace{A_1}_{4 \times m} \underbrace{(1-A_1)}_{4 \times m} \underbrace{[1 \dots 1]^T}_{m \times 1}$$

elementwise multiplication

MLP Learning - Back Propagation



Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Neural Networks & Deep Learning



Table of contents

1. Convolutional Neural Network

2. Dropout

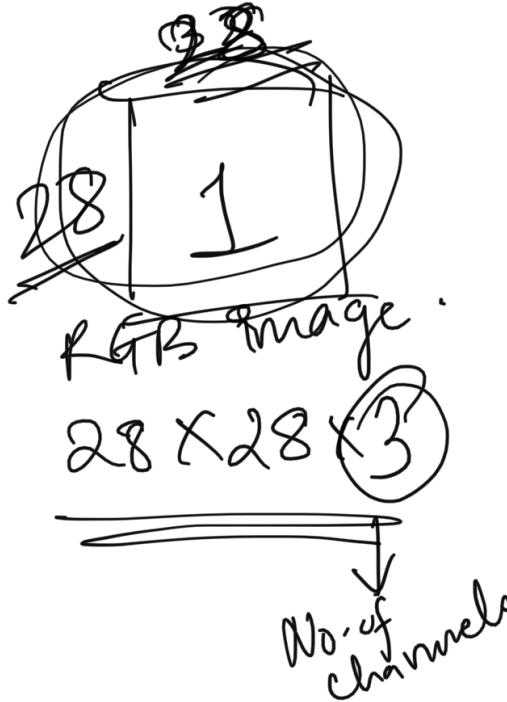
3. Batch normalization

Number of parameters in dense layers

FC (fully connected).

$$\underline{224 \times 224}$$

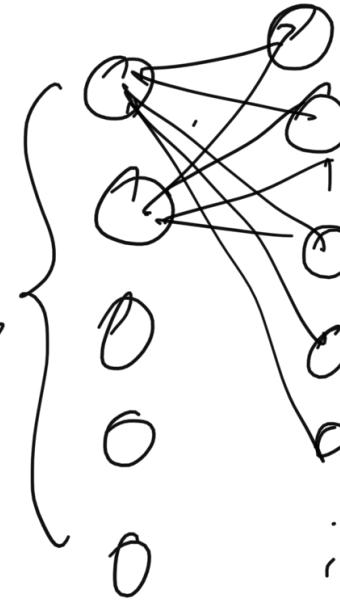
1000.



$$28 \times 28 \times 3 \Rightarrow 784 \times 3$$



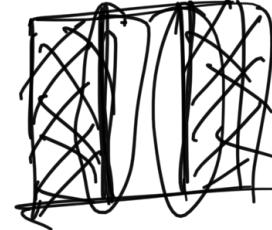
784 x 3
parameters



huge no. of
parameters

$$\underline{784 \times 3 \times 1000}$$

Convolutional Neural Network (CNN)



- Basic building block is the convolution for image processing
- Weight sharing - reduces parameters
- Learn the kernels based on the task.

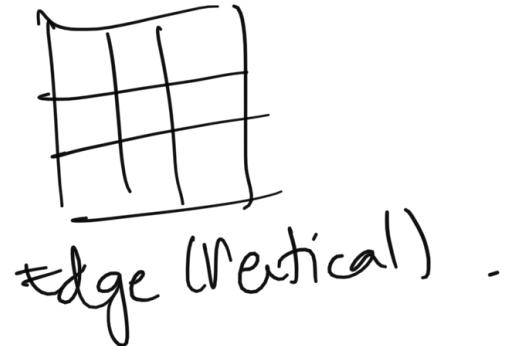
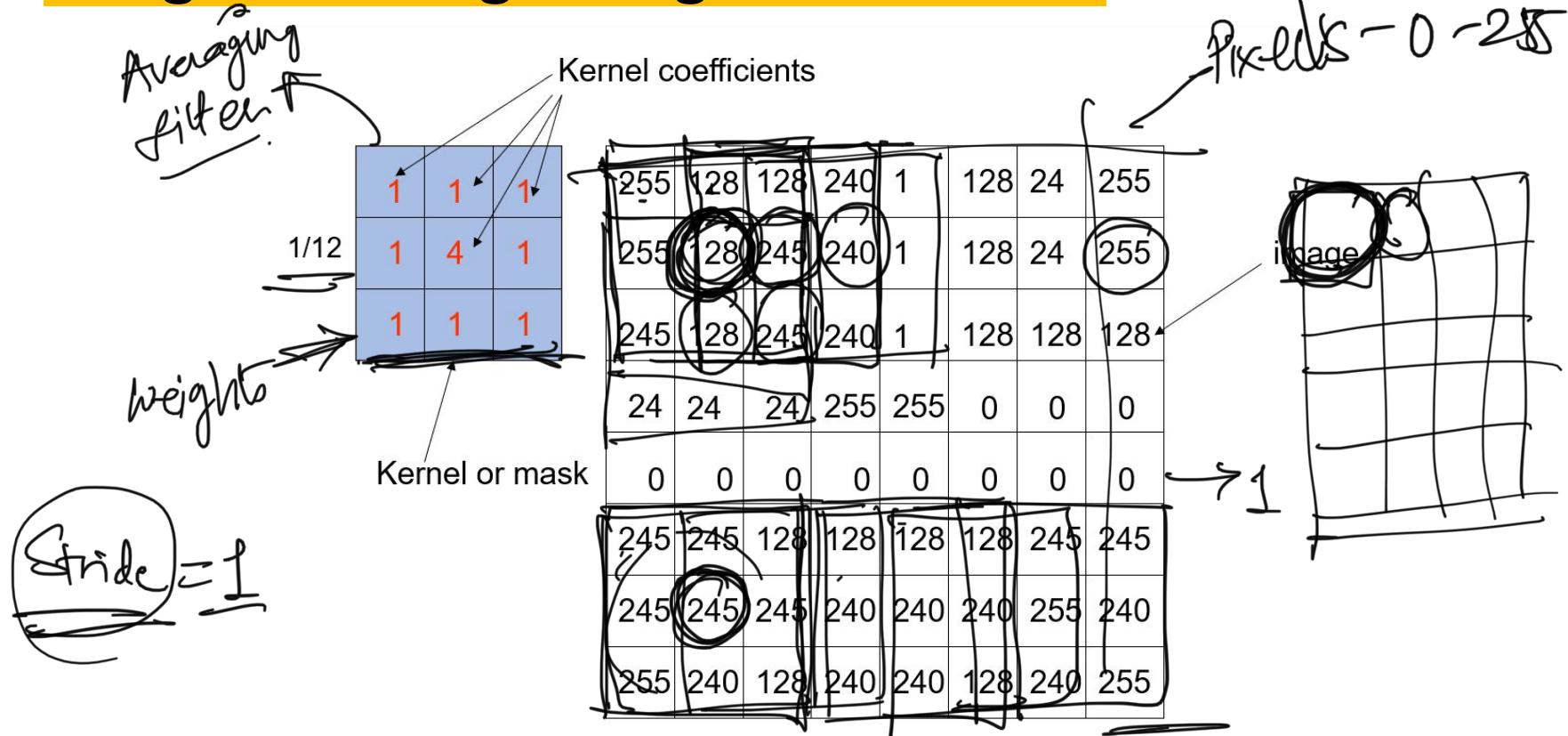
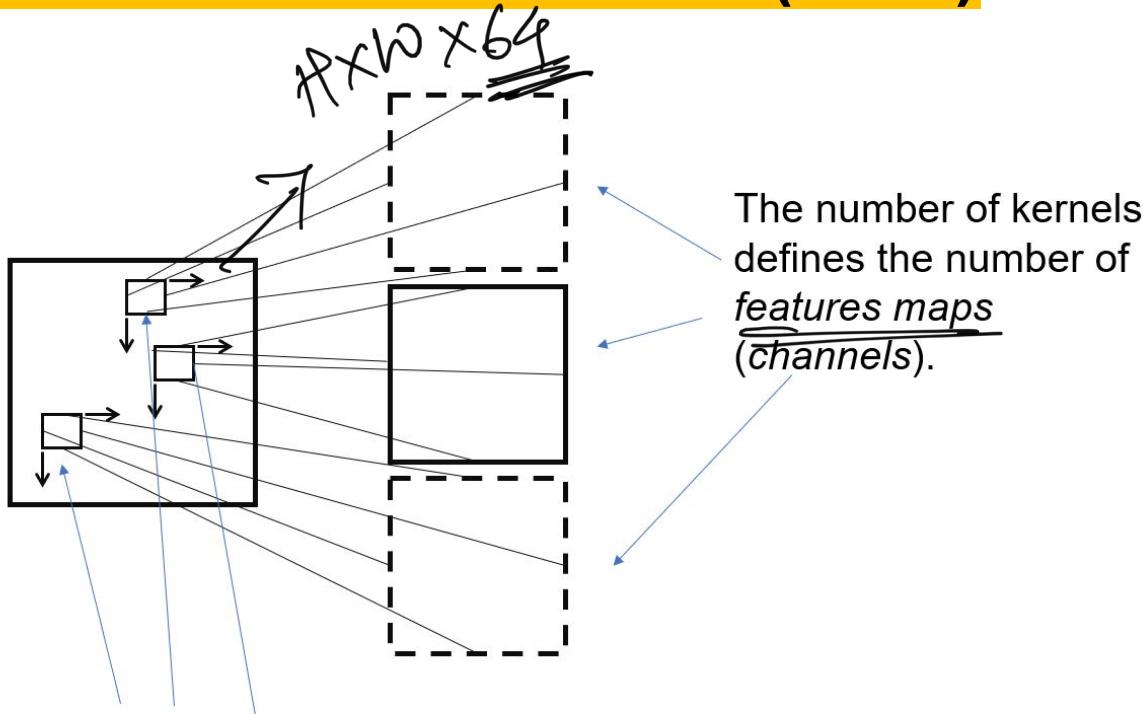


Image Filtering using convolution



Convolutional Neural Network (CNN)



The sizes of the kernels define the *receptive fields*.

Padding in convolutional layer

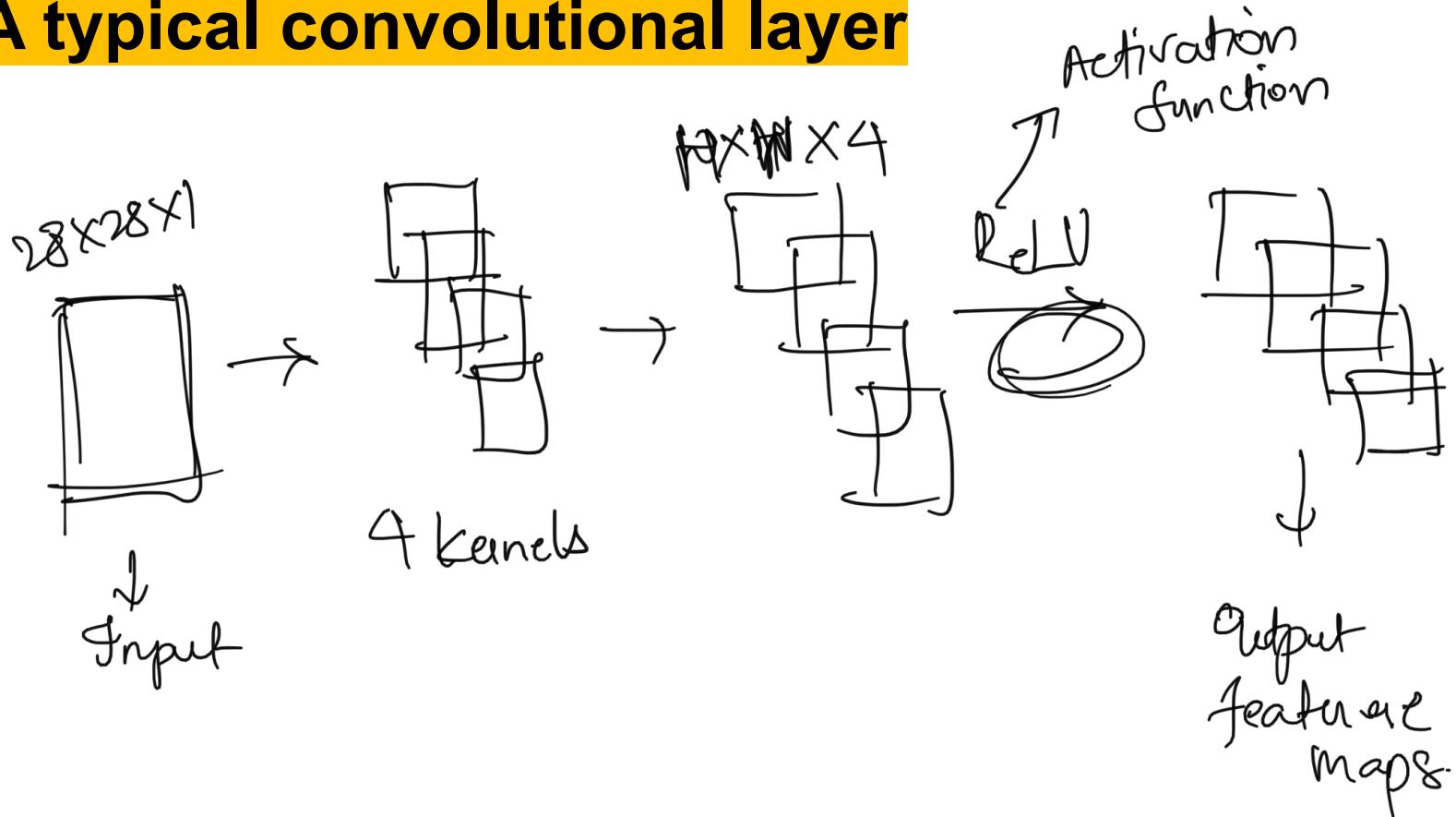
- ## - Zero-padding

-Equally value the boundary pixels -

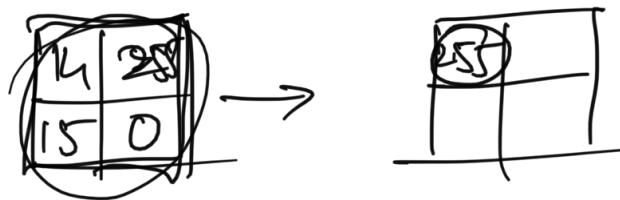
- Same fees reduced
size of the output.

255	128	128	240	1	128	24	255
255	128	245	240	1	128	24	255
245	128	245	240	1	128	128	128
24	24	24	255	255	0	0	0
0	0	0	0	0	0	0	0
245	245	128	128	128	128	245	245
245	245	245	240	240	240	255	240
255	240	128	240	240	128	240	255

A typical convolutional layer



Pooling in CNN

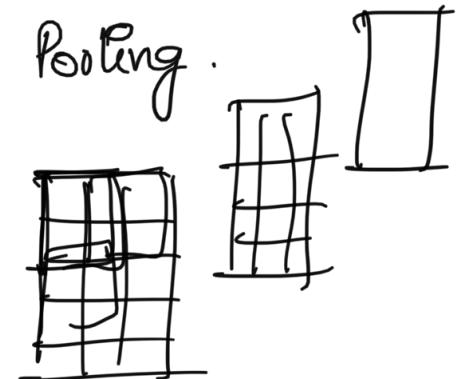


- Used to reduce the size of the output, thereby reducing the number of parameters and preventing overfitting.
- Done on every channel along the height and width of the output. So, the number of channels remain the same.

— Max Pooling

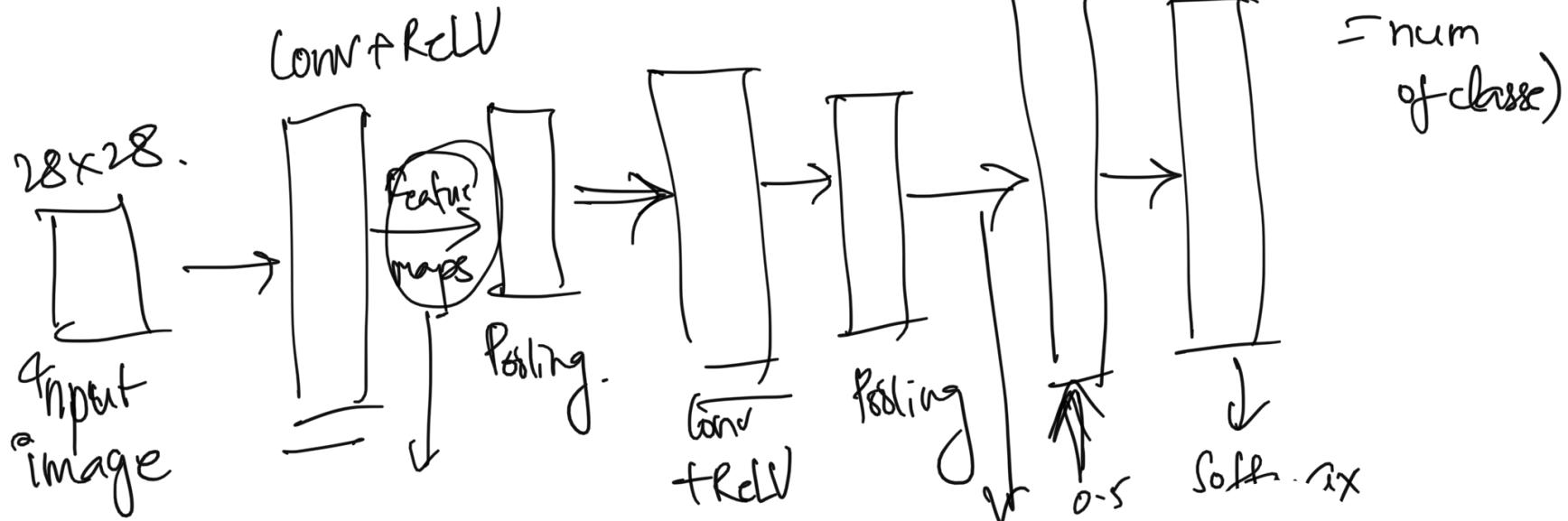
— Average Pooling .

— Kernel size , type of pooling , stride .



$$28 \times 28 \Rightarrow 784$$

A typical CNN



$$224 \times 224$$

$$299 \times 299$$

Batch size \times ~~$B \times H \times W \times C$~~

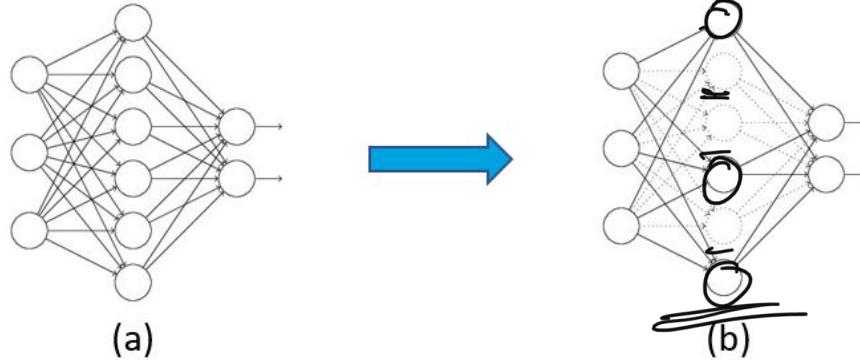
Why is regularization necessary?

- Huge parameter space - poor solutions
- Convergence to local minimum which leads to good performance on training data only.
- To overcome -
 - L2 norm of weights
 - Dropout
 - Batch normalization .

Dropout

- Obtain (b) by randomly deactivate some hidden nodes in (a).
- Reducing co-adaptation of neurons

Probability of dropping a neuron. $(0 - 1) \xrightarrow{\underline{p}}$



Batch Normalization

- Normalizes layer inputs of a batch

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$
$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

variance

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$\hat{x}_i \rightarrow$ mean
 \pm standard deviation

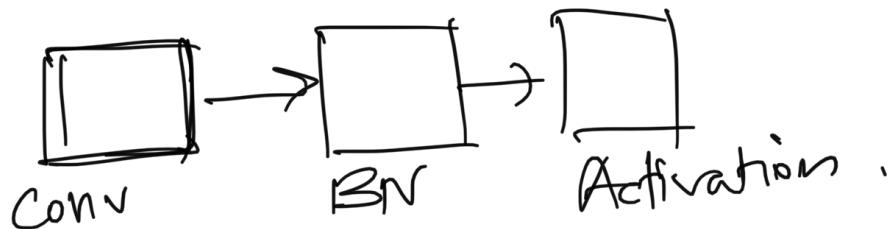
$$y_i \leftarrow \gamma \hat{x}_i + \beta = BN_{\gamma, \beta}(x_i)$$

Learnable parameters

ϵ smoothing parameter

Batch Normalization - Advantages

- Use of higher learning rates
- Covariate shift is reduced
- Avoids vanishing gradients
- Makes the network less sensitive to weight initialization.



Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Neural Networks & Deep Learning



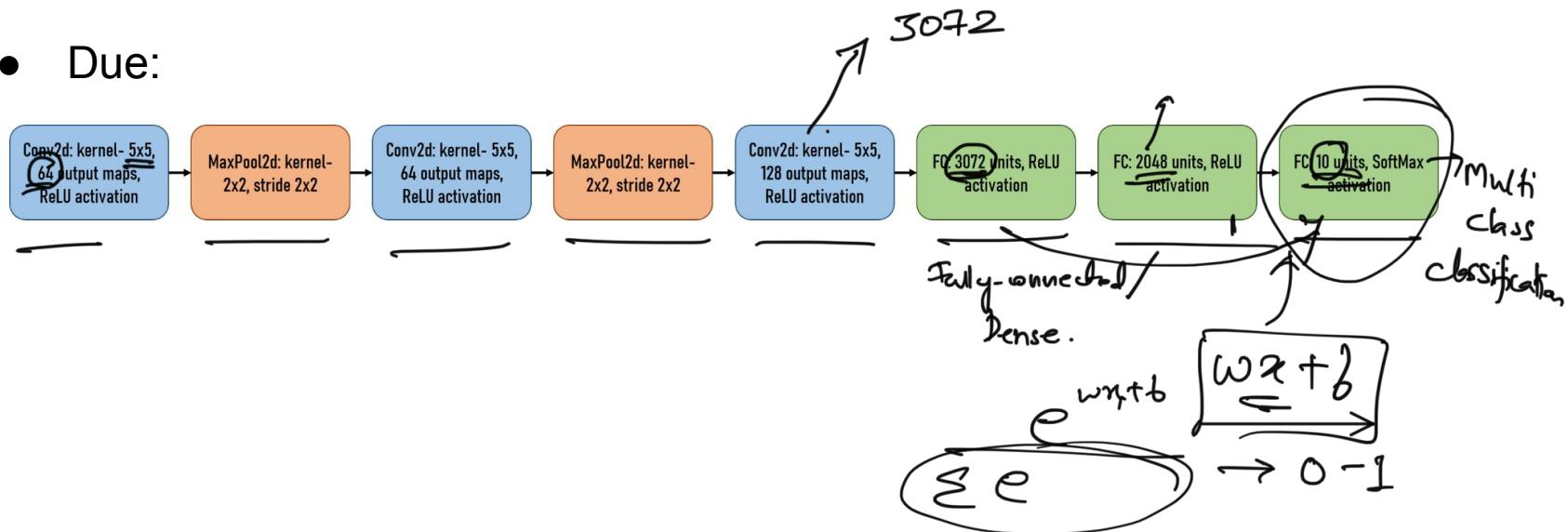
Table of contents

- 1. Project Part 3**
- 2. Autoencoder**
- 3. Recurrent Neural Network**

Project Part 3



- Goal - train a small convolutional neural network for classification of SVHN dataset.
- Due:

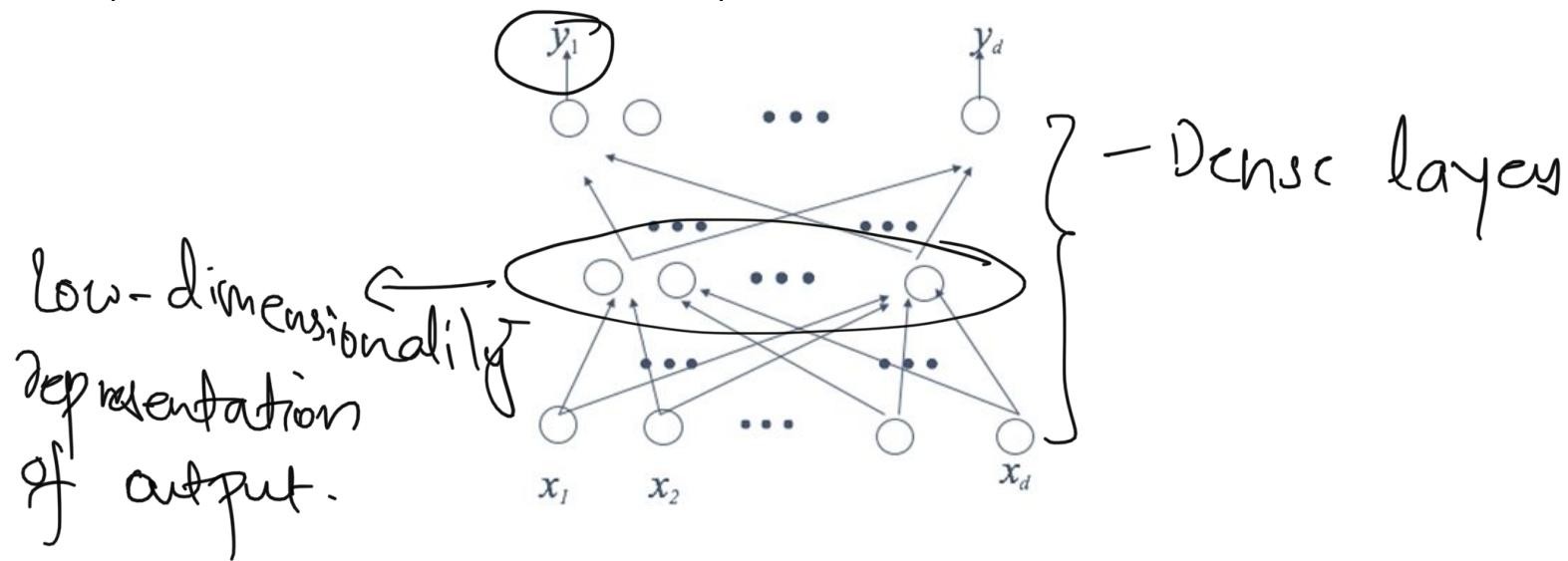


Project Part 3 - Software & Libraries

- Language - Python
- Deep learning libraries - Keras, Tensorflow, PyTorch
- Deliverables-
 - Code with proper comments.
 - A report including the plots for the learning/testing errors and the final classification accuracy.
 - Please submit the code and the report as separate files on Canvas. **Do not zip them.**

Autoencoder

- Train a network without supervision
- y_i being an approximation of x_i

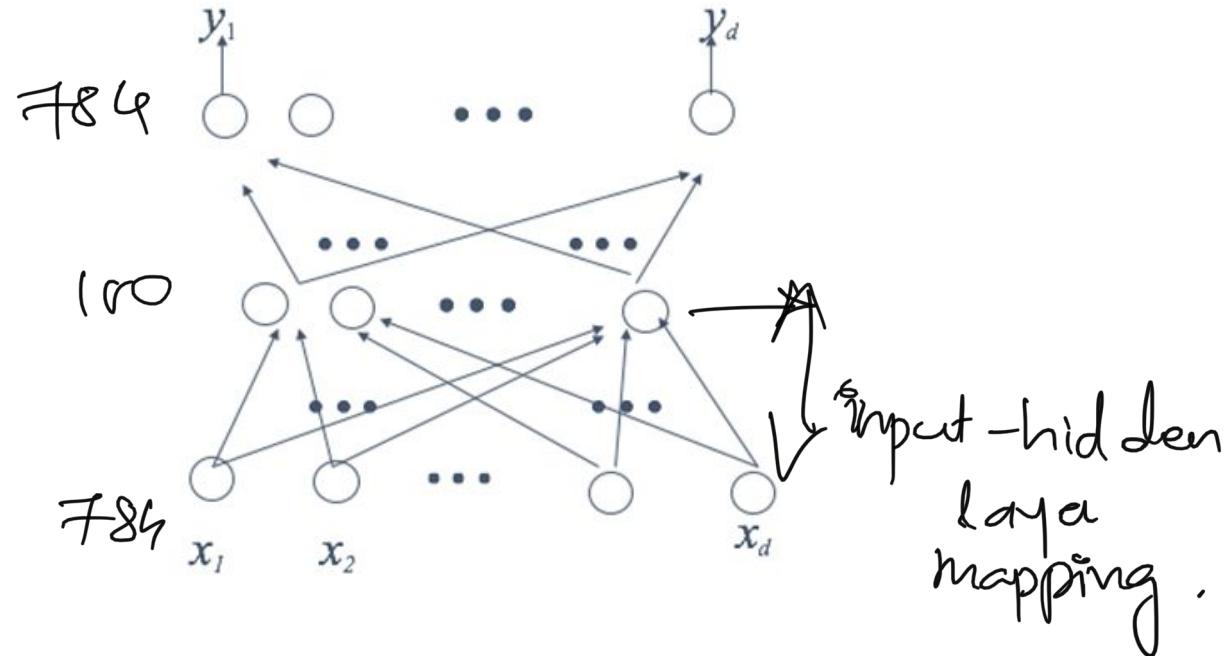


Autoencoder

- Perfect Autoencoder - Reconstructed output will be exactly same as input.
- Practical Scenario -
 - There is error in reconstruction.
 - Learning objective - Learn such that the reconstruction error is minimized.
 - MSE loss.

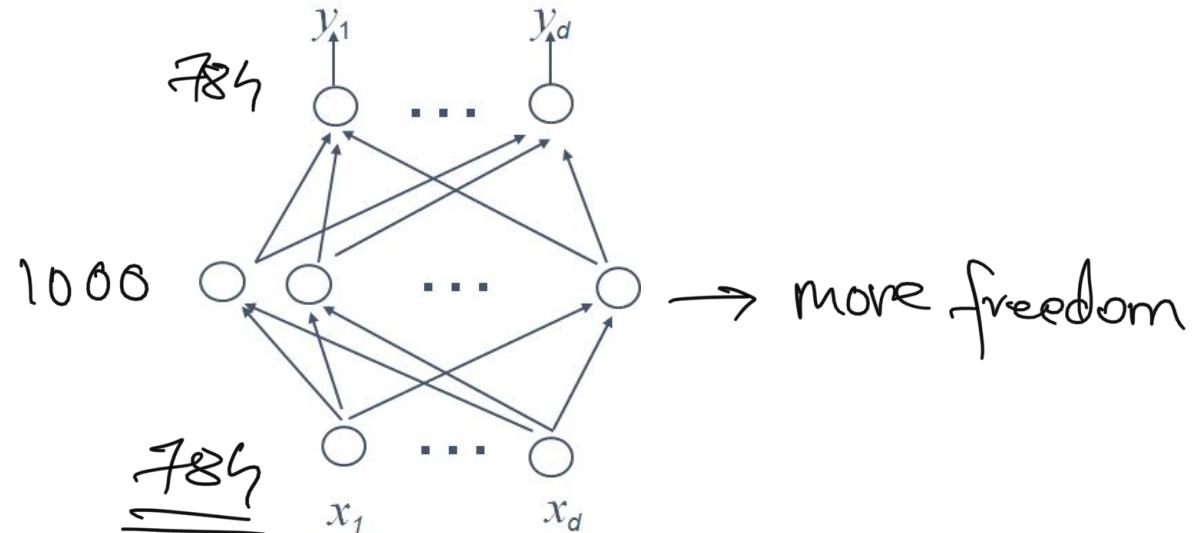
Autoencoder - 2 cases

- Case 1- Much fewer hidden nodes than input nodes



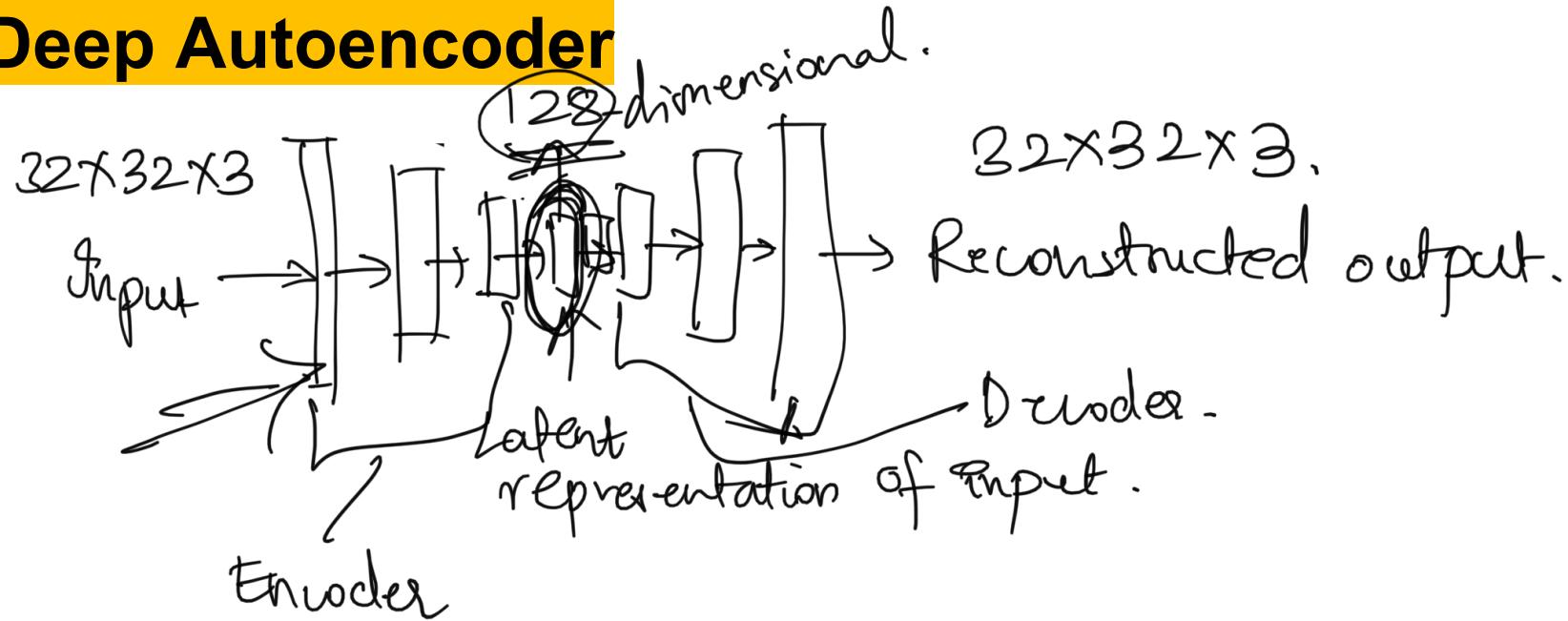
Autoencoder - 2 cases

- Case 1- Allow more hidden nodes than input



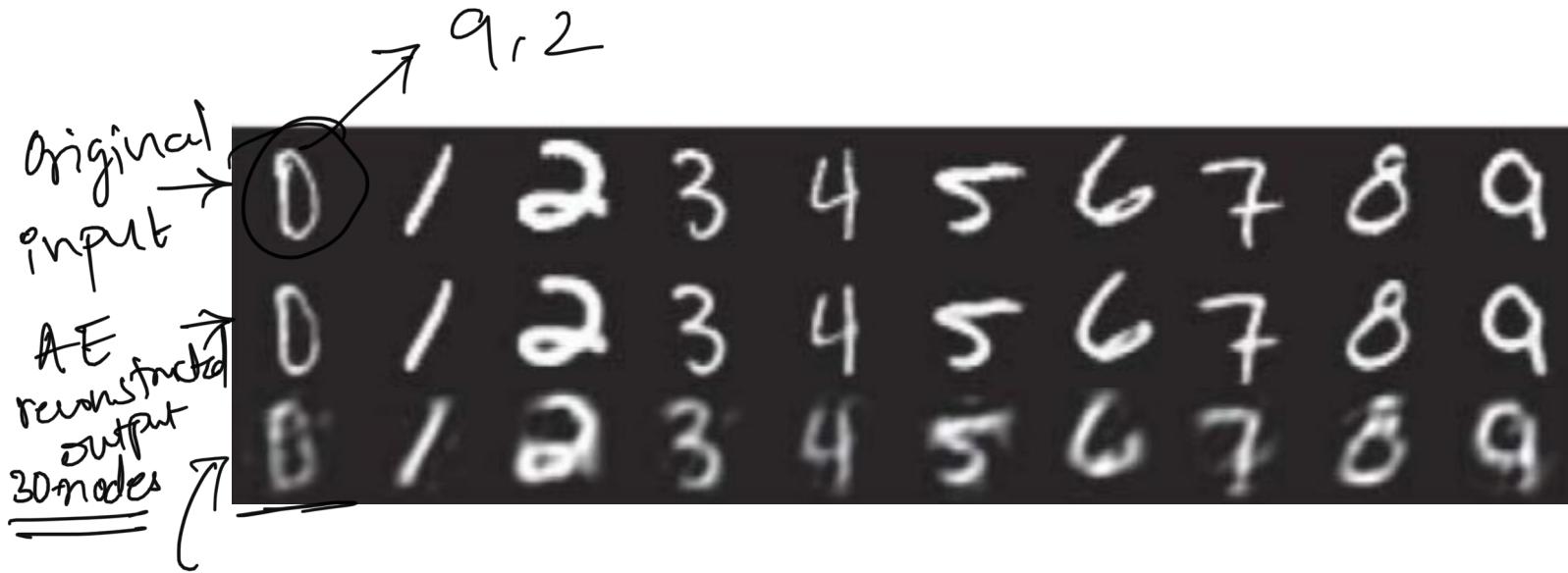
— Regularization .

Deep Autoencoder



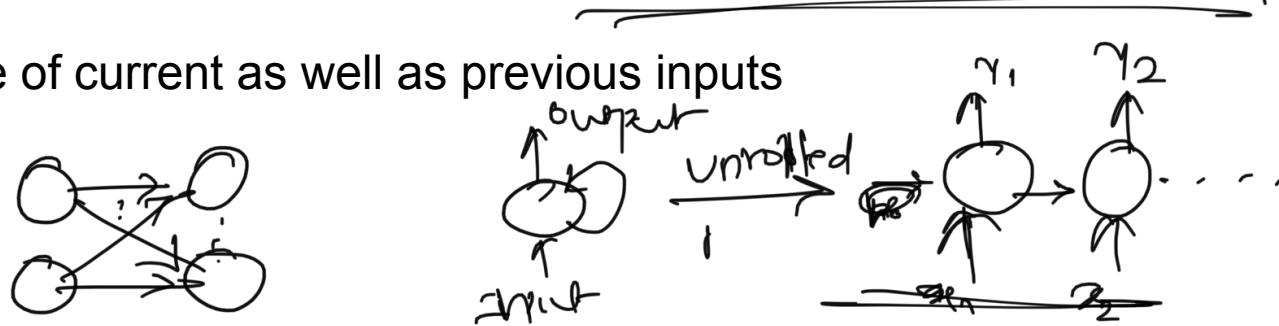
- Encoder \Rightarrow Convolution
- Decoder \Rightarrow Deconvolution

Deep Autoencoder



Recurrent Neural Networks

- Allow directed cycles in connections between neurons
- Such networks could naturally model variable-length sequential data
- Influence of current as well as previous inputs



- Very popular in NLP (Natural Language Processing)
- Vision — to process videos.

LSTM, GRU

Questions?

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

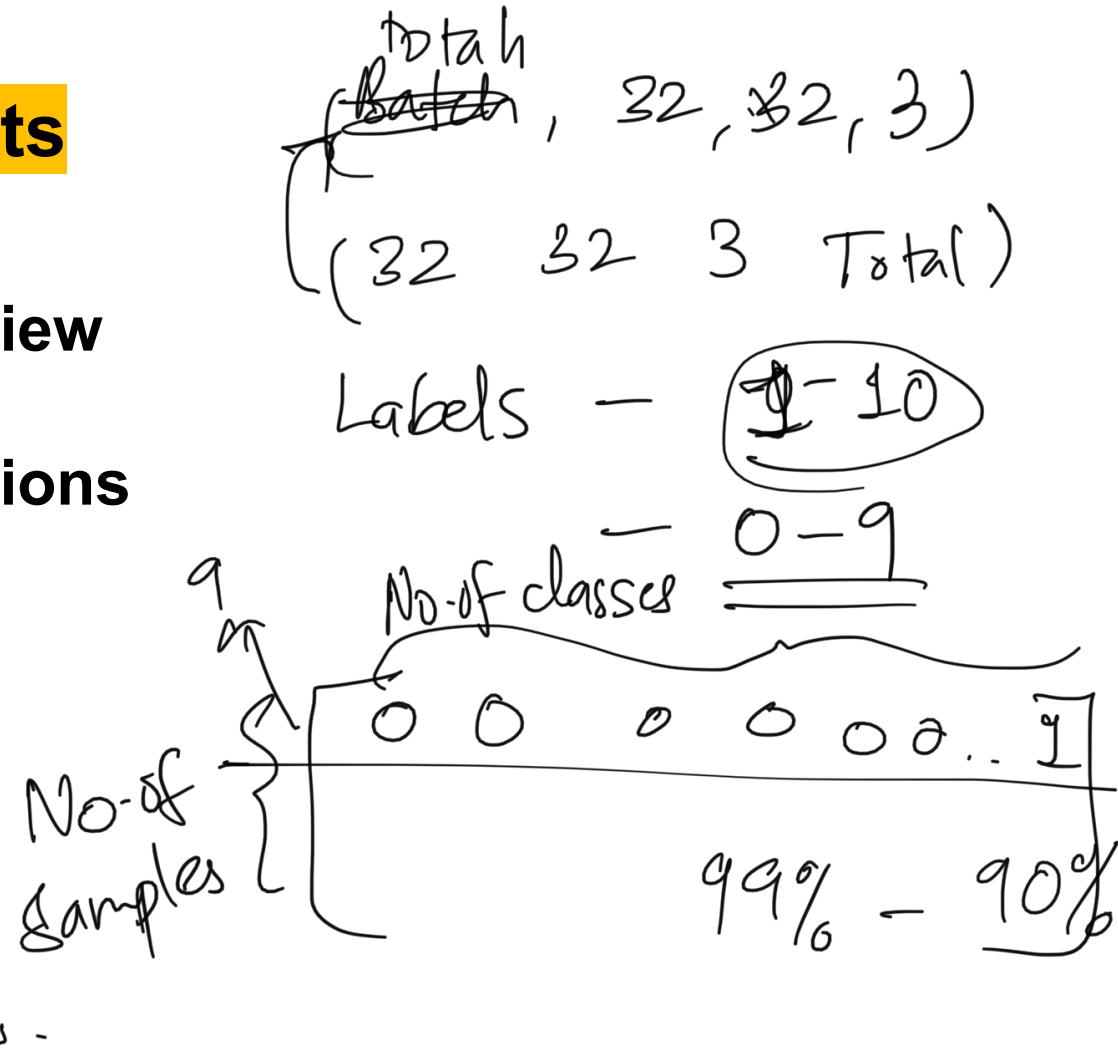
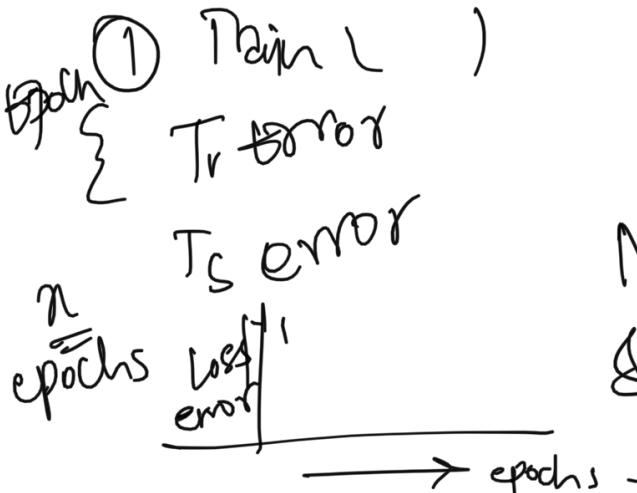
Final Exam Review



Table of contents

1. Concepts Review

2. Sample Questions



Maximum Likelihood Estimation

- Given some training data and assuming a parametric model $p(x|\theta)$; what specific θ will fit/explain the data best?
- To consider all the samples denoted by $D=\{x_1, x_2, \dots, x_n\}$, assume that all the samples are i.i.d - independent and identically distributed.
- So, data likelihood represented by $L(\theta)$ is -

$$L(\theta) = P(D|\theta) = \prod_i P(x_i|\theta)$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} p(x|\theta)$$

Naive Bayes

- The "naive" conditional independence assumption: each feature is (conditionally) independent of every other feature, given the label, i.e.,
 $p(x_i | \{x_j \text{ for any } j \neq i\}, y) = p(x_i | y)$
- The predicted label is given by -

$$\hat{y} = \operatorname{argmax}_y P(y) \prod_{i=1}^d p(x_i | y)$$

Linear Regression

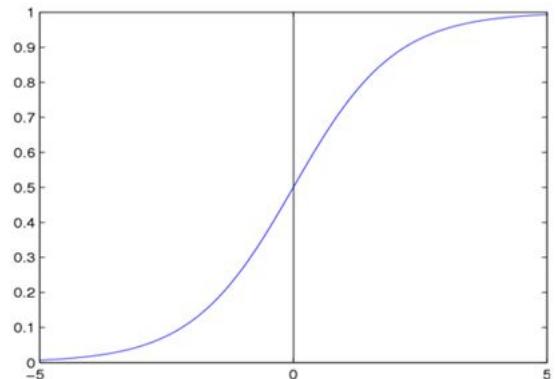
- Regression - A training set of n samples $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$ where $y^{(i)}$ is a continuous “label” (or target value) for $\mathbf{x}^{(i)}$
- Linear regression - modeling the relation between y and x via a linear function $y \approx w_0 + w_1x_1 + \dots + w_dx_d = \mathbf{w}^t\mathbf{x}$
- The error is given as - $\|\mathbf{e}\|^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$

Logistic Regression

- Training set: n labelled samples $\langle \mathbf{x}(i), y(i) \rangle$
- Use the logistic function for modeling $P(y|x)$, considering only the case of $y \in \{0,1\}$

$$P(y = 0|\mathbf{x}) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}$$

$$P(y = 1|\mathbf{x}) = \frac{\exp(w_0 + \sum_{i=1}^d w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}$$



$$\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t}$$

Support Vector Machines (SVM)

- Key idea - To find the decision boundary such that the margin is maximized.
- Data - $\langle x^{(i)}, y^{(i)} \rangle$, $y^{(i)} \in \{-1, 1\}$, $x^{(i)} \in R^d$, for all $i=1, \dots, n$
- Plane equations-
- Margin -

Hidden Markov Model (HMM)

- Dynamic Bayesian network - modeling the process indexed by time.
- Two assumptions -
 - First-order Markov chain - $P(s^t=S_j | s^{t-1}=S_i)$
 - $a_{ij} = P(s^t=S_j | s^{t-1}=S_i), 1 \leq i, j \leq N$, for any t
- Specifying HMM - Θ , Ω , π , A, B

Problems in HMM

- For a given HMM $\Lambda = \{\Theta, \Omega, A, B, \pi\}$
 - Estimation of model parameters
 - Given an observation sequence $O = \{o^1, o^2, \dots, o^k\}$, what is the most likely state sequence $S = \{s^1, s^2, \dots, s^k\}$ that has produced O ?
 - How likely is an observation O ?

K-Means Clustering

Given: n samples, a number k.

Begin

~~Initialize~~ $\mu_1, \mu_2, \dots, \mu_k$ (randomly selected)

do classify n samples according to
nearest μ_i
clusters ←—————
recompute μ_i
until no change in μ_i

return $\mu_1, \mu_2, \dots, \mu_k$

End

EM algorithm for GMM

Expectation
Maximization

Gaussian
Mixture
Models

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \quad (9.23)$$

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\begin{cases} \boldsymbol{\mu}_k^{\text{new}} \\ \boldsymbol{\Sigma}_k^{\text{new}} \\ \pi_k^{\text{new}} \end{cases} = \begin{cases} \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \\ \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T \\ \frac{N_k}{N} \end{cases} \quad (9.24)$$

$$\begin{cases} \boldsymbol{\mu}_k^{\text{new}} \\ \boldsymbol{\Sigma}_k^{\text{new}} \\ \pi_k^{\text{new}} \end{cases} = \begin{cases} \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T \\ \frac{N_k}{N} \end{cases} \quad (9.25)$$

$$\begin{cases} \boldsymbol{\mu}_k^{\text{new}} \\ \boldsymbol{\Sigma}_k^{\text{new}} \\ \pi_k^{\text{new}} \end{cases} = \begin{cases} \frac{N_k}{N} \end{cases} \quad (9.26)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}). \quad (9.27)$$

4. Evaluate the log likelihood

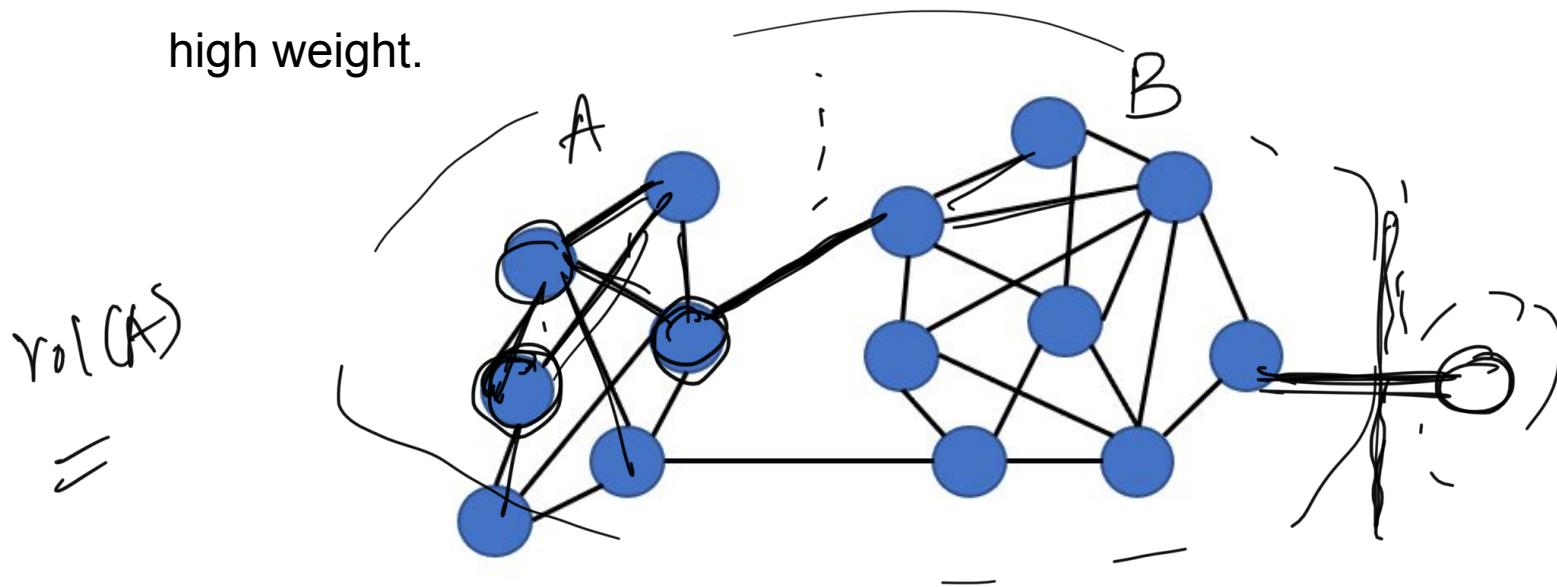
$$\ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (9.28)$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

Bishop's book.

Clustering as a Graph Partition

- Find a partition of a graph such that the edges between different groups have a very low weight while the edges within a group have high weight.



Types of cuts

$$\mathcal{M}\text{in}\text{cut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k W(A_i, \bar{A}_i)$$
$$\text{RatioCut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{|A_i|} = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{|A_i|}$$

nb. of nodes in cluster

$$\text{Ncut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{\text{vol}(A_i)} = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)}$$
$$\text{MinMaxCut}(A_1, \dots, A_k) := \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{W(A_i, A_i)}$$

PCA Algorithm

— Dimensionality Reduction
~~d dimensions~~.

1. Compute the $d \times d$ sample covariance matrix C
2. Find the eigenvalues and corresponding eigenvectors of C
3. Project the original data onto the space spanned by the eigenvectors
 - The projection may be done onto a d' -dimensional subspace spanned by the first d' eigenvectors (ordered by the eigenvalue in descending order)
 - d' is determined by the desired accuracy

Learning in Perceptron

Iterate for t until a stop criterion is met

{

for each sample x_i with label y_i :

{

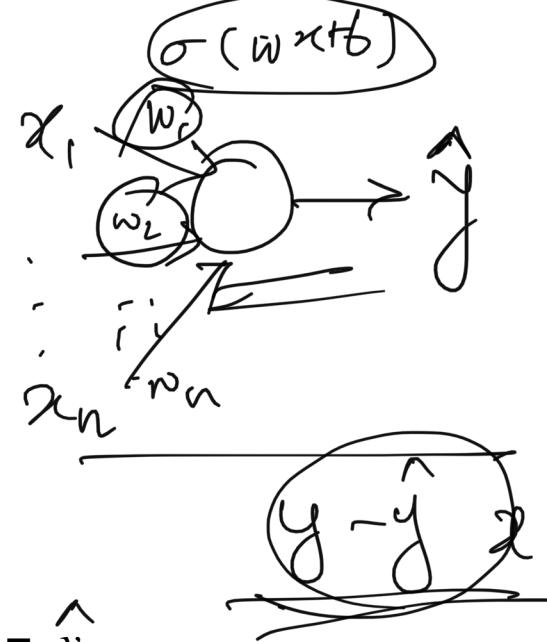
compute the output (y_i) of the network

estimate the error of the network $e(w(t)) = y_i - \hat{y}_i$

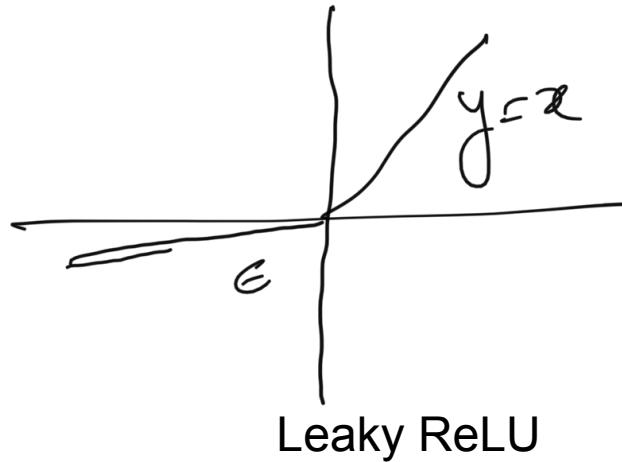
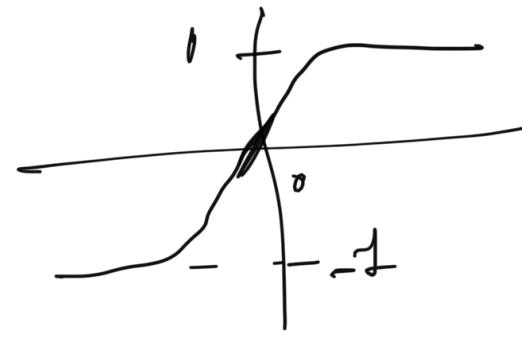
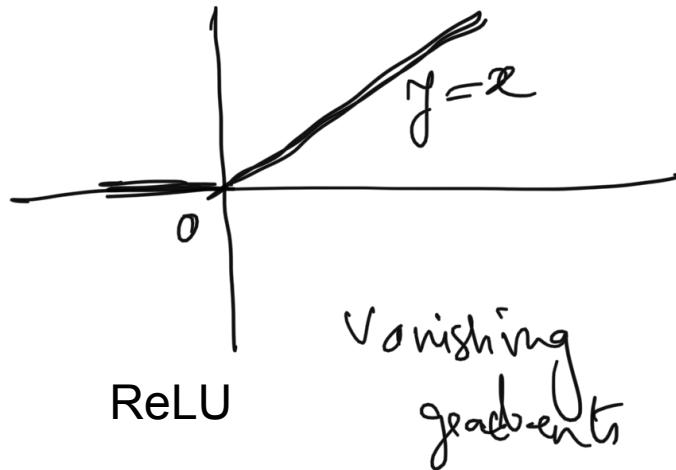
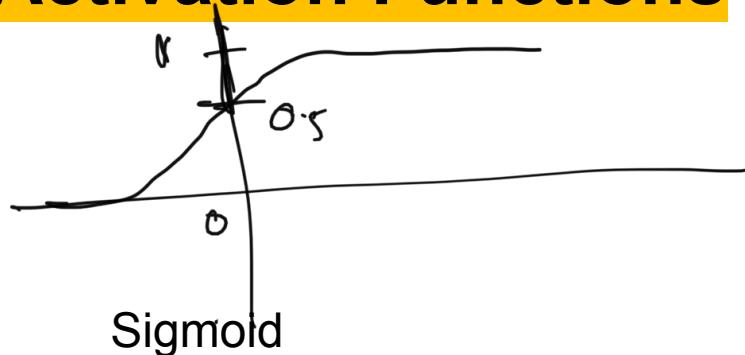
update the weight $w(t+1) = w(t) + e(w(t))x_i$

}

$t++$



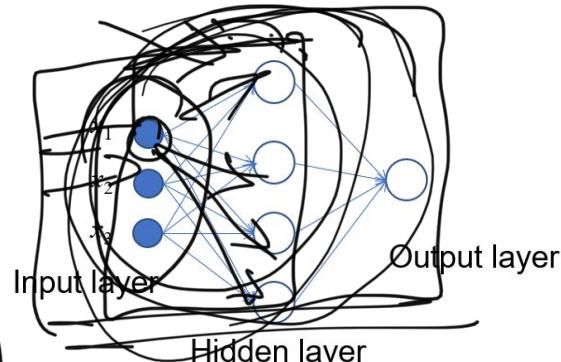
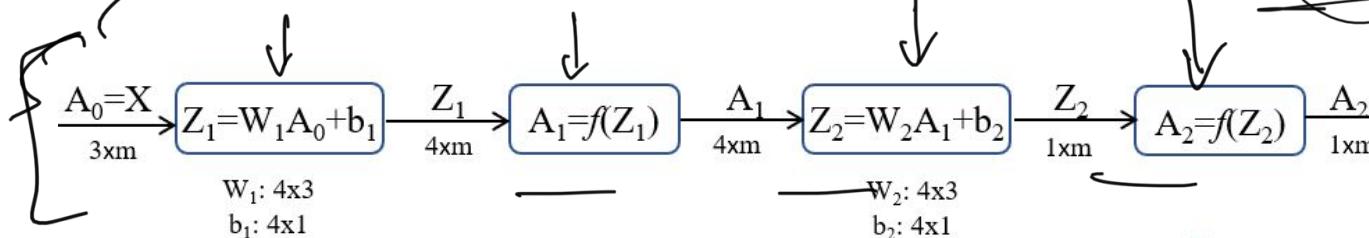
Activation Functions



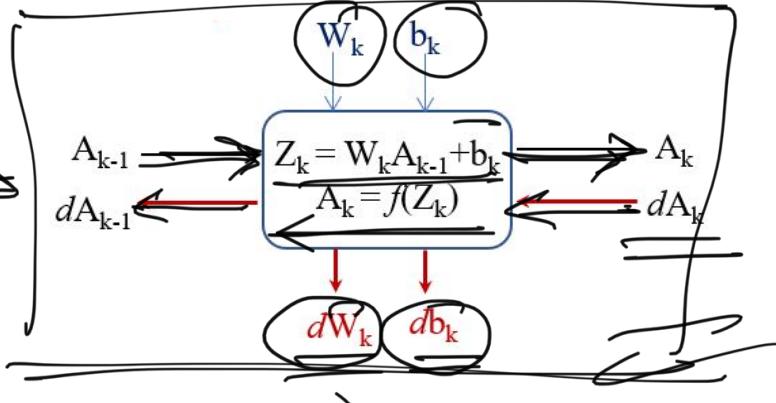
Multi-layer perceptron

MLP Learning

Forward pass



Backpropagation
phase.



3X3X4

Convolutional Neural Network (CNN)

- Basic building block is the convolution for image processing

- Weight sharing - reduces parameters

- Learn the kernels based on the task.

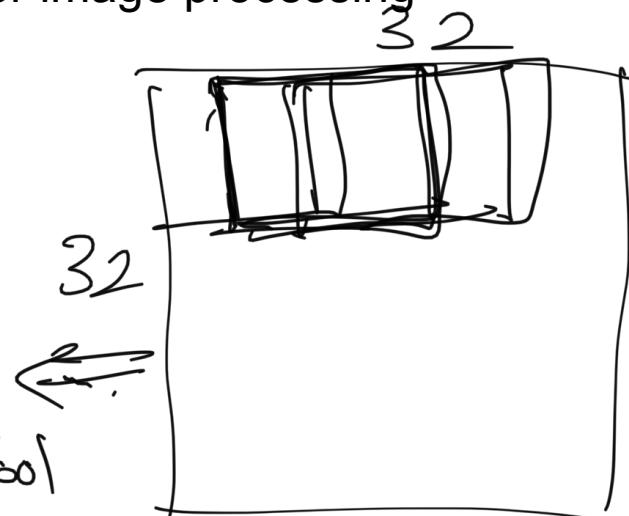
- Layers -

convolutional layers.

Pooling layer. - MaxPool

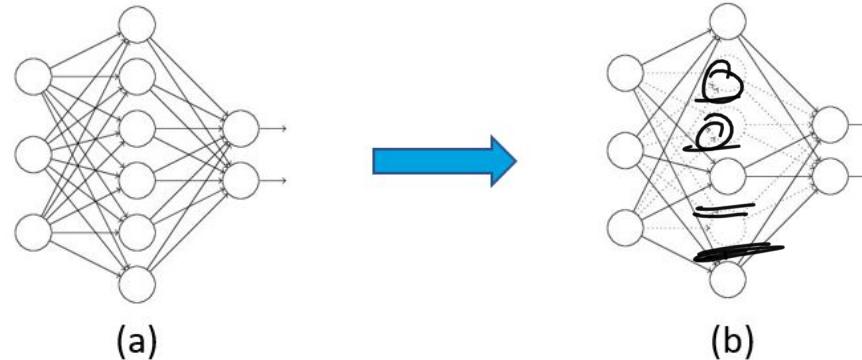
Activation layers - ReLU, Leaky ReLU.

Fully-connected - end of the network



Dropout

- Obtain (b) by randomly deactivate some hidden nodes in (a).
- Reducing co-adaptation of neurons



Batch Normalization

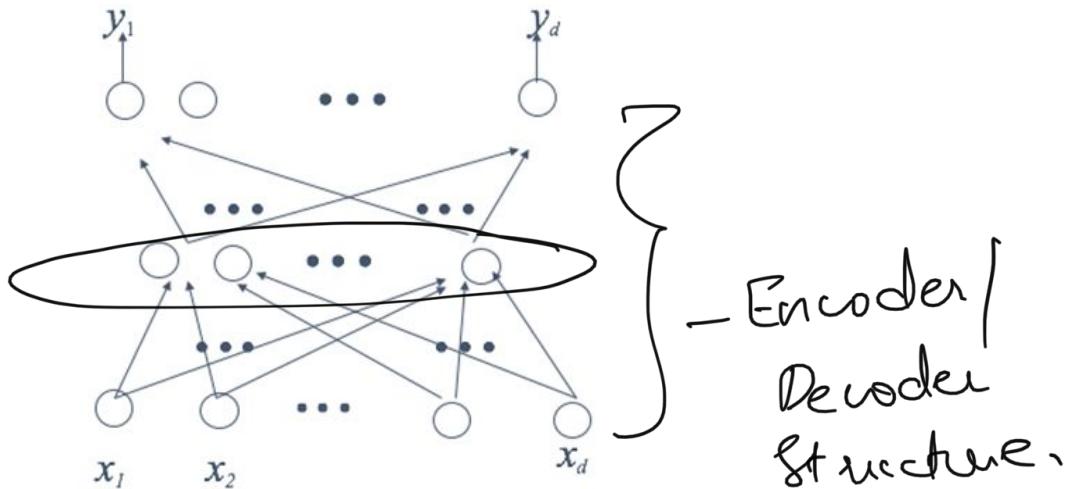
- Normalizes layer inputs of a batch

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad \sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad \hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$
$$y_i \leftarrow \underline{\gamma} \hat{x}_i + \underline{\beta} \equiv \text{BN}_{\gamma, \beta}(x_i)$$

Autoencoder

— Unsupervised.

- Train a network without supervision
- y_i being an approximation of x_i



Questions?