

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

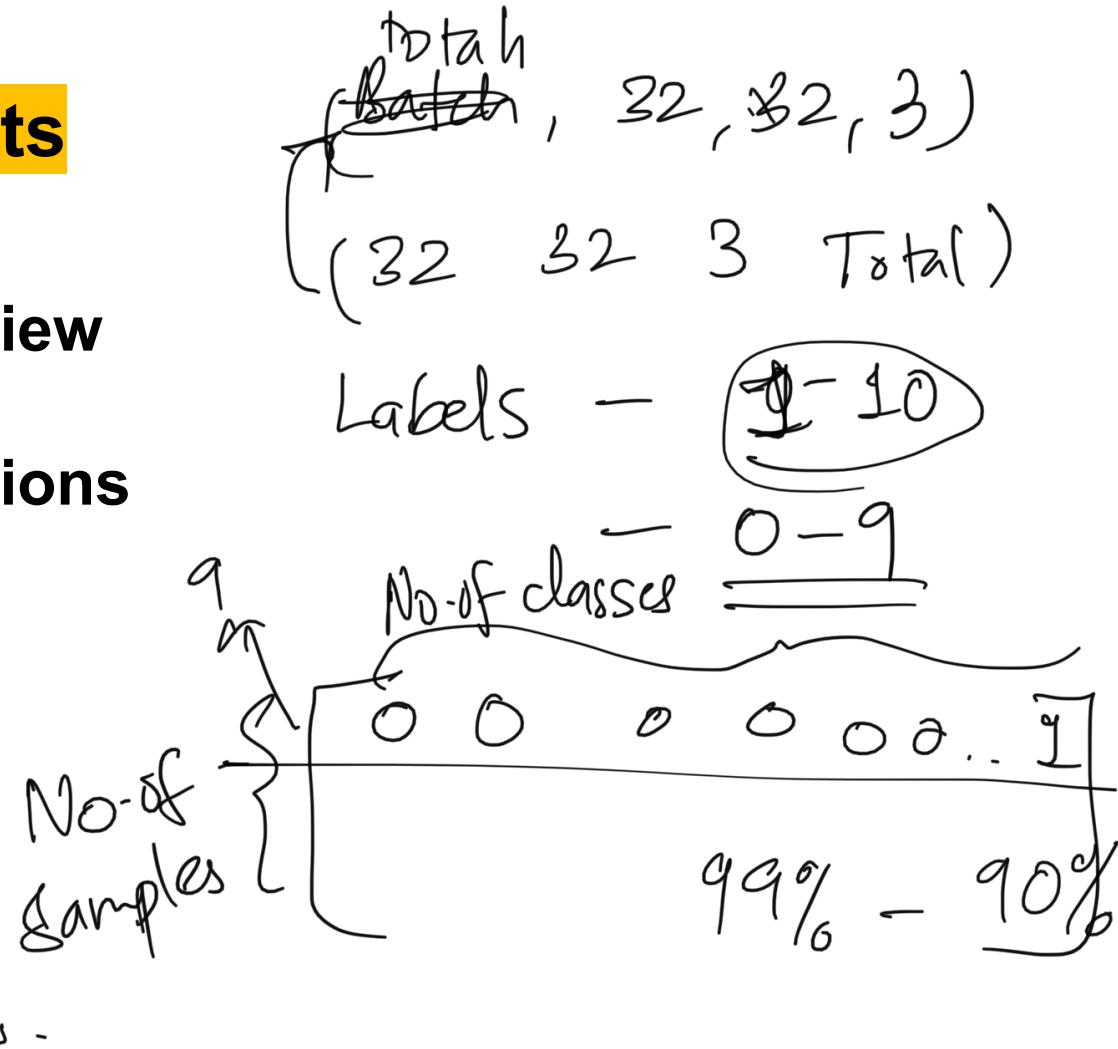
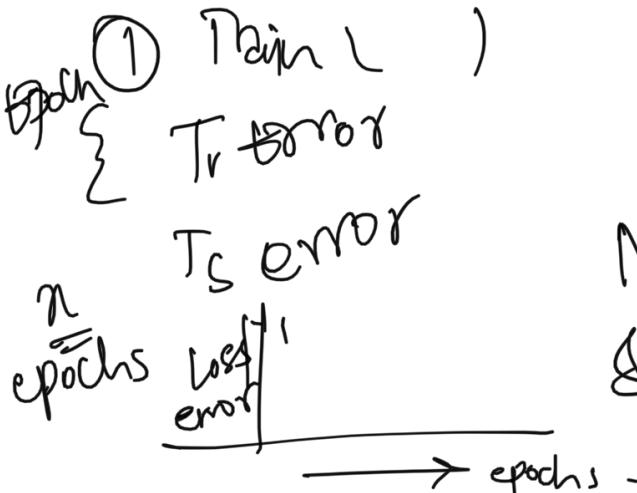
# Final Exam Review



# 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  $\theta$  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

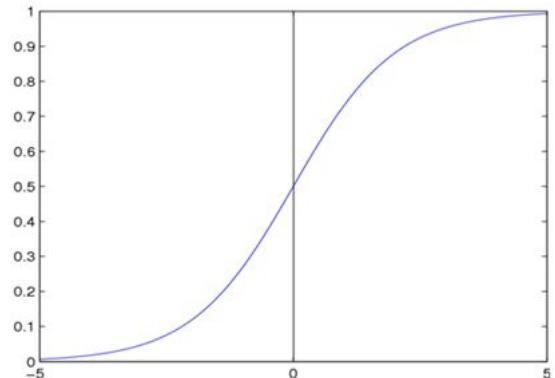
- Regression - A training set of n samples  $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$  where  $y^{(i)}$  is a continuous “label” (or target value) for  $\mathbf{x}^{(i)}$
- Linear regression - modeling the relation between  $y$  and  $x$  via a linear function  $y \approx w_0 + w_1x_1 + \dots + w_dx_d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

# 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  $O$ ?
  - 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  $\mu_i$   
*clusters* ←

recompute  $\mu_i$

until no change in  $\mu_i$

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

End

# EM algorithm for GMM

*Expectation  
Maximization*

Gaussian  
Mixture  
Models

1. Initialize the means  $\mu_k$ , covariances  $\Sigma_k$  and mixing coefficients  $\pi_k$ , 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

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

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

$$\begin{cases} \boldsymbol{\mu}_k^{\text{new}} \\ \boldsymbol{\Sigma}_k^{\text{new}} \\ \pi_k^{\text{new}} \end{cases} = \begin{cases} \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \\ \frac{N_k}{N} \end{cases} \quad (9.26)$$

where

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

Bishop's book.

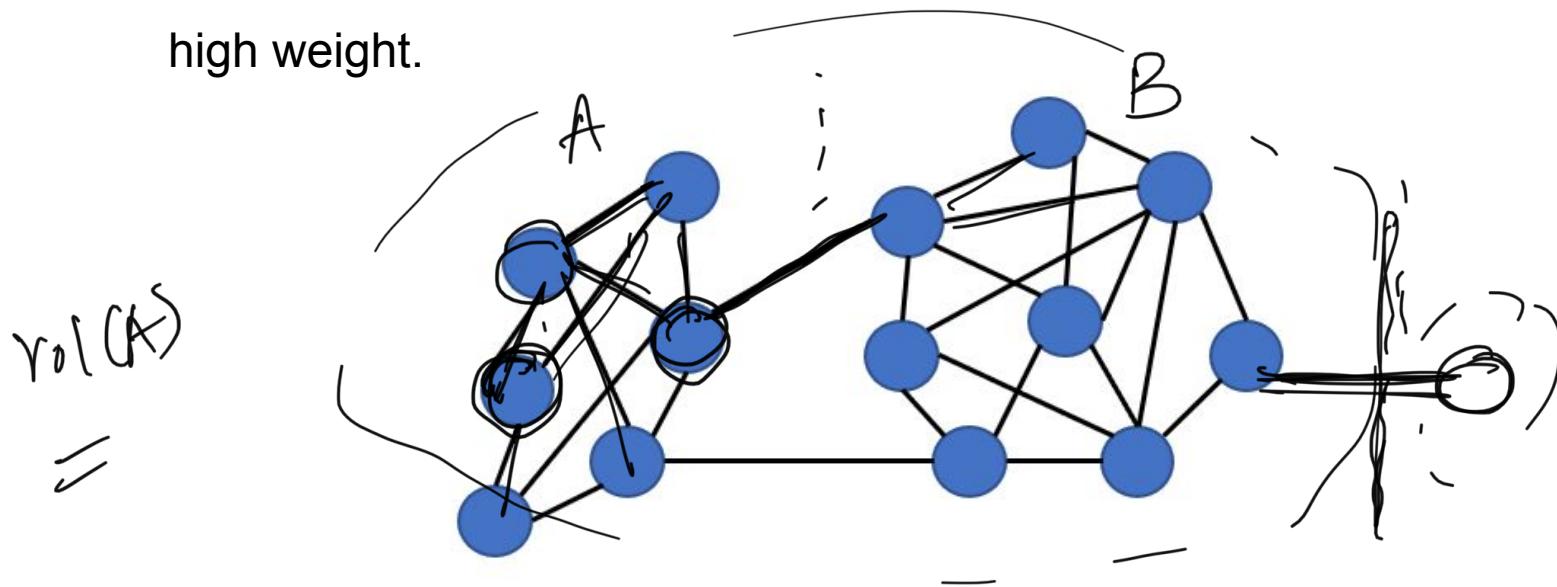
4. Evaluate the log likelihood

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

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

# 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  $y_i$ :

{

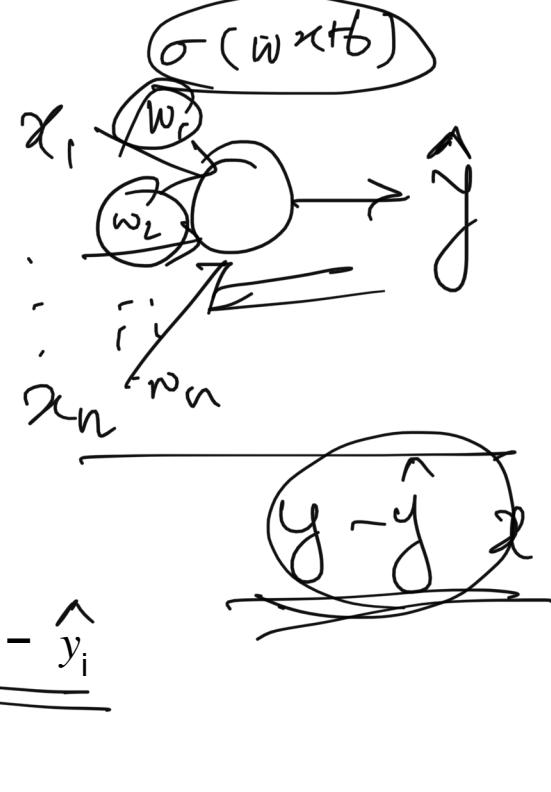
compute the output ( $y_i$ ) of the network

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

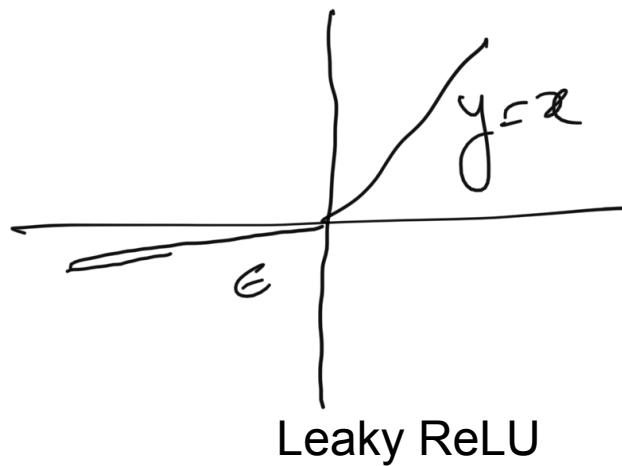
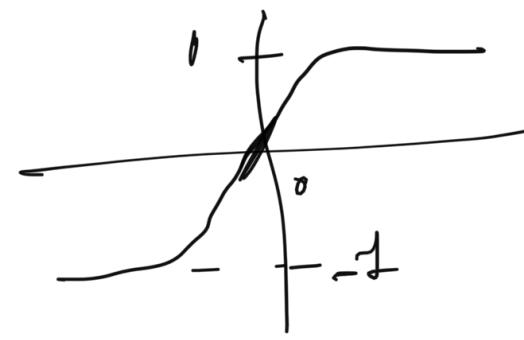
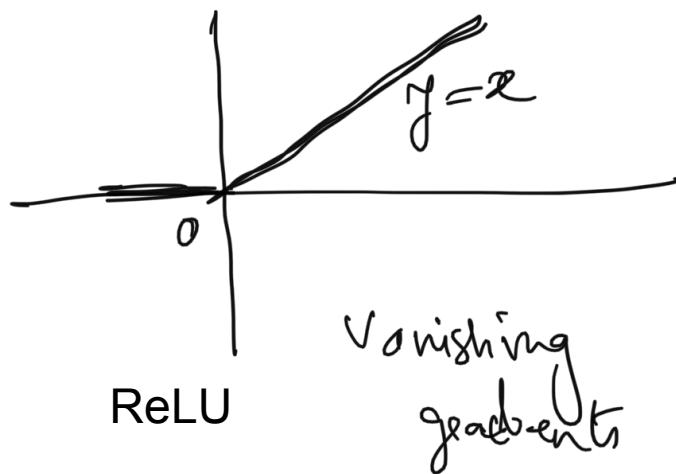
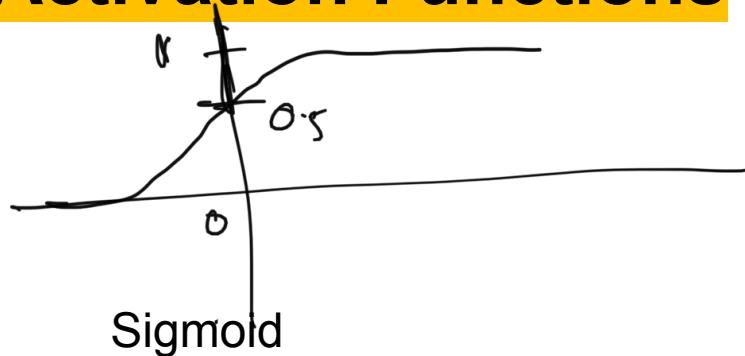
update the weight  $w(t+1) = w(t) + e(w(t))x_i$

}

$t++$



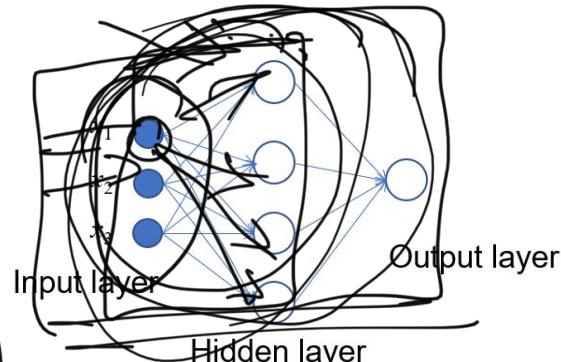
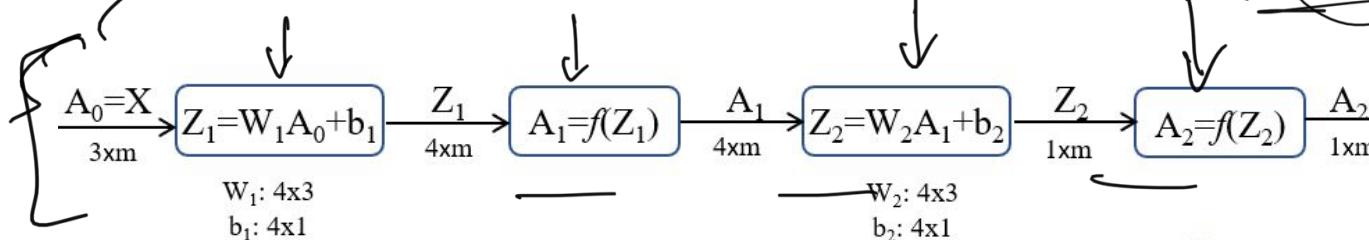
# Activation Functions



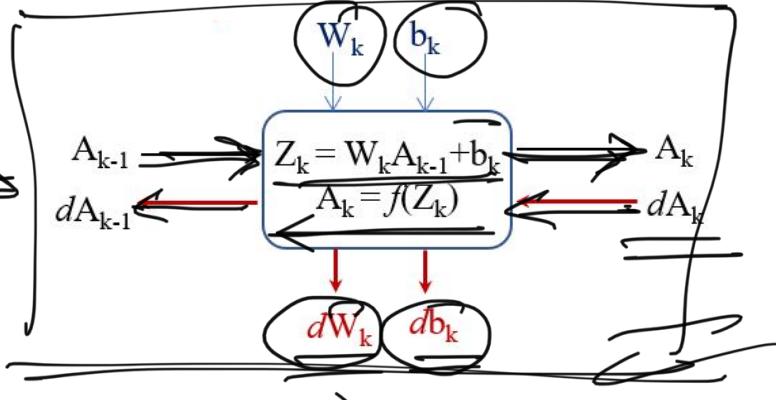
Multi-layer perceptron

## 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.

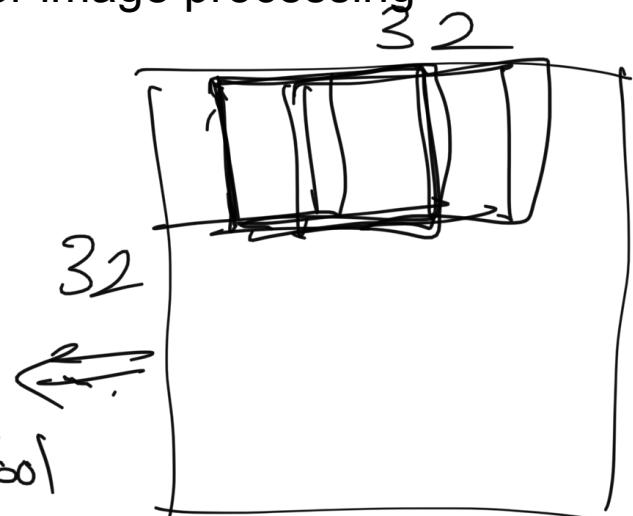
- Layers -

convolutional layers.

Pooling layer. - MaxPool

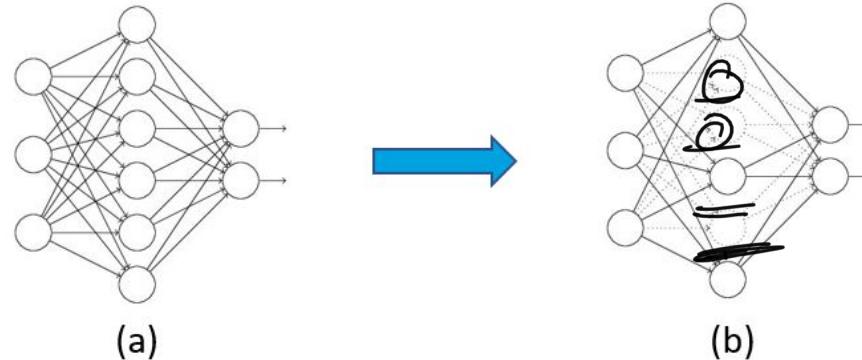
Activation layers - ReLU, Leaky ReLU.

Fully-connected - end of the network



# Dropout

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



# Batch Normalization

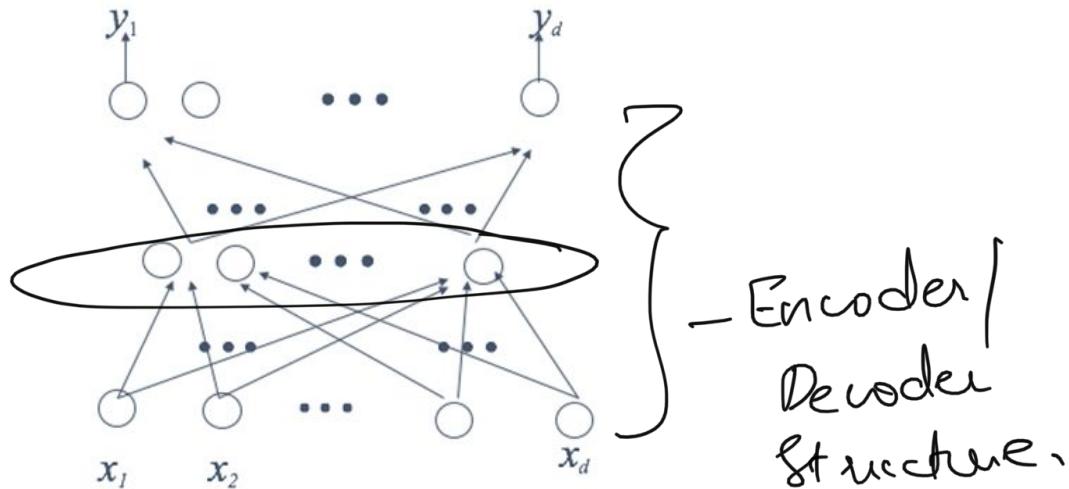
- Normalizes layer inputs of a batch

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

# Autoencoder

— Unsupervised.

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



# Questions?