

# ECE421

# Intro. to Machine Learning

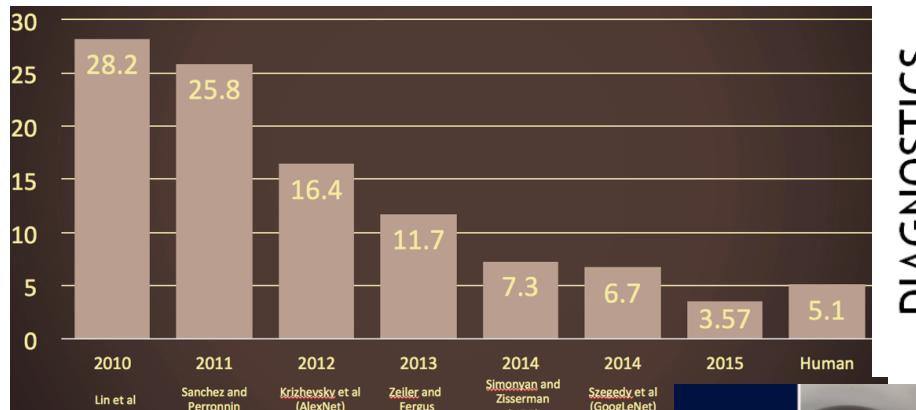
Fall 2019

Ashish Khisti

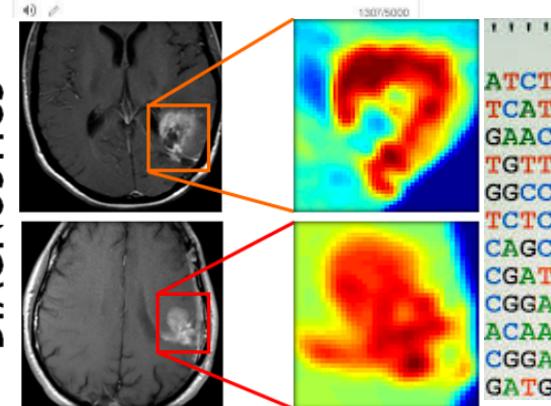
# ImageNet Challenge

IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



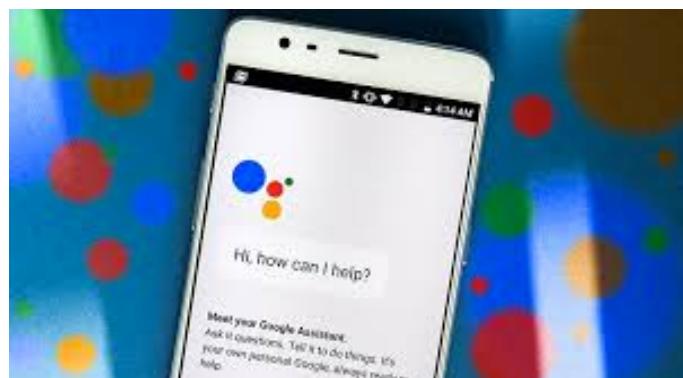
## DIAGNOSTICS



ATCTCTTGGCTCCAGCATCGATGAAGAACGCA  
TCATTTAGAGGAAGTAAAAAGTCGTAACAAGGT  
GAACCTGTCAAAACTTTAAACACGGATCTCTT  
TGTTGCCTTCGGCGCGCCCCAAGGGTGCCCCG  
GGCCTGCCGTGGCAGATCCCCAACGCCGGGCG  
TCTCTTGGCTCCAGCATCGATGAAGAACGCA  
CAGCATCGATGAAGAACGCGAGCAGCAG  
CGATACCTCTGAGTGTTCTTAGCGAACTGTCA  
CGGATCTCTGGCTCCAGCATCGATGAAGAAC  
ACAACGGATCTCTGGCTCCAGCATCGATGA  
CGGATCTCTGGCTCCAGCATCGATGAAGAAC  
GATGAAGAACGCGAGCAGCAGCAGATATGTAAT



LEE SEDOL  
00:01:00



## • Natural language processing

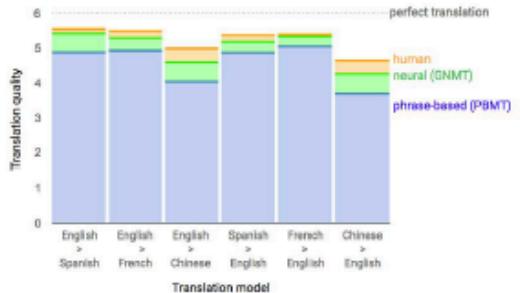
English Spanish French Spanish - detected Translate

Capítulo primero. Que trata de la condición y ejercicio del famoso hidalgo don Quijote de la Mancha

En un lugar de la Mancha, de cuyo nombre no quiero acordarme, no ha mucho tiempo que vivía un hidalgo de los lanza en astillero, adarga antigua, rocin flaco y galgo corredor. Una olla de algo más vaca que carnero, salpicón las más noches, duelos y quebrantos los sábados, lantejas los viernes, algún palomino de arriachura los domingos, consumían las tres partes de su hacienda. El resto della concluian sayo de velante, cañas de velludo para las fiestas, con sus puntiflos de lo mismo, y los días de entresemana se honraba con su vellorí de lo más fino. Tenía en su casa una ama que pasaba de los cuarenta, y una sobrina que no llegaba a los veinte, y un mozo de campo y plaza, que así ensillaba el rocin como tomaba la podadera. Frisaba la edad de nuestro hidalgo con los cincuenta años; era de complejión recha, seco de carne, enjuto de rostro, gran madrugador y amigo de la caza. Quieren decir que tenía el sobrenombre de Quijada, o Quesada, que en esto hay alguna diferencia en los autores que deste caso escriben; aunque, por conjecturas verosímiles, se deja entender que se llamaba Quijana. Pero esto importa poco a nuestro cuento; basta que en la narración dél no se salga un punto de la verdad.

First chapter. Which deals with the condition and exercise of the famous nobleman Don Quixote de la Mancha

In a place of La Mancha, whose name I do not want to remember, there has not been a long time that lived a lord of the lance in shipyard, old pork, thin rocin and greyhound runner. A pot of something more cow than ram, split most nights, duels and breaks on Saturdays, giblets on Fridays, some palomino in addition to Sundays, consumed the three parts of his estate. The rest of the party concluded a velvet dress, hairy tights for the parties, with their slippers of the same, and the days of midweek were honored with their vellorí of the finest. He had a housekeeper in his house who was in his forties, and a niece who was not in his twenties, and a boy in the country and square, who saddled the rocin as he took the pruning. He emphasized the age of our hidalgo at the age of fifty: Was of a hard complexion, dry of flesh, thin of face, great early bird and friend of the hunt. They mean that he had the nickname of Quijada, or Quesada, that in this there is some difference in the authors who in this case write; Although, by plausible conjectures, it is understood that it was called Quejana. But this matters little to our story; It is enough that in the narration of him a point of truth does not come out.

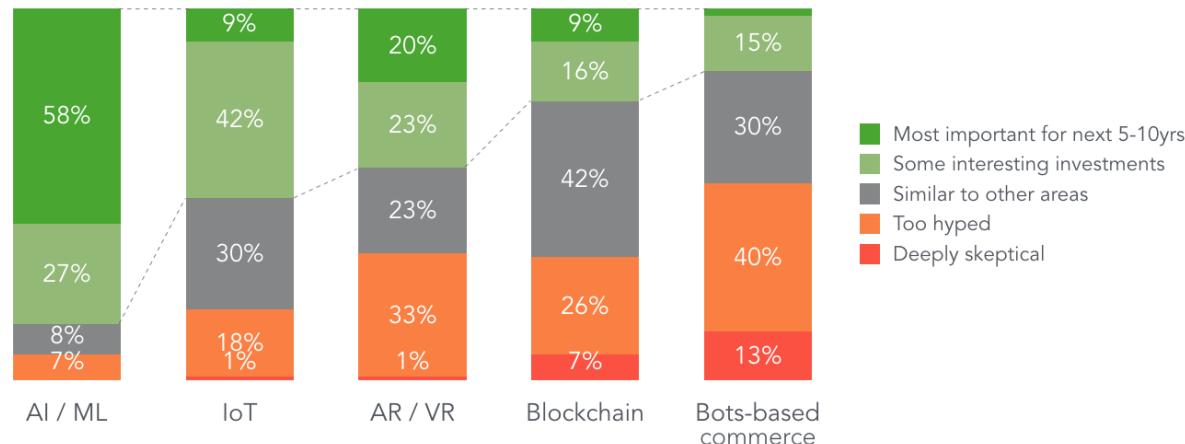


## Google's Neural Machine Translation (Wu et al. 2016)



Most VCs are most excited about AI & Machine Learning as their most important investment theme for the coming 5-10 years.

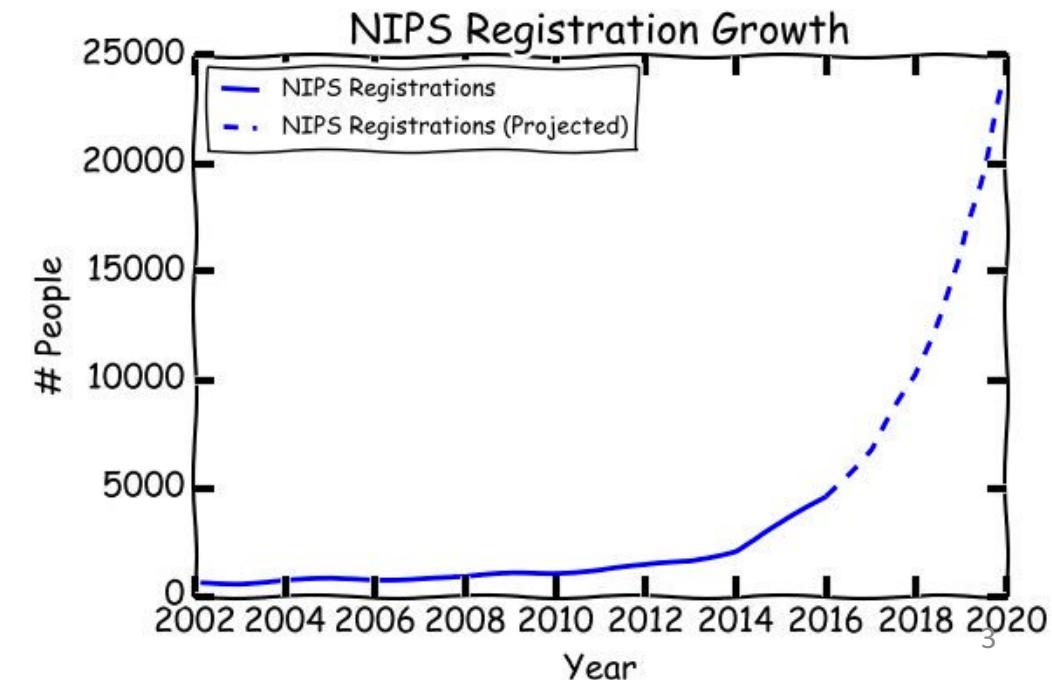
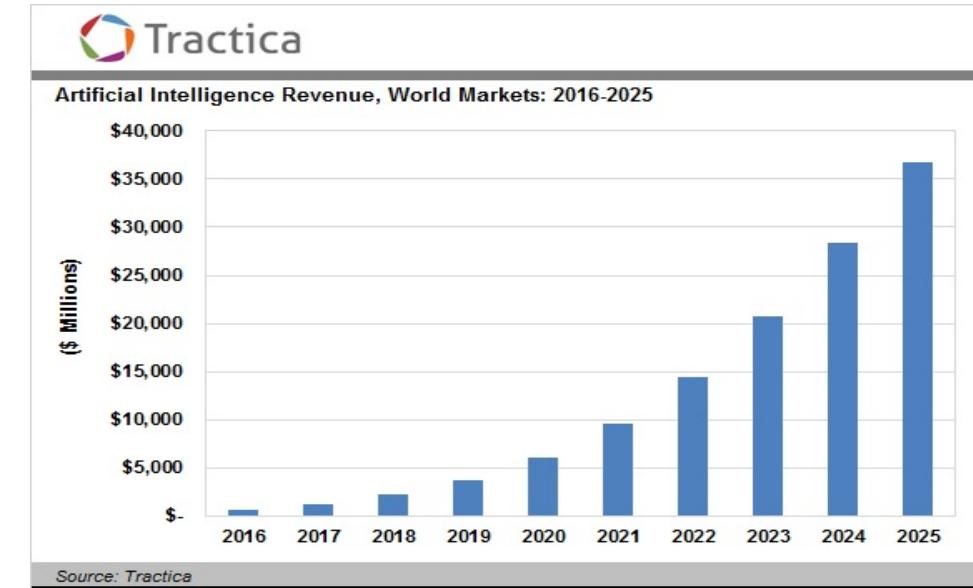
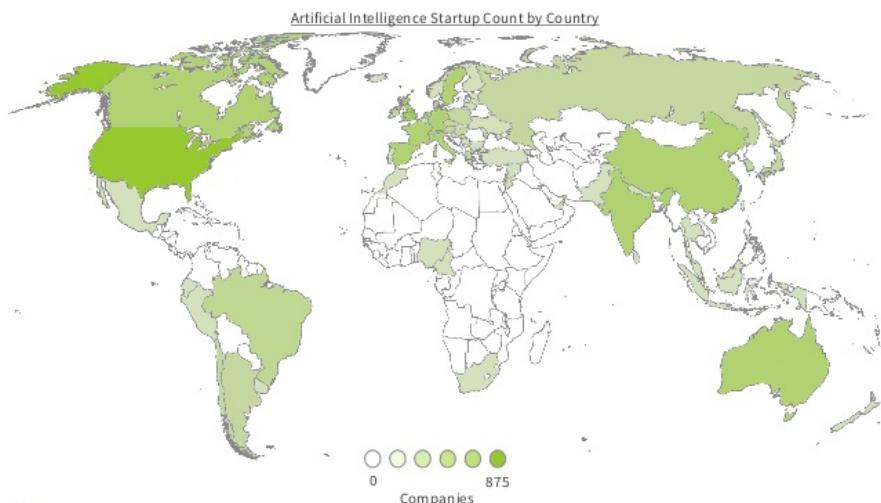
*Q. How do you feel about the following investment areas?*



17 Source: Upfront Ventures survey of VCs (N=115), Jan 2017

**upfront**  
VENTURES

Artificial intelligence startups are a global phenomenon



# Machine Learning Jargon

semi-supervised learning	<b>overfitting</b>	stochastic gradient descent	<b>SVM</b>	<i>Q</i> learning
Gaussian processes	<b>deterministic noise</b>			
<i>distribution-free</i>	<i>linear regression</i>	VC dimension	<b>data snooping</b>	learning curves
collaborative filtering	nonlinear transformation			<b>mixture of experts</b>
<b>decision trees</b>	<i>RBF</i>	<i>training versus testing</i>	<i>neural networks</i>	<i>no free lunch</i>
<b>active learning</b>	linear models	bias-variance tradeoff	noisy targets	<i>Bayesian prior</i>
<i>ordinal regression</i>	cross validation	logistic regression	<b>data contamination</b>	<b>weak learners</b>
<b>ensemble learning</b>	error measures	types of learning	perceptrons	<b>hidden Markov models</b>
<i>exploration versus exploitation</i>		<i>kernel methods</i>		graphical models
<i>clustering</i>	<b>is learning feasible?</b>		<b>soft-order constraint</b>	
	<b>regularization</b>	weight decay	<i>Occam's razor</i>	<i>Boltzmann machine</i>

# Overview

## Learning Methods

### Supervised Learning

- Linear Models
- Neural Networks
- Support Vector Machines

### Unsupervised Learning

- Clustering
- Density Estimation
- Principal Component Anal.

## Learning Theory

- PAC Learning
- VC Dimension
- Bias Variance Tradeoff

## Learning Principles

- Regularization
- Validation

# Linear Classification

- Given Training Samples:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

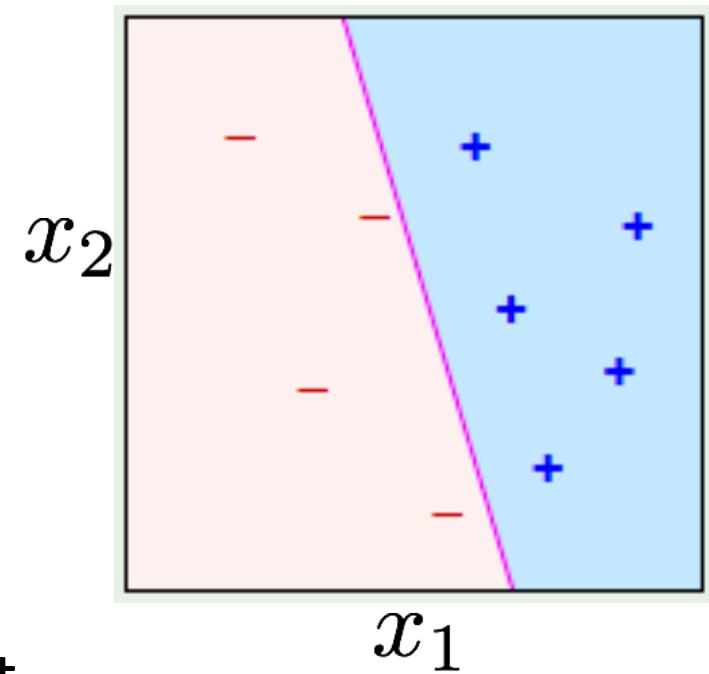
$$\mathbf{x}_i \in \mathbb{R}^d, \quad y_i \in \{\pm 1\}$$

- Determine a *classification rule*

$$y = \text{sign} \left( \sum_{j=0}^d w_j x_j \right)$$

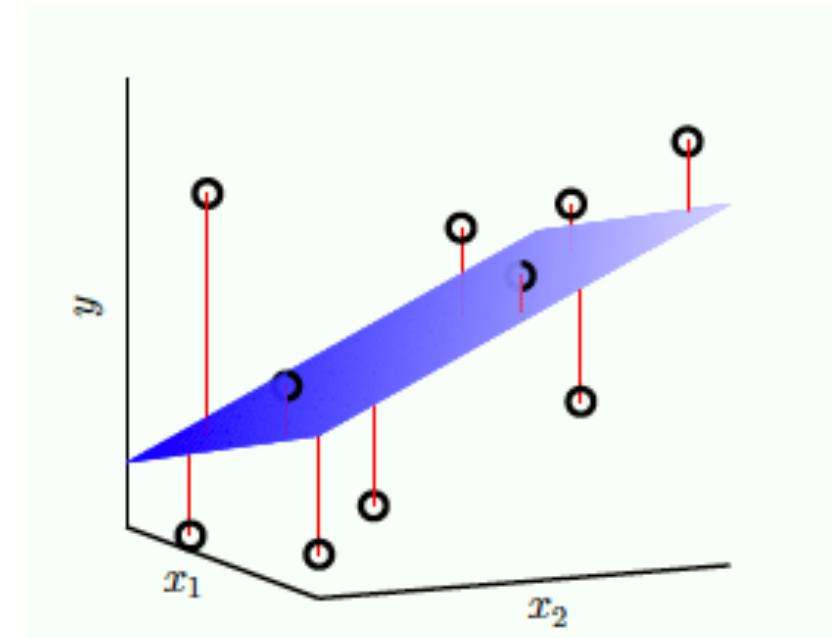
to minimize classification error

- Perceptron Learning Algorithm
- Logistic Regression and Gradient Descent

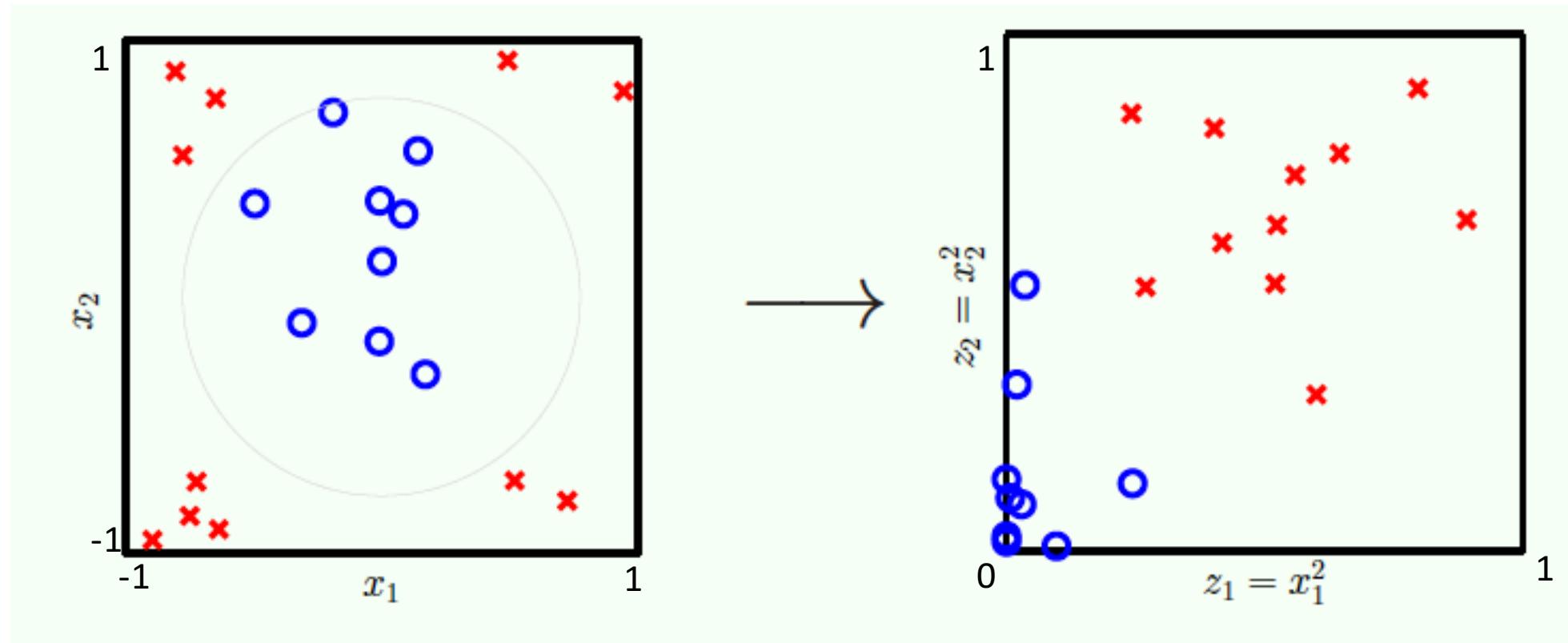


# Linear Regression

- **Given** Training Samples:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$   
 $\mathbf{x}_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R}$
- **Determine** a *regression rule*:  $\hat{y} = \sum_{i=0}^d w_i x_i$   
 $\mathbf{x} = (x_1, \dots, x_d)$
- **Minimize** the prediction error:  
$$\sum_{i=1}^N (y_i - \hat{y}_i)^2$$
- Least Squares and its variations



# Non-Linear Transformations of Data

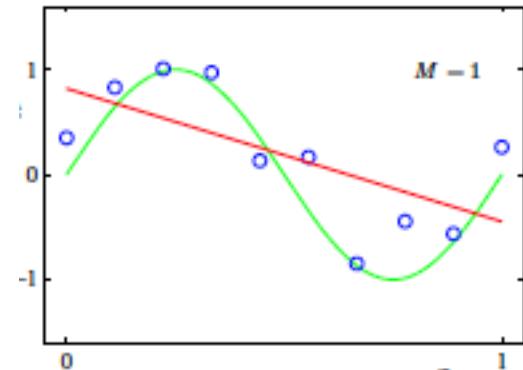


$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \implies \mathbf{z} = \begin{bmatrix} 1 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$

$$\hat{y} = \text{sign} (\mathbf{w}^T \mathbf{z})$$

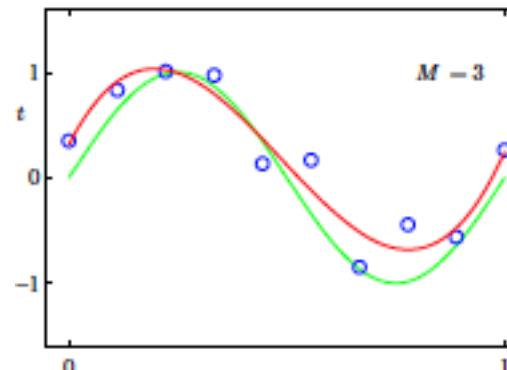
Nonlinear Transforms, Feature Vectors, Kernel Methods

# Non-Linear Transformations of Data



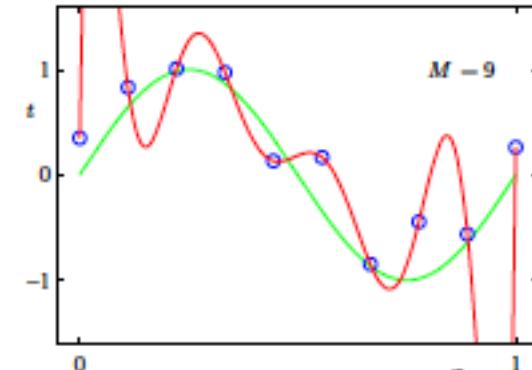
$$\mathbf{z} = [1 \quad x]^T$$

First Order Polynomial



$$\mathbf{z} = [1 \quad x \quad x^2 \quad x^3]^T$$

Third Order Polynomial



$$\mathbf{z} = [1 \quad x \quad x^2 \quad \dots \quad x^9]^T$$

9<sup>th</sup> Order Polynomial

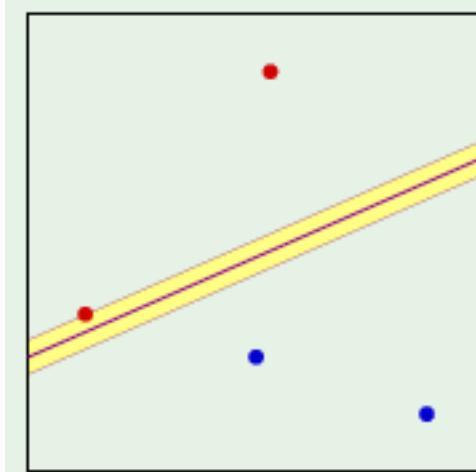
Under-fitting

$$\hat{y} = \mathbf{w}^T \mathbf{z}$$

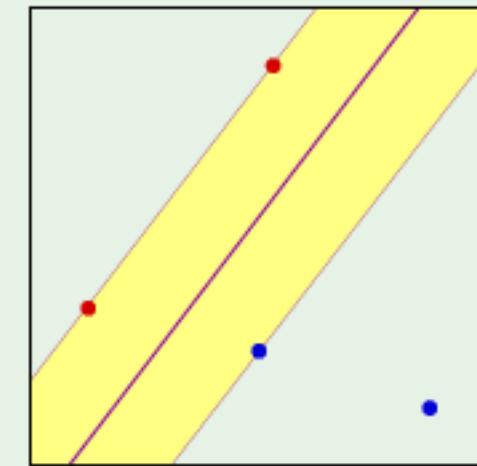
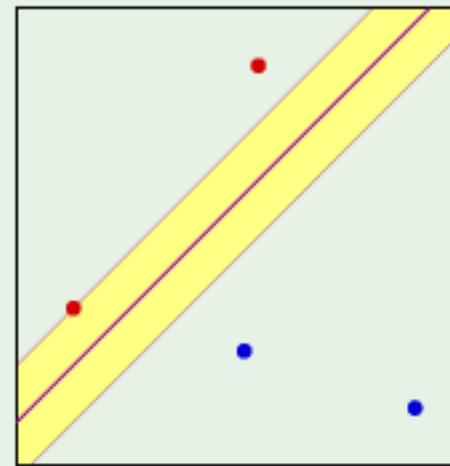
Overfitting

- Higher Order Non-Linearity:
  - Better Fit to Training Data
  - Less Robust to Noise (Overfitting)

# Support Vector Machines



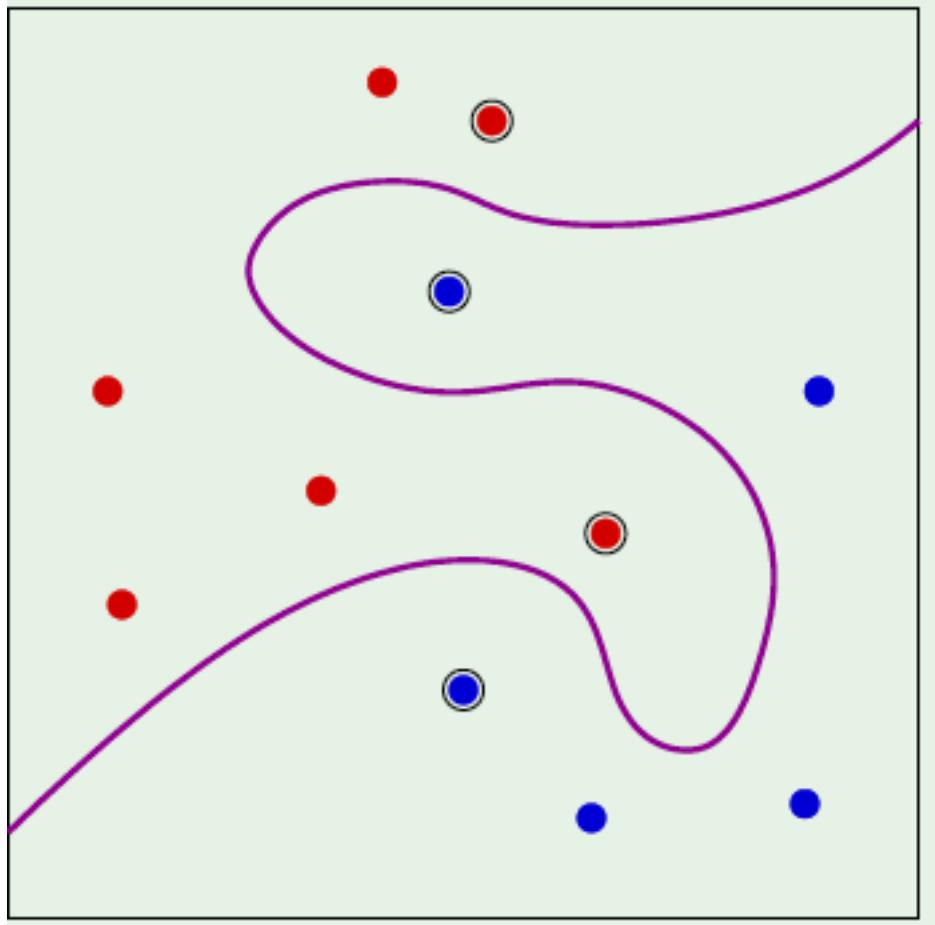
Thin Margin Classifier



Fat Margin Classifier

- Quadratic Programming: Maximizing Margin in Linear Classification
- Lagrange Duality Framework for identifying Support Vectors

# SVMs with Non-Linear Transforms

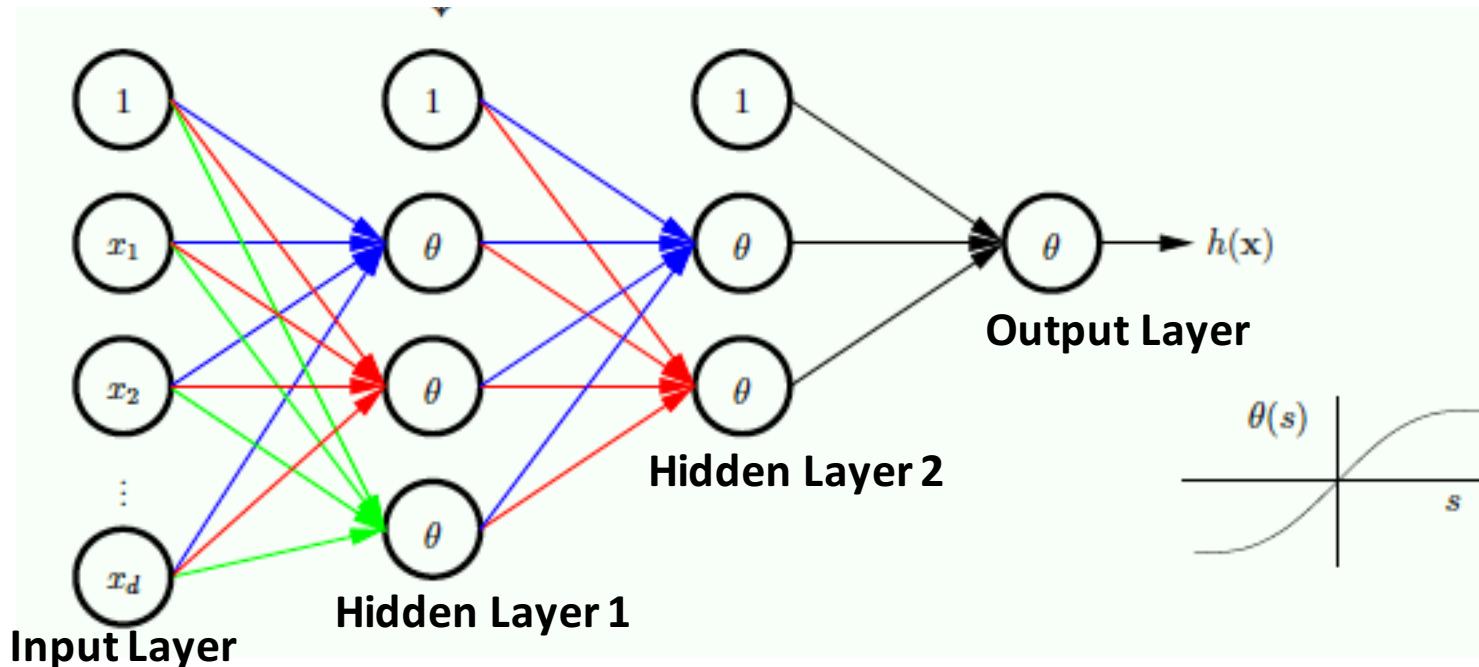


- Radial Basis Functions

$$\Phi(\mathbf{x}) = [\exp \{-\gamma(\|\mathbf{x} - \mathbf{x}_n\|^2)\}]_{1 \leq n \leq N}$$

- Kernel Trick for Efficient Computation
- Bounds on Generalization Error

# Neural Networks - Architecture



$$h_i^{(1)} = \theta \left( \sum_{j=0}^d w_{ji}^{(1)} \cdot x_j \right)$$

↑  
Output 'i'  
(Hidden Layer 1)  
Weights: Layer 1  
Inputs

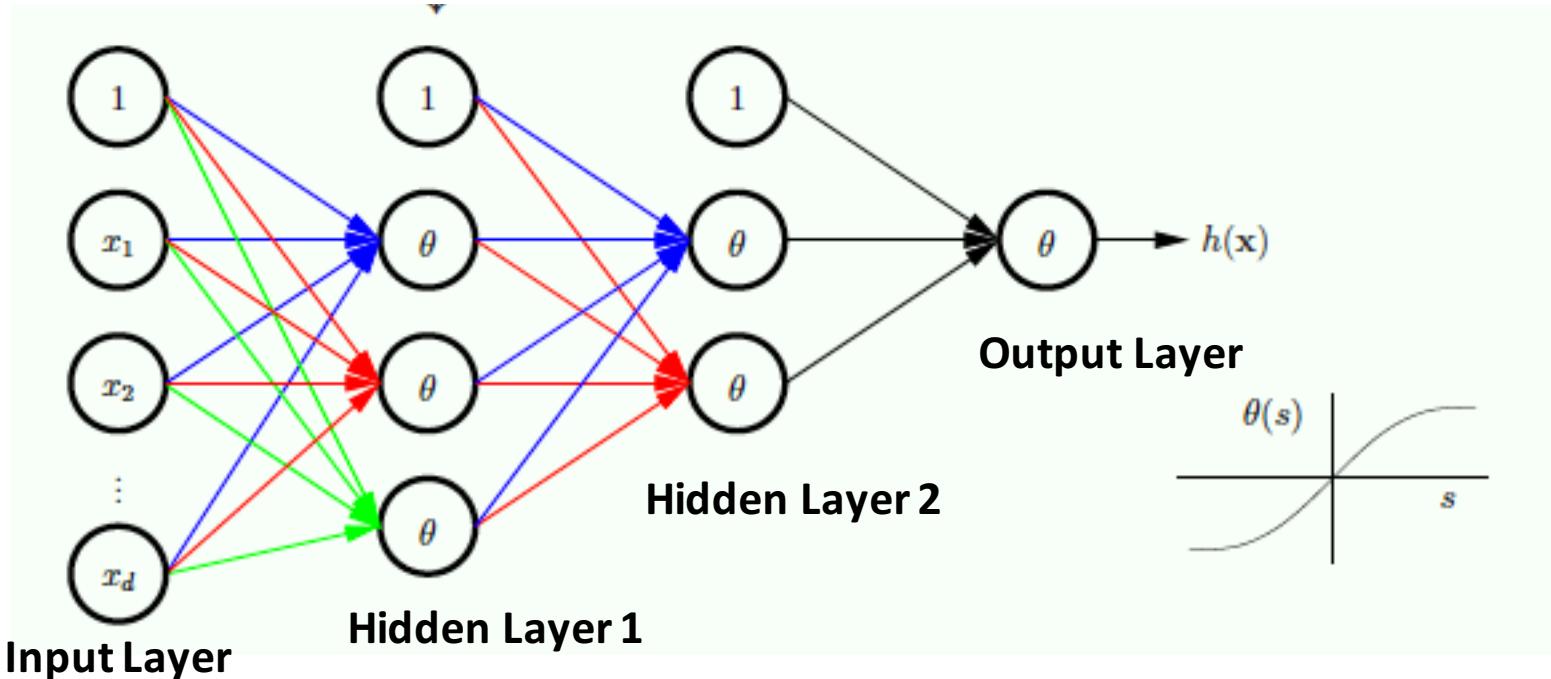
$$h_i^{(2)} = \theta \left( \sum_{j=0}^d w_{ji}^{(2)} \cdot h_j^{(1)} \right)$$

↑  
Output 'i'  
(Hidden Layer 2)  
Weights: Layer 2

$$y = \text{sign} \left( \sum_{j=0}^d w_{ji}^{(3)} \cdot h_j^{(2)} \right)$$

↑  
Output

# Neural Networks - Architecture



$$y = \text{sign} [W_3^T \{ \theta (W_2^T \{ \theta (W_1^T \cdot \mathbf{x}) \}) \}]$$

**Universal Approximation Theorem:** A Neural Network with **one hidden layer** can approximate any “reasonable” non-linear function with sufficiently many hidden units.

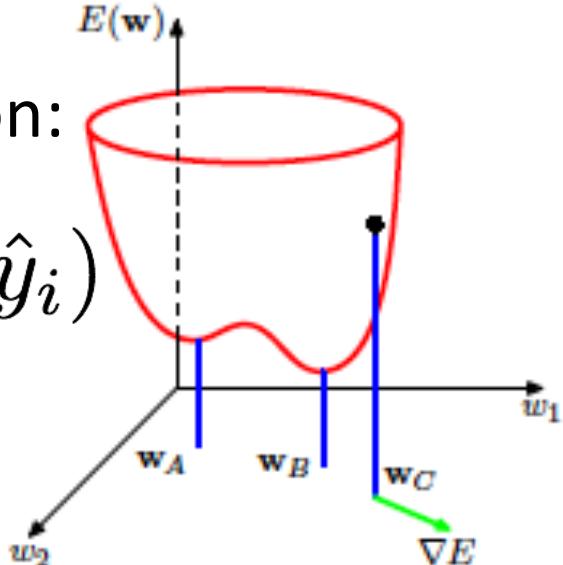
# Neural Network - Learning

- **Given** Training Samples:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- **Fix:**

- Number of Hidden Units per layer
- Nonlinearity:  $\theta$

- Learn: Weights:  $w_{i,j}^k$

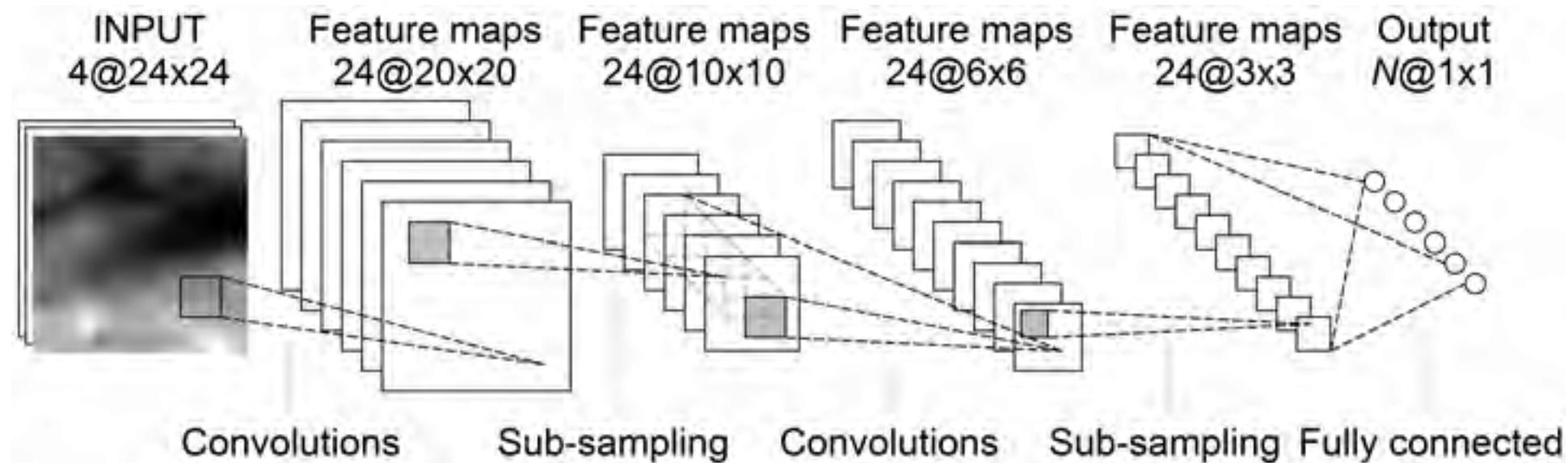
- **Minimize** Loss Function:

$$E(\mathbf{w}) = \sum_{i=1}^n \ell_{\mathbf{w}}(y_i, \hat{y}_i)$$


## Training Procedure

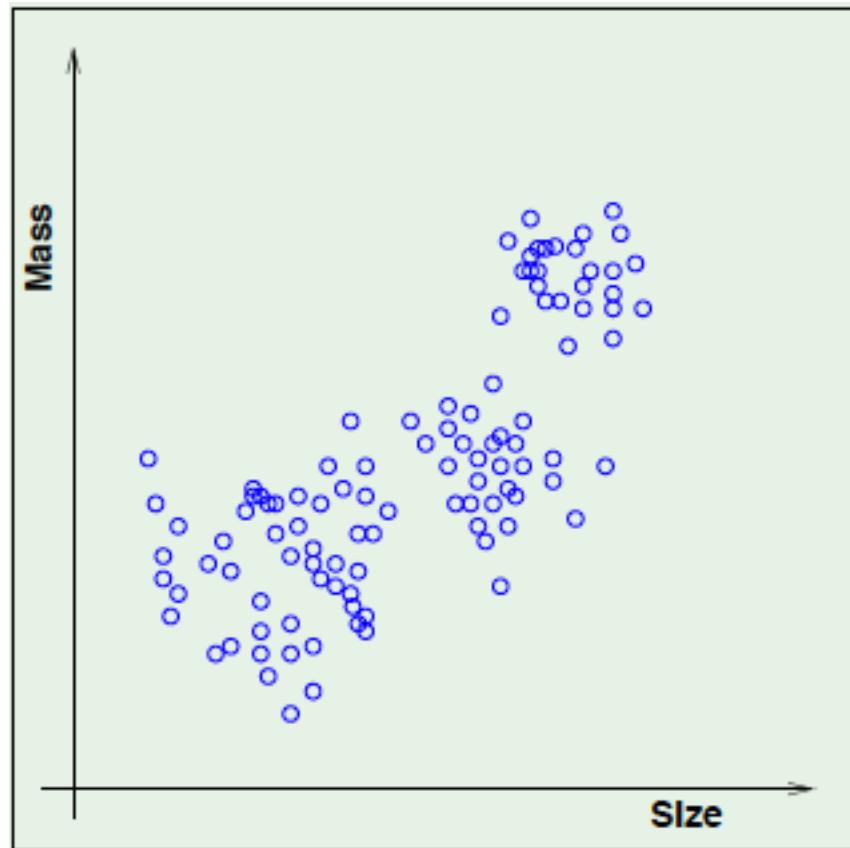
- Stochastic Gradient Descent
- Backpropagation Algorithm
- Early Stopping Rule
- Dropout and Regularization
- *Convolutional Neural Networks.*

# Convolutional Neural Networks



- Image statistics are translation invariant (objects and viewpoint translates)
- Expect low-level features to be local (e.g. edge detector)
- Expect high-level features learned to be coarser

# Unsupervised Learning



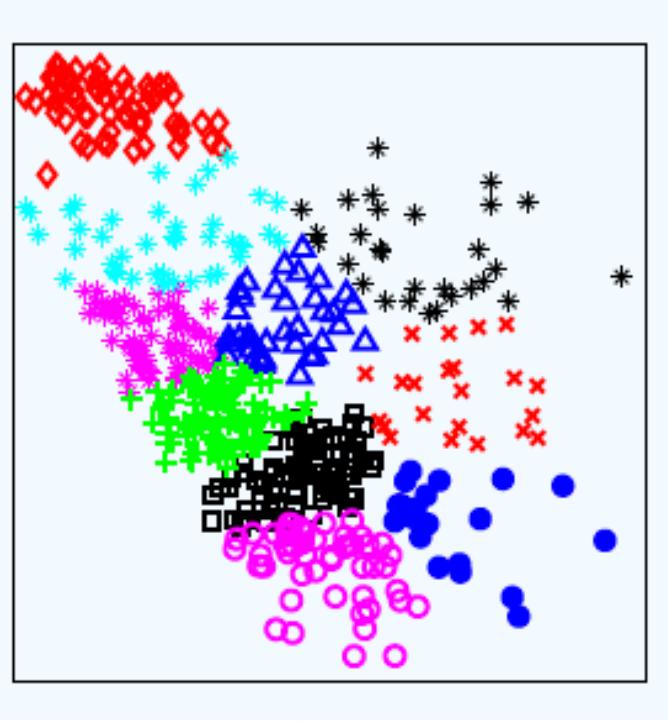
Training Data:

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$$

No Labels!

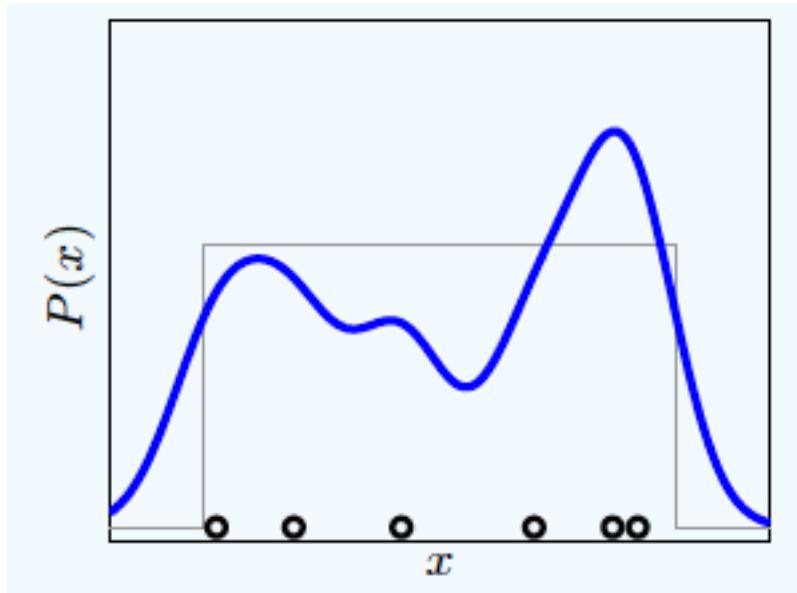
Still Need to do Learning!

# k-Means Clustering



- Cluster Centers:  $\mu_1, \mu_2, \dots, \mu_k$
  - Partitions:  $S_1, S_2, \dots, S_k$
  - Minimize : 
$$\sum_{j=1}^N \|\mathbf{x}_j - \mu(\mathbf{x}_j)\|^2$$
- Lloyd's Algorithm

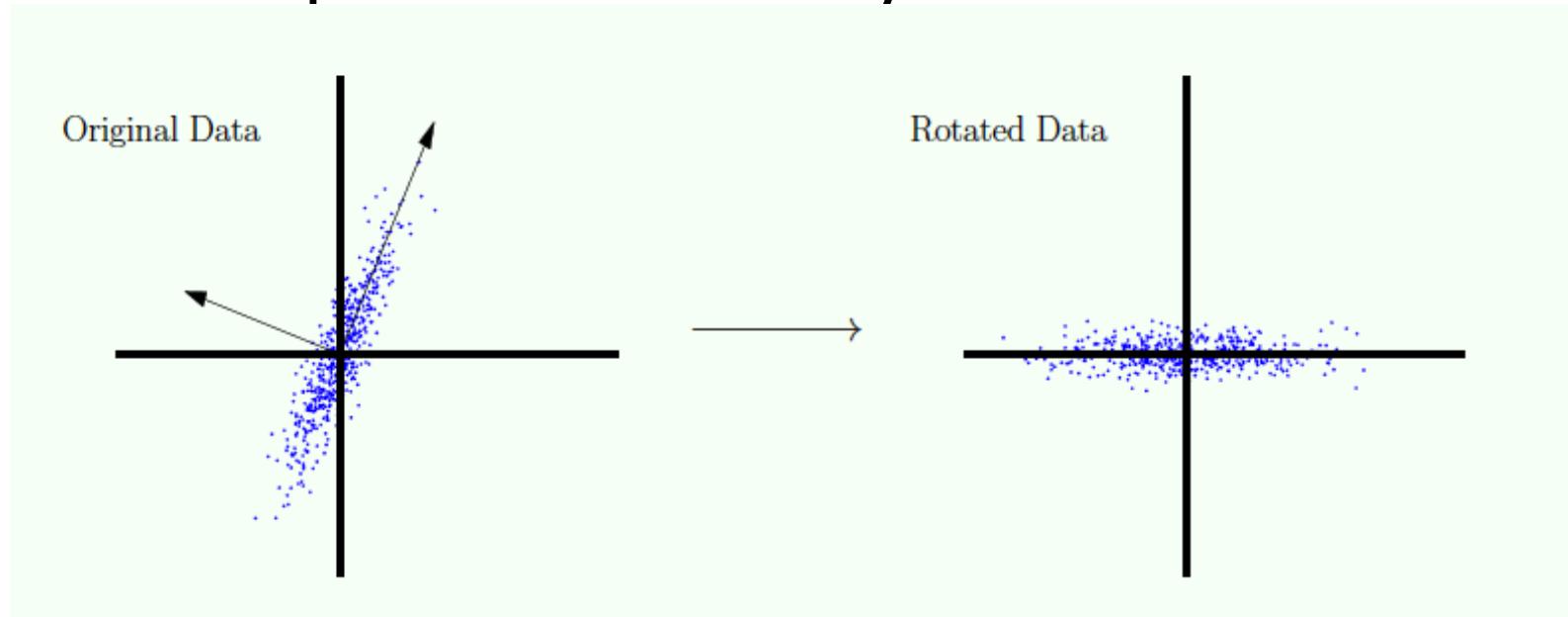
# Density Estimation - Gaussian Mixture Models



$$f(\mathbf{x}) = \sum_{j=1}^k w_j \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

- Estimating Probability Density Function
- Expectation Maximization Algorithm (EM) for GMMs
- General EM (time permitting) for Maximum Likelihood solution

# Principal Component Analysis



- Dimensionality Reduction
- Identify directions with large variance
- Singular Value Decomposition
- Non-Linear Extension: AutoEncoders

# Theory

# Training and Testing Errors

- Training Data:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

- Ground Truth (Not Known):  $y = f(\mathbf{x})$

- Output of Learning:  $\hat{y} = h(\mathbf{x})$

- **Training (In Sample) Error**

$$E_{\text{in}}(h) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y_i \neq h(\mathbf{x}_i)]$$

- **Testing (Out of Sample) Error**

$$E_{\text{out}}(h) = E_{\mathbf{x}, y} [\mathbb{I}[y \neq h(\mathbf{x})]]$$

# Probably Approximately Correct (PAC) Learning

- Fixed Hypothesis Class

$$\mathcal{H} = \{h_1, h_2, \dots, h_M\}, \quad h \in \mathcal{H}$$

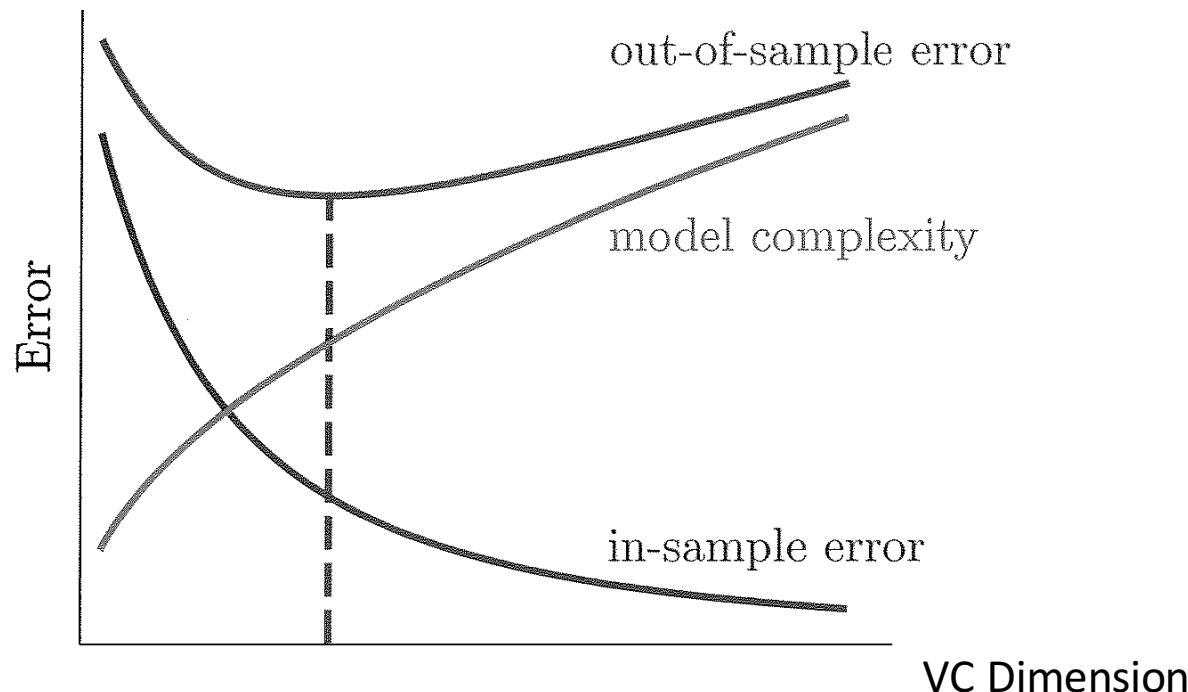
- The following bound holds with probability:  $1 - \delta$

- $$E_{\text{out}}(h) \leq E_{\text{in}}(h) + \sqrt{\frac{1}{2N} \log \frac{2M}{\delta}}$$


$$|\mathcal{H}| = \infty \implies M \leftarrow \text{VC Dimension}(\mathcal{H})$$

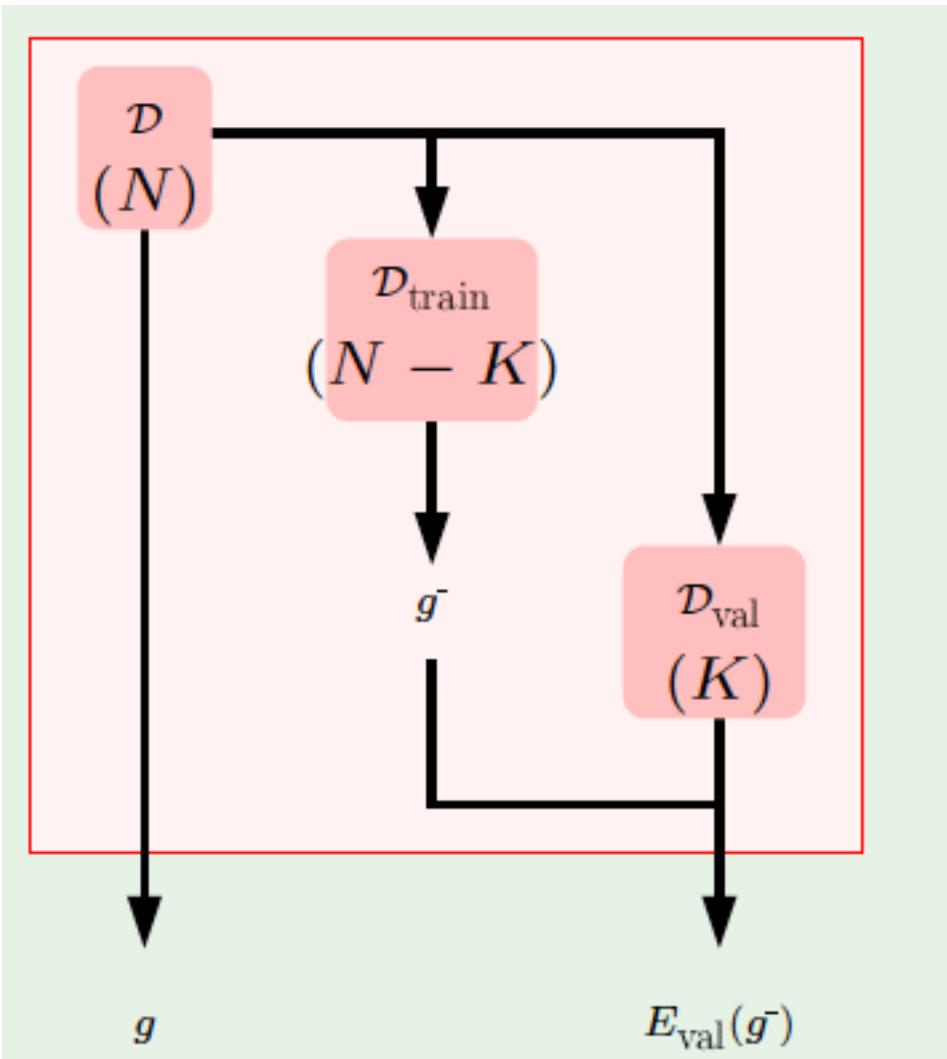
# Vapnik Chervonenkis (VC) dimension

- VC Dimension provides a natural measure for complexity of class
- Linear Models: VC Dimension = dimension + 1
- Neural Networks (roughly) = # of Weight Parameters



# Techniques

# Validation

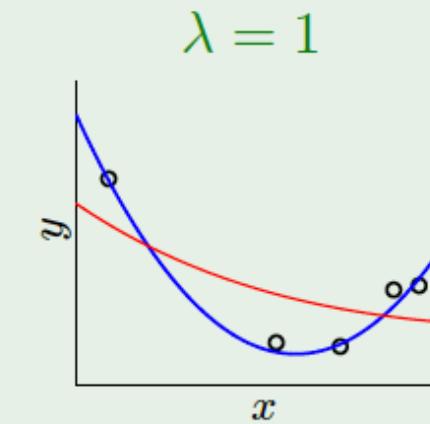
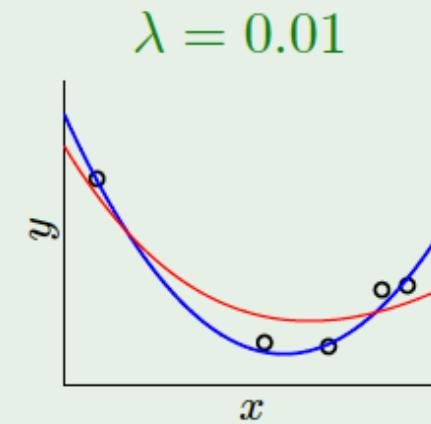
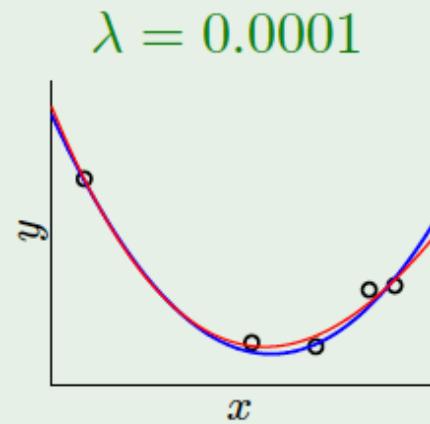
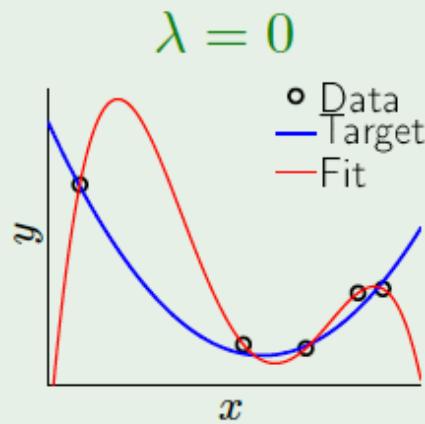


- Divide Training Set into 2 Parts
  - Learning Set
  - Validation Set
- Learning Set: Train Model
- Validation Set: Estimate Test Error
- Applications of Validation
  - Model Selection
  - Selection of Hyperparameters (e.g., regularization coefficient)
  - Number of Training Steps (Early Stopping)
- Cross Validation

# Regularization

$$y = \mathbf{w}^T \cdot \mathbf{z}, \quad \mathbf{z} = [1 \ x \ \dots, x^M]^T$$

Minimizing  $E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$  for different  $\lambda$ 's:



- Weight Decay Method
- Neural Networks: Drop Out Method

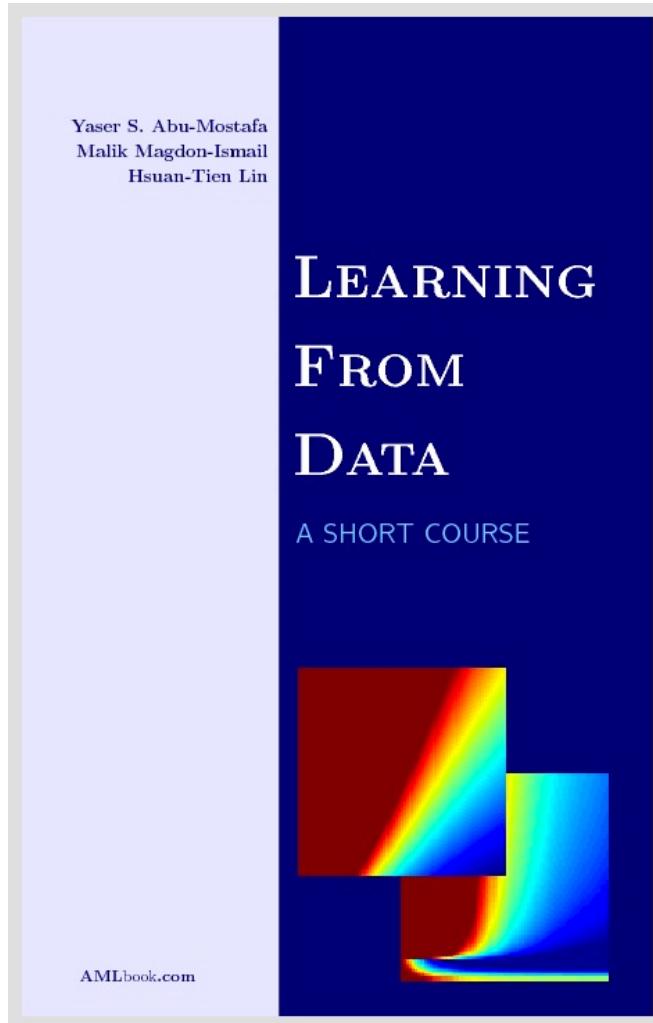
# Logistics

# Course Staff

## Instructor:

- Ashish Khisti
- Lectures
  - Mondays 5-6 (MC252)
  - Wednesdays 4-6 (MC252)
- Tutorials
  - Tutorial 1: Thursdays 1-3pm (GB304)
  - Tutorial 2: Thursdays 4-6pm (BA1230)

# Course Textbook



- **Required Textbook:**
- “Learning from Data” Available at U of T bookstore
  - <http://www.amlbook.com>
  - Supplementary Chapters and Appendices
  - Slides from other courses (Caltech, RPI etc)
  - Video Lectures (Prof. Abu-Mostafa)
  - Discussion Forum
- Deep Learning by Goodfellow et.al (Free Online book)
- Recommended Textbook
  - Machine Learning and Pattern Recognition by Christopher Bishop (Text for CSC411)

# Tentative Schedule

1	September 9-13	<b>Intro/Linear Classification</b>	Tut 1			
2	September 16-20	<b>Linear Regression, Regularization, Gradient Descent</b>	Tut 2	<b>Assnt 1 Posted</b>		
3	September 23-27	<b>Logistic Regression</b>	Tut 3			
4	September 30-Oct 4	<b>Gradient Descent</b>	Tut 4	<b>Assnt 1 due</b>		
5	Oct 7-11	<b>Multilayer Perceptron, Backpropagation</b>	Tut 5	<b>Assnt 2 Posted</b>		
6	Oct 14-18	<b>Thanksgiving and Mid-term Week</b>	Tut 6	<b>Mid-Term Exam</b>	<b>Wed Oct 16th</b>	
7	Oct 21-25	<b>Deep Neural Networks</b>	Tut 7			
8	Oct 28 – Nov 1	<b>Unsupervised Learning: Clustering and Density Estimation</b>	Tut 8	<b>Assnt 2 Due</b>		
9	Nov 4-8	<b>EM Algorithm, PCA</b>	Tut 9	<b>Assnt 3 Posted</b>		March 17th - Drop Date
10	Nov 11-15	<b>PAC Learning - 1 (Finite Hypothesis)</b>	Tut 10			
11	Nov 18-22	<b>PAC Learning - 2 (VC Dim, Bias-Variance)</b>	Tut 11	<b>Assnt 3 Due</b>		
12	Nov 25-29	<b>Validation, Generalization, Bootstrapping</b>	Tut 12			
13	Dec 2-6	<b>Support Vector Machine</b>		Last Day of Class – Dec 4th		

# Pre-requisites

- Undergraduate Course in Probability (Official Pre-Req.)
  - Bayes Theorem, Union Bound, Gaussian Distributions
- Linear Algebra (Strongly Recommended)
  - Vector Space Concepts, Matrices
- Programming
  - We will use Python and Tensor Flow Package for our assignments

# Grade Composition

- Mid Term: Wed Oct 16th (30%)
- Programming Assignment ( $3 \times 10\% = 30\%$ )
- Final Exam: 40% in Exam Week
- Location of Mid-Term: Exam Center
  - MC252 & RW117

# Tensor flow Assignments

- **Python based** ML Library released by Google in 2015
- Automatic Training for Neural Networks
- GPU Support (Not Required for Assignments in this courses)
- Installation through Anaconda Environment is Recommended (See Installation Guide on Course Webpage)
- Tons of Resources!
  - Tensorflow.Org Tutorials
  - CS231n Stanford Tutorial (<http://cs231n.stanford.edu/>)
  - See Course Webpage for a simple tutorial (Updated, Use Chrome Browser)

# Tensor flow Example (<https://www.tensorflow.org>)

```
import tensorflow as tf  
x = tf.placeholder(tf.float32, [None, 784])  
W = tf.Variable(tf.zeros([784, 10]))  
b = tf.Variable(tf.zeros([10]))  
y = tf.nn.softmax(tf.matmul(x, W) + b)
```

Initialize Computational Graph

```
cross_entropy = tf.reduce_mean(-tf.reduce_sum(y_ * tf.log(y), reduction_indices=[1]))  
train_step = tf.train.GradientDescentOptimizer(0.5).minimize(cross_entropy)
```

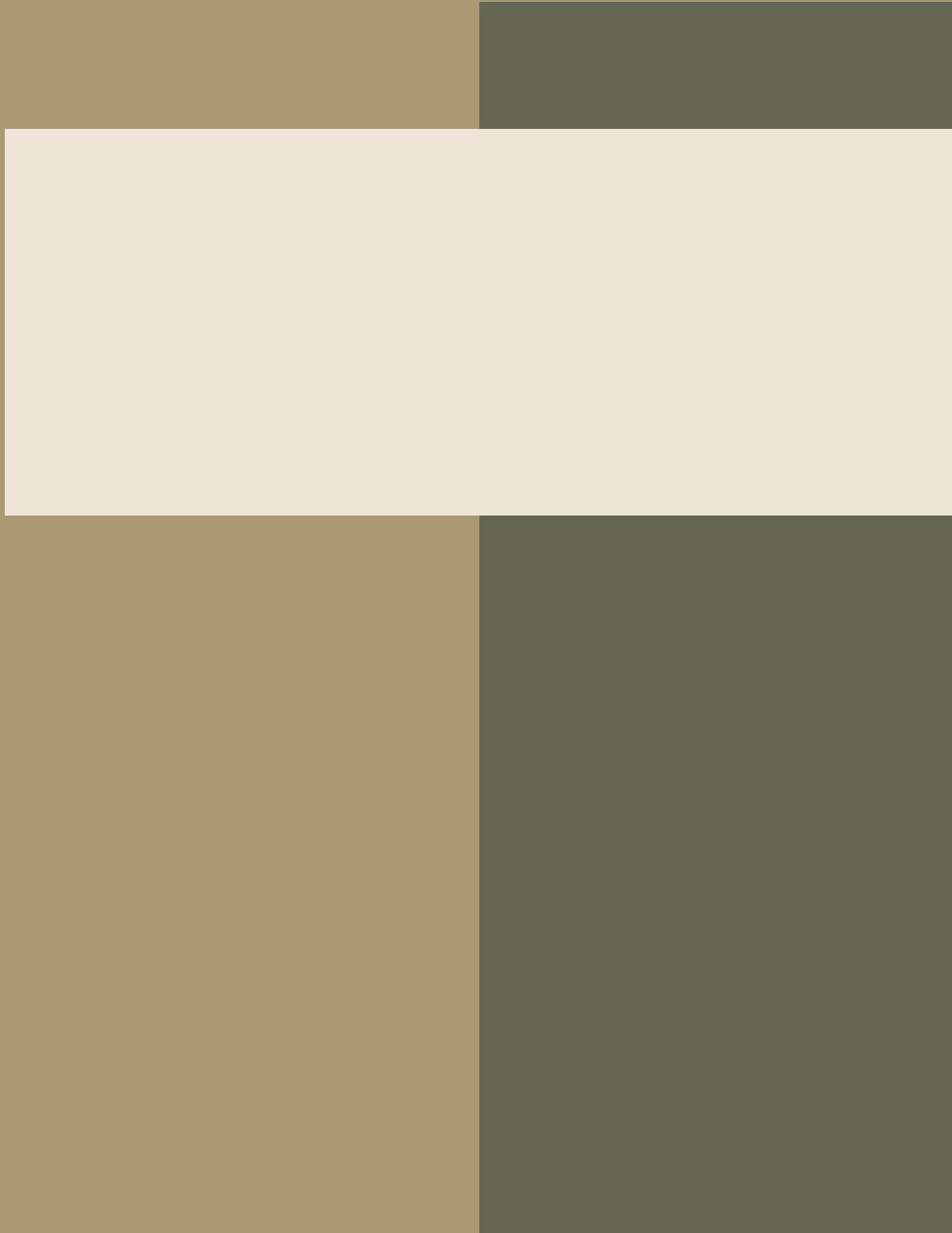
Loss Function  
and Optimizer

```
sess = tf.InteractiveSession()  
  
tf.global_variables_initializer().run()  
  
for _ in range(1000):  
    batch_xs, batch_ys = mnist.train.next_batch(100)  
    sess.run(train_step, feed_dict={x: batch_xs, y_: batch_ys})
```

Training Routing

# ECE421 Course Description

An Introduction to the basic theory, the fundamental algorithms, and the computational toolboxes of machine learning. The focus is on a balanced treatment of the practical and theoretical approaches, along with hands on experience with relevant software packages. Supervised learning methods covered in the course will include: the study of linear models for classification and regression, neural networks and support vector machines. Unsupervised learning methods covered in the course will include: principal component analysis, k-means clustering, and Gaussian mixture models. Theoretical topics will include: bounds on the generalization error, bias-variance tradeoffs and the Vapnik-Chervonenkis (VC) dimension. Techniques to control overfitting, including regularization and validation, will be covered



# ECE 421: Intro. To Machine Learning

## Week #1:

What is Machine learning.

Develop computational systems that adaptively improve performance with experience from accumulated data.

### Example 1: Recommendation Systems.

Problem: Design a system that recommends new movies to users based on their preferences

#### Approach 1: Human Experts.

An expert views each movie & collects various attributes

- Actors ?
- block buster ?
- Action / Comedy ?

Expert Interviews Each user.

- Like comedy ?
- fav. actors ?
- prefer block busters .

⇒ find Suitable recommendations.

Cons: tedious, bias, preferences change etc.

Approach 2: Machine Learning

- data driven
- uses past user ratings.

1) User Preference Vector

$$\underline{x}_i = (x_{i1}, x_{i2}, \dots, x_{id}) \in \mathbb{R}^d$$

for each  $i \in \mathcal{U}$

2) Movie attribute vector

$$\underline{y}_j = (y_{j1}, y_{j2}, \dots, y_{jd}) \in \mathbb{R}^d$$

for each  $j \in \mathcal{M}$

$i$  = user index  
 $\mathcal{U}$  = set of users

$j$  = movie index  
 $\mathcal{M}$  = set of movies.

3) Data Set

$y_{ij}$  = rating by user  $i$  for movie  $j$

		movies	1	2	3	4	-
		users	1	2	3	4	-
1	*	$\gamma_{12}$	$\gamma_{13}$	*			
2	$\gamma_{21}$	$\gamma_{22}$	*				*
3	$\gamma_{31}$	*	*				$\gamma_{34}$
4							

$\mathcal{S}$  = Set of (user, movie) pairs  
for which ratings are available.

=  $\{(i, j) : \gamma_{ij} \text{ is available in our data set}\}$

## 1) Proposed Model

$$\begin{aligned}
 & \text{Rating Estimate} \quad \hat{\gamma}_{ij} = \underline{x}_i^T \underline{y}_j \\
 & \underline{x}_i = (\underline{x}_{i1}, \underline{x}_{i2}, \dots, \underline{x}_{id}) = \sum_{l=1}^d x_{il} y_{jl} \\
 & \underline{y}_j = (\underline{y}_{j1}, \underline{y}_{j2}, \dots, \underline{y}_{jd})
 \end{aligned}$$

How to compute  $\underline{x}_i, \underline{y}_j$ ?

✓ Mode (Parameters)

$$\{S\} = \left\{ \left\{ \underline{x}_i \right\}_{i \in U}, \left\{ \underline{y}_j \right\}_{j \in M} \right\}$$

✓ Loss Function

$$E_{in}(S) = \sum_{(i,j) \in S} (\underline{y}_{ij} - \underbrace{\underline{x}_i^T \underline{y}_j}_{= \hat{y}_{ij}})^2$$

Training Phase

$$\checkmark S^* = \arg \min_{S} E_{in}(S)$$

find model parameters that minimize loss function.

- Gradient descent

- EM Algorithm

Output  $\{x_i\}_{i \in U}, \{y_j\}_{j \in M}$

✓ for all users & movies

Prediction Phase:  $\hat{y}_{ij} = \underline{x}_i^T \underline{y}_j$

Cons : Models are not easily interpretable.

new-movies / users cannot be recommended

Sampling bias: users who have watched too many movies will influence others.

---

## Example 2 : Credit Approval

Bank needs to decide whether or not to approve the credit of user

User attribute vector

$$\underline{x} = (\text{age}, \text{gender}, \text{salary}, \dots) \in \mathbb{R}^d$$

Output

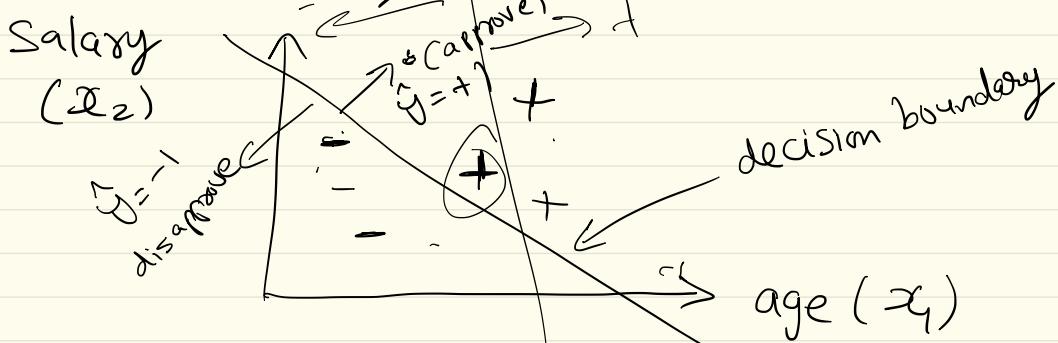
$$y = \begin{cases} +1 & \text{approve} \\ -1 & \text{dis-approve} \end{cases}$$

Proposed Model: Linear classification

$$d=2 \quad \underline{x} = (\text{age}, \text{salary})$$

Historical Data:

$$D = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$



$$\underline{x}_1 = (x_{11}, x_{12}) \text{ model programs}$$

$$(S2) \quad \{(\omega_0, \omega_1, \omega_2)\}$$

$$\omega_0 + \omega_1 x_1 + \omega_2 x_2 = 0$$

# (Binary) Linear Classification

weight vector :  $\underline{w} \in \mathbb{R}^d$   
 Constant  $b \in \mathbb{R}$   $= (\omega_1, \omega_2, \dots, \omega_d)$

Given Input  $\underline{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$

If  $\sum_{i=1}^d \omega_i x_i > b \Rightarrow \hat{y} = +1$   
 else  $\hat{y}_i = -1$

Model Parameters :  $\underline{\Omega} = \{w, b\}$

Training Set :  $\{(x_1, y_1), \dots, (x_N, y_N)\}$

$x_i = (x_{i1}, x_{i2}, \dots, x_{id}) \in \mathbb{R}^d$   
 $y_i \in \{-1, +1\}$

Training Phase: Compute  $\underline{\Omega}$

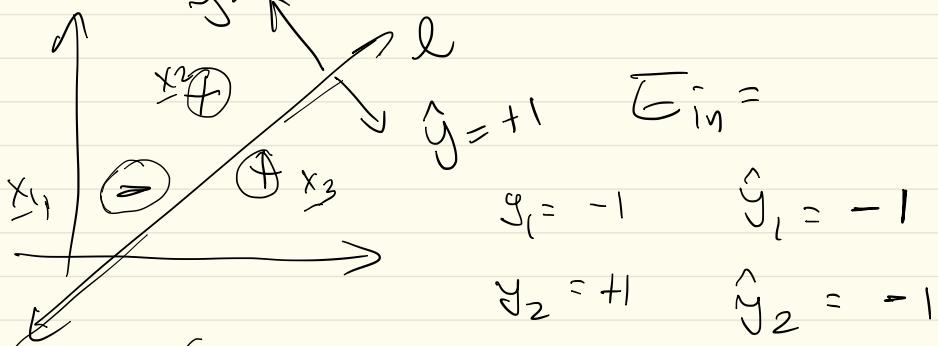
## Loss function:

Average Number of Errors.

$$E_{in}(S) = \frac{1}{N} \sum_{i=1}^N 1 \{ y_i \neq \hat{y}_i \}$$

1(-) = indicator function

$$1(t) = \begin{cases} 1 & t \text{ is true} \\ 0 & t \text{ is false} \end{cases}$$



# Perceptron Learning Algorithm

## Linear Classification

Weight vector  $\underline{w} = (w_1, w_2, \dots, w_d)$   
constant  $b \in \mathbb{R}$

### Change in Notation

#### Expanded dimension

$$\underline{x} = (x_0 = 1, x_1, x_2, \dots, x_d) \in \mathbb{R}^{d+1}$$

$$\underline{w} = (w_0 = -b, w_1, w_2, \dots, w_d) \in \mathbb{R}^{d+1}$$

#### decision Rule

$$w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

$\hat{y} = +1$   
 $\hat{y} = -1$

$$\hat{y} = +1$$

$$\hat{y} = -1$$

$$-b + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

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

Decision Rule:  $\hat{y} = \begin{cases} +1 & \text{if } \underline{w}^\top \underline{x} \geq 0 \\ -1 & \text{if } \underline{w}^\top \underline{x} < 0 \end{cases}$

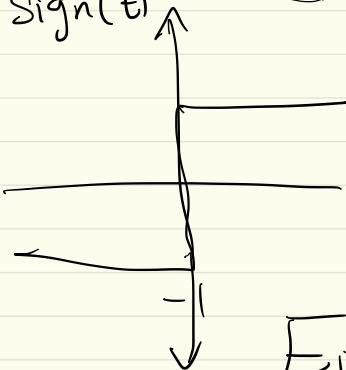
$$\underline{w}^\top \underline{x}$$

$$\begin{cases} +1 & \text{if } \underline{w}^\top \underline{x} \geq 0 \\ -1 & \text{if } \underline{w}^\top \underline{x} < 0 \end{cases}$$

0

$$\hat{y} = h_{\underline{w}}(\underline{x}) = \text{Sign}(\underline{w}^\top \underline{x})$$

$\text{Sign}(t)$



+1

t

Goal:

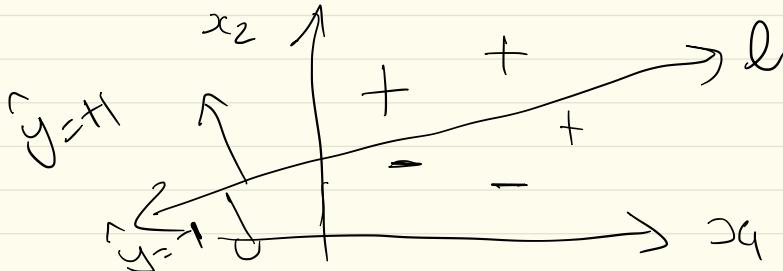
Select  $\underline{w} \in \mathbb{R}^{d+1}$

that minimizes

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{\hat{y}_i \neq \underline{h}_{\underline{w}}(\underline{x}_i)\}}$$

Training Set:

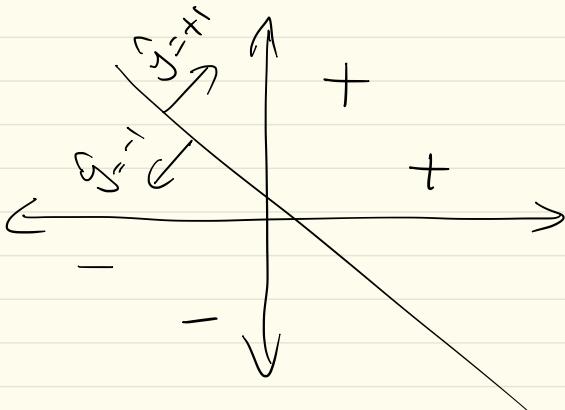
$$\{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$$



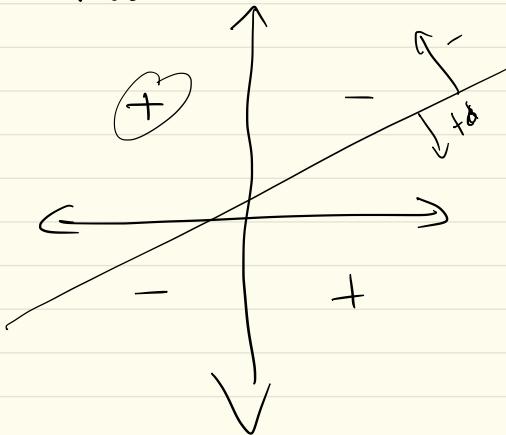
$$E_{in} = \frac{1}{S}$$

## Separable Data Sets:

A data set is separable if there is a line that achieves zero classification error.



(a)  
Separable



(b)  
Not Separable.

If the data set is linearly separable then Perceptron learning Algorithm find  $\underline{w} \in \mathbb{R}^{(J+1)}$  with  $E_{in}(w) = 0$

If the data set is Not Separable

minimizing  $E_{in}(\underline{w})$  is NP-Hard in general. However perceptron algorithm can be extended to non-separable Data sets to attain good Heuristics.

Knowing whether a data set is separable or not is not NP-Hard.

## Perceptron Learning Algorithm

Assume:  $\mathcal{D} = \{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$   
is linearly separable  
 $\underline{x}_i \in \mathbb{R}^{d+1}$        $y_i \in \{-1, +1\}$

Output:  $\underline{w} \in \mathbb{R}^{d+1}$   
that achieves  $E_{in}(\underline{w}) = 0$

Step 1: Initialize  $\underline{w}$  in an arbitrary manner

Step 2: Compute  $E_{in}(\underline{w})$

$$= \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{\{y_i \neq h_w(x_i)\}}$$

If  $E_{in}(\underline{w}) = 0$  Stop.  
else output  $\underline{w}$ .

Step 3 let  $(\underline{x}_n, y_n)$  be any mis-classified point

i.e.  $y_n \neq h_w(\underline{x}_n)$

$$\underline{w} \leftarrow \underline{w} + y_n \underline{x}_n$$
$$\hookrightarrow \{-1, +1\}$$

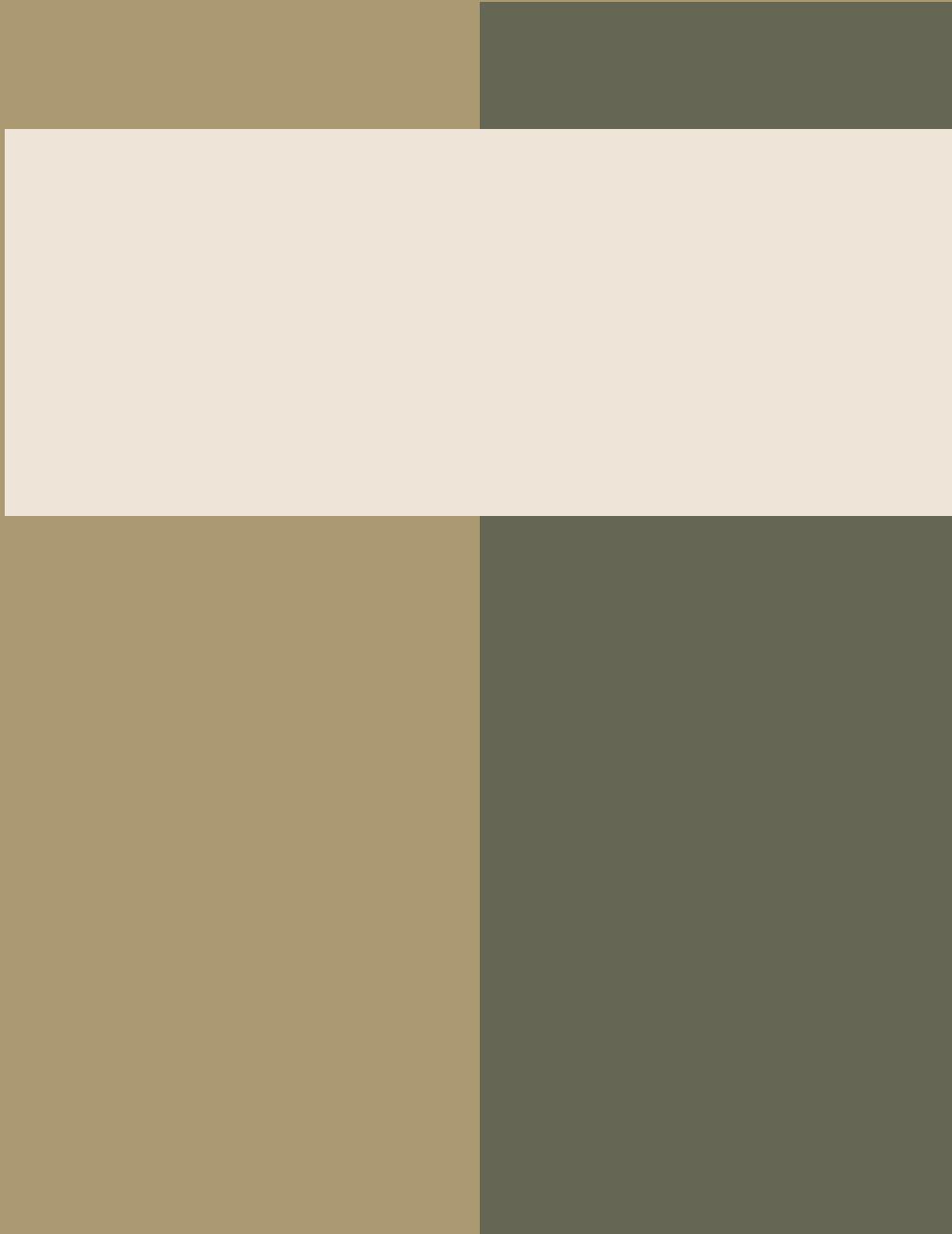
Go to Step 2

Next Lecture

2) Convergence

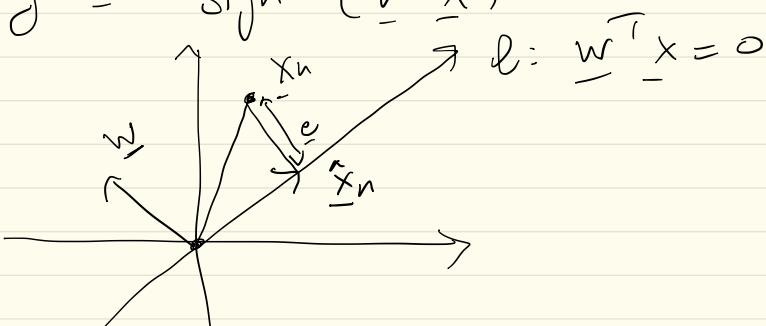
1) Geometric Intuition  
for Perceptron Alg.

3) Extension to non separable data sets



Consider a linear classifier

$$\hat{y} = \text{sign}(\underline{w}^T \underline{x})$$



the decision boundary  $\ell: \underline{w}^T \underline{x} = 0$  is shown above.

Note that  $\underline{w}$  points perpendicular to  $\ell$ .

let  $\underline{x}_n$  be any point, we want to show  $\text{dist}(\underline{x}_n, \ell) = \frac{|\underline{w}^T \underline{x}_n|}{\|\underline{w}\|}$

Let  $\underline{u}$  be a unit vector along  $\underline{w}$

let us decompose

$$\underline{x}_n = \underline{\hat{x}}_n + \underline{e},$$

where  $\hat{\underline{x}}_n$  is a vector along  $l$   
 (i.e. projection of  $\underline{x}_n$  along  $l$ )  
 and  $\underline{e}$  is a vector perpendicular  
 to  $l$  (i.e. parallel to  $u$ )

We can express  $\underline{e} = \|\underline{e}\| \cdot \underline{u}$   
 also note that  
 $\text{dist}(\underline{x}_n, l) = \|\underline{e}\|$

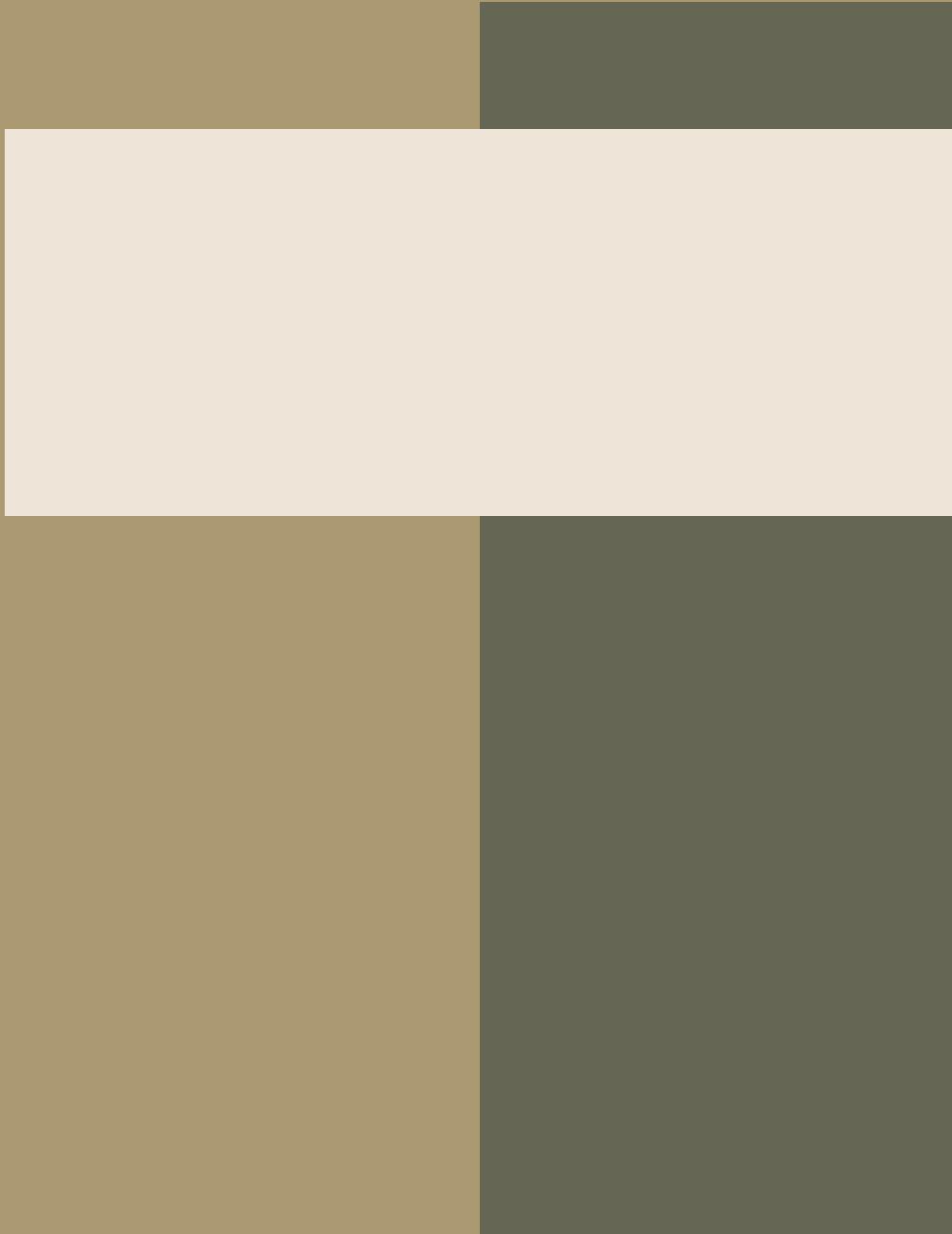
Now consider  
 $\underline{x}_n = \hat{\underline{x}}_n + \underline{e}$   
 $\underline{w}^T \underline{x}_n = \underline{w}^T \hat{\underline{x}}_n + \underline{w}^T \underline{e}$   
 but  $\underline{w}^T \hat{\underline{x}}_n = 0$  as  $\underline{w} \perp \hat{\underline{x}}_n$   
 $\therefore \underline{w}^T \underline{x}_n = \underline{w}^T \underline{e}$   
 Note.  $|\underline{w}^T \underline{x}_n| = |\underline{w}^T \underline{e}|$   
 Since  $\underline{w}$  &  $\underline{e}$  are collinear,

$$|\underline{w}^T \underline{e}| = \|\underline{w}\| \|\underline{e}\|$$

thus  $|\underline{w}^T \underline{x}_n| = \|\underline{w}\| \|\underline{e}\|$

$$\|\underline{e}\| = \frac{|\underline{w}^T \underline{x}_n|}{\|\underline{w}\|}$$

$$\therefore \text{dist}(\underline{x}_n, l) = \frac{|\underline{w}^T \underline{x}_n|}{\|\underline{w}\|}$$



## Week 2

### Binary Linear Classification

$$\mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$

$$\underline{x}_i = (x_{i0}=1, x_{i1}, x_{i2}, \dots, x_{id}) \in \mathbb{R}^{d+1}$$

$$y_i \in \{-1, +1\}$$

### Model Parameters

$$\text{Weight Vector } \underline{w} \in \mathbb{R}^{d+1}$$

### Classification Rule

$$\hat{y} = h_{\underline{w}}(\underline{x}) = \text{Sign}(\underline{w}^T \underline{x})$$

$$\begin{aligned} \hat{y} &= \begin{cases} +1 & \underline{w}^T \underline{x} > 0 \\ -1 & \underline{w}^T \underline{x} < 0 \end{cases} \\ \text{Loss Function} &= \begin{cases} +1 & \underline{w}^T \underline{x} > 0 \\ -1 & \underline{w}^T \underline{x} < 0 \end{cases} \end{aligned}$$

$$| E_{in}(\underline{w}) =$$

Loss function

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N \mathbf{1}\{y_n \neq \hat{y}_n\}$$

Given training set  $\mathcal{D}$ ,  
minimize  $E_{in}(\underline{w})$

---

## Perceptron Learning Algorithm

Input: Training set  $\mathcal{D}$  that is  
linearly separable

Output:  $\underline{w} \in \mathbb{R}^{d+1}$  that achieves  
 $E_{in}(\underline{w}) = 0$

Initialization: Initialize  $\underline{w}$  in  
an arbitrary manner

e.g.  $\underline{w} = (0, 0, \dots, 0) \in \mathbb{R}^{d+1}$

Step 1: Check if  $E_{in}(\underline{w}) = 0$  ?

If Yes STOP & output  $\underline{w}$

Step 2: let  $(\underline{x}_n, y_n)$  be any point in  $D$  that is mis-classified (including any point on boundary i.e.  $\underline{w}^T \underline{x}_n = 0$ )

If  $y_n = +1$  then  $\underline{w} \leftarrow \underline{w} + \underline{x}_n$

If  $y_n = -1$  then  $\underline{w} \leftarrow \underline{w} - \underline{x}_n$

$\underline{w} \leftarrow \underline{w} + y_n \cdot \underline{x}_n$

Go to Step 1

(In 1950's by Rosenblatt)

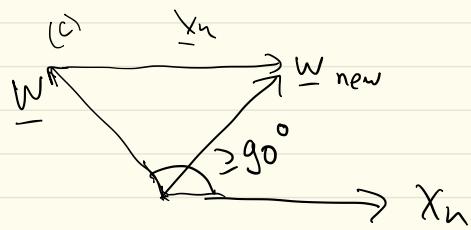
## Update Rule:

### Geometric Intuition

let  $(\underline{x}_n, y_n)$  be a mis-classified Point

Case 1:  $y_n = +1$        $\hat{y}_n = -1$

Recall:  $\hat{y}_n = \text{Sign}(\underline{w}^T \underline{x}_n)$

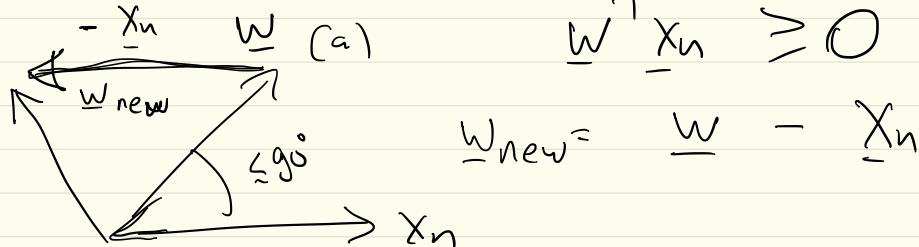


$$\underline{w}^T \underline{x}_n \leq 0$$

$$\underline{w}_{\text{new}} \leftarrow \underline{w} + \underline{x}_n$$

" $\underline{w}_{\text{new}}$ " moves closer to  $\underline{x}_n$

Case 2:  $y_n = -1$ ,  $\hat{y}_n = +1$



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

$$\underline{w}_{\text{new}} = \underline{w} - \underline{x}_n$$

# Algebraic View Point

$y_n$	$w^T x_n$	classification	$y_n \cdot w^T x_n$
+ 1	$\geq 0$	✓	$\geq 0$
+ 1	$< 0$	X	$< 0$
- 1	$\geq 0$	X	$< 0$
- 1	$< 0$	✓	$\geq 0$

If  $(x_n, y_n)$  is correctly classified  $y_n w^T x_n \geq 0$

If  $(x_n, y_n)$  is misclassified  $y_n w^T x_n < 0$

Claim: If  $(x_n, y_n)$  is mis-classified  
in Step 2 and

$$w_{\text{new}} = w + y_n \cdot x_n$$

$$y_n \cdot w_{\text{new}}^T x_n > y_n \cdot w^T x_n$$

Proof:  $w_{\text{new}}^T = w^T + y_n \cdot x_n^T$

$$w_{\text{new}}^T x_n = w^T x_n + y_n \cdot x_n^T x_n$$

$$y_n \cdot w_{\text{new}}^T x_n = y_n \cdot w^T x_n + \boxed{\frac{y_n^2}{||x_n||^2} x_n^T x_n}$$

# Convergence of Perceptron Alg.

let  $\underline{w}^*$  be the output of PLA.  
then PLA will stop in at-most

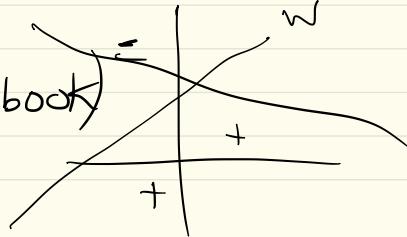
$$T = \frac{R^2}{\rho^2} \text{ steps where}$$

$$R = \max_{1 \leq n \leq N} \|\underline{x}_n\| \leftarrow \begin{matrix} \text{Radius of} \\ \text{D} \end{matrix}$$

$$\rho = \min_{1 \leq n \leq N} \left\{ \frac{|\underline{w}^T \underline{x}_n|}{\|\underline{w}\|} \right\} \text{ margin}$$

Problem 1.3 (textbook)

$w^*$  distance between  
closest point in  
 $D$  to the  
decision  
boundary.



# Pocket Algorithm

Heuristic that applies to non-separable datasets.

Idea: keep the best weight vector upto current iteration aside

If the current perceptron update yields a better weight vector, i.e lower value of  $E_{in}$ , then replace the best weight vector with current one.

Stop after sufficiently many iterations.

## Pocket Algorithm

Initialize: let  $\hat{w} = w(0)$  be an arbitrary weight vector.  
( $T=10,000$ )

for  $t = 0, 1, 2, \dots, T-1$

- Run one update of the perceptron algorithm to obtain  $w(t+1)$

- Evaluate  $E_{in}(w(t+1))$

- If  $E_{in}(w(t+1)) < E_{in}(\hat{w})$   $\leftarrow$  replace  $\hat{w} = w(t+1)$

end for

# Linear Regression

$$\mathcal{D} = \{ (\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N) \}$$

$$\underline{x}_i = (x_{i0}=1, x_{i1}, x_{i2}, \dots, x_{id}) \in \mathbb{R}^{d+1}$$

$$y_i \in \mathbb{R}$$

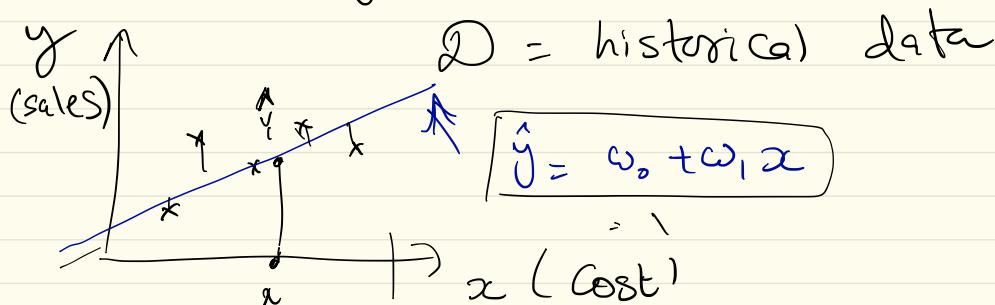
$$\text{Given } \underline{x} = (x_0=1, x_1, x_2, \dots, x_d) \in \mathbb{R}^{d+1}$$

$$\text{Output: } \hat{y} = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_d x_d$$

Model Params

$$\underline{\omega} = (\omega_0, \omega_1, \omega_2, \dots, \omega_d)$$

Example:  $x = \text{advertising cost in 1 week}$   $\in \mathbb{R}^{d+1}$   
 $y = \text{sales in that week}$



$\omega_0 = \text{Sales when } \omega_1 = \text{increase in sales per unit}$   
 $\text{there is no advertising}$   $\omega_1 = \text{increase in advertisement}$

Example 2:

$$\underline{X} = \begin{bmatrix} x_0=1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{cases} \text{cost/week in Internet} \\ \text{cost/week in TV} \\ \text{cost/week in newspaper} \end{cases}$$

adv  
adv  
adv

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

Linear Prediction

$$\hat{y} = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 = \underline{x}^T \underline{\omega}$$

Notation: scalar variable:  $x$

Vector:  $\underline{X} = \begin{bmatrix} x_0=1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^{d+1}$

$$\underline{X}^T = [x_0, x_1, \dots, x_d]$$

Matrix  $\mathcal{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1M} \\ x_{21} & x_{22} & \dots & x_{2M} \\ \vdots & & & \\ x_{N1} & x_{N2} & \dots & x_{NM} \end{bmatrix}$

# Linear Regression : Problem Formulation

Dataset :  $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$\underline{x}_i \in \mathbb{R}^{d+1}, \quad y_i \in \mathbb{R}$$

Data Matrix

$$\mathcal{X} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \in \mathbb{R}^{(N \times d+1)}$$

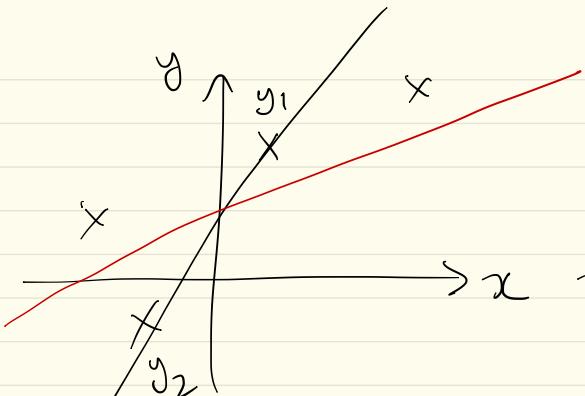
Observation Vector

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Weight vector

$$\underline{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \in \mathbb{R}^{(d+1)} \quad \hat{\underline{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} \underline{x}_1^T \underline{w} \\ \underline{x}_2^T \underline{w} \\ \vdots \\ \underline{x}_N^T \underline{w} \end{bmatrix}$$

Given  $\mathcal{X}, \underline{y}$  find  $\underline{w}$



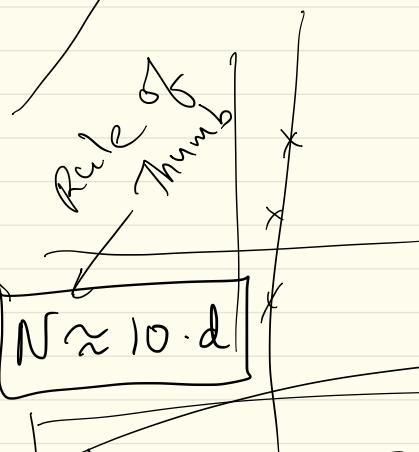
Ideally

$$y_i = \hat{y}_i$$

$\forall i=1, 2, \dots, N$

$$y_i = \underline{x}_i^T \underline{w}$$

for  $i=1, 2, \dots, N$



$$y_1 = w_0 + x_{11} w_1 + x_{12} w_2 \dots + x_{1d} w_d$$

$$y_2 = w_0 + x_{21} w_1 + x_{22} w_2 + \dots + x_{2d} w_d$$

.

$$y_N = w_0 + x_{N1} w_1 + x_{N2} w_2 + \dots + x_{Nd} w_d$$

Unknowns:  $w_0, w_1 \dots w_d$

$$\boxed{N \leq d+1}$$

# of eqs  $\leq$  # of variables.

In practice  $N \gg d$ , so  
 Setting  $\underline{Y} = \hat{\underline{Y}}$  is not feasible

## Loss Function

$$\begin{aligned} E_{in}(\underline{w}) &= \frac{1}{N} \|\underline{Y} - \hat{\underline{Y}}\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - \underline{x}_i^\top \underline{w})^2 \end{aligned}$$

(Squared Error Loss function)

$$\underline{w}^* = \arg \min_{\underline{w} \in \mathbb{R}^{d+1}} E_{in}(\underline{w})$$

$$= \arg \min_{\underline{w} \in \mathbb{R}^{d+1}} \|\underline{Y} - \mathcal{X} \underline{w}\|^2$$

$$= \hat{\underline{Y}}$$

Analytic Solution

# Least Squares Solution

## Loss Function

$$f(\underline{w}) = \left\| \underline{y} - \underline{x} \underline{w} \right\|^2$$

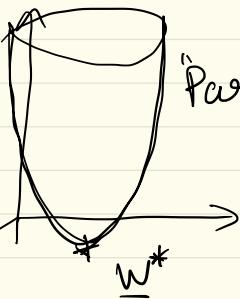
↓                      ↑                      ↑  
 weight vector    obs. vector    data Matrix

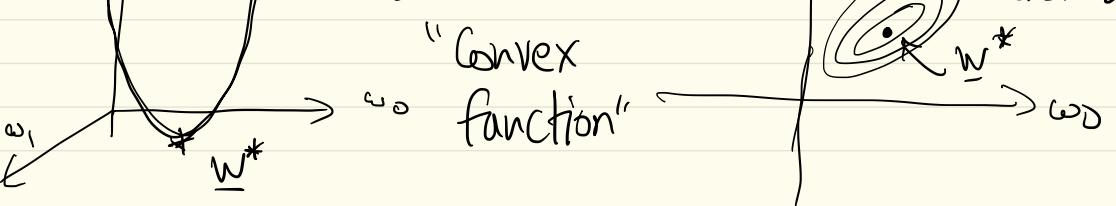
Special Case:  $d = 1$

$$\underline{x}_i = \begin{pmatrix} x_{i0} = 1 \\ x_{i1} \end{pmatrix} \quad \underline{w} = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix}$$

$$f(\underline{w}) = \sum_{i=1}^N (y_i - \underline{x}_i^\top \underline{w})^2$$

$$f(w_0, w_1) = \sum_{i=1}^N (y_i - w_0 - x_{i1} w_1)^2$$

↓                      Given                       $f(w_0, w_1)$   

 "Paraboloid"              "Convex function"



$$\underset{\underline{w} \in \mathbb{R}^{d+1}}{\text{minimize}} \quad f(\underline{w})$$

$$\nabla f(\underline{w}) = \begin{bmatrix} \frac{\partial f}{\partial w_0} \\ \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix} = \underline{0}$$

$$f(\underline{w}) = \| \underline{y} - \mathcal{X} \underline{w} \|^2$$

Gradient:

$$g(\underline{w}) = \underline{w}^T \underline{v} = w_0 v_0 + w_1 v_1 + \dots + w_d v_d$$

$$= \underline{v}^T \underline{w}$$

$$\nabla g(\underline{w}) = \underline{v} \rightarrow \begin{bmatrix} \frac{\partial g}{\partial w_0} \\ \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_d} \end{bmatrix} = \begin{bmatrix} v_0 \\ v_1 \\ \vdots \\ v_d \end{bmatrix} = \underline{v}$$

$$g(\underline{w}) = \underline{w}^T A \underline{w}$$

$$\nabla_{\underline{w}} g(\underline{w}) = \nabla_{\underline{w}} (\underline{w}^T A \underline{w}) = (A + A^T) \underline{w}$$

$$(\text{exercise}) \quad = 2A\underline{w} \quad \text{if } A = A^T$$

1.  $\Leftrightarrow A$  is symmetric

$$f(\underline{w}) = \|\underline{y} - \mathcal{X}\underline{w}\|^2$$

$$\nabla f(\underline{w}) = 0$$

$$f(\underline{w}) = (\underline{y} - \mathcal{X}\underline{w})^T (\underline{y} - \mathcal{X}\underline{w})$$

$$= \underline{y}^T \underline{y} - \underline{y}^T \underbrace{\mathcal{X}\underline{w}}_v - \underbrace{\underline{w}^T \mathcal{X}^T \underline{y}}_{(\mathcal{X}\underline{w})^T (\mathcal{X}\underline{w})} +$$

$$= \|\underline{y}\|^2 - 2 \underline{w}^T \mathcal{X}^T \underline{y} + \underline{w}^T \mathcal{X}^T \mathcal{X}\underline{w}$$

$$\begin{aligned} \nabla f(\underline{w}) &= \nabla_{\underline{w}} \left( \|\underline{y}\|^2 - 2 \underline{w}^T \mathcal{X}^T \underline{y} + \underline{w}^T \underbrace{\mathcal{X}^T \mathcal{X}\underline{w}}_A \right) \\ &= -2 \mathcal{X}^T \underline{y} \end{aligned}$$

$$\boxed{A = (\mathcal{X}^T \mathcal{X}) \quad A^T = \mathcal{X}^T \mathcal{X}}$$

$$+ 2 \mathcal{X}^T \mathcal{X} \underline{w} = 0$$

$$\mathcal{X} : \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$$

$$\underline{y} \in \mathbb{R}^{N \times 1} \quad \mathcal{X}^T \mathcal{X} : \in \mathbb{R}^{(d+1) \times (d+1)}$$

Solve for w

$$\boxed{\mathcal{X}^T \mathcal{X}} w = \mathcal{X}^T y$$

$$M : (d+1) \times (d+1)$$

$$\mathcal{X} = \begin{bmatrix} & \\ & \\ N & \end{bmatrix}^{d+1}$$
$$\mathcal{X}^T = \begin{bmatrix} & \\ & \\ d+1 & \end{bmatrix}^N$$

$$\mathcal{X}^T \mathcal{X} = \begin{bmatrix} & \\ & \\ d+1 & \end{bmatrix}^N$$

full rank Matrix (usually)

$$\underline{w} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$

$\hat{\mathbf{y}} = \mathbf{x} \underline{w} = \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$

Least Square Solution in  
Linear Regression.

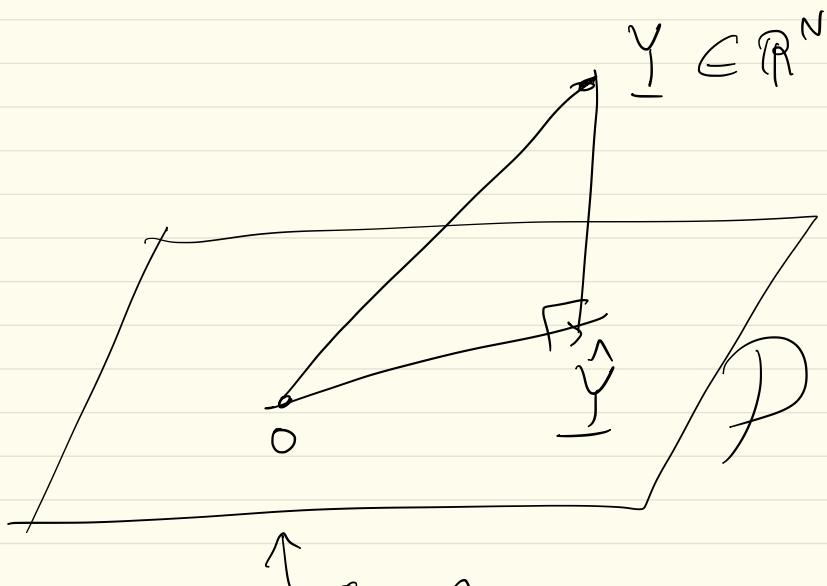
$\mathbf{x}$ :  $N \times (d+1)$  Matrix

Projection Matrix

Geometric Interpretation

$\mathbf{x} = [\underbrace{\mathbf{q}_0 \ \mathbf{q}_1 \ \dots \ \mathbf{q}_d}_{\text{columns of } \mathbf{x}}]$

$\hat{\mathbf{y}} = \mathbf{x} \underline{w} = \mathbf{q}_0 w_0 + \mathbf{q}_1 w_1 + \dots + \mathbf{q}_d w_d$   
 $\in \text{Span}\{\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_d\}$



$$\text{Col-Span } \{x\}$$

$$= \text{Span } \{g_0, g_1, \dots, g_d\}$$

$$\hat{y} = \underbrace{\mathcal{X} (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T y}_{\text{projection matrix}}$$

# Regularized Least Squares

$$\underline{w}_\lambda^* = \arg \min_{\underline{w} \in \mathbb{R}^{d+1}} \left\{ \| \underline{Y} - \mathcal{X} \underline{w} \|^2 + \lambda \| \underline{w} \|^2 \right\}$$

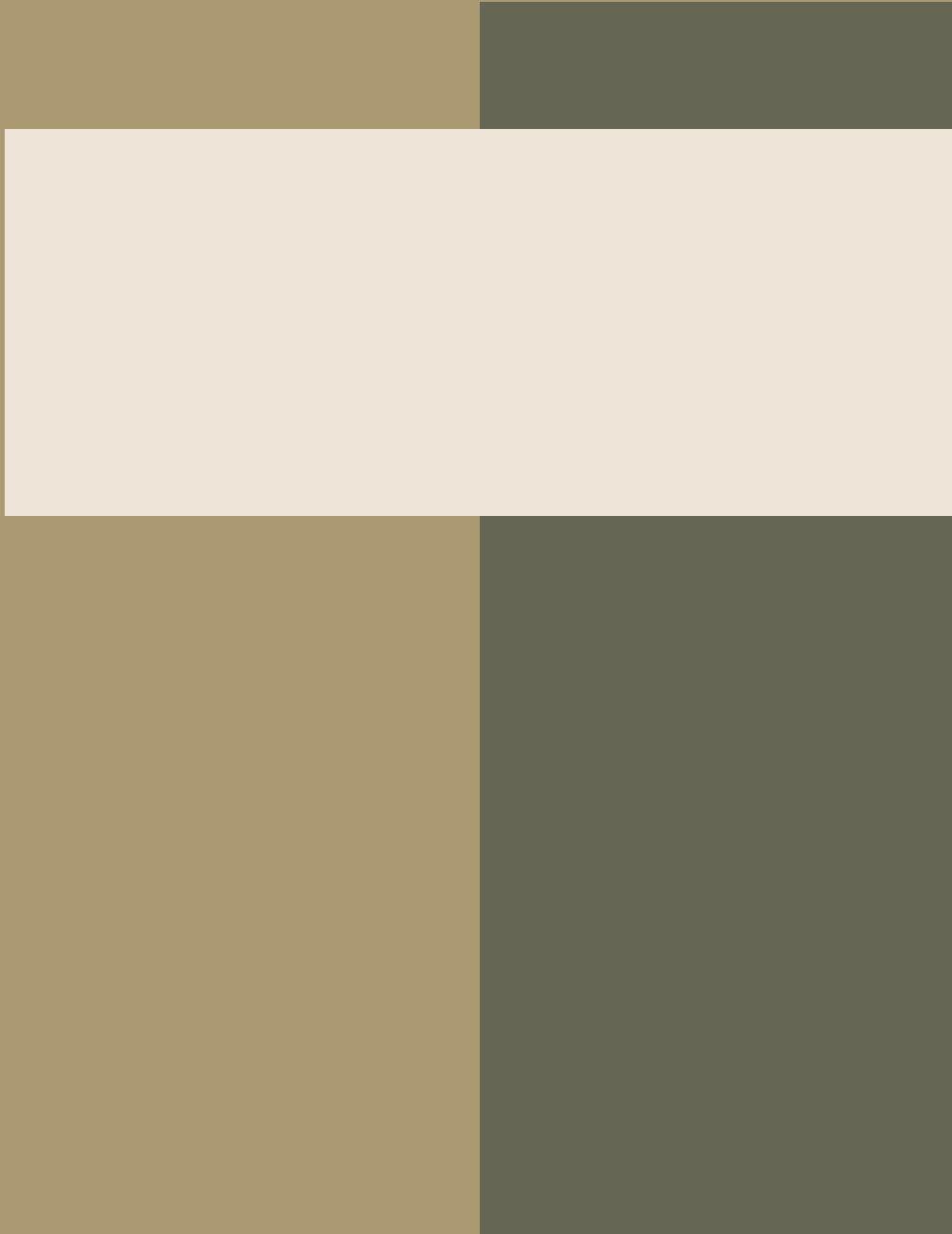
$\lambda$  = Regularization coefficient

$\lambda$  can be selected empirically through cross-validation

$$f(\underline{w}) = \| \underline{Y} - \mathcal{X} \underline{w} \|^2 + \lambda \| \underline{w} \|^2$$

$$\nabla f(\underline{w}) = 0$$

$$\underline{w}_\lambda^* = (\mathcal{X}^T \mathcal{X} + \lambda I)^{-1} \mathcal{X}^T \underline{Y}$$

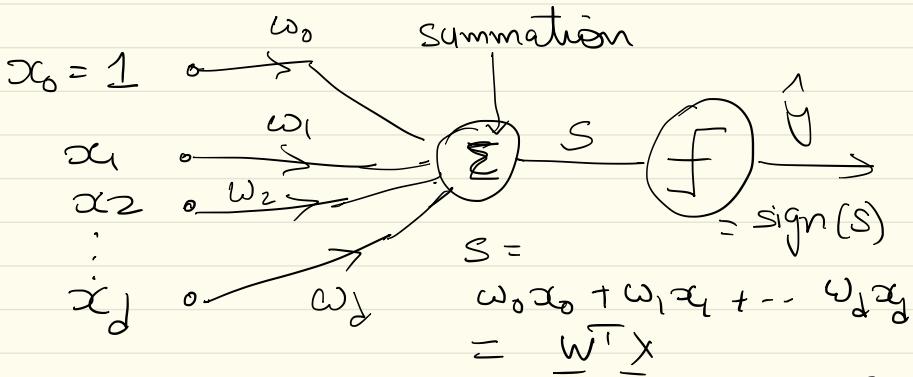


# Logistic Regression

## Linear Classification

Data Vector  $\underline{x} \in \mathbb{R}^{d+1}$   $y \in \{-1, +1\}$

$$\hat{y} = \text{sign}(\underline{w}^T \underline{x})$$



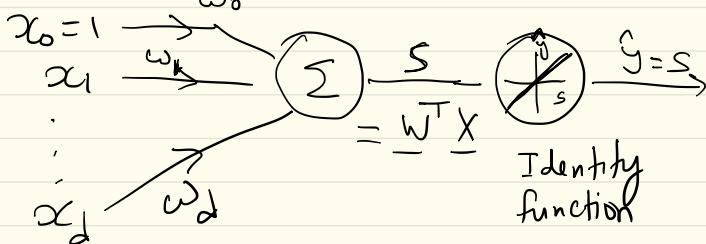
$$\text{Loss Function} = 1 \{ y \neq \text{sign}(S) \}$$

## Linear Regression

Data Vector  $\underline{x} \in \mathbb{R}^{d+1}$   $y \in \mathbb{R}$

$$\hat{y} = \underline{w}^T \underline{x}$$

$$\text{Loss Function} = (y - \hat{y})^2$$

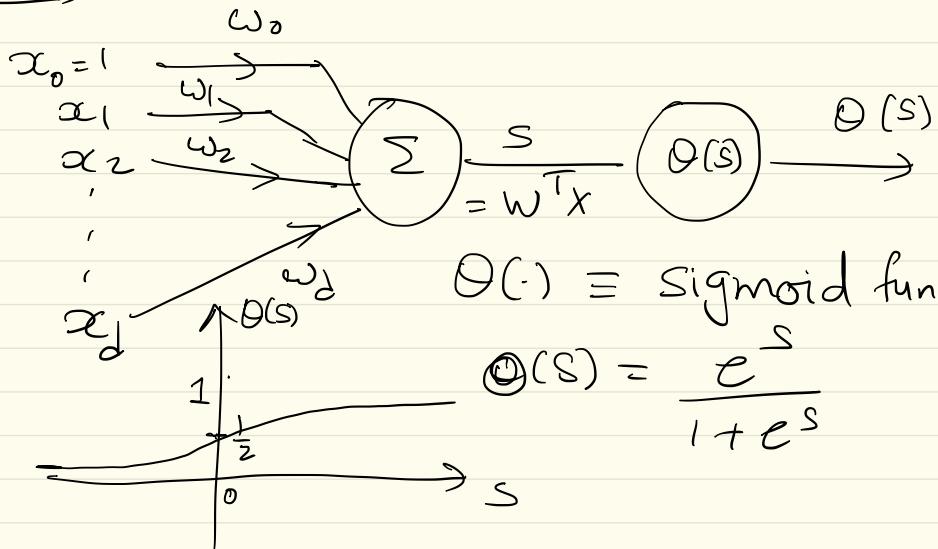


# Binary Classification

$$x \in \mathbb{R}^{d+1}$$

$$y \in \{-1, +1\}$$

Model



$\Theta(\cdot) \equiv$  Sigmoid function

$$\Theta(s) = \frac{e^s}{1 + e^s}$$

Interpretation

$$\begin{aligned} \Theta(s) &= \\ \hookrightarrow &= w^T x \end{aligned}$$

$$\Pr(y=1|x) = \hat{P}_w(1|x)$$

$$\hat{P}_w(1|x) = \Theta(w^T x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

$$\hat{P}_w(-1|x) = 1 - \hat{P}_w(1|x) = \frac{1}{1 + e^{w^T x}} = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$$

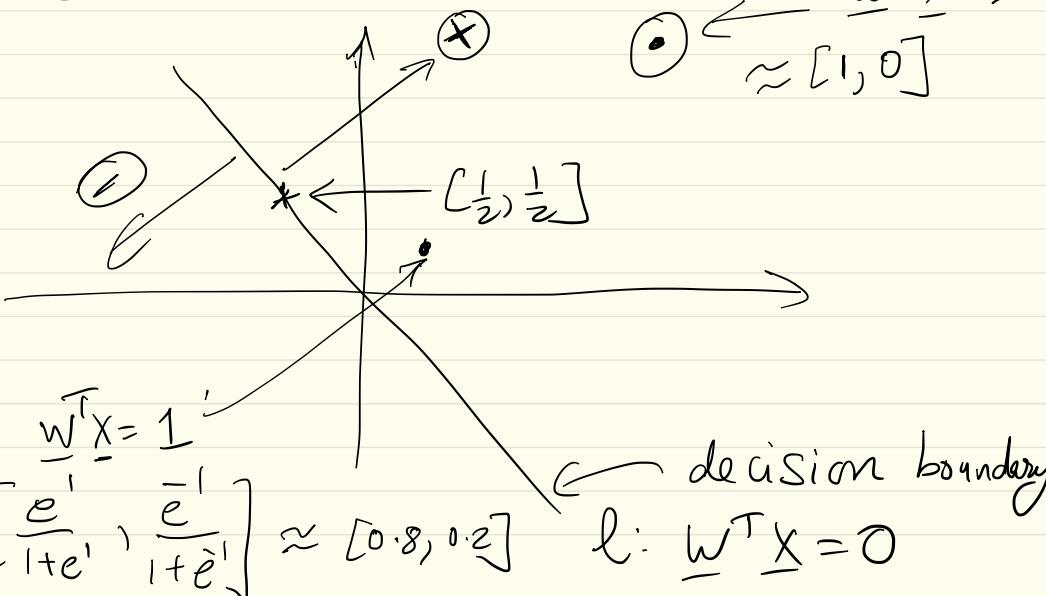
$$\hat{P}_w(y|x) = \frac{e^{y \cdot \underline{w}^T \underline{x}}}{1 + e^{y \cdot \underline{w}^T \underline{x}}}$$

$y \in \{-1, +1\}$

Recap: Input  $\underline{x} \in \mathbb{R}^{d+1}$   
Model Param:  $\underline{w} \in \mathbb{R}^{d+1}$

Output: Probability Vector  
 $[\hat{P}_w(+1|\underline{x}), \hat{P}_w(-1|\underline{x})]$

"Soft decision"



# Training Set

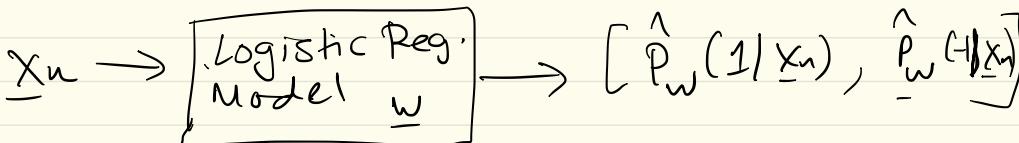
$$\mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$

Model Parameter:  $\underline{w} \in \mathbb{R}^{d+1}$

Output:  $[\hat{P}_w(1|\underline{x}), \hat{P}_w(-1|\underline{x})]$

~~Given~~ Given  $\mathcal{D}$  how to select  $\underline{w}$ ?  
which loss function to select?

Loss Function:  $-\log \hat{P}_w(y_n | \underline{x}_n)$



## Log-loss Function

Example: A given input  $\underline{x} \in \mathbb{R}^{d+1}$

generates the following output

$$[\hat{P}_w(1|\underline{x}), \hat{P}_w(-1|\underline{x})] = [0.8, 0.2]$$

$$\begin{cases} y=+1 & \text{loss} = -\log 0.8 \approx 0.22 \\ y=-1 & \text{loss} = -\log 0.2 \approx 1.39 \end{cases}$$

$$\text{Output Vector} = [0.001, 0.999]$$

$$= [\hat{P}_w(1|x), \hat{P}_w(-1|x)]$$

$$y=+1 \quad \text{loss} = -\log 0.001 \approx 10$$

$$y=-1 \quad \text{loss} \approx -\log 0.999 \approx 10^{-4}$$


---

$E_{in}(\underline{w})$  = training error

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w})$$

where  $e_n(\underline{w})$  = loss associated  
with  $(\underline{x}_n, y_n)$

$$e_n(\underline{w}) = -\log \hat{P}_w(y_n | \underline{x}_n)$$

$$\hat{P}_w(y_n | \underline{x}_n) = \frac{e^{y_n \underline{w}^\top \underline{x}_n}}{1 + e^{y_n \underline{w}^\top \underline{x}_n}} = \frac{1}{1 + e^{-y_n \underline{w}^\top \underline{x}_n}}$$

$$e_n(\underline{w}) = \log \frac{1}{\hat{P}_w(y_n | \underline{x}_n)} = \log \left( 1 + e^{-y_n \underline{w}^\top \underline{x}_n} \right)$$

$$e_n(\underline{w}) = \log \left( 1 + e^{-y_n \underline{w}^T \underline{x}_n} \right)$$

$y_n$	$\underline{w}^T \underline{x}_n$	$y_n \underline{w}^T \underline{x}_n$	$e_n(\underline{w})$
+1	$\gg 0$	$\gg 0$	Small
+1	$\ll 0$	$\ll 0$	large
-1	$\gg 0$	$\ll 0$	large
-1	$\ll 0$	$\gg 0$	small

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w})$$

$$= \frac{1}{N} \sum_{n=1}^N \log \left( 1 + e^{-y_n \underline{w}^T \underline{x}_n} \right)$$

$$\underline{w}^* = \arg \min_{\underline{w} \in \mathbb{R}^{d+1}} E_{in}(\underline{w})$$

# Logistic Regression

$$\mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$
$$\underline{x}_i \in \mathbb{R}^{d+1} \quad y_i \in \{-1, +1\}$$

Model Parameters :  $\underline{w} \in \mathbb{R}^{d+1}$

Given  $\underline{x} \in \mathbb{R}^{d+1}$

outputs  $[\hat{P}_{\underline{w}}(1|\underline{x}), \hat{P}_{\underline{w}}(-1|\underline{x})]$

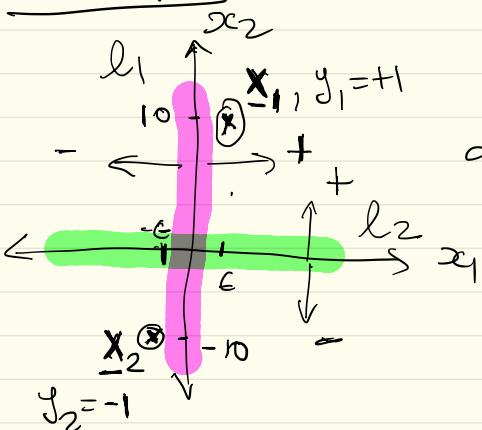
$$\hat{P}_{\underline{w}}(y|\underline{x}) = \frac{e^{y \cdot \underline{w}^T \underline{x}}}{1 + e^{y \cdot \underline{w}^T \underline{x}}} \quad y \in \{-1, +1\}$$

## Loss Function

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n)$$

$$= \frac{1}{N} \sum_{n=1}^N \log \left( 1 + \frac{-y_n}{e} \underline{w}^T \underline{x}_n \right)$$

Example :  $N = 2$   $d = 2$



$$\epsilon \approx 10^{-3}$$

$$\mathcal{D} = \{(x_1, y_1), (x_2, y_2)\}$$

$$x_1 = (1, \epsilon, 10) \quad y_1 = +1$$

$$x_2 = (1, -\epsilon, -10) \quad y_2 = -1$$

$$l_1 : \omega_0 = 0$$

$$l_1 : \underline{w}^T \underline{x} = 0 \Rightarrow \omega_0 + \omega_1 x_1 + \omega_2 x_2 = 0$$

$$[\omega_0 \ \omega_1 \ \omega_2] = [0, 1, 0]$$

evaluate loss function  $E_{in}(w)$   
for (a) linear classification

(b) logistic Regression

$$(a) E_{in}(w) = \frac{1}{2} \sum_{n=1}^2 \mathbb{1} \{ y_n \neq \text{sign}(\underline{w}^T \underline{x}_n) \}$$

$$(b) E_{in}(w) = \frac{1}{2} \left( \log \left( 1 + \frac{y_1}{e} \underline{w}^T \underline{x}_1 \right) + \log \left( 1 + \frac{y_2}{e} \underline{w}^T \underline{x}_2 \right) \right) \\ = \frac{1}{2} \left[ \log \left( 1 + e^{1-\epsilon} \right) + \log \left( 1 + e^{(-1)-\epsilon} \right) \right] \approx 0.693$$

# Regularization

$$l_2 : \quad x_2 = 0$$

$$\underline{w} = (0, 0, 1)$$

$$\underline{x}_1 = (1, \epsilon, 10)$$

$$y_1 = +1$$

$$\underline{x}_2 = (1, -\epsilon, -10)$$

$$y_2 = -1$$

## linear classification

$$E_{in}(\underline{w}) = \frac{1}{2} \sum_{n=1}^2 1 \{ y_n \neq \text{sign}(\underline{w}^\top \underline{x}_n) \} \\ = 0$$

## logistic Regression

$$E_{in}(\underline{w}) = \frac{1}{2} \left( \log \left( 1 + \frac{-y_1 \underline{w}^\top \underline{x}_1}{e} \right) + \log \left( 1 + \frac{-y_2 \underline{w}^\top \underline{x}_2}{e} \right) \right) \\ = \frac{1}{2} \left[ \log \left( 1 + \bar{e}^{1+10} \right) + \log \left( 1 + \bar{e}^{-(-1)-(-10)} \right) \right] \\ = \log \left( 1 + \bar{e}^{10} \right) \approx 5 \times 10^{-4}$$

$$\underline{w}' = (0, 0, 100) \Rightarrow$$

$$100 \quad x_2 = 0$$

$$\Rightarrow x_2 = 0 \Leftrightarrow l_2$$

$$E_{in}(\underline{w}') = E_{in}(\underline{w}) ?$$

$$\min_{\underline{w}} \left[ E_{in}(\underline{w}) + \lambda \|\underline{w}\|^2 \right]$$

$$\frac{1}{N} \sum_{n=1}^N \log \left( 1 + e^{-y_n \underline{w}^\top \underline{x}_n} \right)$$

$\underline{w}$

## Maximum-Likelihood Viewpoint

Given: Training Set

$$\mathcal{D} = \{ (\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N) \}$$

$$\underline{x}_i \in \mathbb{R}^{d+1} \quad y_i \in \{-1, +1\}$$

Data vector sequence:  $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$

label sequence:  $y_1, y_2, \dots, y_N$

logistic Regression Model assigns a probability to each  $y_i$  given  $\underline{x}_i$

Select  $\underline{w} \in \mathbb{R}^{d+1}$  that maximizes

$$P(y_1, y_2, \dots, y_N | \underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)$$

Assume data samples are generated independently

$$P_{\underline{w}}(y_1, y_2, \dots, y_N | \underline{x}_1, \underline{x}_2, \dots, \underline{x}_N) = \prod_{n=1}^N \hat{P}_{\underline{w}}(y_n | \underline{x}_n)$$

Select  $\underline{w}$  that maximized

$$\prod_{n=1}^N \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = \frac{1}{1 + e^{-y_n \underline{w}^\top \underline{x}_n}}$$

$$\underline{w}^* = \arg \max_{\underline{w}} \frac{1}{N} \prod_{n=1}^N \hat{P}_{\underline{w}}(y_n | \underline{x}_n)$$

$$= \arg \max_{\underline{w}} \frac{1}{N} \log \prod_{n=1}^N \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = \arg \min_{\underline{w}} E_{in}(\underline{w})$$

$$= \arg \max_{\underline{w}} \frac{1}{N} \sum_{n=1}^N \log \hat{P}_{\underline{w}}(y_n | \underline{x}_n) - E_{in}(\underline{w})$$

# Cross- Entropy Viewpoint

let  $S = \{S_1, S_2, \dots, S_M\}$  be some discrete alphabets

let  $P$  &  $Q$  be two distributions on  $S$

$$P = (P(S_1), P(S_2), \dots, P(S_M))$$

$$Q = (q(S_1), q(S_2), \dots, q(S_M))$$

Cross- Entropy function

$$CE(P, Q) = - \sum_{i=1}^M P(S_i) \log q(S_i)$$

log-loss function can be viewed as a cross-entropy.

Consider

$$e_n(w) = - \log \hat{P}_w(y_n | \underline{x}_n)$$

$$= \left\{ \begin{array}{l} -1 \{y_n = +1\} \log \hat{P}_w(1 | \underline{x}_n) \\ -1 \{y_n = -1\} \cdot \log \hat{P}_w(-1 | \underline{x}_n) \end{array} \right\} \quad \begin{array}{l} y_n \in \{-1, +1\} \end{array}$$

$$\mathcal{S} = \{-1, +1\}$$

$$CE(P, Q) =$$

$$-P(-1) \log Q(-1) - P(1) \log Q(1)$$

$$= 1 \{y_n = -1\} \log \hat{P}_w(-1 | x_n)$$

$$- 1 \{y_n = +1\} \log \hat{P}_w(+1 | x_n)$$

$$Q = [\hat{P}_w(-1 | x_n), \hat{P}_w(+1 | x_n)]$$

$$P = [1 \{y_n = -1\}, 1 \{y_n = +1\}]$$

$$= \begin{cases} [1, 0] & y_n = -1 \\ [0, 1] & y_n = +1 \end{cases}$$

ground truth probability

④ "distance" between Output (Q)  
↔ ground truth (P)

Training Set

$$\{(x_1, q_1), (x_2, q_2), \dots, (x_n, q_n)\}$$

$$q_i = (q_i(-1), q_i(+1))$$

$$\text{Output: } (\hat{P}_w(-1|x_i), \hat{P}_w(1|x_i))$$

Loss function

$$= -q_i(-1) \log \hat{P}_w(-1|x_i) \\ - q_i(+1) \log \hat{P}_w(1|x_i)$$

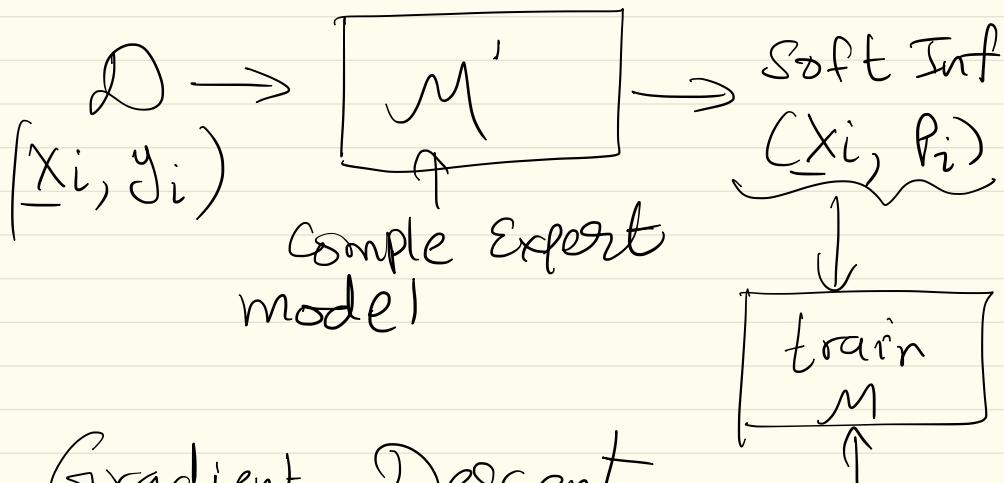
$$= CE(q_i, \hat{P})$$

Knowledge Distillation (2016)

- limited Dataset  $\mathcal{D}$
- train simple Model  $M$
- large & complex Model  $M'$

(Hinton 2016)

Pretrained Model :  $M'$



## Gradient Descent

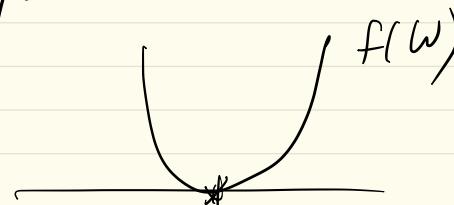
$$\underline{w}^* = \arg \min_{\underline{w} \in \mathbb{R}^{d+1}} \frac{1}{N} \sum_{n=1}^N \log \left( \frac{y_n}{1 + e^{-\underline{w}^T \underline{x}_n}} \right)$$

numerical methods

Convex in  $\underline{w}$

Linear Regression

$$\underline{w}^* = \arg \min_{\underline{w} \in \mathbb{R}^{d+1}} \frac{1}{N} \sum_{i=1}^N (y_i - \underline{w}^T \underline{x}_i)^2$$

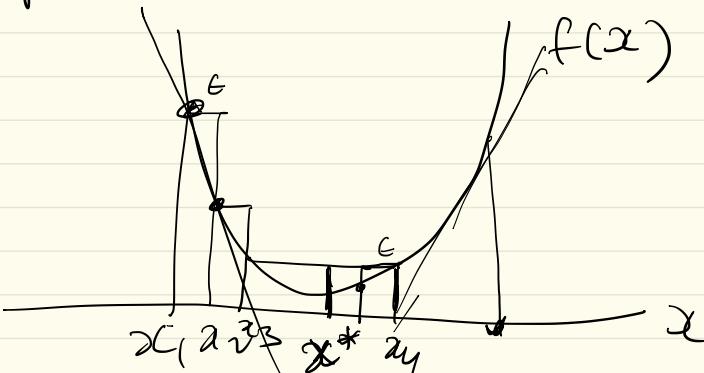


# Gradient Descent

Given some function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   
minimize  $f(\underline{x})$   
 $\underline{x} \in \mathbb{R}^n$

Assume  $f(\underline{x})$  is a differentiable function, but a closed-form solution will not exist in general.

special case  $n = 1$   $f: \mathbb{R} \rightarrow \mathbb{R}$



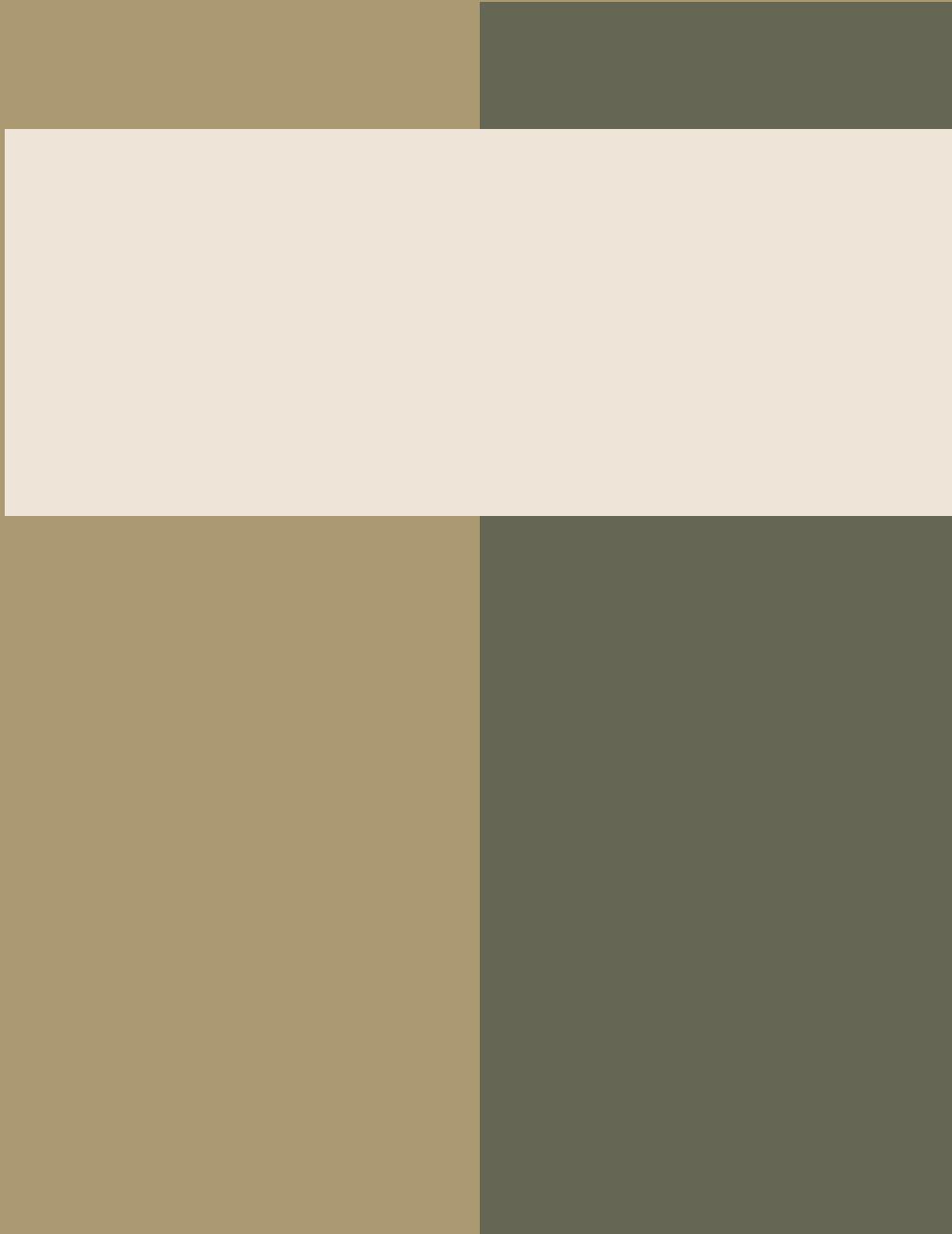
$f(x)$  is  
a convex  
function

Compute  $x^*$  numerically.

If  $x = x^*$   $f'(x) = 0$  Consider  $f'(x) = \frac{\partial f}{\partial x}$

if  $x > x^*$   $f'(x) > 0$

if  $x < x^*$   $f'(x) < 0$



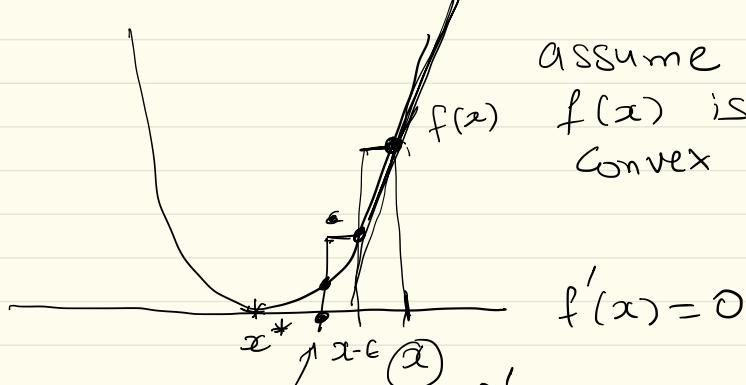
# Gradient Descent

Given a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

minimize  $f(\underline{x})$   
 $\underline{x} \in \mathbb{R}^n$

Assume  $f(\underline{x})$  is differentiable  
but a closed form does not exist

Special Case  $n=1$   $f: \mathbb{R} \rightarrow \mathbb{R}$

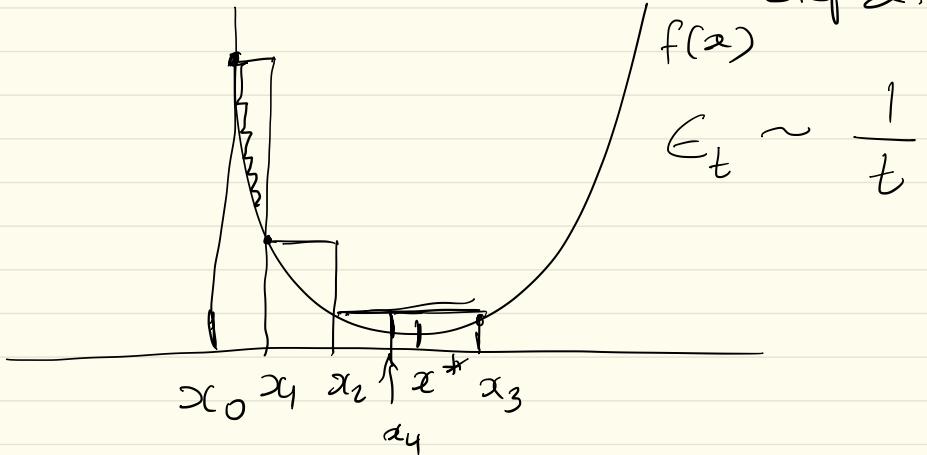


If  $x = x^*$   $x - x^* \Rightarrow f'(x) = 0$

If  $x > x^* \Rightarrow f'(x) > 0$  [is increasing]

If  $x < x^* \Rightarrow \underbrace{f'(x) < 0}_{f(x) \text{ is decreasing}}$

1. Initialize  $x = x_0$
2. If  $f'(x) \approx 0$  then stop & output  $x$
3. If  $f'(x) > 0$   
then  $x \leftarrow x - \epsilon$  goto step 2
4. If  $f'(x) < 0$   
then  $x \leftarrow x + \epsilon$  goto step 2.



typically decrease  $\epsilon$  as the iteration progressed.

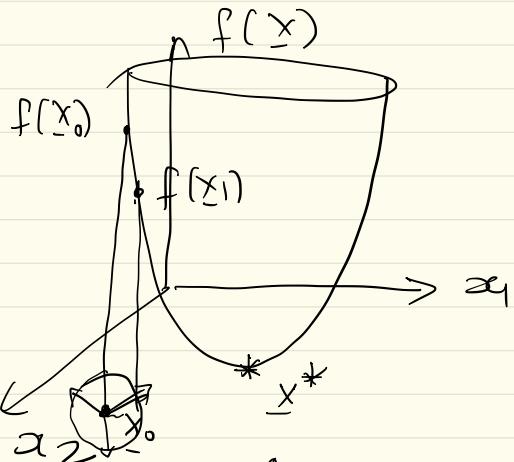
large  $\epsilon \Rightarrow$  less accuracy      small  $\epsilon \Rightarrow$  slow convergence

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

$$n = 2$$

assume  $f(\underline{x})$  is convex.

$$f(x_1, x_2)$$



Given  $\underline{x}_0$  is starting point

$$\underline{x}_1 = \underline{x}_0 + \underline{\epsilon} \cdot \underline{u}$$

$\underline{\epsilon} \equiv$  step size

$$\underline{u} \equiv \text{unit vector } \|\underline{u}\|=1$$

Assume  $\epsilon$  is some small constant (e.g.  $10^{-3}$ )

How to find  $\underline{u}$ ?

$$f(\underline{x}_1) = f(\underline{x}_0 + \epsilon \cdot \underline{u})$$

Select  $\underline{u}$  that minimizes  $f(\underline{x}_0 + \epsilon \cdot \underline{u})$

Taylor Series Approximation

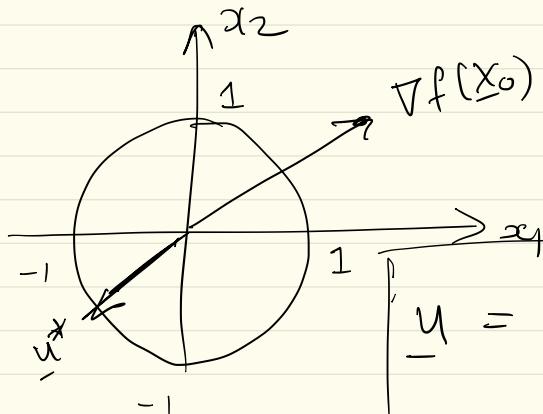
$$f(\underline{x}_1) \approx f(\underline{x}_0) + \epsilon \cdot \underline{u}^T \nabla f(\underline{x}_0) + O(\epsilon^2)$$

$$\underline{u} = \arg \min_{\underline{u}: \|\underline{u}\|=1} \left\{ f(\underline{x}_0) + \epsilon \cdot \underline{u}^T \nabla f(\underline{x}_0) \right\}$$

2<sup>nd</sup> order terms.

$$\underline{u} = \arg \min \left( \underline{u}^T \cdot \underbrace{\nabla f(\underline{x}_0)}_{\begin{bmatrix} \frac{\partial f(\underline{x}_0)}{\partial x_1} \\ \frac{\partial f(\underline{x}_0)}{\partial x_2} \end{bmatrix}} \right)$$

$\text{U: } \|\underline{u}\| = 1$



$$\boxed{\underline{u} = -\frac{\nabla f(\underline{x}_0)}{\|\nabla f(\underline{x}_0)\|}}$$

Given  $\underline{x}_0 \rightarrow \underline{x}_1$

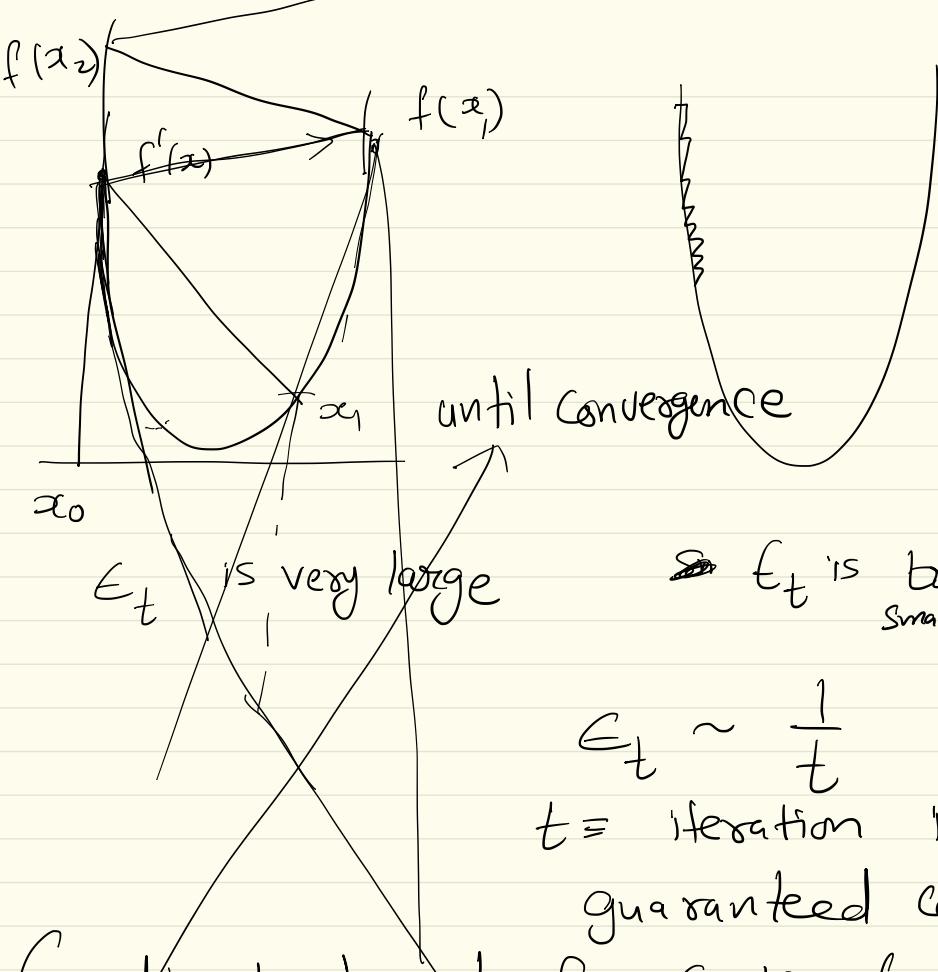
$$\underline{x}_{t+1} = \underline{x}_t + (\text{Step-size})_t \cdot \underline{u}_t$$

$$\underline{u}_t = -\frac{\nabla f(\underline{x}_t)}{\|\nabla f(\underline{x}_t)\|}$$

$$(\text{Step-size})_t = \|\nabla f(\underline{x}_t)\| \cdot \epsilon_t$$

$$\boxed{\underline{x}_{t+1} = \underline{x}_t - \epsilon_t \nabla f(\underline{x}_t)}$$

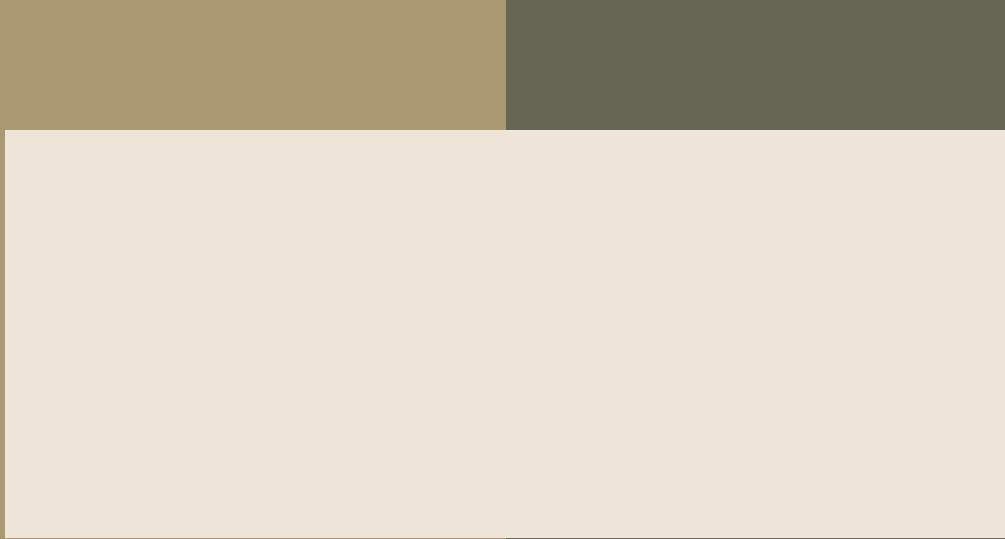
learning rate



## Gradient descent for Convex functions

Initialize  $\underline{x}_0$  in some arbitrary fashion  
for  $t = 0, 1, 2, \dots$

- compute  $\underline{g}_t = \nabla f(\underline{x}_t)$
  - Set  $\underline{v}_t = -\underline{g}_t$
  - update  $\underline{x}_{t+1} = \underline{x}_t + E_t \cdot \underline{v}_t$
- Continue



## Announcements

- 1) Assignment 1: due Oct 9<sup>th</sup> in class  
& online (on quercus)
  - 2) Midterm: Oct 16<sup>th</sup> (locations TBA)
- 

## Gradient Descent for Linear Regression

Loss function

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w})$$

$$e_n(\underline{w}) = (\underline{w}^T \underline{x}_n - y_n)^2$$

$$\underline{w}^* = \arg \min_{\underline{w} \in \mathbb{R}^{d+1}} E_{in}(\underline{w})$$

$\Rightarrow$  find a closed-form (analytical)  
solution.

$$\underline{w}_{LS}^* = (\underbrace{\underline{x}^T \underline{x}}_{d+1 \times N}^{-1}) \underline{x}^T \underline{y}$$

Apply Gradient descent to  
compute  $\underline{w}^*$

$$\nabla_{\underline{w}} f(\underline{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix}$$

Update Rule       $\forall t \geq 0$

$$\underline{w}_{t+1} = \underline{w}_t - \epsilon_t \underbrace{\nabla E_{in}(\underline{w}_t)}_{\substack{\sum_{n=1}^N (\underline{w}_t^T \underline{x}_n - y_n)^2 \\ = e_n(\underline{w})}}$$

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N (\underline{w}^T \underline{x}_n - y_n)^2$$

$= e_n(\underline{w})$

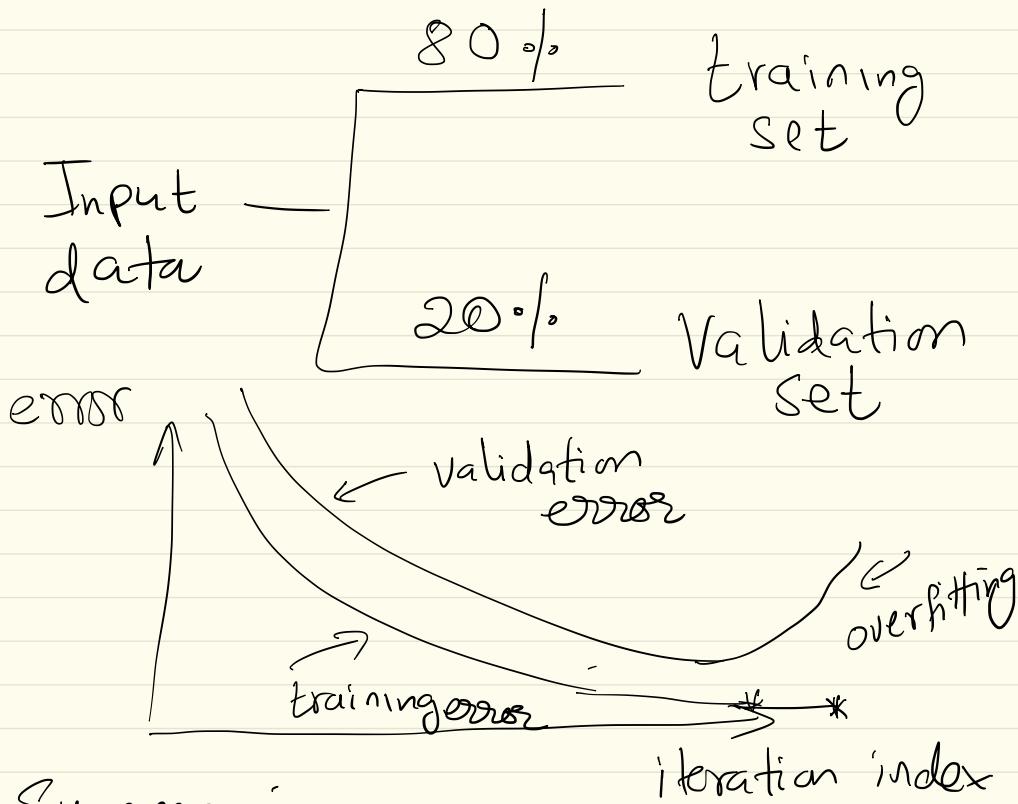
$$\begin{aligned} \nabla E_{in}(\underline{w}_t) &= \nabla_{\underline{w}_t} \left( \frac{1}{N} \sum_{n=1}^N (\underline{w}_t^T \underline{x}_n - y_n)^2 \right) \\ &= \frac{1}{N} \sum_{n=1}^N \nabla_{\underline{w}_t} (\underline{w}_t^T \underline{x}_n - y_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^N 2 (\underline{w}_t^T \underline{x}_n - y_n) \nabla_{\underline{w}_t} (\underline{w}_t^T \underline{x}_n - y_n) \\ &= \frac{1}{N} \sum_{n=1}^N 2 (\underline{w}_t^T \underline{x}_n - y_n) \underline{x}_n \end{aligned}$$

---


$$\underline{w}_{t+1} = \underline{w}_t - \frac{2\epsilon_t}{N} \sum_{n=1}^N (\underline{w}_t^T \underline{x}_n - y_n) \underline{x}_n$$


---

as  $t \rightarrow \infty$      $\underline{w}_t \rightarrow \underline{w}_{LS}$



### Summarize

Least - Squares Solution

- (1) Can computationally expensive
- (2) Cause overfitting
- (3) Can be an overkill

## Gradient Descent update for Linear Regression

$$\underline{w}_{t+1} = \underline{w}_t - \epsilon_t \frac{1}{N} \sum_{n=1}^N 2 \underline{x}_n (\underline{w}_t^T \underline{x}_n - y_n)$$

For each iteration require  
 $\Theta(Nd)$  computations

## Stochastic Gradient Descent

at each iteration select one  
training example at random  
(uniformly) and compute

$$\underline{w}_{t+1} = \underline{w}_t - \epsilon_t \nabla C_n(\underline{w}_t)$$
$$n \sim \text{Unif}\{1, 2, \dots, N\}$$

$$\underline{w}_{t+1} = \underline{w}_t - 2\epsilon_t (\underline{w}_t^T \underline{x}_n - y_n) \cdot \underline{x}_n$$

# Justification for SGD

Consider

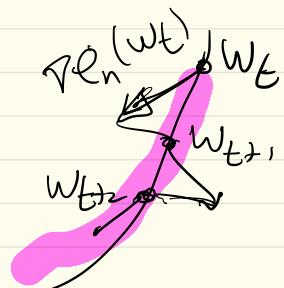
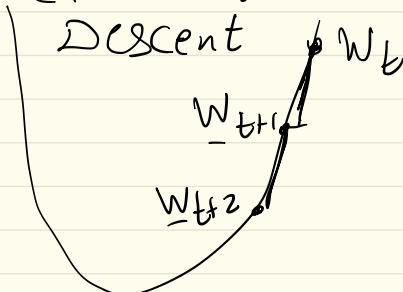
$$\mathbb{E}_n \left[ \nabla e_{(n)}(w_t) \right] = \nabla E_{in}(w_t)$$

$$n \in \text{Unif } \{1, 2, \dots, N\}$$

$$\Pr(n=1) \nabla e_1(w_t) + \Pr(n=2) \nabla e_2(w_t) + \dots + \Pr(n=N) \nabla e_N(w_t)$$

$$= \frac{1}{N} (\nabla e_1(w_t) + \nabla e_2(w_t) + \dots + \nabla e_N(w_t))$$

Exact Gradient



# Mini-Batch Gradient Descent

at each iteration draw  $M$  examples from the dataset at random (without replacement)

let  $S_k = \{j_1, j_2, \dots, j_M\}$  denote the set of examples randomly selected

## Update Rule

$$\underline{w}_{k+1} = \underline{w}_k - \frac{\epsilon_k}{M} \sum_{n=1}^M \nabla e_{(\underline{w}_n)}(w_k)$$

$M = 1 \Rightarrow$  recover SGD

$M = N \Rightarrow$  Exact Gradient descent

$M = 32, 64$  etc.

# Logistic Regression

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w})$$

$$e_n(\underline{w}) = \log \left( 1 + \frac{-y_n \underline{w}^T \underline{x}_n}{e} \right)$$

SGD update

$$\underline{w}_{k+1} = \underline{w}_k - \epsilon_k \nabla e_n(\underline{w}_k)$$

$$n \in \text{Unif}\{1, 2, \dots, N\}$$

$$\nabla e_n(\underline{w}_k) = \nabla_{\underline{w}_k} \left( \log \left( 1 + \frac{-y_n \underline{w}_k^T \underline{x}_n}{e} \right) \right)$$

$$= \frac{1}{1 + \frac{-y_n \underline{w}_k^T \underline{x}_n}{e}} \nabla_{\underline{w}_k} \left( 1 + \frac{-y_n \underline{w}_k^T \underline{x}_n}{e} \right) = -y_n \underline{x}_n$$

$$= \frac{\nabla_{\underline{w}_k} \left( e^{-y_n \underline{w}_k^T \underline{x}_n} \right)}{1 + \frac{-y_n \underline{w}_k^T \underline{x}_n}{e}} = \frac{-y_n \underline{w}_k^T \underline{x}_n}{e} \nabla_{\underline{w}_k} \left( -y_n \underline{w}_k^T \underline{x}_n \right)$$

Update Rule :  $\underline{w}_{k+1} = \underline{w}_k + \frac{\epsilon_k y_n \underline{x}_n}{1 + e^{y_n \underline{w}_k^T \underline{x}_n}}$

$$\underline{w}_{k+1} = \underline{w}_k + \frac{\epsilon_k y_n \underline{x}_n}{1 + e^{y_n \underline{w}_k^T \underline{x}_n}}$$

$n \in \{1, 2, \dots, N\}$

### Special Cases

- Select a highly misclassified point

$$y_n \underline{w}_k^T \underline{x}_n \ll 0$$

$$e^{y_n \underline{w}_k^T \underline{x}_n} \approx 0$$

$$\underline{w}_{k+1} = \underline{w}_k + \epsilon_k y_n \underline{x}_n$$

Set  $\epsilon_k = 1 \Rightarrow$  Perceptron Alg

$\rightarrow$  Alternately  $y_n \underline{w}_k^T \underline{x}_n \gg 0 \Rightarrow \underline{w}_{k+1} \approx \underline{w}_k$

# Multi-class Logistic Regression

Input Dataset

$$\mathcal{D} = \{ (\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_n, y_n) \}$$
$$\underline{x}_i \in \mathbb{R}^{d+1} \quad y_i \in \{1, 2, \dots, C\}$$

We will consider  $C=3$

let  $\underline{w}(1)$ ,  $\underline{w}(2)$ ,  $\underline{w}(3) \in \mathbb{R}^{d+1}$   
be weight vectors for classes  
 $c_1$ ,  $c_2$  &  $c_3$ .

$$\mathcal{S} = \{ \underline{w}(1), \underline{w}(2), \underline{w}(3) \}$$

(Model Parameters)

Output:  $\left[ \hat{P}_{\mathcal{S}}(1|\underline{x}), \hat{P}_{\mathcal{S}}(2|\underline{x}), \hat{P}_{\mathcal{S}}(3|\underline{x}) \right]$   
given  $\underline{x} \in \mathbb{R}^{d+1}$

# Soft Max - classification

$$\hat{P}_{S2}(i | \underline{x}) = \frac{e^{\underline{w}^T(i) \underline{x}}}{\sum e^{\underline{w}^T(1) \underline{x}} + e^{\underline{w}^T(2) \underline{x}} + e^{\underline{w}^T(3) \underline{x}}}$$

$i \in \{1, 2, 3\}$

## Loss Function

Log-loss function  $(\underline{x}_n, y_n)$

$$y_n \in \{1, 2, 3\}$$

Output  $\left[ \hat{P}_{S2}(1 | \underline{x}_n), \hat{P}_{S2}(2 | \underline{x}_n), \hat{P}_{S2}(3 | \underline{x}_n) \right]$

$$e_n(S2) = -\log \hat{P}_{S2}(y_n | \underline{x}_n)$$

$$\text{find } S2^* = \arg \min_{S2} \frac{1}{N} \sum_{n=1}^N e_n(S2)$$

$$S2^* = \left( \underline{w}^*(1), \underline{w}^*(2), \underline{w}^*(3) \right) \quad \begin{cases} \text{using} \\ \text{SGD} \end{cases}$$

## SGD update

Iteration # $k$

$$\mathcal{S}_k = (\underline{w}_k(1), \underline{w}_k(2), \underline{w}_k(3))$$

$$\underline{w}_{k+1}(1) = \underline{w}_k(1) - \epsilon_k \nabla_{\underline{w}_k(1)} e_n(\omega)$$

$$n \in \text{Unif } \{1, 2, \dots, N\}$$

$$\ell \in \{1, 2, 3\}$$

$$\underline{w}_{k+1}(1) = \underline{w}_k(1) - \epsilon_k \nabla_{\underline{w}_k(1)} e_n(\omega)$$

$$\underline{w}_{k+1}(2) = \underline{w}_k(2) - \epsilon_k \nabla_{\underline{w}_k(2)} e_n(\omega)$$

$$\underline{w}_{k+1}(3) = \underline{w}_k(3) - \epsilon_k \nabla_{\underline{w}_k(3)} e_n(\omega)$$

Standard SGD

Compute

$$\nabla_{w_k(1)} e_n (\mathcal{D}_k)$$

$$= \nabla_{w_k(1)} \left( -\log \hat{P}_{\mathcal{D}_k} (y_n | \underline{x}_n) \right)$$

$$= \nabla_{w_k(1)} \left[ -\log \frac{\underline{w}_k^T (y_n) \cdot \underline{x}_n}{e^{\underline{w}_k^T (1) \cdot \underline{x}_n} + e^{\underline{w}_k^T (2) \cdot \underline{x}_n} + e^{\underline{w}_k^T (3) \cdot \underline{x}_n}} \right]$$

$$= \nabla_{w_k(1)} \left[ -\underline{w}_k^T (y_n) \cdot \underline{x}_n + \log \left( e^{\underline{w}_k^T (1) \cdot \underline{x}_n} + e^{\underline{w}_k^T (2) \cdot \underline{x}_n} + e^{\underline{w}_k^T (3) \cdot \underline{x}_n} \right) \right]$$

$y_n \in \{1, 2, 3\}$

$$y_n = 1 : \quad \nabla_{w_k(1)} \left( \underline{w}_k^T (1) \cdot \underline{x}_n \right) = \underline{x}_n$$

$$y_n = 2, 3 \quad \nabla_{w_k(1)} \left( \underline{w}_k^T (y_n) \cdot \underline{x}_n \right) = 0$$

Verify

$$\nabla_{w_k(l)} \left( \log \left( e^{\underline{w}_k^\top(1) \cdot \underline{x}_n} + e^{\underline{w}_k^\top(2) \cdot \underline{x}_n} + e^{\underline{w}_k^\top(3) \cdot \underline{x}_n} \right) \right)$$

$$= \frac{e^{\underline{w}_k^\top(l) \cdot \underline{x}_n}}{\sum_{j=1}^3 e^{\underline{w}_k^\top(j) \cdot \underline{x}_n}} \cdot \underline{x}_n$$

Show that

$$\nabla_{w_k(l)} \left[ -\log \hat{P}_{\mathcal{D}_K}(\underline{x}_n) \right]$$

$$= \begin{cases} -\underline{x}_n + \frac{e^{\underline{w}_k^\top(l) \cdot \underline{x}_n}}{\sum_{j=1}^3 e^{\underline{w}_k^\top(j) \cdot \underline{x}_n}} & \text{if } y_n = l \\ \frac{e^{\underline{w}_k^\top(l) \cdot \underline{x}_n}}{\sum_{j=1}^3 e^{\underline{w}_k^\top(j) \cdot \underline{x}_n}} & \text{if } y_n \neq l \end{cases}$$

# Momentum Based Methods

Basic SGD update

$$\underline{g}_t = \nabla \ell_n(\underline{w}_t)$$

$$\underline{w}_{t+1} = \underline{w}_t - \epsilon_t \underline{g}_t$$

SGD + Momentum  $\underline{v}_0 = 0$

$$\underline{g}_t = \nabla \ell_n(\underline{w}_t)$$

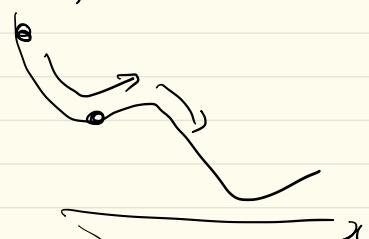
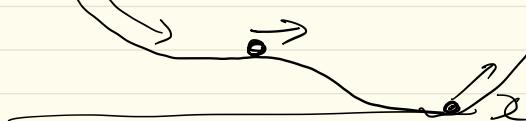
$$\underline{v}_t = -\epsilon_t \underline{g}_t + \eta \cdot \underline{v}_{t-1}$$

typically  $\eta \approx 0.9$

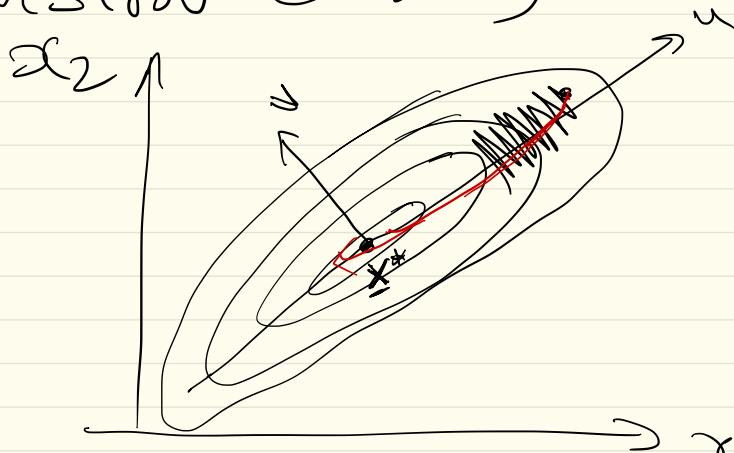
$$\underline{w}_{t+1} = \underline{w}_t + \underline{v}_{t_f(x)}$$

$\underline{v} \equiv$  velocity vector

$f(x)$

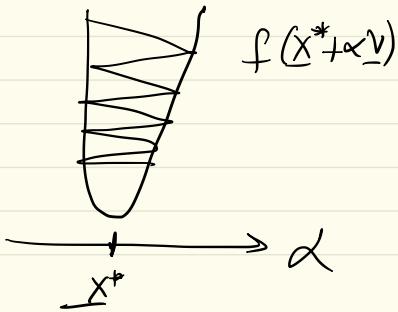
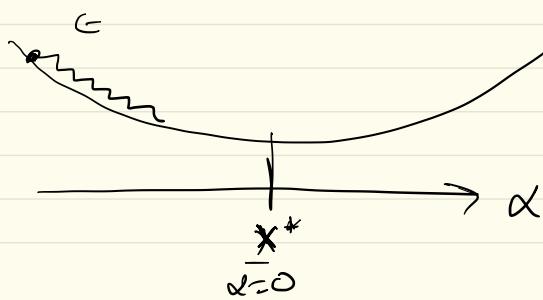


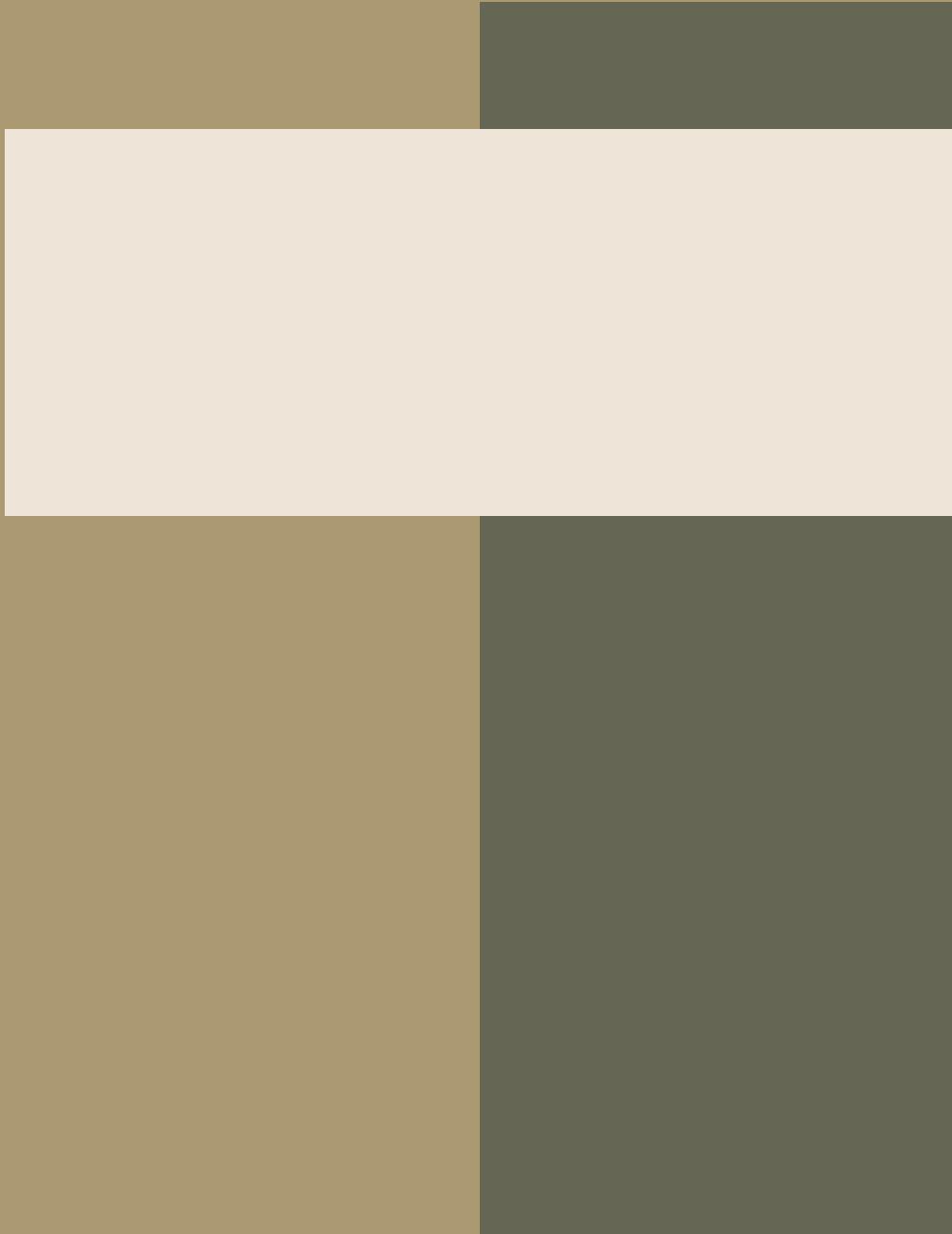
Even for convex function  
 momentum based methods have  
 provably faster convergence  
 than gradient descent  
 Nesterov (1983)



slice along "u" direction  
 $f(x^* + \alpha u)$

slice along "v" direction  
 $f(x^* + \alpha v)$





# Announcements

- 1) Assignment # 1: Due on Wed.
  - hard copy in class
  - Code submission online
- 2) Mid-term - Next Wednesday
  - will cover upto momentum methods (first part of today's lecture)
- 3) Wed Oct 9<sup>th</sup> lecture by TA  
NO office hour on wed.
  - Extra off : TBA

# Momentum Method.

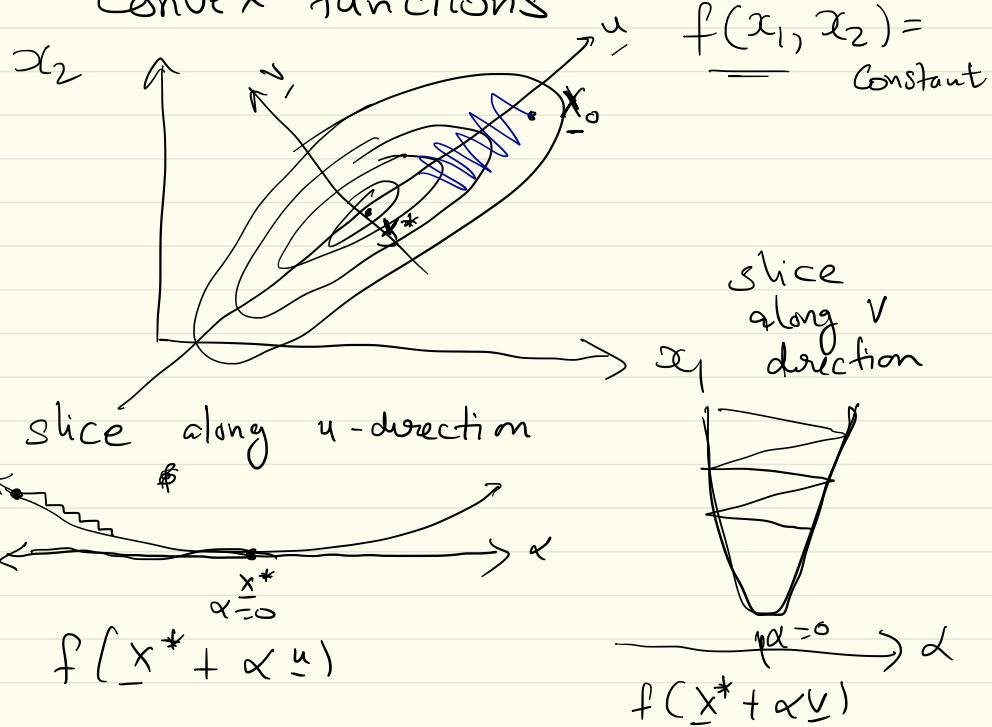
gradient vector:  $\underline{g}_t = \nabla e_n(\underline{w}_t)$

velocity vector:  $\underline{v}_t = -\epsilon_t \underline{g}_t$

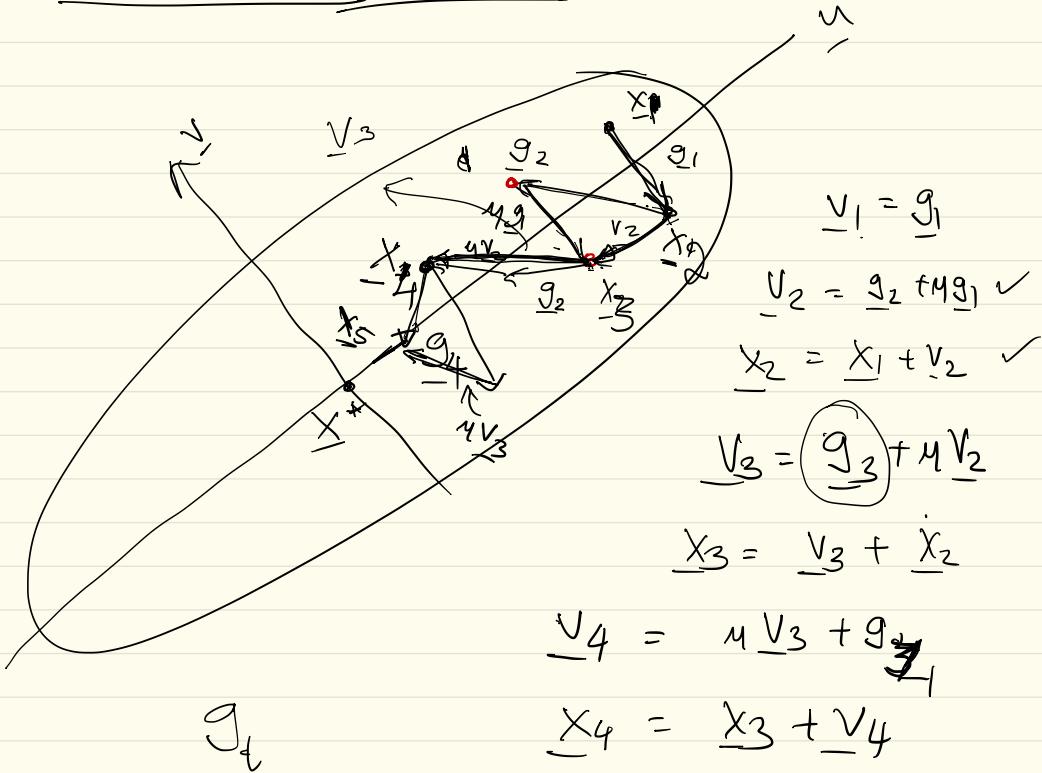
$$\underline{v}_0 = 0$$

$$\underline{w}_{t+1} = \underline{w}_t + \underline{v}_t$$

Momentum methods lead to faster convergence even for convex functions



## Momentum Method



## Nestrov Momentum

$$v_t = \gamma v_{t-1} - \epsilon_t \nabla c_n(w_t + \gamma v_{t-1})$$

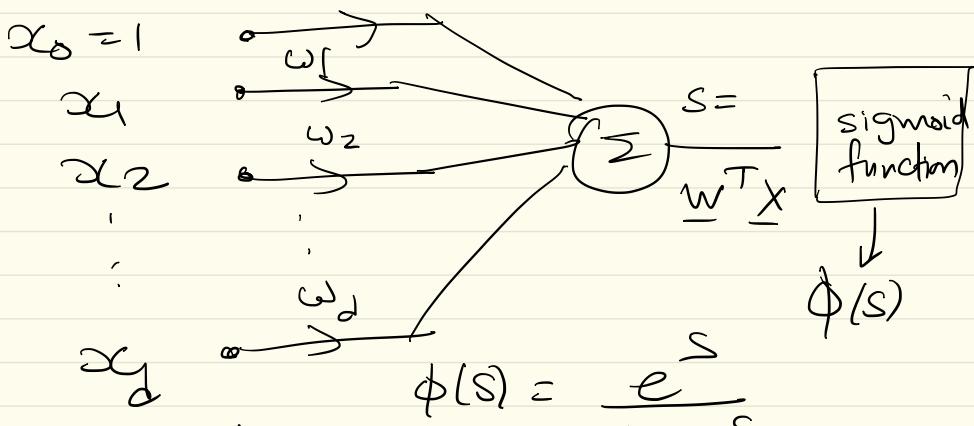
$$w_{t+1} = w_t + v_t$$

ADAM Method (<sup>built into</sup> Tensorflow)

# Neural Networks

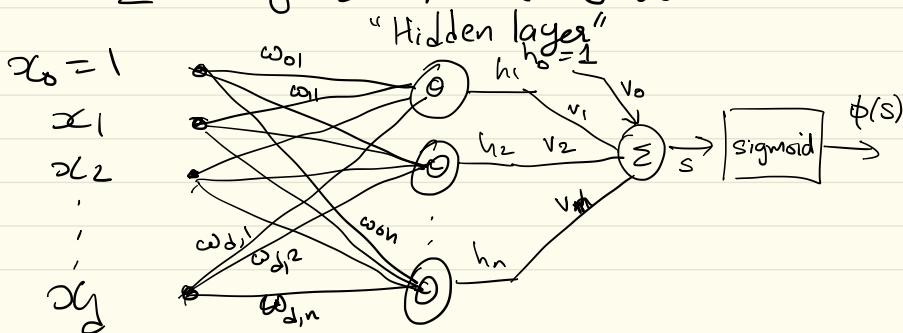
Recall: Binary logistic regression  
"single layer network"

Input layer  $\omega_0$



$$\phi(s) = \hat{P}_w(1|x)$$

2-layer Network



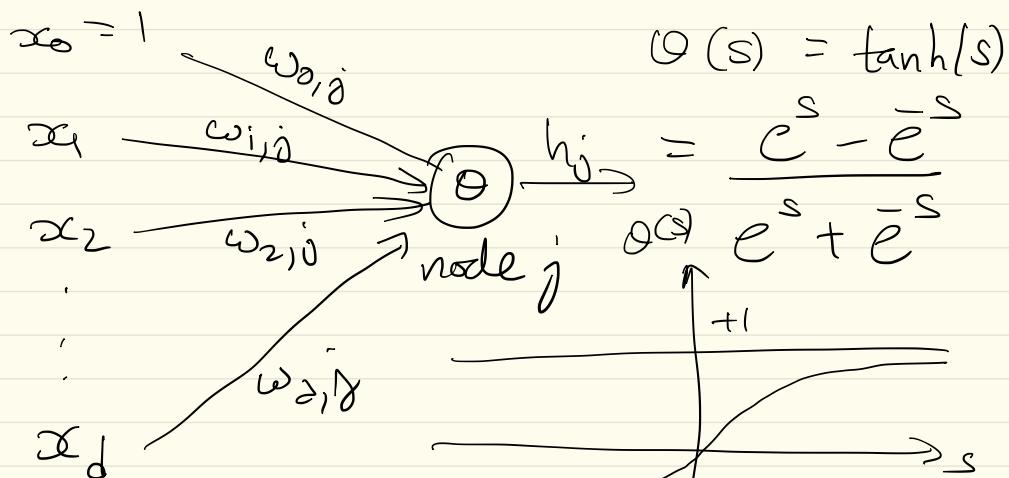
$\Theta$  is a non linear function.

$\omega_{i,j}$  = weight between input node i & hidden layer node j

$v_i$  = weight on the edge between hidden node i & output node

$h_j$  = output of hidden node j

$$h_j = \Theta(\omega_{0j} + \sum_{i=1}^d \omega_{ij} x_i)$$



$$\Theta(s) \approx s \text{ if } |s| \ll 1$$

$$\Theta'(s) = 1 - \Theta^2(s)$$

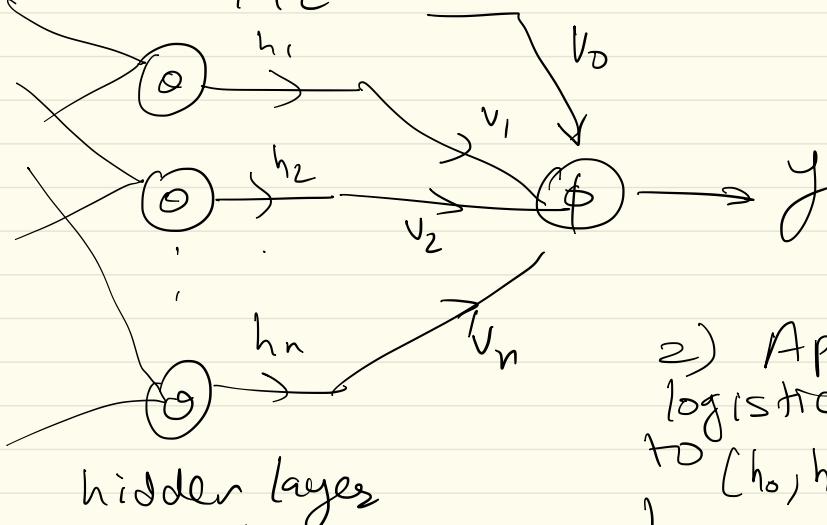
$$\text{ReLU}(s) = \max(0, s)$$

$$h_j = \Theta(\underline{w}_j^T \underline{x})$$

$$\underline{w}_j = (\omega_{0,j} \quad \omega_{1,j} \quad \dots \quad \omega_{d,j})^T$$

$$y = \phi(v_0 + \sum_{j=1}^n v_j h_j)$$

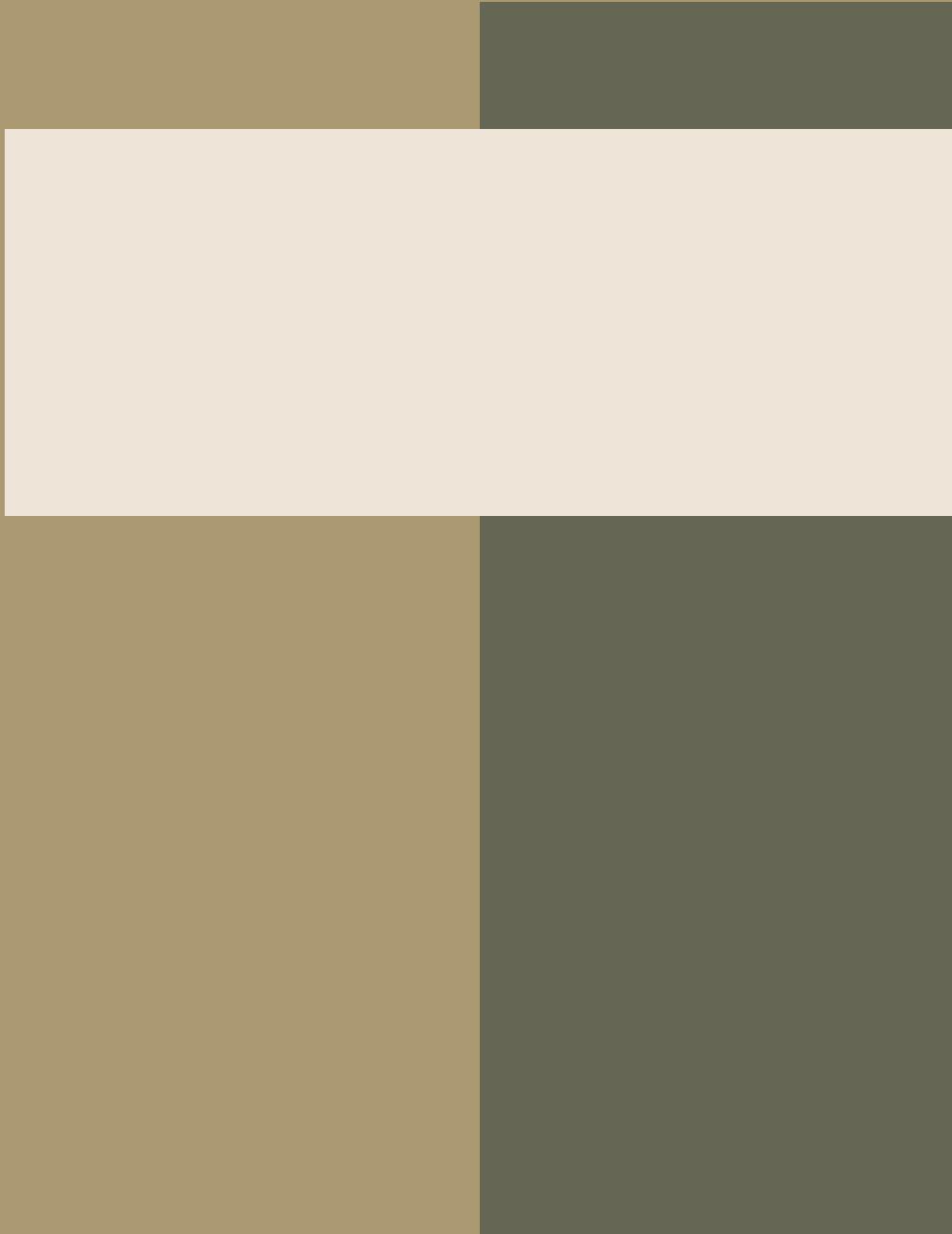
$$\phi(s) = \frac{e^s}{1+e^s} \quad h_0 = 1$$



2) Apply logistic regression  
to  $(h_0, h_1, \dots, h_n)$

Two-step approach

i) Apply non-linear transforms  
 $h_j = \Theta(\underline{w}_j^T \underline{x})$



# Training of Neural Networks

Notation :  $L$  layers.

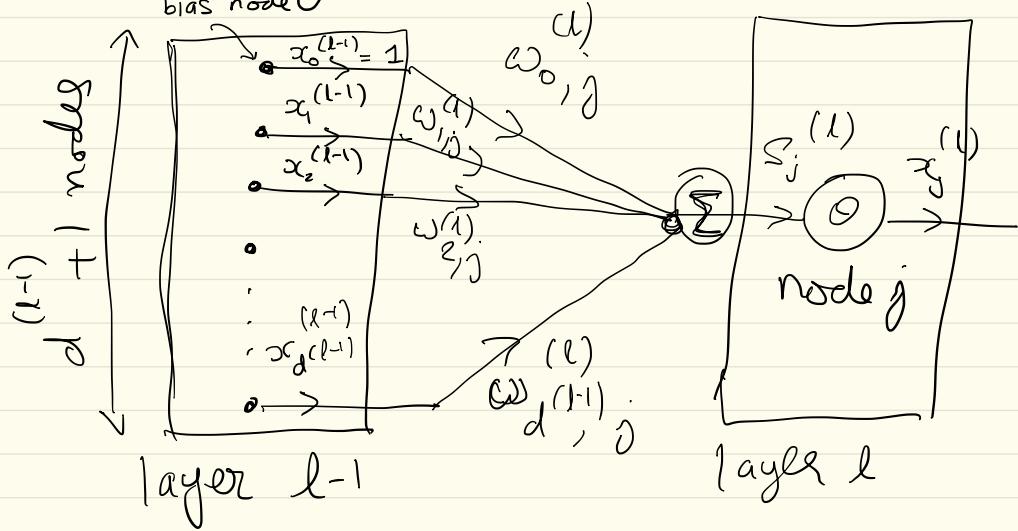
Input layer :  $l = 0$

Hidden layers :  $1 \leq l \leq L-1$

Output layer :  $l = L$

$w_{i,j}^{(l)}$  : weight connecting node  $i$  in layer  $(l-1)$  & node  $j$  in layer  $l$

$d^{(l)}$  : number of nodes in layer  $l$   
(excluding bias node)



$x_j^{(l)}$  = output from node  $j$  in layer  $l$

$s_j^{(l)}$  = input into node  $j$  in layer  $l$

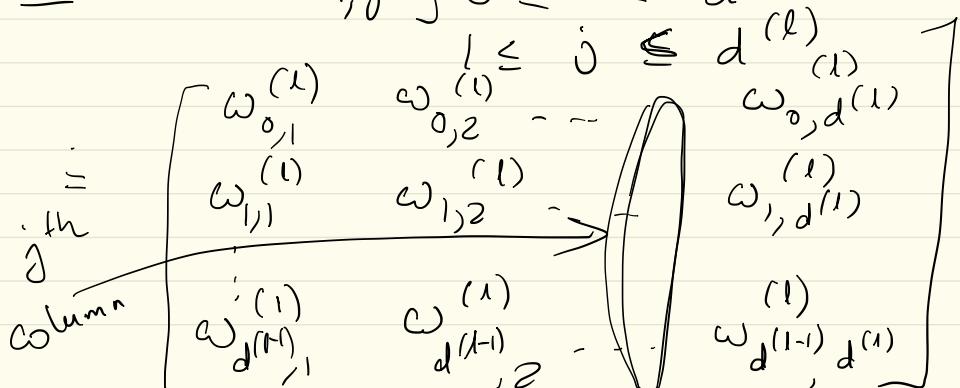
$$s_j^{(l)} = \omega_{0,j}^{(l)} + \sum_{i=1}^{d^{(l-1)}} \omega_{i,j}^{(l)} \cdot x_i^{(l-1)}$$

$$x_j^{(l)} = \Theta(s_j^{(l)})$$

## Vector Notation

$$\underline{s}^{(l)} = \begin{bmatrix} s_1^{(l)} \\ s_2^{(l)} \\ \vdots \\ s_d^{(l)} \end{bmatrix} \quad \Theta(\underline{s}^{(l)}) = \begin{bmatrix} \Theta(s_1^{(l)}) \\ \Theta(s_2^{(l)}) \\ \vdots \\ \Theta(s_d^{(l)}) \end{bmatrix}$$

$$\underline{w}^{(l)} = \{ \omega_{i,j}^{(l)} \} \quad 0 \leq i \leq d^{(l-1)}$$

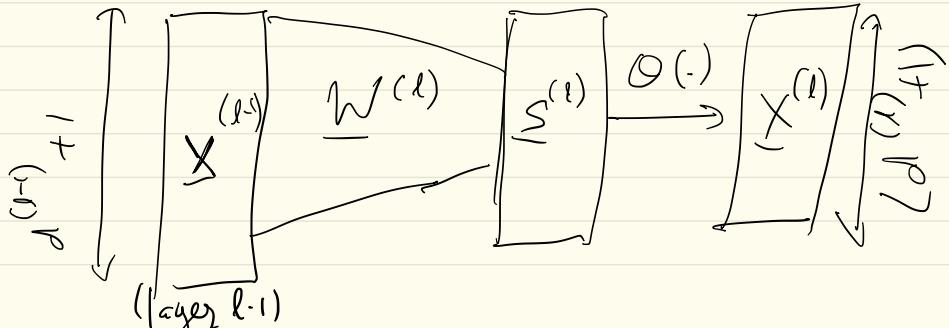
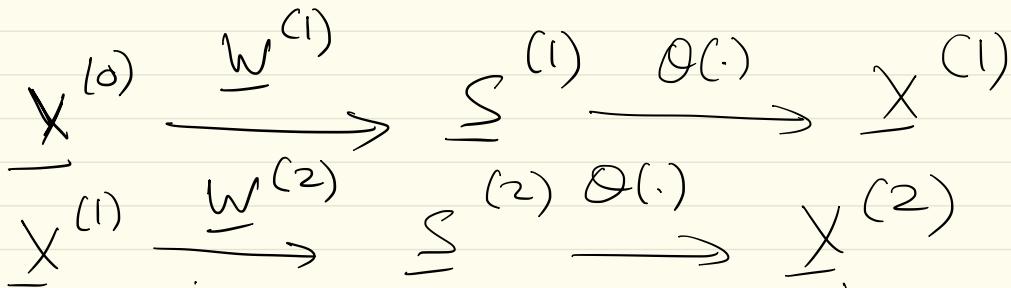


# Forward Propagation Alg.

$j^{\text{th}}$  column of  $\underline{w}^{(l)}$  is the collection of weights incident on node  $j$  in layer  $l$ .

$$\underline{s}^{(l)} = [\underline{w}^{(l)}]^T \underline{x}^{(l-1)}$$

Input Vector:  $\underline{x}^{(0)} = [1, x_1, x_2 \dots x_d]$



## Computational Complexity

$$\underline{X}^{(l-1)} \xrightarrow{\underline{W}^{(l)}} \underline{S}^{(l)} \xrightarrow{\underline{Q}^{(l)}} \underline{X}^{(l)}$$

= total number of edges  
between layer  $(l-1)$  & layer  $l$

$$= (d^{(l-1)} + 1) \cdot d^{(l)}$$

## Total Computational Complexity

$$= \Theta \left( \sum_{l=1}^L (d^{(l-1)} + 1) d^{(l)} \right)$$

= # of edges in the  
whole network.

Typically a few millions

# Introduced: Neural Networks.

## Model Parameters

$$\mathcal{D} = \left\{ \omega_{i,j}^{(l)} \right\}_{1 \leq l \leq L}$$

Input:  $\underline{x}^{(0)} = [1, x_1, x_2, \dots, x_d]$

Output:  $\underline{x}^{(L)}$  (computed using forward pass)

## Loss Function

Regression  $\rightarrow$  Squared loss

$$\underline{x}^{(L)} \rightarrow x^{(L)} \in \mathbb{R}$$

loss function  $g(x^{(L)}, y) = (x^{(L)} - y)^2$

Classification: log loss

$$\underline{x}^{(L)} = (\hat{P}_1, \hat{P}_2, \dots, \hat{P}_c) \quad g(\underline{x}^{(L)}, y=i) = -\log \hat{P}_i$$

let  $(\underline{x}_n, y_n)$  be a point in the training set

$$e_n(\mathcal{S}) = g(\underline{x}_n^{(L)}, y_n)$$

$$e.g = (\underline{x}_n^{(L)} - y_n)^2$$

(for squared loss)

Select  $\mathcal{S}$  to minimize the training error =  $\frac{1}{N} \sum_{n=1}^N e_n(\mathcal{S})$

---

## Stochastic Gradient Descent

In iteration # $k$ , let  $\mathcal{S}_k$  be the current choice of model params.

- Select a training point, say  $(\underline{x}_n, y_n)$  at random.
- use forward pass to compute

$$e_n(\mathcal{S}_k) = g(\underline{x}_n^{(L)}, y_n)$$

Compute

$$\frac{\partial \epsilon_n(\omega)}{\partial \omega_{i,j}^{(l)}} = \frac{e_n(\omega') - e_n(\omega)}{\Delta}$$

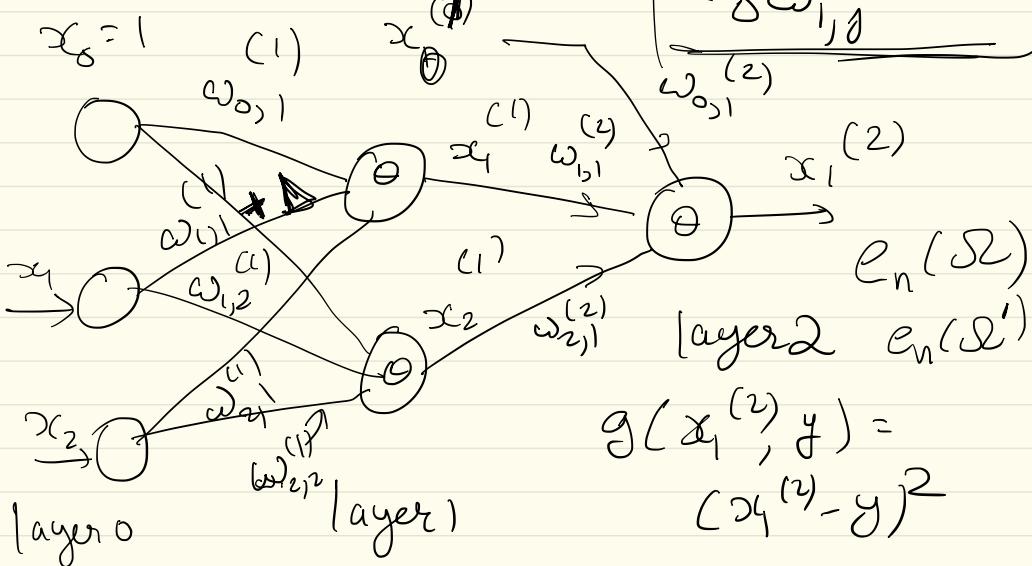
$$\frac{\partial \epsilon_n(\omega_k)}{\partial \omega_{i,j}^{(l)}}$$

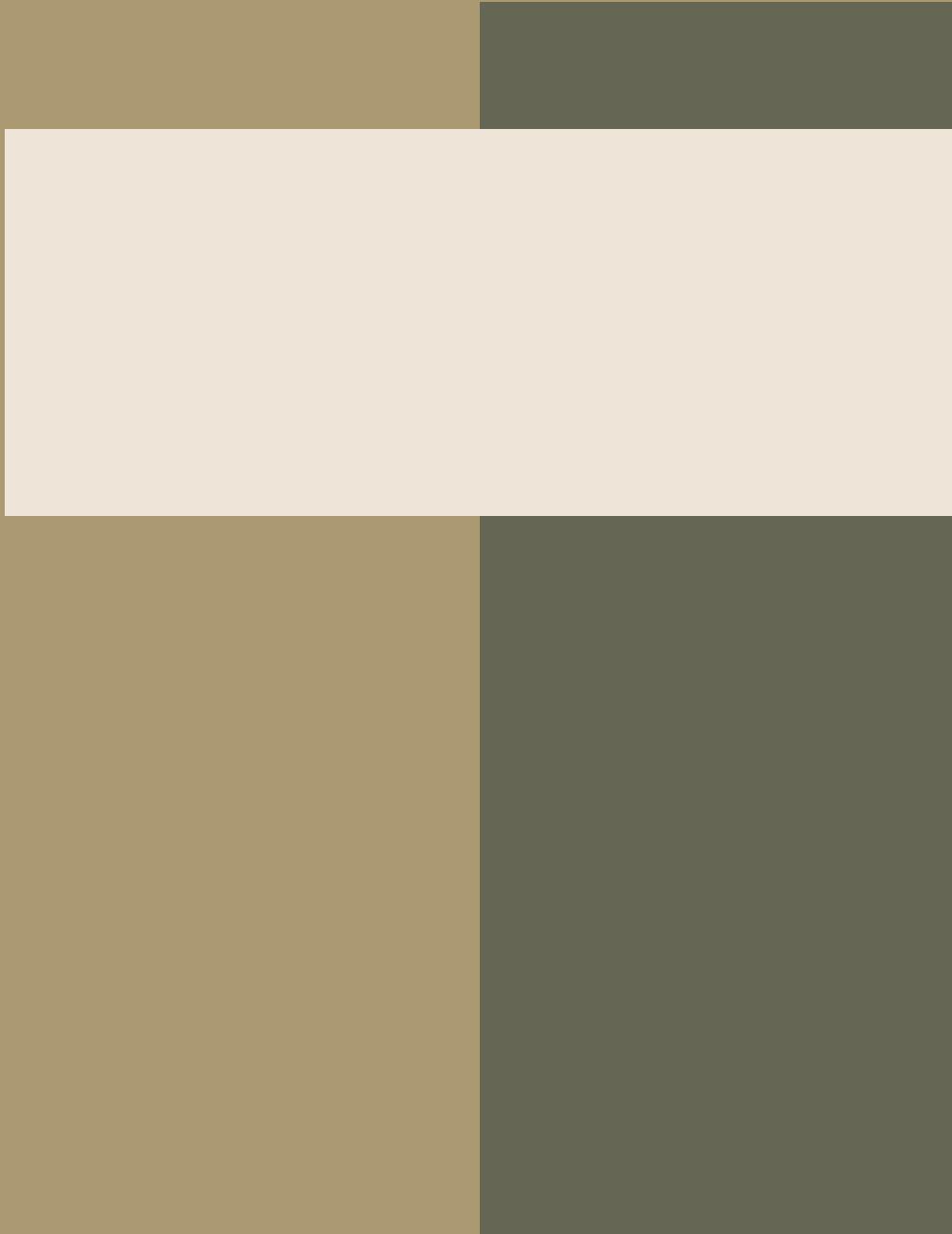
$\omega_{i,j}^{(l)}$

$$\omega_{i,j}^{(l)} \leftarrow \omega_{i,j}^{(l)} - \epsilon \frac{\partial \epsilon_n(\omega_k)}{\partial \omega_{i,j}^{(l)}}$$

- continue until convergence.

How to Compute





# Training for Neural Networks

## Last Lecture: Stochastic Gradient Descent

In each iteration

- select one training example  $(\underline{x}, y)$  at random
- forward Pass : output  $\underline{x}^L$
- loss:  $e(S2) = g(\underline{x}^L, y)$
- weight update

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \epsilon \frac{\partial e(S2)}{\partial w_{ij}^{(l)}}$$

(1986 Rumelhart, Hinton - Williams)

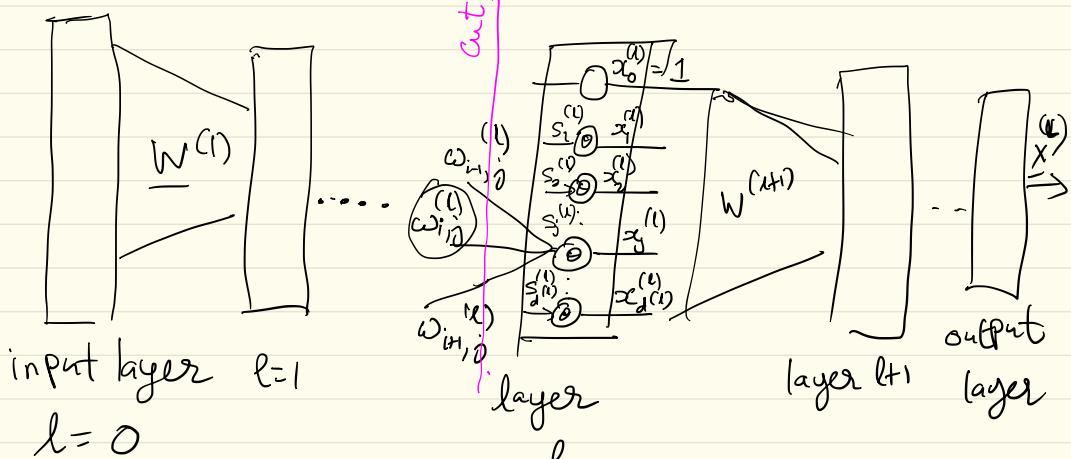
f Continue until convergence.

## Back Propagation Algorithm

- exact computation of all derivatives
- Complexity is linear<sup>in</sup> size of neural network

$$\frac{\partial e}{\partial w_{ij}^{(l)}} = \frac{\partial e}{\partial s_j^{(l)}} \cdot \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}}$$

L-layer neural Network



$$e(S) = g(x^{(l)}, y) \\ = (x^{(l)} - y)^2$$

Compute  $\frac{\partial e(S)}{\partial w_{ij}^{(l)}}$   $\forall i, j, l$

Given  $S^{(l)} = (s_1^{(l)}, s_2^{(l)}, \dots, s_d^{(l)})$   
&  $w^{(l+1)}, w^{(l+2)}, \dots, w^{(L)}$  the  
output vector can be computed.

$e(S) = e(s_1^{(1)}, s_2^{(1)}, \dots, s_j^{(N)}, \dots, s_d^{(N)})$ ,  
 $w^{(l+1)}, w^{(l+2)}, \dots, w^{(L)}$   
depends on  $w_{ij}^{(l)}$

"Composition trick"

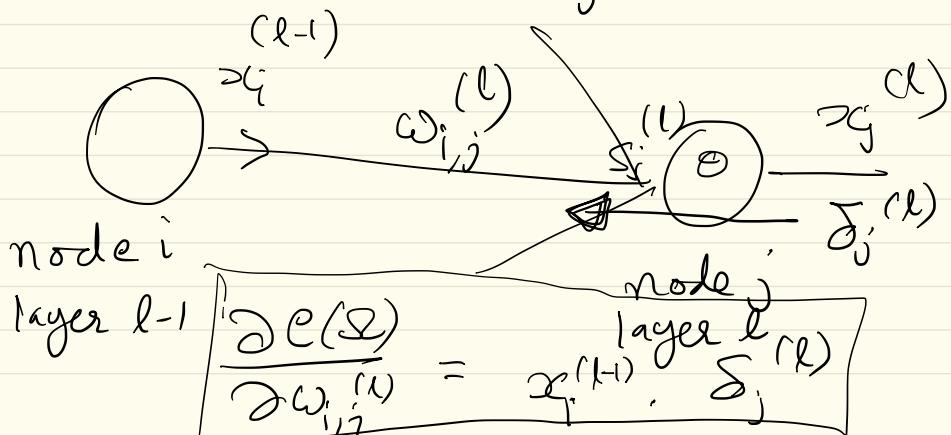
$$\frac{\partial e(\varphi)}{\partial \omega_{i,j}^{(l)}} = \frac{\partial e(\varphi)}{\partial s_j^{(l)}} \cdot \left[ \frac{\partial s_j^{(l)}}{\partial \omega_{i,j}^{(l)}} \right]$$

$$s_j^{(l)} = \sum_{k=0}^d x_k^{(l-1)} \cdot \omega_{k,j}^{(l)}$$

$$\frac{\partial s_j^{(l)}}{\partial \omega_{i,j}^{(l)}} = \sum_{k=0}^d \frac{\partial}{\partial \omega_{i,j}^{(l)}} (x_k^{(l-1)} \cdot \omega_{k,j}^{(l)})$$

$$= x_i^{(l-1)}$$

define  $\delta_j^{(l)} = \frac{\partial e(\varphi)}{\partial s_j^{(l)}}$  Backward Message



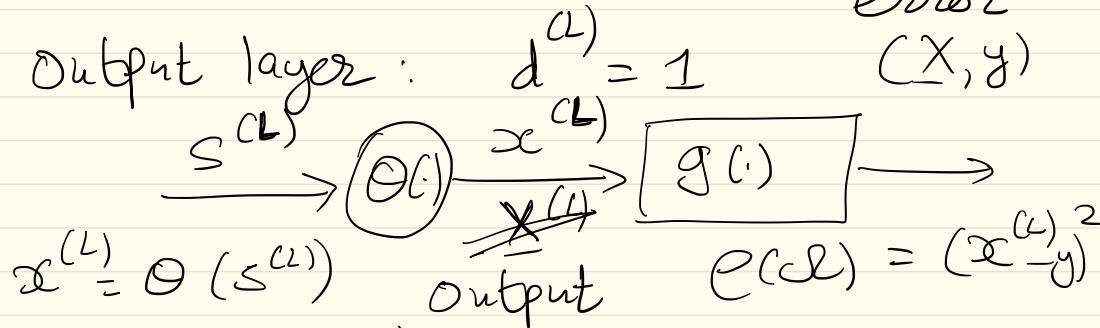
use backward pass  $\rightarrow$

Compute  $(\delta_j^{(L)})$   $j=1, 2, \dots, d^{(L)}$

$$\underline{\delta}^{(L)} \rightarrow \underline{\delta}^{(L-1)} \rightarrow \underline{\delta}^{(L-2)} - \underline{\delta}^{(L)} \rightarrow \underline{\delta}^{(1)}$$

special case: Regression, squared error

Output layer:  $d^{(L)} = 1$   $(x, y)$



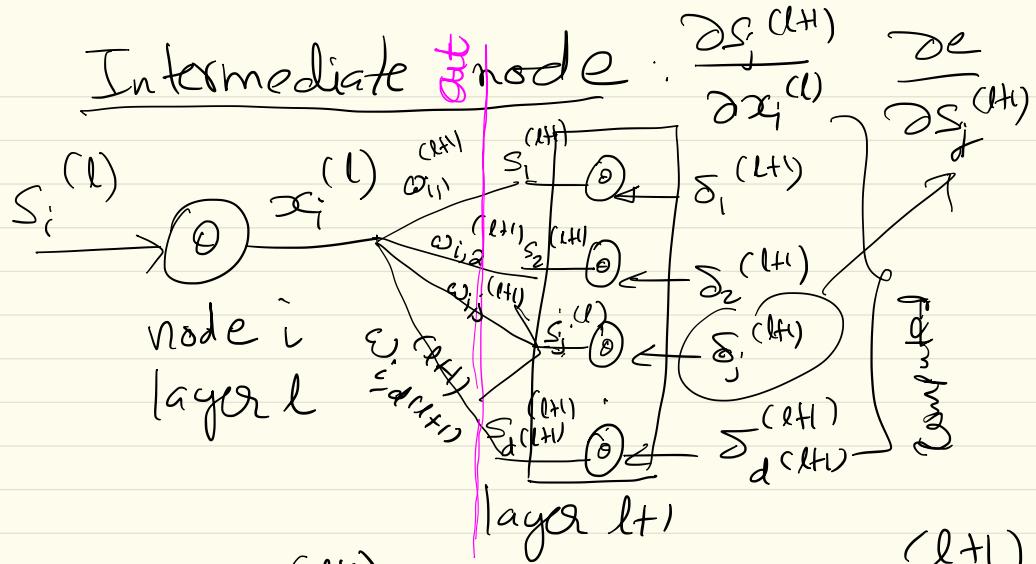
Compute:  $\underline{\delta}^{(L)} = \frac{\partial e(x)}{\partial s^{(L)}}$

$$= \frac{\partial e}{\partial x^{(L)}} \cdot \frac{\partial x^{(L)}}{\partial s^{(L)}} = 2(x^{(L)} - y) \underbrace{\Theta'(s^{(L)})}_{}$$

e.g.  $\theta(s) = \tanh(s)$

$$\text{check } \theta'(s) = \frac{\partial \theta}{\partial s} = 1 - \theta^2(s)$$

$\delta^{(L)}$  at output layer  $\rightarrow$  easy computation.



Given  $S_1^{(l+1)}, S_2^{(l+1)}, \dots, S_d^{(l+1)}$   
& wish to compute  $S_i^{(l)} = \underbrace{\frac{\partial e(S)}{\partial S_i^{(l)}}}_{\partial S_i^{(l)}}$   
for  $i=1, 2, \dots, d(l)$

$$\frac{\partial e(S)}{\partial S_i^{(l)}} = \frac{\partial e(d)}{\partial x_i^{(l)}} \cdot \underbrace{\frac{\partial x_i^{(l)}}{\partial S_i^{(l)}}}_{\theta'(S_i^{(l)})}$$

$$x_i^{(l)} = \theta(S_i^{(l)}) \quad \theta'(S_i^{(l)})$$

$$\frac{\partial x_i^{(l)}}{\partial S_i^{(l)}} = \theta'(S_i^{(l)})$$

$$e(S) = e(S_1^{(l+1)}, S_2^{(l+1)}, \dots, S_d^{(l+1)}, W^{(l+2)}, W^{(l)})$$

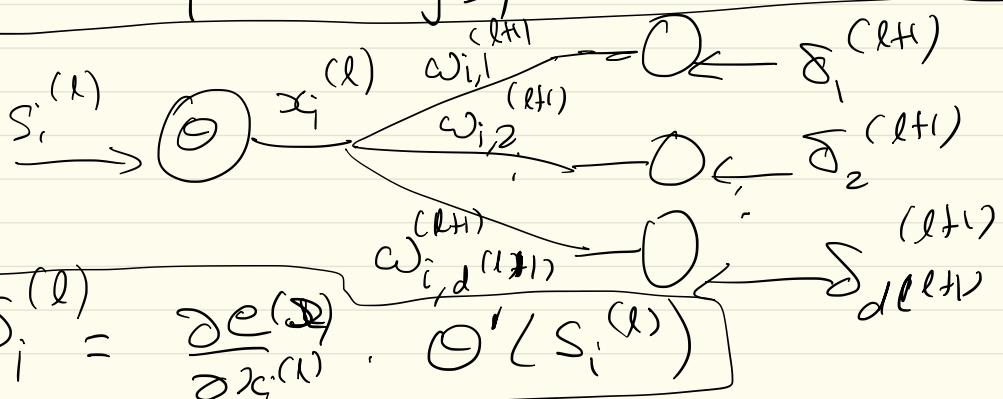
depend on  $\partial x_i^{(l)}$

$$\frac{\partial e(\alpha)}{\partial x_i^{(l)}} = \sum_{j=1}^d \underbrace{\frac{\partial e(\alpha)}{\partial s_j^{(l+1)}}}_{\delta_j^{(l+1)}} \cdot \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}}$$

$$s_j^{(l+1)} = \sum_{k=0}^d x_k^{(l)} \cdot \omega_{k,j}^{(l+1)}$$

$$\frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} = \omega_{i,j}^{(l+1)}$$

$$\boxed{\frac{\partial e(\alpha)}{\partial x_i^{(l)}} = \sum_{j=1}^d \delta_j^{(l+1)} \cdot \omega_{i,j}^{(l+1)}}$$



# Backpropagation Algorithm

Input data point:  $(\underline{x}, y)$

$$\mathcal{S} = \{\underline{w}^{(1)}, \underline{w}^{(2)}, \dots, \underline{w}^{(L)}\}$$

Output:  $\frac{\partial e(\mathcal{S})}{\partial w_i^{(l)}}$

$$\underline{\delta}^{(L)} = (\delta_1^{(L)}, \delta_2^{(L)}, \dots, \delta_d^{(L)})^T$$

for  $l = L-1, L-2, \dots, 2, 1$ .

1. Run Forward Pass to compute

$$\{\underline{s}^{(l)}, \underline{x}^{(l)}\} \quad l=1, 2, \dots, L$$

&  $e(\mathcal{S}) = g(\underline{x}^{(L)}, y)$  element wise multiplication

2.

$$\underline{\delta}^{(L)} = \left[ \nabla_{x^{(L)}} g(\underline{x}^{(L)}, y) \right] \otimes \Theta'(\underline{s}^{(L)})$$

$$\underline{s}^{(L)} = \begin{bmatrix} s_1^{(L)} \\ s_2^{(L)} \\ \vdots \\ s_d^{(L)} \end{bmatrix} \quad \underline{x}^{(L)} = \begin{bmatrix} 1 \\ \Theta(\underline{s}^{(L)}) \end{bmatrix} \quad \Theta'(\underline{s}^{(L)}) = \begin{bmatrix} \Theta'(s_1^{(L)}) \\ \Theta'(s_2^{(L)}) \\ \vdots \\ \Theta'(s_d^{(L)}) \end{bmatrix}$$

3. for  $\ell = L-1, L-2, \dots, 1$ :

$$\underline{\delta}^{(\ell)} = \left[ \underline{w}^{(\ell+1)} \cdot \underline{\delta}^{(\ell+1)} \right] \otimes \Theta'(\underline{s}^{(\ell)})$$

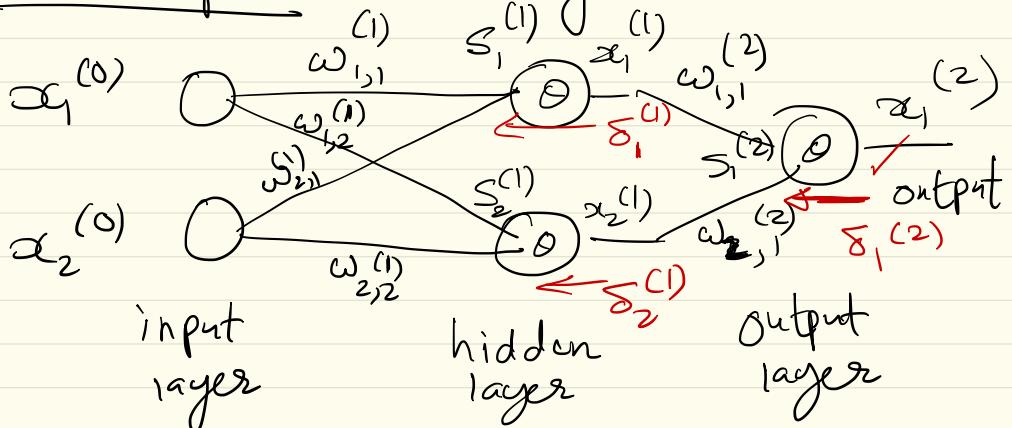
$$\underline{w}^{(\ell+1)} = \begin{bmatrix} \omega_{1,1}^{(\ell+1)} & \omega_{1,2}^{(\ell+1)} & \cdots & \omega_{1,d^{(\ell+1)}}^{(\ell+1)} \\ \omega_{2,1}^{(\ell+1)} & \omega_{2,2}^{(\ell+1)} & \cdots & \omega_{2,d^{(\ell+1)}}^{(\ell+1)} \\ \vdots & \vdots & & \vdots \\ \omega_{d^{(\ell)},1}^{(\ell+1)} & \cdots & \cdots & \omega_{d^{(\ell)},d^{(\ell+1)}}^{(\ell+1)} \end{bmatrix}$$

end.

$$\underline{\delta}^{(L)} \rightarrow \underline{\delta}^{(L-1)} \rightarrow \underline{\delta}^{(L-2)} \rightarrow \dots \rightarrow \underline{\delta}^{(1)}$$

$$\frac{\partial e(\underline{\delta})}{\partial \omega_{i,j}^{(\ell)}} = x_i^{(\ell-1)} \cdot \underline{\delta}_j^{(\ell)}$$

# Example : 2-layer network



$\ell = 0$

$\ell = 1$

$\ell = 2$

Compute

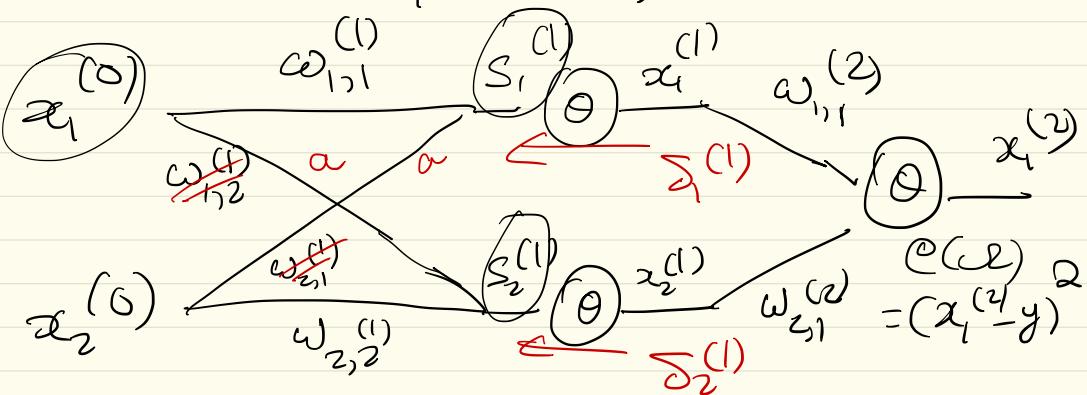
$$\frac{\partial e(\omega)}{\partial \omega_{1,1}^{(2)}} =$$

$$x_1^{(1)} \cdot \delta_1^{(2)} = \frac{\partial e}{\partial x_1^{(2)}} \cdot \frac{\partial x_1^{(2)}}{\partial \omega_{1,1}^{(1)}} \quad ?$$

$$\delta_1^{(2)} = \frac{\partial e(\omega)}{\partial \delta_1^{(2)}} = \frac{\partial e}{\partial x_1^{(2)}} \cdot \frac{\partial x_1^{(2)}}{\partial \delta_1^{(2)}} \\ = z(x_1^{(2)} - y) \cdot \theta'(S_1^{(2)})$$

$$\delta_1^{(1)} = \frac{\partial e(\omega)}{\partial x_1^{(1)}} \cdot \frac{\partial x_1^{(1)}}{\partial \delta_1^{(1)}} = \delta_1^{(2)} \cdot \omega_{1,1}^{(2)} \cdot \theta'(S_1^{(1)})$$

$$\delta_2^{(1)} = \delta_1^{(2)} \cdot \omega_{2,1}^{(2)} \cdot \theta^1(\zeta_2^{(1)})$$



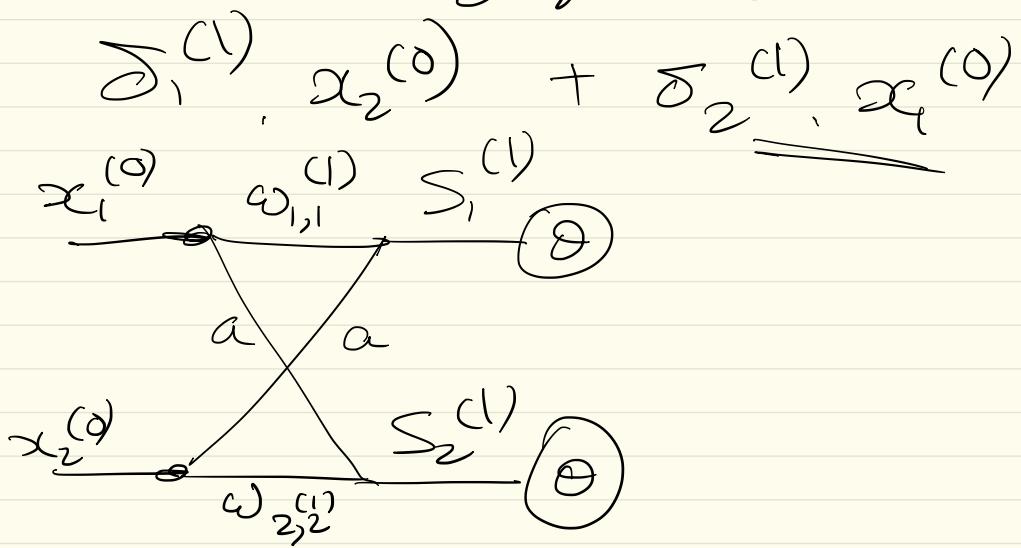
$$\begin{aligned} \frac{\partial e(\zeta)}{\partial x_1^{(0)}} &= \frac{\partial e(\zeta)}{\partial \delta_1^{(1)}} \cdot \frac{\partial \delta_1^{(1)}}{\partial x_1^{(0)}} \\ &\quad + \frac{\partial e(\zeta)}{\partial \delta_2^{(1)}} \cdot \frac{\partial \delta_2^{(1)}}{\partial x_1^{(0)}} \\ &= \delta_1^{(1)} \cdot \omega_{1,1}^{(1)} + \delta_2^{(1)} \cdot \omega_{1,2}^{(1)} \end{aligned}$$

$$\begin{aligned} \frac{\partial e(\zeta)}{\partial x_2^{(0)}} &= \delta_1^{(1)} \cdot \omega_{2,1}^{(1)} \\ &\quad + \delta_2^{(1)} \cdot \omega_{2,2}^{(1)} \end{aligned}$$

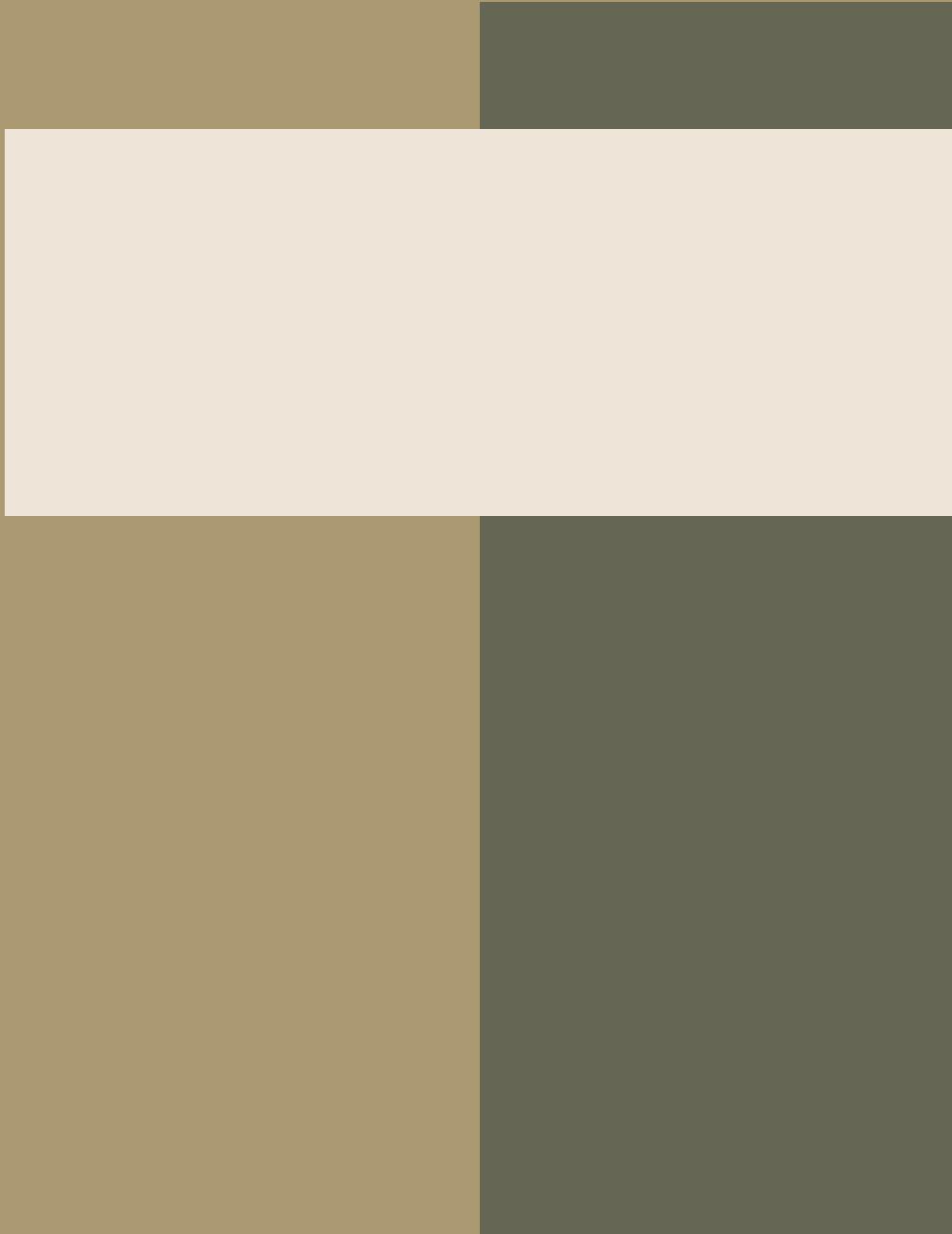
Compute  $\frac{\partial e(\zeta)}{\partial \alpha}$

$\boxed{\omega_{1,2}^{(1)} = \omega_{2,1}^{(1)} = \alpha}$

$$\frac{\partial e(r)}{\partial a} = \underline{\underline{\frac{\partial e}{\partial S_1^{(1)}}}} \cdot \underline{\underline{\frac{\partial S_1^{(1)}}{\partial a}}} + \underline{\underline{\frac{\partial e}{\partial S_2^{(1)}}}} \cdot \underline{\underline{\frac{\partial S_2^{(1)}}{\partial a}}}$$



$$S_1^{(1)} = a \cdot \underline{\underline{x_2^{(0)}}} + \underline{\underline{\omega_{1,1}^{(1)}}} \cdot \underline{\underline{x_1^{(0)}}}$$



# Implementation of Neural Networks

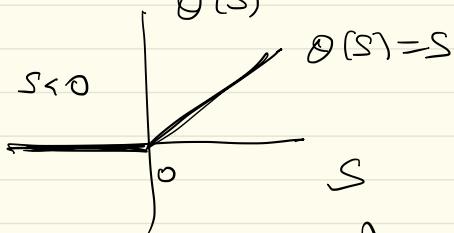
- Mid-term Avg.: 61.5/100

## i) choice of Activation functions

Sigmoid  $\Theta(s) = \frac{e^s}{1+e^s}$

tanh  $\Theta(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$

ReLU.  $\Theta(s) = \max\{0, s\}$



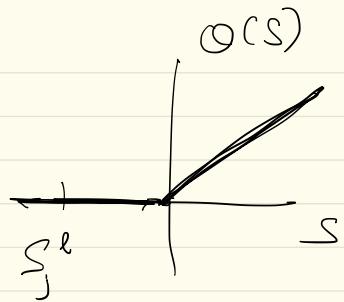
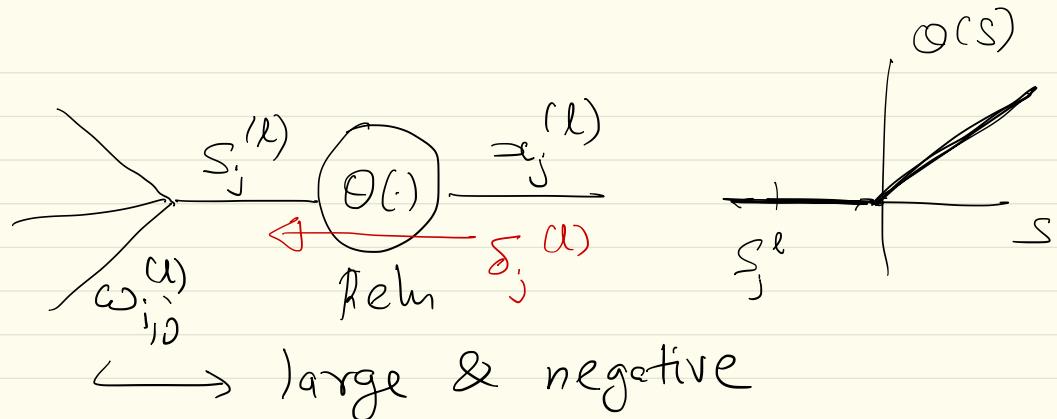
- modern NNs lead to 6x improvement in speed of training

- Simple in computation
- Avoids non-linear effects

Issues: (Dead neurons)

~ 40% of hidden units are

"dead - neurons"



$\hookrightarrow$  large & negative

$$s_j^{(l)} < 0 \Rightarrow a_j^{(l)} = 0$$

$$\text{also } \theta'(s_j^{(l)}) = 0$$

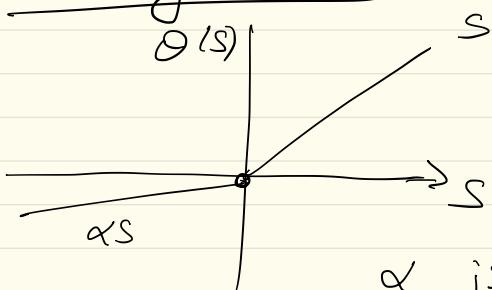
$$\bar{\omega}_j^{(l)} = 0$$

$$\frac{\partial e}{\partial \omega_{i,j}^{(l)}} = 0$$

$$\omega_{i,j}^{(l)} \leftarrow \omega_{i,j}^{(l)} - \epsilon \underbrace{\frac{\partial e}{\partial \omega_{i,j}^{(l)}}}_{=0}$$

$\sim 40\%$  are "dead neurons"

## Leaky Relu



$$O(S) = \begin{cases} S & S > 0 \\ \alpha S & S \leq 0 \end{cases}$$

$$\alpha \approx 0.1$$

$\alpha$  is set to a small constant to carry over -ve inputs forward.

Parametric Relu:  $\alpha$  is not fixed but updated during training

## 2) Input Pre-processing

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

$$x_i \in \mathbb{R}^d$$

$$x_i = (x_{i(1)}, x_{i(2)}, \dots, x_{i(d)})$$

Zero mean - unit variance normalization

$$\bar{x}_i = \frac{1}{N} \sum_{n=1}^N x_{i(n)} \quad (\text{Sample mean})$$

$$\sigma_{(j)} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_{i(j)} - \bar{x}_{i(j)})^2} \quad (\text{Sample standard deviation})$$

$$\underline{x}_i(j) \leftarrow \frac{\underline{x}_i(j) - \bar{y}(j)}{\sigma(j)}$$

for each  $j = 1, 2, \dots, d$  &  
each  $i = 1, 2, \dots, N$

Batch Norm: normalization  
at the output of each hidden  
layer.

### Data-Augmentation

Given  $(\underline{x}, y)$  in training set,  
also include  $(F(\underline{x}), y)$   
 $F = \text{scaling, cropping, flipping}$   
 $\text{etc.}$

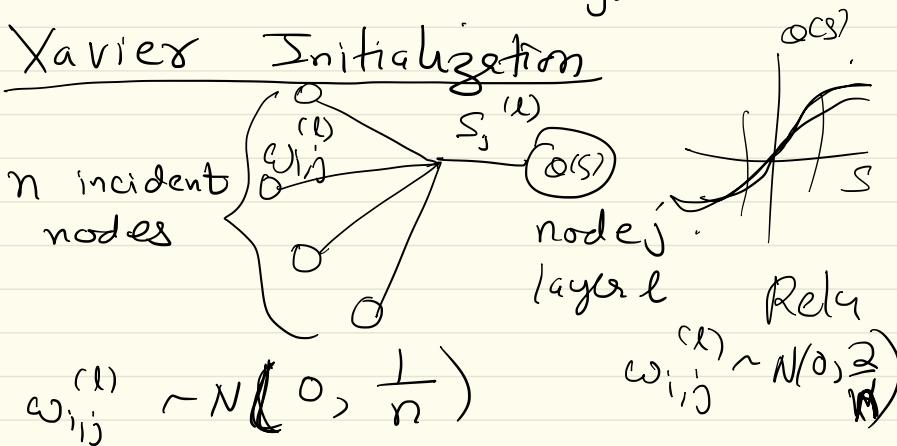
(Section 8.4  
in Deep-Learning  
Book)

## Weight Initialization

traditionally  $\omega_{i,j}^{(l)} \sim N(0, \sigma^2)$   
 $\sigma^2 \approx \frac{3}{10}$

- this works on "Shallow networks"
- for deep networks with many layers,  $\rightarrow$  Vanishing gradients

## Xavier Initialization

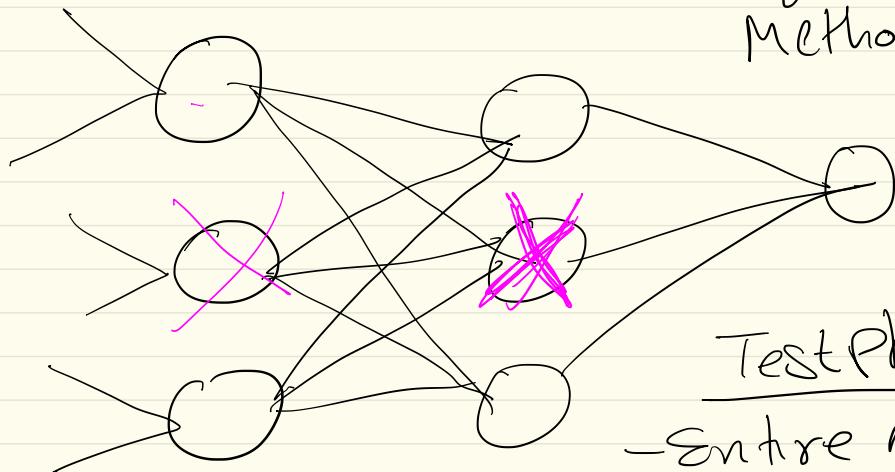


$$S_j^{(l)} = \sum_{i=1}^n \underbrace{x_i^{(l-1)}}_{\substack{\text{unit} \\ \text{Variance}}} \cdot \underbrace{\omega_{i,j}^{(l)}}_{\substack{\text{Variance} \approx \frac{1}{n} \\ \times \text{independent}}}$$

$$\text{Variance}(S_j^{(l)}) \approx 1$$

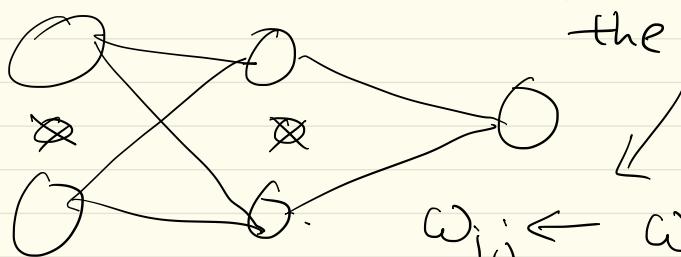
# Drop out (D.L. book Sec. 7.12)

- with probability  $P$ , set output of a given node to zero in the forward pass. "Regularization Method"

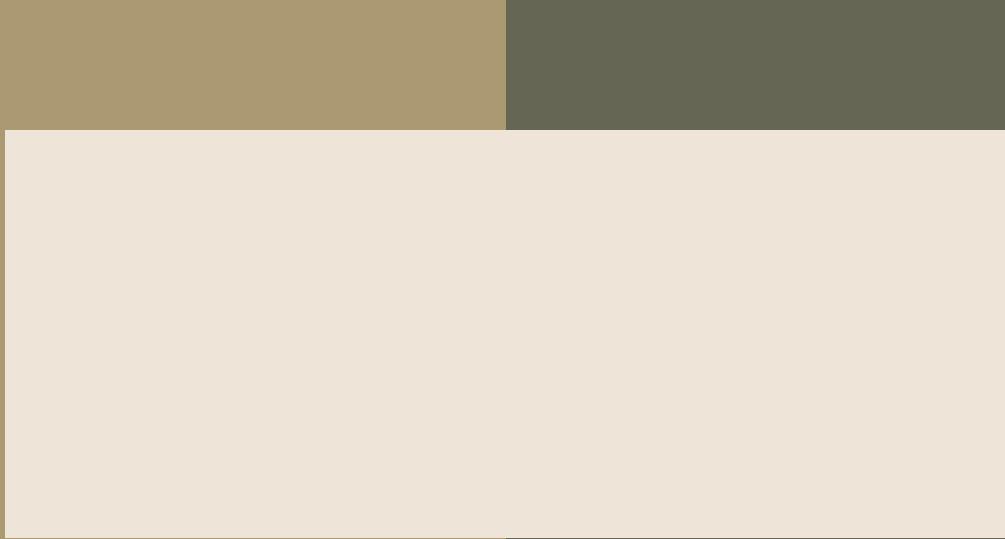


Test Phase:

- entire network  
but scale  
the weights



$$\omega_{i,j} \leftarrow \omega_{i,j} (1 - P)$$

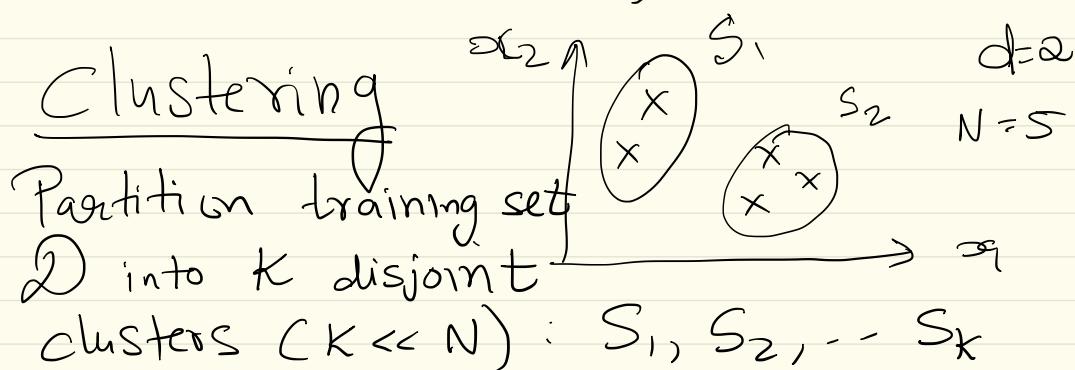


# Unsupervised Learning

Training Set

$$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\} \quad \underline{x}_i \in \mathbb{R}^d$$

- 1) Clustering.
  - 2) Density Estimation.
  - 3) Dimensionality Reduction
- PCA → ECE367 } X



such that elements within each cluster are close to each other

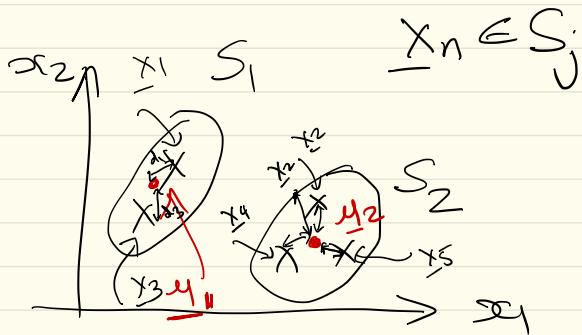
$$S_i \subseteq \mathcal{D}, \quad S_i \cap S_j = \emptyset \text{ (disjoint)} \\ S_1 \cup S_2 \cup \dots \cup S_k = \mathcal{D}$$

# Loss function for clustering.

For cluster  $S_i$  let  $\underline{y}_i$  be the cluster center

Approximation Error for  $S_j$

$$E_j = \sum \| \underline{x}_n - \underline{y}_j \|^2$$

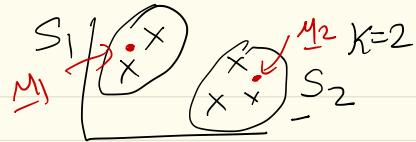


$$E_{in} (S_1, S_2, \dots, S_k, \underline{y}_1, \underline{y}_2, \dots, \underline{y}_k)$$

$$= \sum_{j=1}^K E_j = \sum_{j=1}^K \sum_{\underline{x}_n \in S_j} \| \underline{x}_n - \underline{y}_j \|^2$$

$$= \sum_{n=1}^N \| \underline{x}_n - \underbrace{\underline{y}(\underline{x}_n)}_{\text{cluster center associated with } \underline{x}_n} \|^2$$

## Optimal Clustering



Given  $\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n\}$

minimize  $E_{in}(S_1, \dots, S_k, \underline{M}_1, \dots, \underline{M}_k)$

$$\begin{aligned} S_1 - S_K & \quad \underbrace{\qquad \qquad \qquad}_{n} \\ \underline{M}_1, \dots, \underline{M}_K & = \sum_{n=1}^n \|\underline{x}_n - \underline{M}(\underline{x}_n)\|^2 \end{aligned}$$

N-P Hard

Good News:

Efficient heuristics  
that are widely  
used in practice.

K-means Algorithm

Applications

1) Topic Discovery : N documents  
Partition into  $k$  clusters

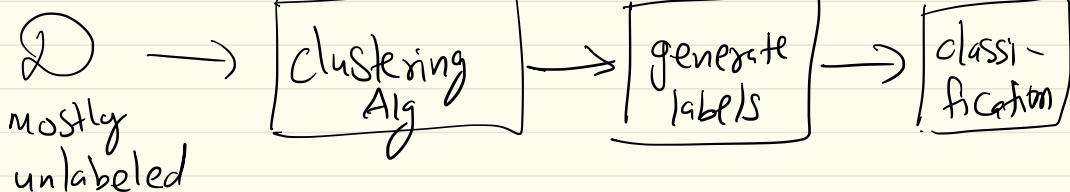
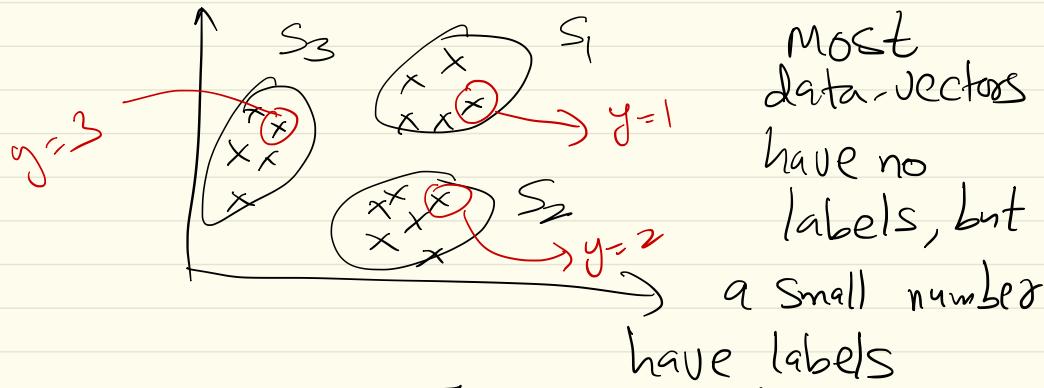
$\underline{x}_i$  = histogram of words in document  $i$

$\underline{M}_j$  = mean histogram in cluster  $j$

## 2) Recommendation Systems:

- cluster consumers based on some common preferences to provide new recommendations.

## 3) Semi-Supervised Learning

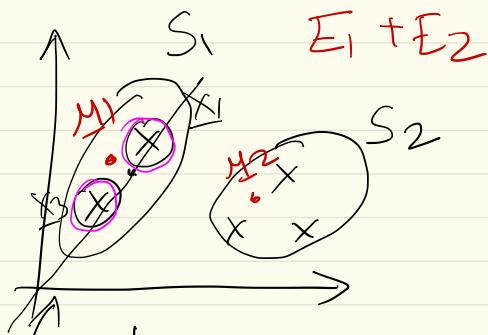


K-Means Algorithm

$$\text{minimize}_{\underline{\mu}_1, \dots, \underline{\mu}_K} E_{in}(S_1, \dots, S_K, \underline{\mu}_1, \dots, \underline{\mu}_K)$$

$$S_1 - S_K = \sum_{j=1}^K E_j$$

Problem 1: Given  $S_1 - S_K$   
find  $\underline{\mu}_1, \dots, \underline{\mu}_K$  to minimize  $E_{in}$



Select  $\underline{\mu}_j$  in  
order to minimize

$$E_j$$

$$\underline{\mu}_j^* = \arg \min_{\underline{\mu}_j \in \mathbb{R}^d} \sum_{x_n \in S_j} \|x_n - \underline{\mu}_j\|^2$$

Example:

$$E_1(\underline{\mu}_1) = \|\underline{x}_1 - \underline{\mu}_1\|^2 + \|\underline{x}_3 - \underline{\mu}_1\|^2$$

$$\nabla E_1(\underline{\mu}_1) = 0 \quad \underline{\mu}_1^* = \frac{1}{2} (\underline{x}_1 + \underline{x}_3)$$

$$(E_2)(\underline{\mu}_2) = \|\underline{x}_2 - \underline{\mu}_2\|^2 + \|\underline{x}_4 - \underline{\mu}_2\|^2 + \|\underline{x}_5 - \underline{\mu}_2\|^2$$

$$\nabla E_2(\underline{M}_2) = 0 \quad \underline{M}_2^* = \frac{1}{3} (\underline{x}_2 + \underline{x}_4 + \underline{x}_5)$$

General Solution

$$S_1 = \{ \underline{x}_1 \\ S_2 = \{$$

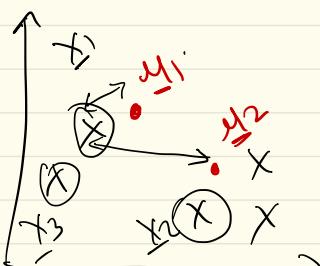
$$\underline{M}_j^* = \frac{1}{|S_j|} \sum_{\underline{x}_n \in S_j} \underline{x}_n$$

will minimize  $E_j(\underline{M}_j)$

Problem 2: Given  $\underline{M}_1, \underline{M}_2, \dots, \underline{M}_K$

→ find  $S_1, S_2, \dots, S_K$  to

minimize  $E_{in}(S_1, S_2, \dots, S_K, \underline{M}_1, \underline{M}_2, \dots, \underline{M}_K)$



$$= \underbrace{\sum_{n=1}^N \| \underline{x}_n - \underline{M}(\underline{x}_n) \|^2}_{\text{For each point } \underline{x}_n}$$

we simply select the cluster center closest to  $\underline{x}_n$  & assign  $\underline{x}_n$  to that cluster

# K-Means Alg (Lloyd 1956)

1. Initialize  $\underline{m}_1, \underline{m}_2, \dots, \underline{m}_k$  in

some fashion.

e.g. (K-means ++ Initialization)

2. Construct  $S_1, \dots, S_k$  by  
solving sub-problem #2

3. update  $\underline{m}_1, \dots, \underline{m}_k$  by  
solving sub-problem #1

4. Repeat steps 2 & 3  
until convergence.

$(E_{in}(S_1 \dots S_k, \underline{m}_1 \dots \underline{m}_k))$

is sufficiently small  $\rightarrow$  stop

Alg. is known to converge to a  
local minima.

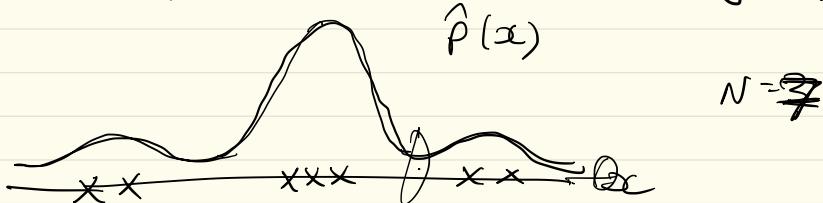
# Density Estimation

$$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\} \quad \underline{x}_i \stackrel{iid}{\sim} p(x)$$

$p(\underline{x})$  is unknown to us.

Output  $\hat{p}(\underline{x})$  an estimate of  $p(\underline{x})$

$$d=1$$



Applications : Anomaly detection.

Given Training set  $\mathcal{D}$

Compute  $\hat{p}(\underline{x})$

given query point  $q$  if  $\hat{p}(q) < \epsilon$

declare  $q$  to be an *anomaly*.

1) Histogram Method

2) Nearest neighbour

3) Gaussian mixture model ( $EM$  algorithm)

# Histogram Method

assume that each  $\underline{x}$  lies in a

bounded region e.g.  $\underline{x} \in [0, 1]^d$

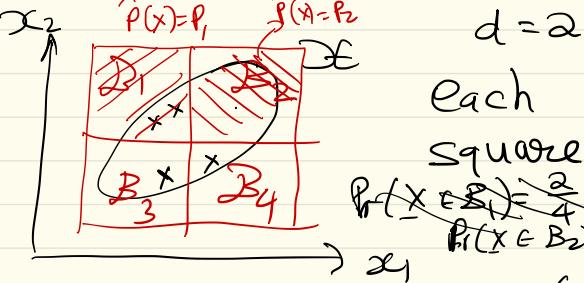
i.e.  $P(\underline{x}) = 0$  if  $\underline{x} \notin \mathcal{X} \subseteq \mathbb{R}^d$

Cover  $\mathcal{X}$  using uniform

hypercubes  $B_1, B_2, \dots, B_k$  of

constant volume  $V$ .

$N = 4$



each  $B_i$  is a square with area of  $V$ .

$$P_i(\underline{x} \in B_i) = \frac{2}{4}$$

$$P_i(\underline{x} \in B_2) = 0$$

Step 1: Compute  $P_i(\underline{x} \in B_i) = \frac{N_i}{N}$

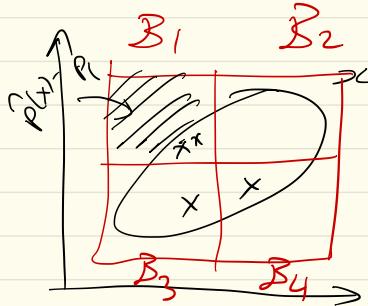
$$i = 1, 2, 3, 4$$

$N_i \equiv$  total # of Points from  $\mathcal{D}$  in  $B_i$

$$N = N_1 + N_2 + \dots + N_k = \text{total # of training pts}$$

$$P_i(\underline{x} \in B_1) = \frac{2}{4} \quad P_i(\underline{x} \in B_2) = 0$$

$$P_i(\underline{x} \in B_3) = P_i(\underline{x} \in B_4) = \frac{1}{4}$$



$$\Pr(x \in B_1) = \frac{2}{4}$$

$$\Pr(x \in B_2) = 0$$

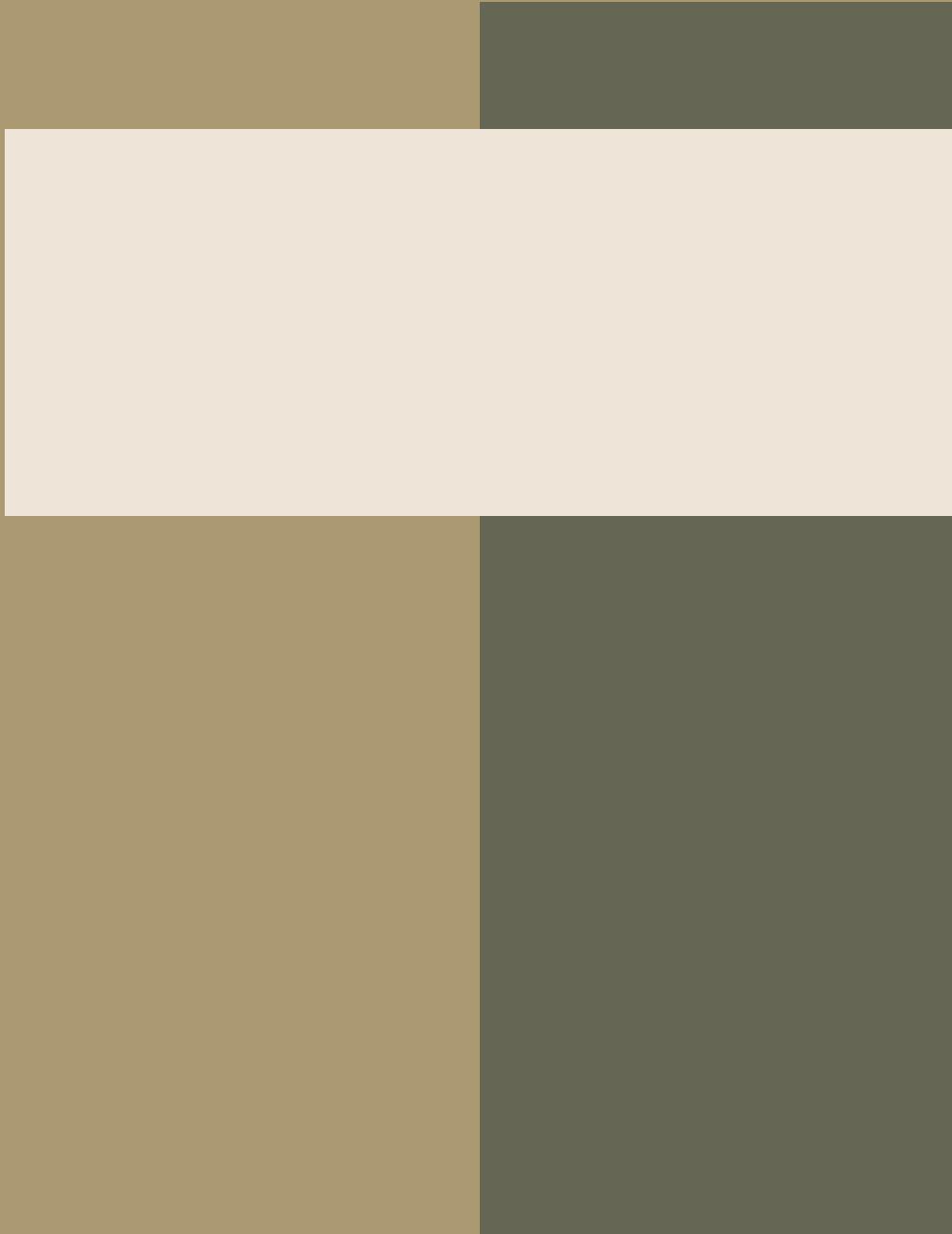
$$\Pr(x \in B_3) = \Pr(x \in B_4) = \frac{1}{4}$$

Assume density function  $\hat{P}(x)$  is constant on each  $B_i$

$$1. e \quad \hat{P}(x) = \begin{cases} P_1 & \forall x \in B_1 \\ P_2 & \forall x \in B_2 \\ P_3 & \forall x \in B_3 \\ P_4 & \forall x \in B_4 \end{cases}$$

$$\text{find } P_1, P_2, P_3, P_4 \rightarrow P_1 = \frac{2}{4V}$$

$$\int_{x \in B_1} \hat{P}(x) dx = \Pr(x \in B_1) = \frac{2}{4} \Rightarrow \text{area of } B_1$$



# Density Estimation

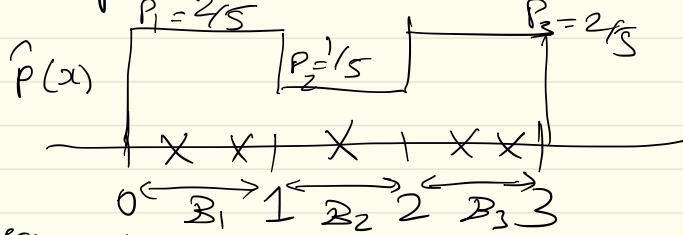
$$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n\} \quad p(x) > q$$

$$\underline{x}_i \in \mathbb{R}^d \quad \underline{x}_i \stackrel{iid}{\sim} p(\underline{x})$$

Output an estimate, say  $\hat{p}(\underline{x})$   
of  $p(\underline{x})$

## Histogram Method

$$\text{Example: } d=1 \quad N=5$$



Assume:

$\mathcal{X} \not\subset \mathcal{X} = [0, 3]$  then  $p(x)$

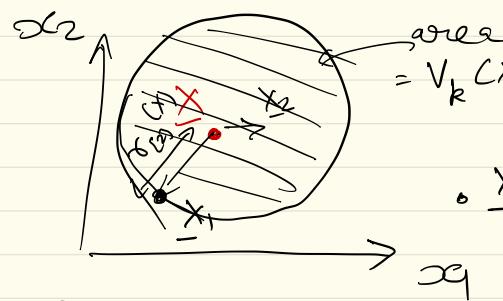
$$\mathcal{B}_1 = [0, 1] \quad \mathcal{B}_2 = [1, 2] \quad \mathcal{B}_3 = [2, 3]$$

$$P_r(x \in \mathcal{B}_1) = \frac{2}{5} \quad P_r(x \in \mathcal{B}_2) = \frac{1}{5} \quad P_r(x \in \mathcal{B}_3) = \frac{2}{5}$$

$$\hat{p}(x) = \frac{2}{5} \quad x \in \mathcal{B}_1 \quad \hat{p}(x) = \frac{1}{5} \quad x \in \mathcal{B}_2 \quad \hat{p}(x) = \frac{2}{5} \quad x \in \mathcal{B}_3$$

## 2) Nearest Neighbour Estimation

$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$  &  
estimate  $\hat{p}(\underline{x})$  for any given  $\underline{x} \in \mathbb{R}^d$



$$\text{area} = \pi r_k^2 \\ = V_k(x)$$

$$d=2$$

Given  $\underline{x}$   
output  $\hat{p}(\underline{x})$   
 $k=2$

$$d \geq 1$$

Given  $\underline{x}$  let

$\underline{x}_{[1]}, \underline{x}_{[2]}, \dots, \underline{x}_{[k]}$  be the  $k$  nearest neighbours in  $\mathcal{D}$  to  $\underline{x}$

$$r_k(\underline{x}) = \|\underline{x} - \underline{x}_{[k]}\| \quad \& \quad V_k(\underline{x})$$

be the volume of a sphere in  $\mathbb{R}^d$  with radius  $r_k(\underline{x})$

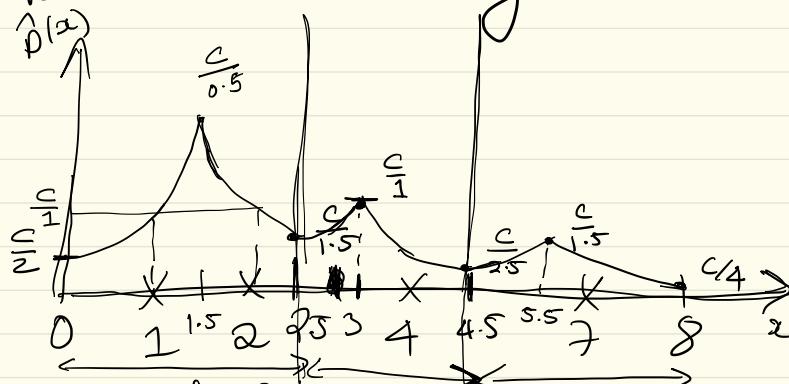
$$\hat{p}(\underline{x}) \propto \frac{1}{V_k(\underline{x})} \Rightarrow \hat{p}(\underline{x}) = \frac{C}{V_k(\underline{x})}$$

where  $C$  = normalizing constant

$$d=1$$

Example:  $\mathcal{D} = \{x_1=1, x_2=2, x_3=4,$

$k=2 \quad \text{NN-density estimation } x_4=7\}$

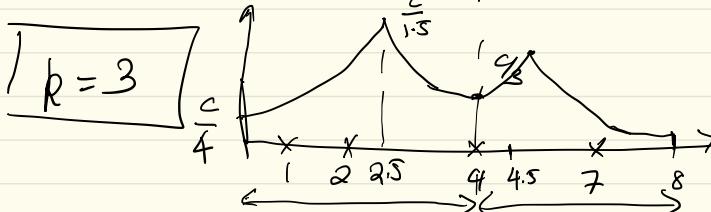


$$p(x) = 0 \quad \text{if } x \notin [0, 8]$$

$$x \in [0, 2.5] : NN = \{1, 2\}$$

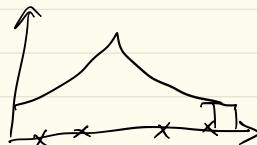
$$x \in [2.5, 4.5] : NN = \{2, 4\}$$

$$x \in [4.5, 8] : NN = \{4, 7\}$$



~~$k=1?$~~

$$\mathcal{P} = \{1, 2, 4, 7\}$$



$$K = E[(\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})^T]$$

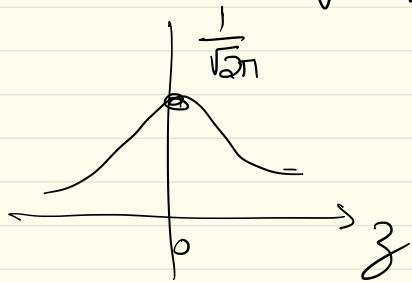
So far . . .

1) Histogram Method

2) K-Nearest Neighbour.

- Gaussian Mixture Models.

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{z^2}{2} \right\}$$



mean : 0

Variance : 1

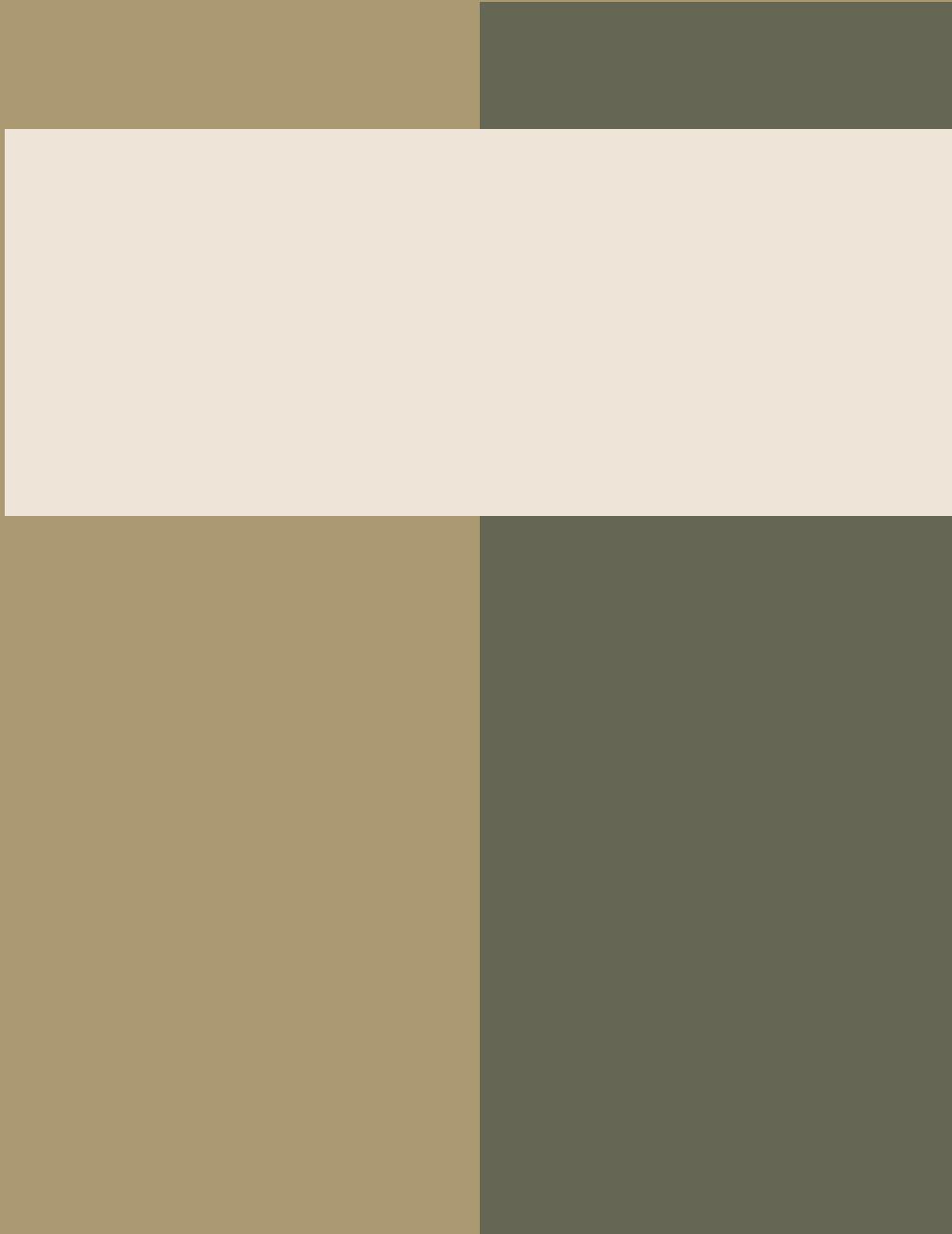
$$z \sim N(0, 1)$$

Two Jointly Gaussian Random Variables.

$\underline{x}_1, \underline{x}_2$  : Jointly Gaussian

Covariance Matrix





# Density Estimation

Sec 6.3.3 & Sec 6.4

Training Set

$$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$$

$$\underline{x}_i \in \mathbb{R}^d \quad \underline{x}_i \stackrel{iid}{\sim} P(\underline{x})$$

output  $\hat{P}(\underline{x})$  an estimate of  $P(\underline{x})$

## Previous Methods

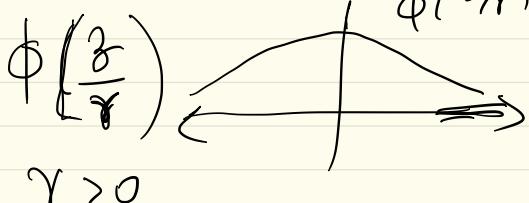
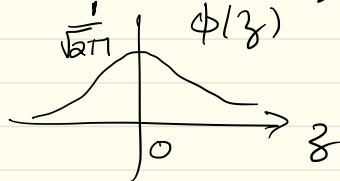
1) Histogram

2) K-Nearest Neighbour

Parzen Window Estimation using Gaussian Kernels.

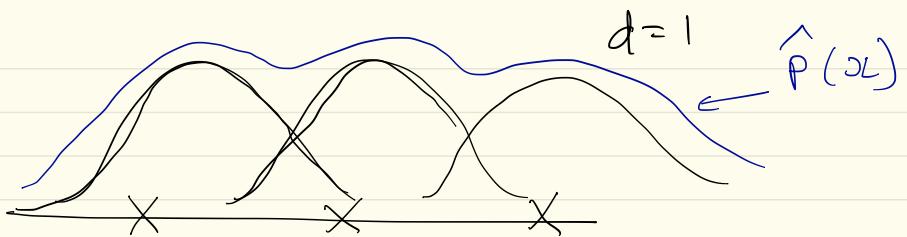
Gaussian Density function  $Z \sim N(0, 1)$

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$



"width =  $r$ "

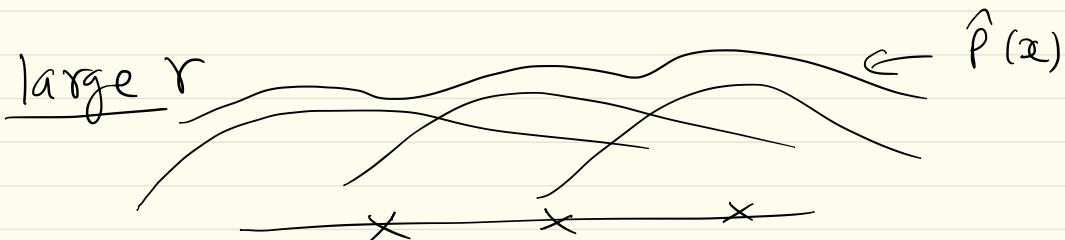
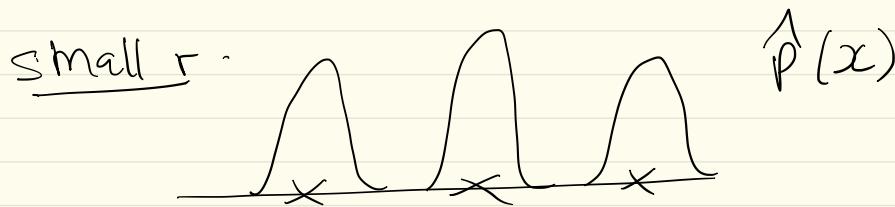
Gaussian Kernel



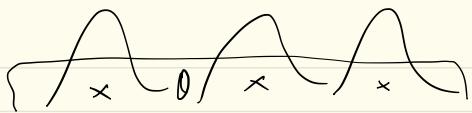
$$\hat{P}(x) = \frac{1}{K} \sum_{n=1}^N \phi\left(\frac{x - x_i}{\gamma}\right)$$

$K$  = normalization constant

$\gamma$  = kernel width.



choice of  $\gamma$  can be important  
in practice.  $\therefore$  Validation to  
select "r"



Validation set

$\mathcal{V} = \{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_M\}$   
not included in  $\mathcal{D}$ .

for different values of  $\gamma$ ,

compute  $\hat{p}_r(x) = \frac{1}{K} \sum_{n=1}^N \phi\left(\frac{x-x_i}{\gamma}\right)$

evaluate  $\hat{p}_r(x)$  on  $\mathcal{V}$  for  
each choice of  $\gamma$ .

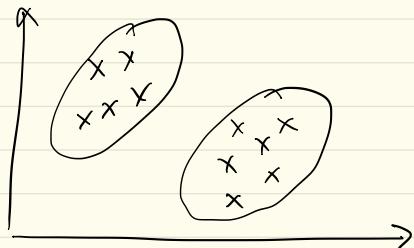
Now evaluate

$$\begin{aligned} -\log \hat{p}_r(\hat{x}_1, \hat{x}_2, \dots, \hat{x}_M) &= \log \prod_{m=1}^M \hat{p}_r(\hat{x}_m) \\ &= -\underbrace{\sum_{m=1}^M \log \hat{p}_r(\hat{x}_m)}_{\text{Loss function}} \end{aligned}$$

$$\text{Loss function} = -\sum_{m=1}^M \log \hat{p}_r(\hat{x}_m)$$

Select  $\gamma$  to minimize this loss function

# Gaussian Mixture Model



$K=2$  Gaussian density functions.

Multi variate Gaussian Density function

Univariate Gaussian:  $Z \sim N(\mu, \sigma^2)$

$$f_Z(z) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(z-\mu)^2}{2\sigma^2}\right\}$$

$\mu = E[Z]$  (mean)

$\sigma^2 = E[(Z-\mu)^2]$  (variance)

Bivariate Gaussian Density function

$$X_1 \sim N(\mu_1, \sigma_1^2)$$

$$E[X_1] = \mu_1$$

$$\text{Var}(X_1) = \sigma_1^2$$

$$X_2 \sim N(\mu_2, \sigma_2^2)$$

$$E[X_2] = \mu_2$$

$$\text{Var}(X_2) = \sigma_2^2$$

## Covariance:

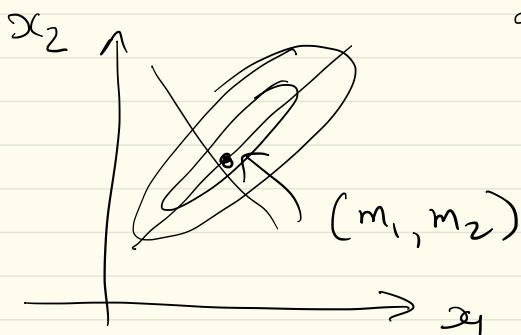
$$\sigma_{12} = E \{ (X_1 - m_1) (X_2 - m_2) \}$$

auto correlation coefficient

$$p = \frac{\sigma_{12}}{\sigma_1 \cdot \sigma_2}$$

$$-1 \leq p \leq 1$$

$$f_{X_1, X_2}(x_1, x_2) = \frac{1}{2\pi \sigma_1 \sigma_2 \sqrt{1-p^2}} \exp \left[ -\frac{\left( \frac{(x_1 - m_1)^2}{\sigma_1^2} + \frac{(x_2 - m_2)^2}{\sigma_2^2} - 2p \left( \frac{x_1 - m_1}{\sigma_1} \right) \left( \frac{x_2 - m_2}{\sigma_2} \right) \right)}{2(1-p^2)} \right]$$



$$f(x_1, x_2) = \text{constant}$$

$$\left( \frac{(x_1 - m_1)^2}{\sigma_1^2} + \frac{(x_2 - m_2)^2}{\sigma_2^2} - 2p \left( \frac{x_1 - m_1}{\sigma_1} \right) \left( \frac{x_2 - m_2}{\sigma_2} \right) \right) = \text{constant}$$

Multi-variate Gaussian Density

$$\underline{X} = [x_1, x_2, \dots, x_d]^T$$

$x_i$  is Gaussian

$$E[x_i] = m_i$$

$$\underline{\mu} = [m_1, m_2, \dots, m_d]^T$$

Covariance Matrix

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2d} \\ \vdots & & & \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{bmatrix}$$

$$\begin{aligned} \sigma_{ij} &= E\{(x_i - m_i)(x_j - m_j)\} \\ &= \sigma_{ji} \quad \underline{x} \sim N(\underline{\mu}, \Sigma) \end{aligned}$$

Density function

$$f(x) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}}$$

$$\exp \left\{ -\frac{(x - \mu)^T \Sigma^{-1} (x - \mu)}{2} \right\}$$

(Pg 240 - Leon-Garcia Text)

Example:

$$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$$

$\underline{x}_i \stackrel{iid}{\sim} p(x)$  it is known that  $p(x)$  is Gaussian.

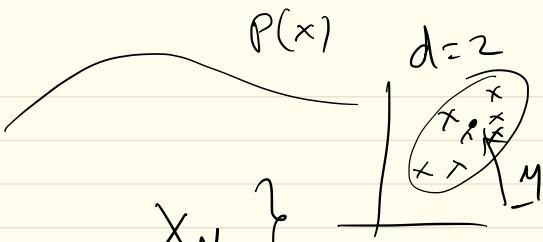
$p(x)$  is of the form  $N(\underline{\mu}, \Sigma)$

We only need to estimate mean  $\underline{\mu}$ , & covariance matrix  $\Sigma$  from  $\mathcal{D}$ .

$$\hat{\underline{\mu}} = \frac{1}{N} \sum_{n=1}^N \underline{x}_n \quad (\text{Empirical mean})$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\underline{x}_n - \underline{\mu})(\underline{x}_n - \underline{\mu})^\top$$

$$\Sigma = E[(\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})^\top]$$



Given  $P(x)$

# Gaussian Mixture Model

$$d=1$$

$$N(\mu_1, \sigma_1^2)$$

$$N(\mu_2, \sigma_2^2)$$

$$N(\mu_3, \sigma_3^2)$$



$$\hat{P}(x) = \frac{1}{4} \cdot \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left\{ -\frac{(x-\mu_1)^2}{2\sigma_1^2} \right\}$$

$$+ \frac{1}{2} \cdot \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ -\frac{(x-\mu_2)^2}{2\sigma_2^2} \right\}$$

$$+ \frac{1}{4} \cdot \frac{1}{\sqrt{2\pi\sigma_3^2}} \exp \left\{ -\frac{(x-\mu_3)^2}{2\sigma_3^2} \right\}$$

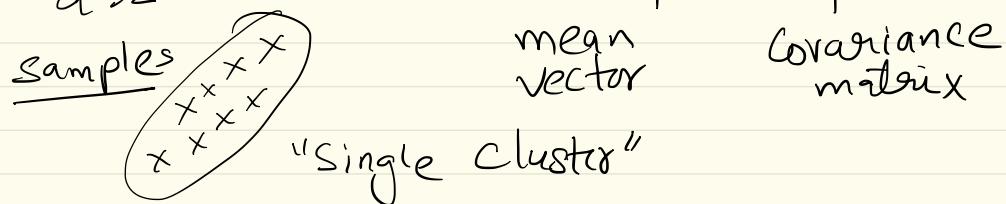
$$= \frac{1}{4} N(x; \mu_1, \sigma_1^2) + \frac{1}{2} N(x; \mu_2, \sigma_2^2) \\ + \frac{1}{4} N(x; \mu_3, \sigma_3^2)$$

# Gaussian Mixture Model

Multivariate Gaussian

$$P_j(\underline{x}) = N(\underline{x}; \underline{\mu}_j, \Sigma_j)$$

$d=2$



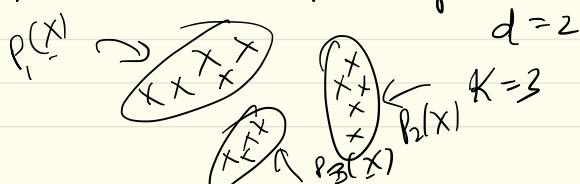
GMM

$$\hat{P}(\underline{x}) = \omega_1 \cdot P_1(\underline{x}) + \omega_2 \cdot P_2(\underline{x}) + \dots + \omega_k \cdot P_k(\underline{x})$$

$$P_j(\underline{x}) = N(\underline{x}; \underline{\mu}_j, \Sigma_j)$$

$\omega_1, \omega_2, \dots, \omega_K$  = weights  
associated with each Gaussian

$$\omega_i \geq 0, \quad \sum_{i=1}^K \omega_i = 1$$



Given

$$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}, \underline{x}_i \in \mathbb{R}^d$$

Estimate Parameters of Gaussian Mixture Model.

$$\hat{P}_{\mathcal{D}}(\underline{x}) = \sum_{j=1}^K \omega_j \cdot N(\underline{x}_j | \underline{\mu}_j, \Sigma_j)$$

Model Parameters:

$$\mathcal{D} = \{\omega_j, \underline{\mu}_j, \Sigma_j\}_{j=1}^K$$

find  $\mathcal{D}$  that is the "best-fit" for  $\mathcal{D}$ .

Maximum-Likelihood Solution

find  $\mathcal{D}$  that maximized

$$\hat{P}_{\mathcal{D}}(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N) = \prod_{n=1}^N \hat{P}_{\mathcal{D}}(\underline{x}_n)$$
$$|\mathcal{D}^* = \arg \max \left[ \prod_{n=1}^N \hat{P}_{\mathcal{D}}(\underline{x}_n) \right]^{1/N}$$

$$\begin{aligned}
 \mathcal{S}^* &= \arg \max_{\mathcal{S}} \log \prod_{n=1}^N \hat{P}_{\mathcal{S}}(\underline{x}_n) \\
 &= \arg \max_{\mathcal{S}} \sum_{n=1}^N \log \hat{P}_{\mathcal{S}}(\underline{x}_n) \\
 &= \arg \min_{\mathcal{S}} \frac{1}{N} \underbrace{\sum_{n=1}^N -\log \hat{P}_{\mathcal{S}}(\underline{x}_n)}_{e_n(\mathcal{S})} \\
 e_n(\mathcal{S}) &= -\log \hat{P}_{\mathcal{S}}(\underline{x}_n) \\
 &= -\log \left( \sum_{j=1}^K \omega_j N(\underline{x}_n; \underline{m}_j, \bar{\Sigma}_j) \right) \\
 \mathcal{S}^* &= \arg \min_{\mathcal{S}} \frac{1}{N} \sum_{n=1}^N e_n(\mathcal{S})
 \end{aligned}$$

Gradient Descent for  $\omega_j \geq 0$

$$\mathcal{S} = \left\{ (\omega_j), \underline{m}_j, \bar{\Sigma}_j \right\}_{j=1}^K$$

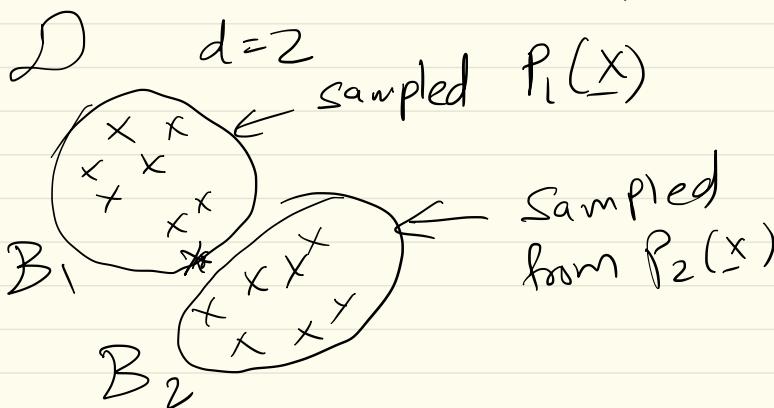
# Expectation - Maximization

$$\hat{\Sigma}^* = \underset{\Sigma}{\operatorname{arg\,min}} \frac{1}{N} \sum_{n=1}^N \ell_n(\Sigma)$$

## Auxiliary Variables

$$B_1, B_2, \dots, B_K \subseteq \mathcal{D}$$

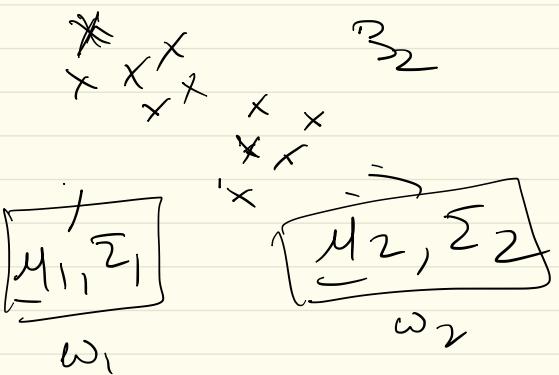
$B_i$  = set of points in  $\mathcal{D}$   
associated with cluster #  $i$   
and hence sampled from  $p_i(\underline{x})$



## Sub - Problem 1 :

Suppose  $B_1, B_2, \dots, B_K$  are given 1-e for each  $\underline{x}_n \in \mathbb{D}$  we know the cluster associated with it

Find  $\omega_j = \{\omega_j, \underline{\mu}_j, \Sigma_j\}_{j=1}^K$



use  
empirical  
estimated

$N_j = \# \text{ of Points}$   
in  $B_j$

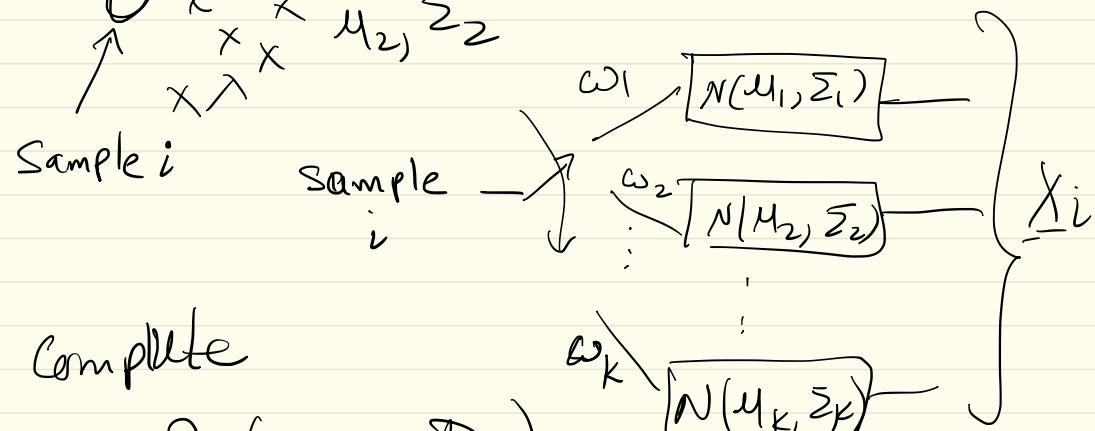
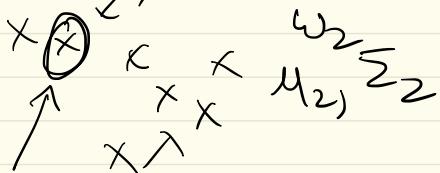
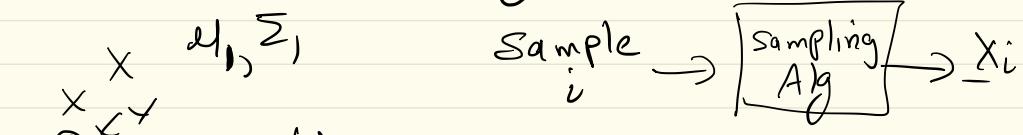
$$\omega_j = \frac{N_j}{N} \quad | \quad \underline{\mu}_j = \frac{1}{N_j} \sum_{\underline{x}_n \in B_j} \underline{x}_n$$

$$\Sigma_j = \frac{1}{N_j} \sum_{\underline{x}_n \in B_j} (\underline{x}_n - \underline{\mu}_j)(\underline{x}_n - \underline{\mu}_j)^T$$

## Sub-Problem #2

Suppose  $\Omega = \{\omega_j, \mu_j, \Sigma_j\}_{j=1}^K$  is known. find  $\beta_1, \dots, \beta_K$

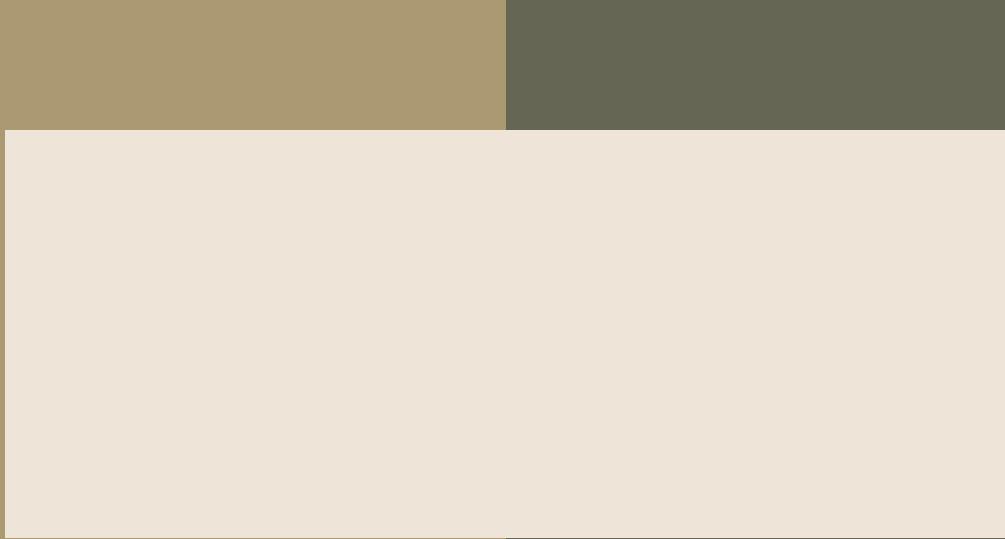
i.e for each  $x_n$ , determine the cluster  $\omega_i$  it belongs to.



Complete

$\Pr(x_n \in \beta_j)$  for each  $j = 1, 2, \dots, K$

& select the  $\beta_j$  with largest value



# Gaussian Mixture Model & EM Algorithm

$$\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}, \underline{x}_i \in \mathbb{R}^d$$

$$\hat{P}_{\mathcal{D}}(\underline{x}) = \sum_{j=1}^K \omega_j \cdot \underbrace{N(\underline{x}; \underline{\mu}_j, \underline{\Sigma}_j)}$$

$$\omega_j \geq 0 \quad \sum_{j=1}^K \omega_j = 1$$

multivariate Gaussian density

Given  $\mathcal{D}$ : find  $\hat{P}_{\mathcal{D}}(\underline{x})$  that is "best-fit" for  $\mathcal{D}$ .

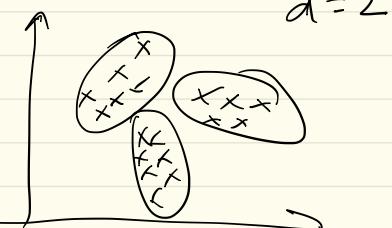
## Loss - Function

Maximum - Likelihood Solution

$$E_{in}(\Sigma) = \frac{1}{N} \sum_{n=1}^N -\log \hat{P}_{\mathcal{D}}(\underline{x}_n)$$

$$\Sigma = \{\omega_j, \underline{\mu}_j, \underline{\Sigma}_j\}_{j=1}^K \times \text{SGD.}$$

$\times$  EM Algorithm.



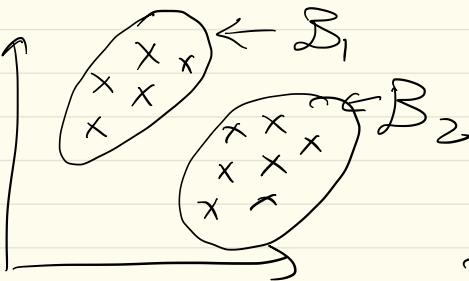
# Expectation-Maximization Algorithm

## Auxiliary Variables

$$B_1, B_2, \dots, B_k \subseteq D$$

$B_i$  = Set of Points in  $D$   
belonging to "cluster =  $i$ "  
and hence sampled from

$$d=2 \quad P_i(\underline{x}) = \mathcal{N}(\underline{x}; \underline{\mu}_i, \Sigma_i)$$



Assume:

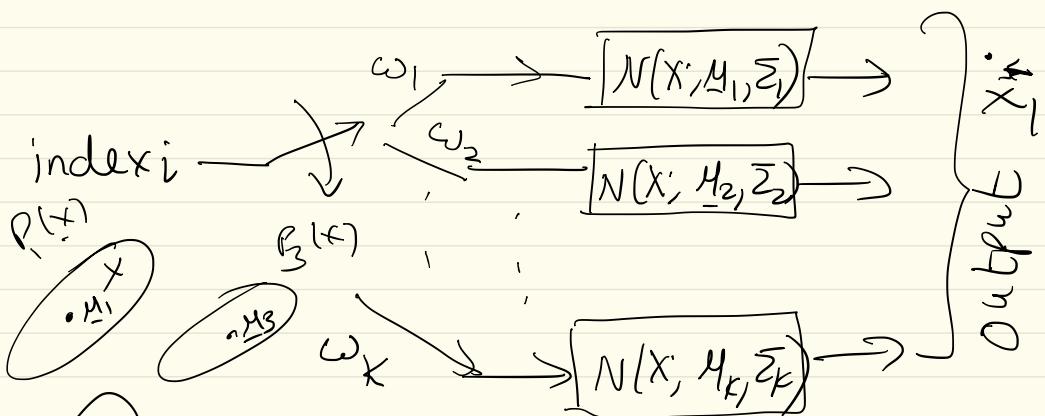
$$\begin{aligned} & k \quad B_i \quad \text{are disjoint} \\ & \bigcup_{i=1}^k B_i = D \end{aligned}$$

$$f(x_i) \approx \hat{P}_{S2}(x_i)$$

Thought Experiment <sup>density  
of  $x_i$</sup>

Assuming each  $\underline{x}_i \in D$  is sampled independently from

$$\hat{P}_{S2}(x) = \sum_{j=1}^k \omega_j N(x; \underline{\mu}_j, \underline{\Sigma}_j)$$



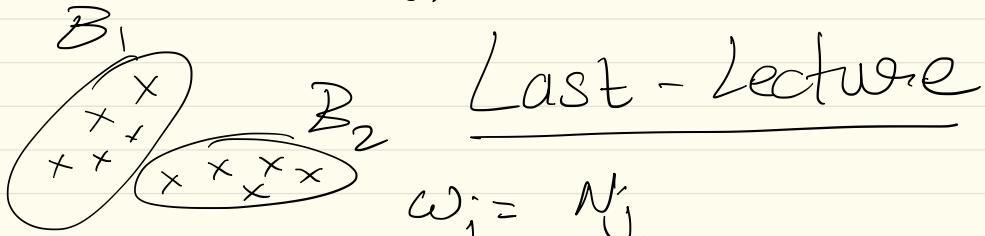
$$\Pr(\underline{x}_i \text{ belongs to cluster } j) = \omega_j$$

$$f(\underline{x}_i | \underline{x}_i \in B_j) = N(\underline{x}_i; \underline{\mu}_j, \underline{\Sigma}_j)$$

# Sub-Problem 1

Suppose  $B_1, B_2, \dots, B_K$  are given. Identify

$$S = \{\omega_j, \underline{y}_j, \Sigma_j\}_{j=1}^K$$



$$\omega_j = \frac{N_j}{N}$$

$N_j$  = # of Points in  $B_j$

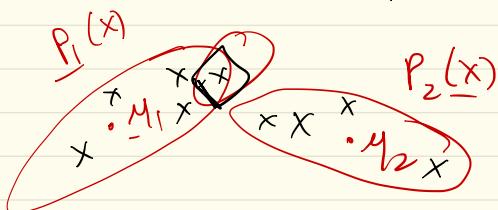
$$\underline{y}_j = \frac{1}{N_j} \sum_{x_n \in B_j} x_n$$

$$\Sigma_j = \frac{1}{N_j} \sum_{x_n \in B_j} (x_n - \underline{y}_j)(x_n - \underline{y}_j)^T$$

(Empirical Estimates)

## Sub problem 2

Given  $\mathcal{D} = \{\underline{\omega}_j, \underline{\mu}_j, \Sigma_j\}_{j=1}^K$   
 find  $\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_K$



$$P_1(\underline{x}) = N(\underline{x}; \underline{\mu}_1, \Sigma_1)$$

$$P_2(\underline{x}) = N(\underline{x}; \underline{\mu}_2, \Sigma_2)$$

$$\vdots$$

$$P_K(\underline{x}) = N(\underline{x}; \underline{\mu}_K, \Sigma_K)$$

Equivalently

Given any  $\underline{x}_n \in \mathcal{D}$  find the  
 "most-likely" cluster associated  
 with  $\underline{x}_n$  "compute"  
 Posterior Probability

$$\Pr(\text{cluster} = j \mid \underline{x}_n)$$

for  
 each  $j = 1, 2, \dots, K$

$\underline{x}_n \in \mathcal{B}_j$  if  $j^* = \arg \max_{j \in \{1, 2, \dots, K\}} \Pr(\text{cluster} = j \mid \underline{x}_n)$

~~0.8, 0.2~~

# Baye's Rule

$$P(\text{cluster} = j | \underline{x}_n)$$

$$= \frac{\hat{P}(\underline{x}_n | \text{cluster} = j)}{P_0(\text{cluster} = j)}$$

$$= \frac{\hat{P}(\underline{x}_n)}{N(\underline{x}_n; \underline{\mu}_j, \Sigma_j) \cdot \omega_j}$$

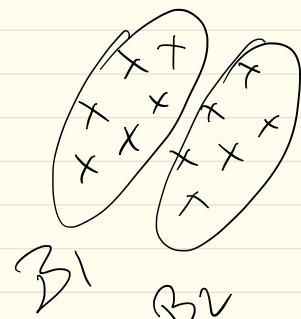
$$\underbrace{\hat{P}_{\Sigma}(\underline{x}_n)}_{\text{independent of } j} = \sum_{j=1}^K \omega_j N(\underline{x}_n; \underline{\mu}_j, \Sigma_j)$$

$$j^* = \arg \max_{j \in \{1, 2, \dots, K\}} \omega_j \cdot N(\underline{x}_n; \underline{\mu}_j, \Sigma_j)$$

$$\underline{x}_n \in \mathcal{D}$$

# EM - Algorithm

- start with some arbitrary choice of  $B_1, B_2, \dots, B_K$



Step 1:

Given  $B_1, \dots, B_K$

use sub-problem # 1

to compute

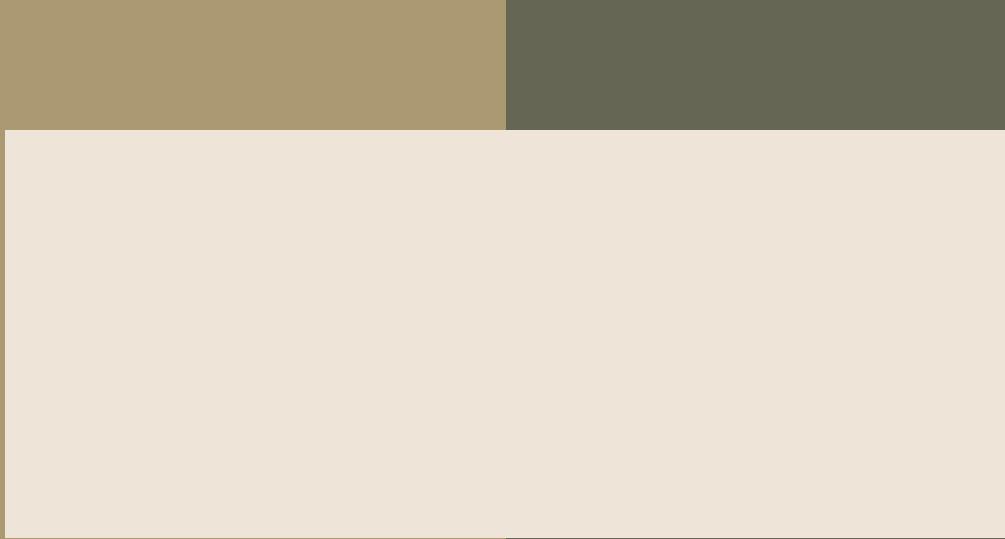
$$\mathcal{D} = \{w_j, \mu_j, \Sigma_j\}_{j=1}^K$$

Step 2: Given  $\mathcal{D} = \{w_j, \mu_j, \Sigma_j\}_{j=1}^K$

use sub-problem # 2 to compute

$B_1, \dots, B_K$

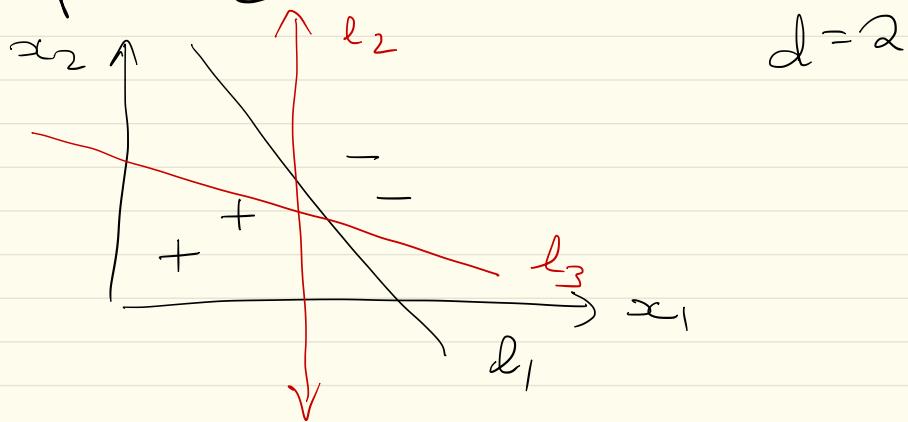
- Alternate between Step # 1 & # 2 until convergence.



# Support Vector Machines.

## - Binary Linear Classification

Separable data set



SVM finds a linear classifier with largest "margin"

## Training Set

$$\mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$

$$\underline{x}_i \in \mathbb{R}^d, \quad y_i \in \{-1, +1\}$$

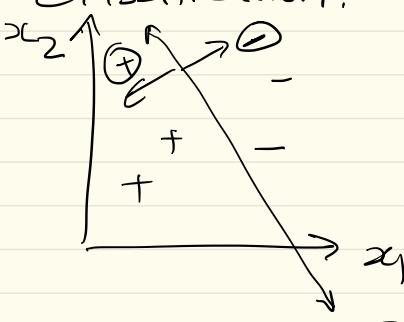
## Decision Rule:

$$\hat{y} = \text{Sign}(\underline{w}^T \underline{x} + b)$$

$$l: \underline{w}^T \underline{x} + b, \quad \underline{w} \in \mathbb{R}^d, \quad b \in \mathbb{R}$$

assume  $\mathcal{D}$  is linearly separable

data-set &  $l$  achieves perfect classification.



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

Given any  $\underline{x}_n \in \mathbb{R}^d$

$$\text{if } y_n = +1 \Rightarrow \underline{w}^T \underline{x}_n + b > 0$$

$$\text{if } y_n = -1 \Rightarrow \underline{w}^T \underline{x}_n + b < 0$$

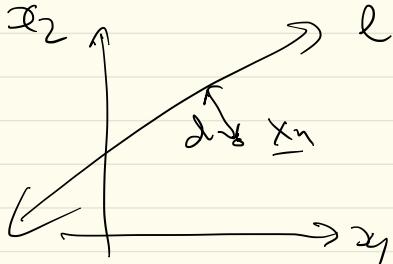
$$\underbrace{y_n (\underline{w}^T \underline{x}_n + b)}_{> 0}$$

$$\forall n = 1, 2, \dots, N$$

Margin of  $\ell$ :  $\underline{w}^T \underline{x} + b$

distance between  $\underline{x}_n$  &  $\ell$

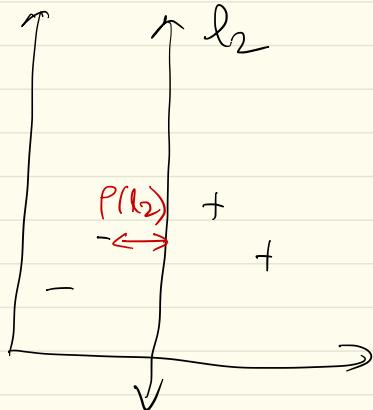
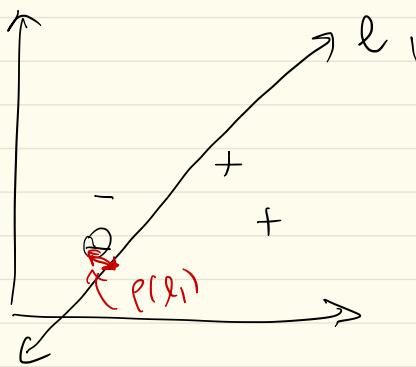
$$\text{dist}(\underline{x}_n, \ell) = \frac{|\underline{w}^T \underline{x}_n + b|}{\|\underline{w}\|}$$



$$= \frac{y_n (\underline{w}^T \underline{x}_n + b)}{\|\underline{w}\|}$$

margin of  $\ell$ :

$$\rho(\ell) = \min_{1 \leq n \leq N} \text{dist}(\underline{x}_n, \ell)$$



Idea in SVM:

find  $\ell$  that maximizes  $\rho(\ell)$

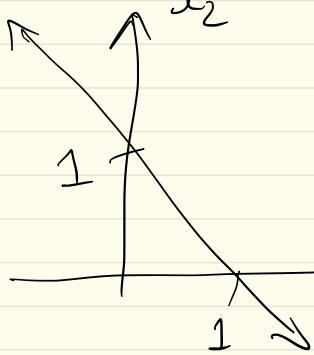
# Normalization of $(\underline{w}, b)$

$$l_1: x_1 + x_2 - 1 = 0$$

$$\underline{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad b = -1$$

$$l_2: 5x_1 + 5x_2 - 5 = 0$$

$$\underline{w} = \begin{bmatrix} 5 \\ 5 \end{bmatrix} \quad b = -5$$



$(\underline{w}, b)$  &

$(\alpha \underline{w}, \alpha b)$

$\alpha l_1 = l_2$  represent the  
same decision boundary for  
any non-zero  $\alpha$ .

Given  $\ell: (\underline{w}^T \underline{x} + b)$

define

$$\delta = \min_{1 \leq n \leq N} y_n (\underline{w}^T \underline{x}_n + b)$$

Without loss of generality  
we can select  $(\underline{w}, b)$  such that  
 $\delta = 1$ .

Proof: Note that since  $\ell$  is  
a perfect classifier,  $y_n (\underline{w}^T \underline{x}_n + b) > 0$   
 $\forall n$ . Hence  $\delta > 0$

Suppose  $\delta \neq 1$ . Define

$$\begin{aligned} \tilde{\underline{w}} &= \frac{\underline{w}}{\delta}, & \tilde{b} &= \frac{b}{\delta} \\ \tilde{\ell}: (\tilde{\underline{w}}^T \tilde{\underline{x}} + \tilde{b}) & \end{aligned}$$

$\delta > 0$   
 $\delta \neq 1$   
 $\tilde{\ell} = \ell$ .

Claim:  $\min_{1 \leq n \leq N} y_n (\tilde{\underline{w}}^T \tilde{\underline{x}}_n + \tilde{b}) = 1$

$$\min_{1 \leq n \leq N} y_n (\underline{w}^T \underline{x}_n + \underline{b})$$

$$= \min_{1 \leq n \leq N} y_n \left( \frac{\underline{w}^T \underline{x}_n}{S} + \frac{b}{S} \right)$$

$$= \underbrace{\min_{1 \leq n \leq N} y_n (\underline{w}^T \underline{x}_n + b)}_S = 1$$

---


$$\text{but } S = \min_{1 \leq n \leq N} y_n (\underline{w}^T \underline{x}_n + b)$$

From now assume that  
for any  $\ell: (\underline{w}^T \underline{x} + b)$   
 $\underline{w}$  &  $b$  are such that

$$\min_{1 \leq n \leq N} y_n (\underline{w}^T \underline{x}_n + b) = 1$$

Margin

$$\rho(l) =$$

$$l: \underline{w}^T \underline{x} + b$$

$$\min_{1 \leq n \leq N} \frac{y_n (\underline{w}^T \underline{x}_n + b)}{\|\underline{w}\|}$$

$$= \frac{\min_{1 \leq n \leq N} y_n (\underline{w}^T \underline{x}_n + b)}{\|\underline{w}\|}$$

$$= \frac{1}{\|\underline{w}\|} = \frac{1}{\sqrt{\underline{w}^T \underline{w}}}$$

Problem A

maximize

$$\frac{1}{\sqrt{\underline{w}^T \underline{w}}} \Rightarrow \text{minimize } \frac{1}{2} \underline{w}^T \underline{w}$$

subject to

$$\min_{1 \leq n \leq N} y_n (\underline{w}^T \underline{x}_n + b) = 1$$

Variables:  $\underline{w} \in \mathbb{R}^d, b \in \mathbb{R}$

$$\min_{1 \leq n \leq N} y_n (\underbrace{w^T x_n + b}_{=a_n}) = 1$$

$$\text{minimize } \frac{1}{2} \underline{w}^T \underline{w}$$

subject to

$$y_n (\underline{w}^T \underline{x}_n + b) \geq 1$$

$$\forall n = 1, 2, \dots, N$$

with equality for at-least one  $n \in \{1, 2, \dots, N\}$

$$\min \{a_1, a_2, \dots, a_N\} = 1$$

original constraint

$$a_1 \geq 1, a_2 \geq 1, \dots, a_N \geq 1$$

with equality for at-least

some  $a_n$

Problem

Variables:  $\underline{w}, b$

minimize  $\frac{1}{2} \underline{w}^T \underline{w}$

$\forall n = 1, 2, \dots$

$N$

$$y_n (\underline{w}^T \underline{x}_n + b) \geq 1$$

with ~~Equality for at least one value of  $n \in \{1, 2, \dots, N\}$~~

Solution of Problem A:  $(\underline{w}^*, b^*)$

Objective:  $\frac{1}{2} \underline{w}^* \underline{w}^*$

Solution to Problem B:  $(\hat{\underline{w}}, \hat{b})$

Objective:  $\frac{1}{2} \hat{\underline{w}}^T \hat{\underline{w}}$

is smaller than

$$\frac{1}{2} \underline{w}^* \underline{w}^*$$

Problem  
B

## Problem A:

Objective: minimize  $\frac{1}{2} \underline{w}^T \underline{w}$   
 such that  $y_n (\underline{w}^T \underline{x}_n + b) \geq 1$

$\forall n = 1, 2, \dots N$ , Equality for at least one  $n$   
 $\Rightarrow$  optimal  $(\underline{w}^*, b)$  given SVM.

## Problem B:

minimize  $\frac{1}{2} \underline{w}^T \underline{w}$   
 Such that  $y_n (\underline{w}^T \underline{x}_n + b) \geq 1$   
 $\forall n = 1, 2, \dots N$

optimal  $\boxed{(\hat{\underline{w}}, \hat{b})}$

then  $\frac{1}{2} \hat{\underline{w}}^T \hat{\underline{w}} \leq \frac{1}{2} \underline{w}^{*T} \underline{w}^*$   
 (in general)

Relaxation of Problem A  
to Problem B is Exact.

Proof: let  $(\hat{\underline{w}}, \hat{b})$  be  
a solution for Problem B

Suppose:  $y_n(\hat{\underline{w}}^T \underline{x}_n + \hat{b}) > 1$

$\forall n = 1, 2, \dots, N$

$$\text{i.e. } \Theta = \min_{1 \leq n \leq N} [y_n(\hat{\underline{w}}^T \underline{x}_n + \hat{b})]$$

$\Theta > 1$ .

Consider  $\underline{w}^* = \frac{\hat{\underline{w}}}{\Theta}$ ,  $b^* = \frac{\hat{b}}{\Theta}$  ✓

Claim: (1)  $\frac{1}{2} \underline{w}^{*T} \underline{w}^* < \frac{1}{2} \hat{\underline{w}}^T \hat{\underline{w}}$

2)  $(\underline{w}^*, b^*)$  is feasible in Problem A

$$\begin{aligned}
 & \frac{1}{2} \underline{\underline{w}}^T \underline{\underline{w}} \\
 = & \frac{1}{2} \left( \frac{\underline{\underline{w}}}{\theta} \right)^T \left( \frac{\underline{\underline{w}}}{\theta} \right) \\
 = & \frac{1}{2} \frac{\underline{\underline{w}}^T \underline{\underline{w}}}{\theta^2} \\
 < & \frac{1}{2} \underline{\underline{w}}^T \underline{\underline{w}}
 \end{aligned}$$

Thus:  
 optimal  $(\underline{\underline{w}}, b)$   
 in problem B  
 is also the  
 optimal  $(\underline{\underline{w}}, b)$   
 in problem A

$$\min_{1 \leq n \leq N} y_n (\underline{\underline{w}}^*{}^T x_n + b^*)$$

$$= \min_{1 \leq n \leq N} y_n \left( \left( \frac{\hat{\underline{\underline{w}}}}{\theta} \right)^T x_n + \frac{\hat{b}}{\theta} \right)$$

$$= \min_{1 \leq n \leq N} y_n \left( \hat{\underline{\underline{w}}}^T x_n + \hat{b} \right)$$

$$= \frac{\min_{1 \leq n \leq N} y_n (\hat{\underline{\underline{w}}}^T x_n + \hat{b})}{\min_{1 \leq n \leq N} y_n (\hat{\underline{\underline{w}}}^T x_n + \hat{b})} = 1$$

$$l: (\underline{w}^T \underline{x} + b)$$

Given  $\mathcal{D} = \{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$

$$\text{minimize } \frac{1}{2} \underline{w}^T \underline{w}$$

such that  $y_n (\underline{w}^T \underline{x}_n + b) \geq 1$   
 $\forall n \in \{1, 2, \dots, N\}$

$\frac{1}{2} \underline{w}^T \underline{w} \leftarrow$  quadratic in  $\underline{w}$   
 $y_n (\underline{w}^T \underline{x}_n + b)$  is linear in  $(\underline{w}, b)$

Quadratic Programming  
Problem.

# Quadratic Programming

variable:  $\underline{u} \in \mathbb{R}^L$

objective:  $\left[ \frac{1}{2} \underline{u}^T Q \underline{u} + \underline{p}^T \underline{u} \right]$

minimize  $Q \in \mathbb{R}^{L \times L}$  matrix

$\underline{p} \in \mathbb{R}^L$  vector

Constraints:  $l_i \leq u_i \leq h_i$   
 $i = 1, 2, \dots, L$

$$\underline{a}_m^T \underline{u} \geq c_m$$

$$m = 1, 2, \dots, M$$

Quadratic Solver

Input:  $Q, \underline{p}, \{l_i, h_i\}_{i=1}^d, \{\underline{a}_m\}_{m=1}^M$

SVM is a Quadratic Programming Problem

SVM : minimize  $\frac{1}{2} \underline{w}^T \underline{w}$

such that  $y_n (\underline{w}^T \underline{x}_n + b) \geq 1$

variables:  $(\underline{w}, b) \quad \forall n \in \{1, 2, \dots, N\}$

Quadratic Programming

minimize  $\frac{1}{2} \underline{u}^T Q \underline{u} + \underline{p}^T \underline{u}$

such that  $\underline{a}_m^T \underline{u} \geq c_m \quad m=1, 2, \dots, M$

Identify  $\underline{u}, Q, P, \{a_m, c_m\}$

$\underline{u} = \begin{bmatrix} \underline{w} \\ b \end{bmatrix} \in \mathbb{R}^{d+1}$        $a_m = [y_m \underline{x}_m^T \underline{y}_m]^T$

$Q = \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix}$        $P = 0$

# Example

$$N = 4$$

$$x_1 = (0, 0)$$

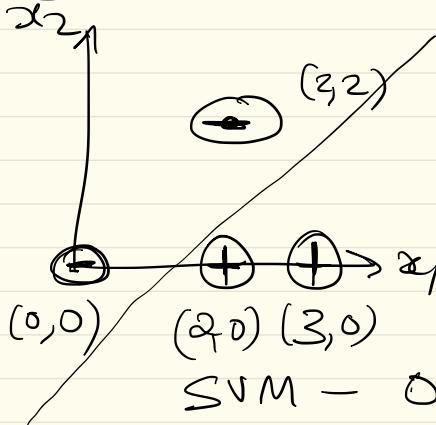
$$y_1 = -1$$

$$x_3 = (2, 0) \quad y_3 = +1$$

$$x_2 = (2, 2)$$

$$y_2 = -1$$

$$x_4 = (3, 0) \quad y_4 = +1$$



find a linear classification boundary with max margin.

## SVM - Optimization

$$\text{minimize} : \frac{1}{2} (w_1^2 + w_2^2)$$

$$\text{such that } y_n (w_1 x_{1n} + w_2 x_{2n} + b) \geq 1 \quad \forall n \in \{1, 2, 3, 4\}$$

$$n=1 : -1 (0 \cdot w_1 + 0 \cdot w_2 + b) \geq 1 \quad \therefore -b \geq 1$$

$$n=2 : -1 (2 w_1 + 2 w_2 + b) \geq 1$$

$$n=3 : 1 (2 w_1 + 0 \cdot w_2 + b) \geq 1 \quad | 3 w_1 + b \geq 1$$

$$\text{minimize}_{\omega_1, \omega_2} \frac{1}{2} (\omega_1^2 + \omega_2^2)$$

$$n=1 : -b \geq 1' \quad \textcircled{1}$$

$$n=2 : -\underbrace{(2\omega_1 + 2\omega_2 + b)}_{\text{---}} \geq 1 \quad \textcircled{2}$$

$$n=3 : 2\omega_1 + b \geq 1 \quad \textcircled{3}$$

$$n=4 : 3\omega_1 + b \geq 1 \quad \textcircled{4}$$

$$(1) \oplus (3) \quad 2\omega_1 \geq 2 \Rightarrow \omega_1 \geq 1$$

$$(2) \oplus (3) \quad -2\omega_2 \geq 2 \Rightarrow \omega_2 \leq -1$$

$$\text{Select } \omega_1 = +1 \quad \& \quad \omega_2 = -1$$

$$\boxed{b \leq -1} \quad b \leq -1 \quad b \geq -1$$

$$\boxed{b = -1} \quad \text{satisfies all constraints} \quad b \geq -2$$

$$\omega_1^* = +1, \quad \omega_2^* = -1, \quad b^* = -1 \quad \begin{matrix} \text{is} \\ \text{solution} \end{matrix}$$