

Mixture models & EM

Sunday, September 29, 2024

9:38 PM

See Chapter 9, Pattern Recognition & Machine Learning

Supervised vs Unsupervised Learning:

Supervised: Input x . Predict output y .

Unsupervised: Learn about x .

- Which values of x are "good".
- One version of "good"
 - = Likely to be seen.
- $P(x)$ is high for training cases, $\therefore P(x)$ is low elsewhere.

K-means clustering

Sunday, September 29, 2024

9:41 PM

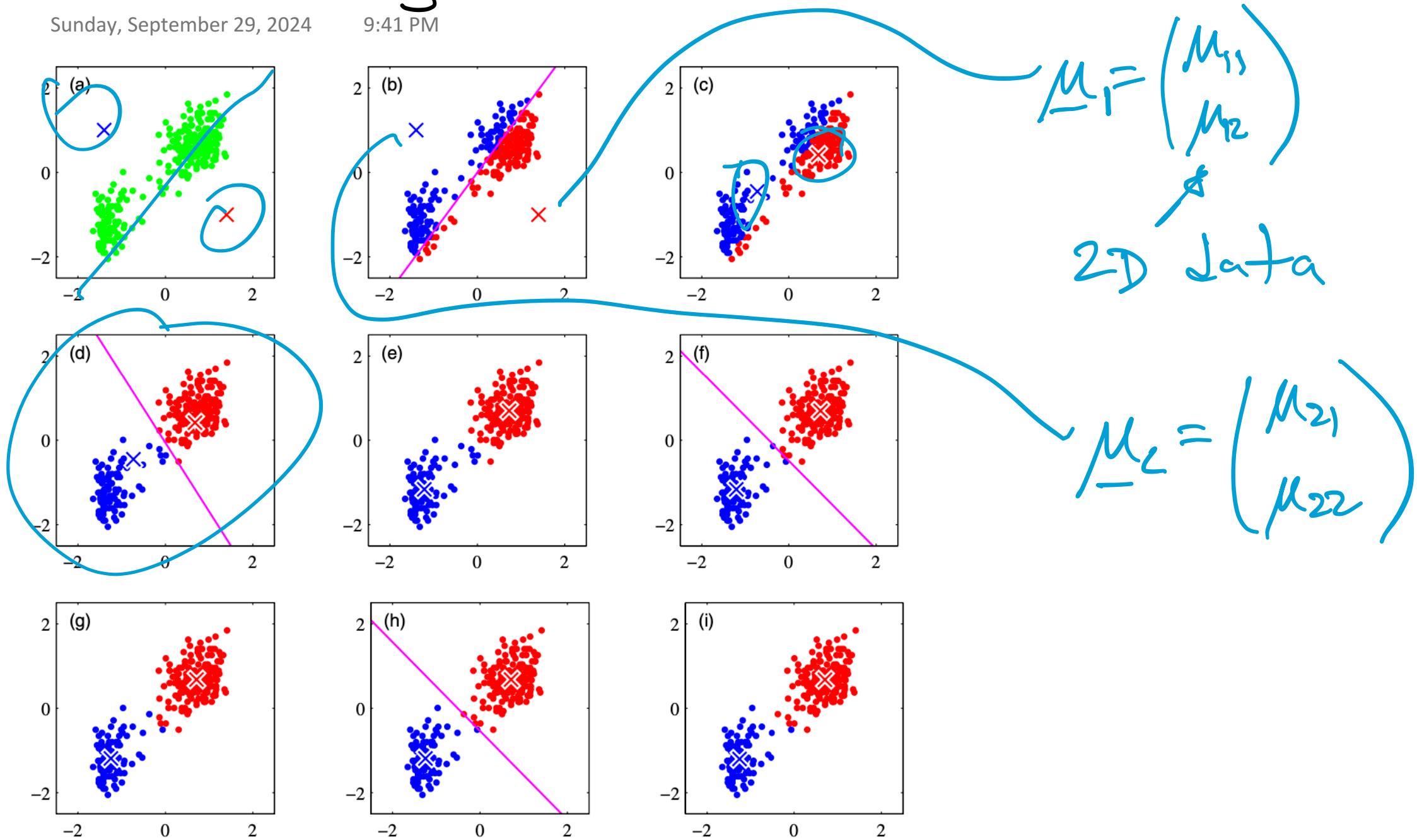


Figure 9.1 Illustration of the K -means algorithm using the re-scaled Old Faithful data set. (a) Green points denote the data set in a two-dimensional Euclidean space. The initial choices for centres μ_1 and μ_2 are shown by the red and blue crosses, respectively. (b) In the initial E step, each data point is assigned either to the red cluster or to the blue cluster, according to which cluster centre is nearer. This is equivalent to classifying the points according to which side of the perpendicular bisector of the two cluster centres, shown by the magenta line, they lie on. (c) In the subsequent M step, each cluster centre is re-computed to be the mean of the points assigned to the corresponding cluster. (d)–(i) show successive E and M steps through to final convergence of the algorithm.

k-means error function

Sunday, September 29, 2024

9:44 PM

$$E(X) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

Annotations:

- A blue circle encloses the variable K in the summation.
- An arrow points from the circled K to the text "K clusters (means)".
- An arrow points from the term μ_k to the text "kth mean (vector)".
- An arrow points from the term x_n to the text "vector datapoint".

$r_{n1} \dots r_{nK}$ is a "one hot encoding" of the cluster assignment of \underline{x}_n

Eg, $r_{n1} \dots r_{nK} = (0, 0, 0, \textcolor{blue}{1}, 0)$ $\leftarrow \underline{x}_n \text{ is in cluster 4}$

Algorithm

Sunday, September 29, 2024

9:50 PM

Initialization

Method 1: Randomly initialize r_{nk} 's.

Proceed to update means.

Method 2: Randomly pick K \bar{x} 's, set $\underline{\mu}$'s to those

Iterate to convergence:

For $n=1 \dots N$ set $r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_j \|x_n - \underline{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases}$

For $k=1 \dots K$ set

$$\underline{\mu}_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$$

Proof Part 1

Sunday, September 29, 2024

9:51 PM

$$E = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\underline{x}_n - \mu_k\|^2$$

$$\min E = \sum_{k=1}^K r_{mk} \|\underline{x}_n - \mu_k\|^2$$

r_{mj}
Don't use r_{nk}

$$= \begin{cases} 1 & \text{if } j = \arg \min_k \|\underline{x}_n - \mu_k\|^2 \\ 0 & \text{otherwise} \end{cases}$$

Proof Part 2

Sunday, September 29, 2024

9:52 PM

$$E = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

$$\min E: \text{ Set } \frac{\partial E}{\partial \mu_j} = 0$$

$$\begin{cases} \frac{\partial E}{\partial \mu_{j1}} \\ \frac{\partial E}{\partial \mu_{j2}} \\ \vdots \\ \frac{\partial E}{\partial \mu_{jD}} \end{cases}$$

$$\frac{\partial E}{\partial \mu_j} = \frac{\partial}{\partial \mu_j} \sum_{n=1}^N \sum_{k=1}^K r_{nk} (\underline{x}_n - \underline{\mu}_k)^T (\underline{x}_n - \underline{\mu}_k)$$

$\frac{\partial r_{nk}(\underline{x}_n - \dots)}{\partial \mu_j} \rightarrow \infty \text{ if } k \neq j$

$$\underline{x}_n^T \underline{x}_n - 2 \underline{x}_n^T \underline{\mu}_k + \underline{\mu}_k^T \underline{\mu}_k$$

$$= \sum_{n=1}^N r_{nj} \frac{\partial}{\partial \mu_j} (\underline{x}_n^T \underline{x}_n - 2 \underline{x}_n^T \underline{\mu}_j + \underline{\mu}_j^T \underline{\mu}_j)$$

$$= \sum_{n=1}^N r_{nj} (-2 \underline{x}_n + 2 \underline{\mu}_j) = 0$$

$$\sum_{n=1}^N r_{nj} \underline{x}_n = \sum_{n=1}^N r_{nj} \cancel{\underline{\mu}_j} \Rightarrow \underline{\mu}_j = \frac{\sum_{n=1}^N r_{nj} \underline{x}_n}{\sum_{n=1}^N r_{nj}}$$

Error minimization example

Sunday, September 29, 2024

9:52 PM

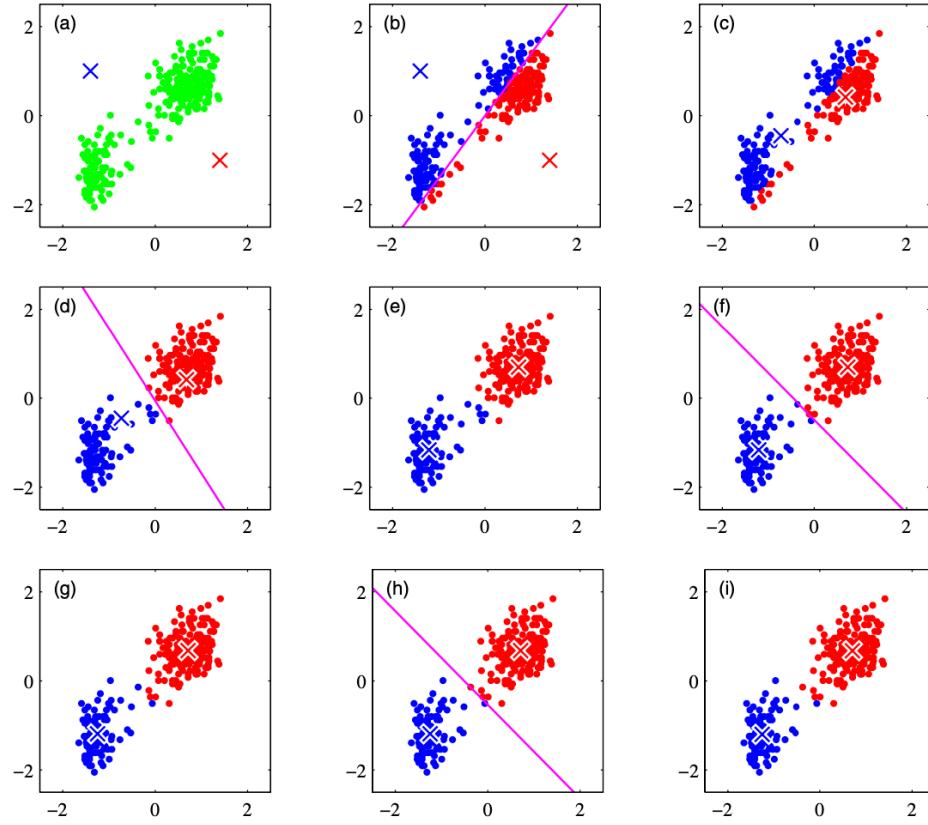
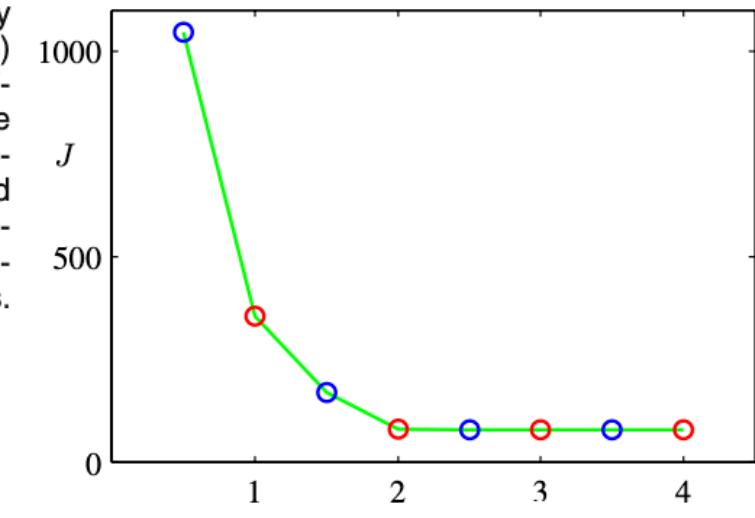


Figure 9.1 Illustration of the K -means algorithm using the re-scaled Old Faithful data set. (a) Green points denote the data set in a two-dimensional Euclidean space. The initial choices for centres μ_1 and μ_2 are shown by the red and blue crosses, respectively. (b) In the initial E step, each data point is assigned either to the red cluster or to the blue cluster, according to which cluster centre is nearer. This is equivalent to classifying the points according to which side of the perpendicular bisector of the two cluster centres, shown by the magenta line, they lie on. (c) In the subsequent M step, each cluster centre is re-computed to be the mean of the points assigned to the corresponding cluster. (d)–(i) show successive E and M steps through to final convergence of the algorithm.

Plot of the cost function J given by (9.1) after each E step (blue points) and M step (red points) of the K -means algorithm for the example shown in Figure 9.1. The algorithm has converged after the third M step, and the final EM cycle produces no changes in either the assignments or the prototype vectors.



Application

Sunday, September 29, 2024

9:55 PM

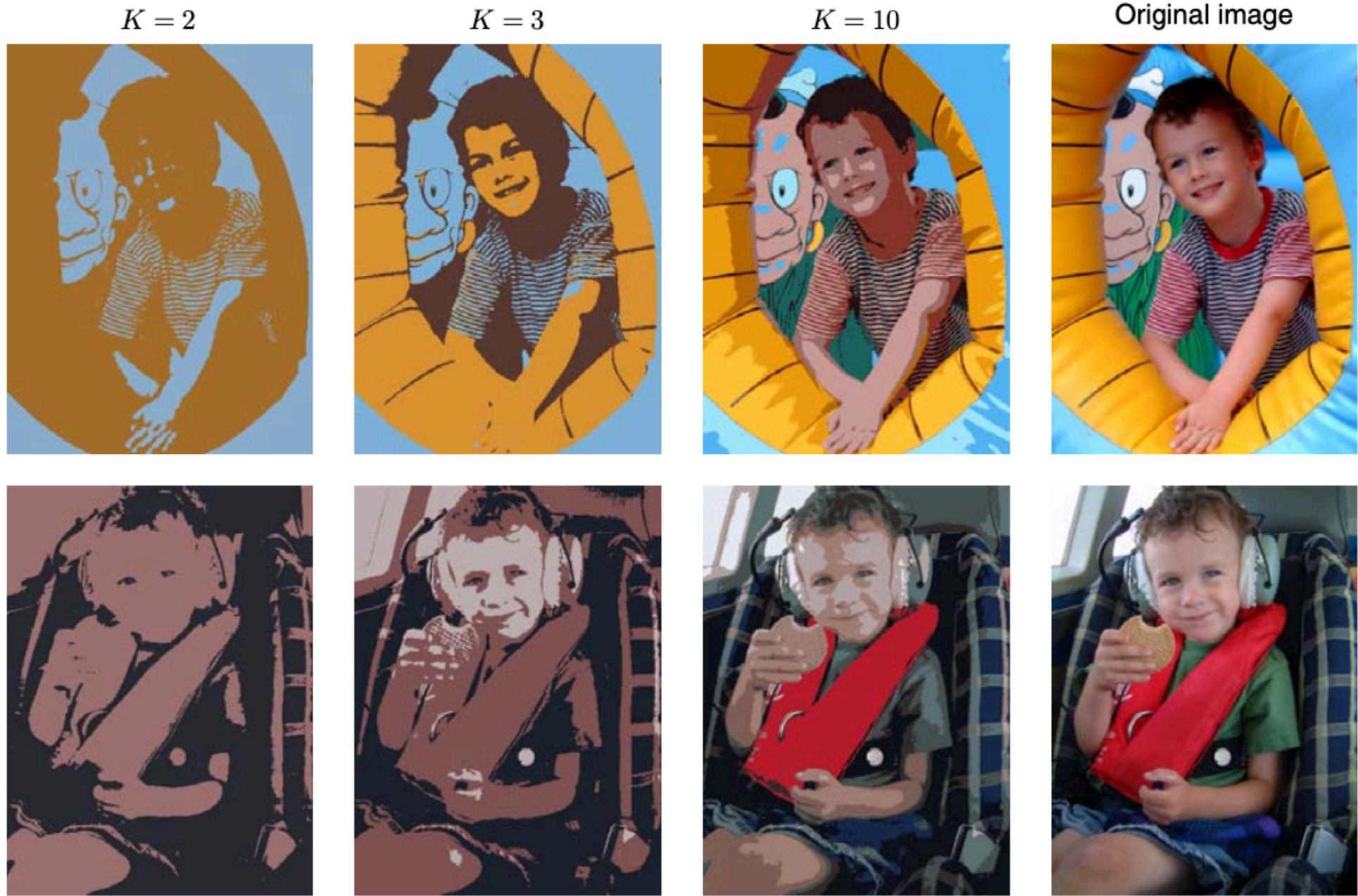
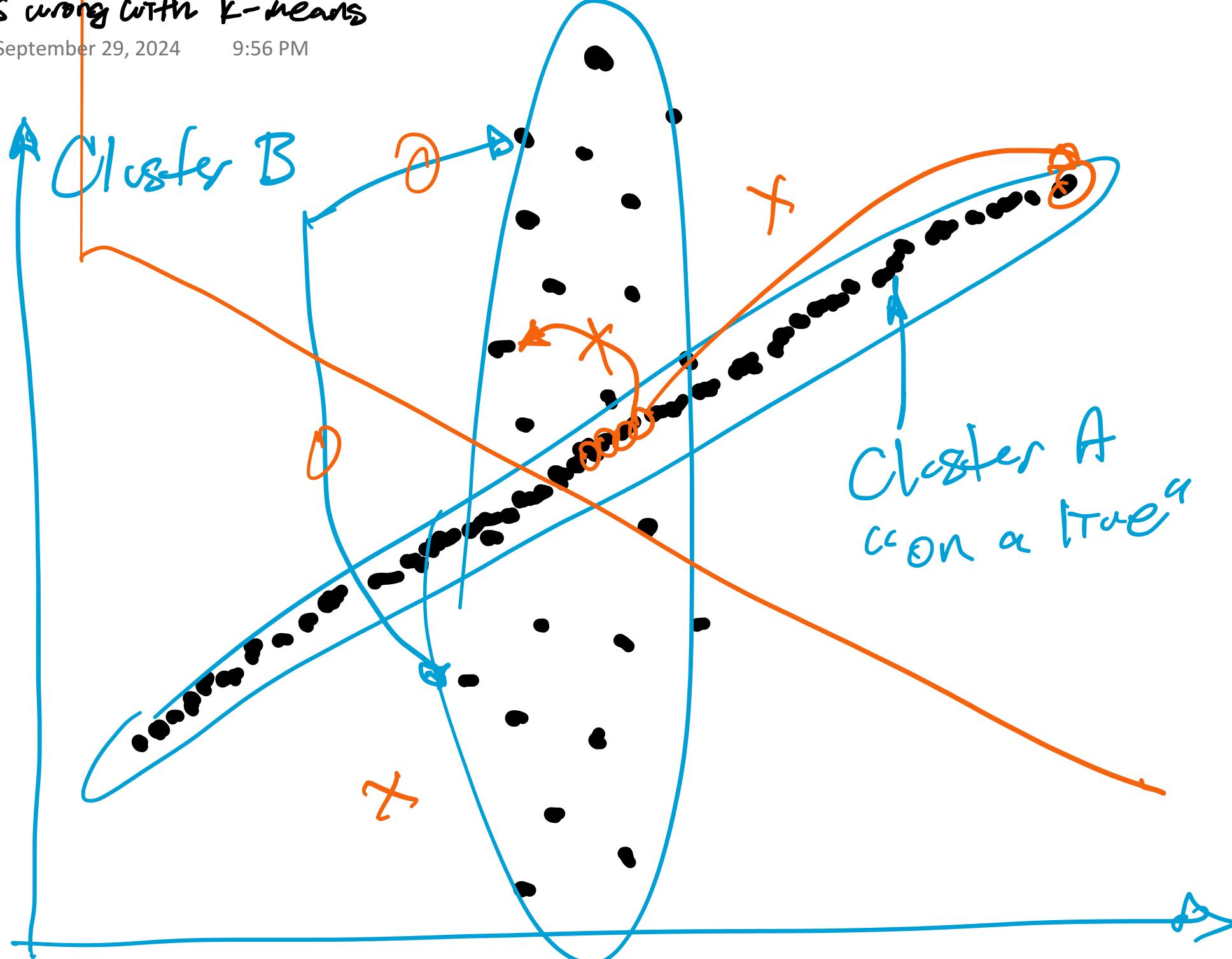


Figure 9.3 Two examples of the application of the K -means clustering algorithm to image segmentation showing the initial images together with their K -means segmentations obtained using various values of K . This also illustrates of the use of vector quantization for data compression, in which smaller values of K give higher compression at the expense of poorer image quality.

What's wrong with k-means

Sunday, September 29, 2024

9:56 PM



Mixture of Gaussians

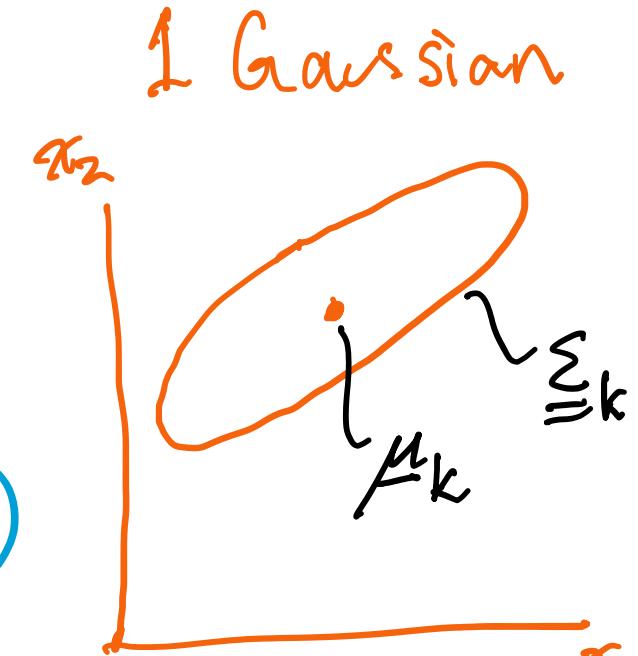
Sunday, September 29, 2024

9:57 PM

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

N(μ, Σ)

parameter
not π or Σ

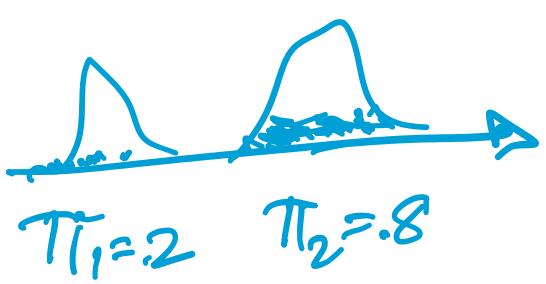


μ_k = mean of cluster k

Σ_k indicates a matrix

Σ_k = covariance matrix of Gaussian k

π_k = mixing proportion of Gaussian k



$$\Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_{22}^2 \end{pmatrix}$$

Example: 2D Gaussian

Tuesday, October 1, 2024

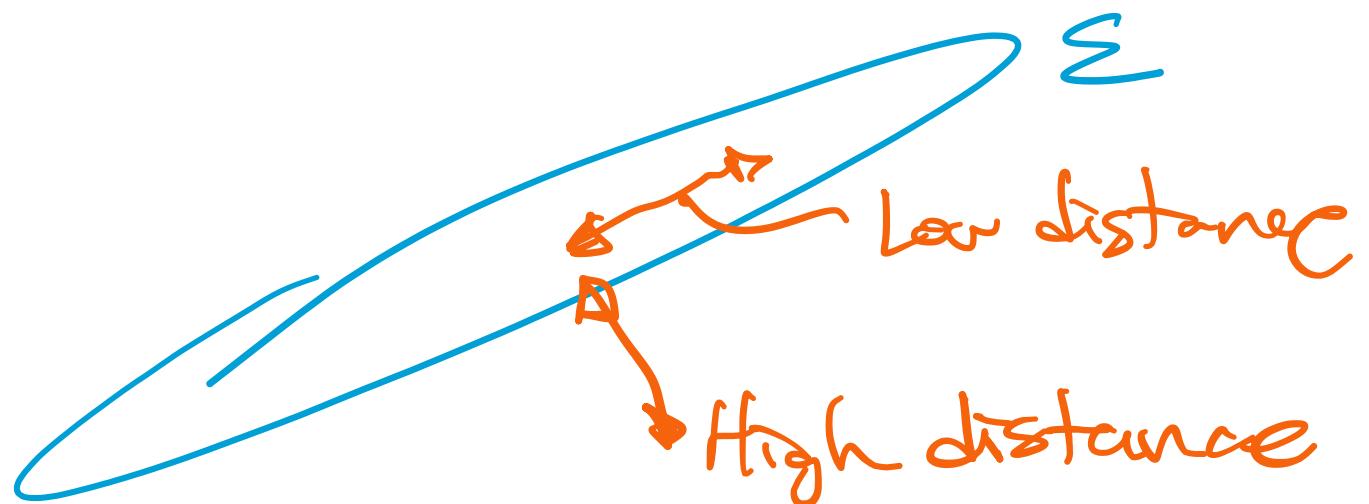
10:45 PM

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_{22}^2 \end{pmatrix}$$

$$N(\underline{x} | \mu, \Sigma) = \frac{1}{\sqrt{2\pi \det \Sigma}} \exp\left(-\frac{1}{2} (\underline{x} - \mu)^T \Sigma^{-1} (\underline{x} - \mu)\right)$$

Determinant

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



$D=5 = \text{Dimensionality of data}$

$$\Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \sigma_{13}^2 & \sigma_{14}^2 & \sigma_{15}^2 \\ \sigma_{12}^2 & \sigma_{22}^2 & \sigma_{23}^2 & \sigma_{24}^2 & \sigma_{25}^2 \\ \vdots & \vdots & \ddots & \ddots & \ddots \\ \sigma_{15}^2 & \sigma_{25}^2 & - & - & \sigma_{55}^2 \end{pmatrix}$$

K is separate, like 20

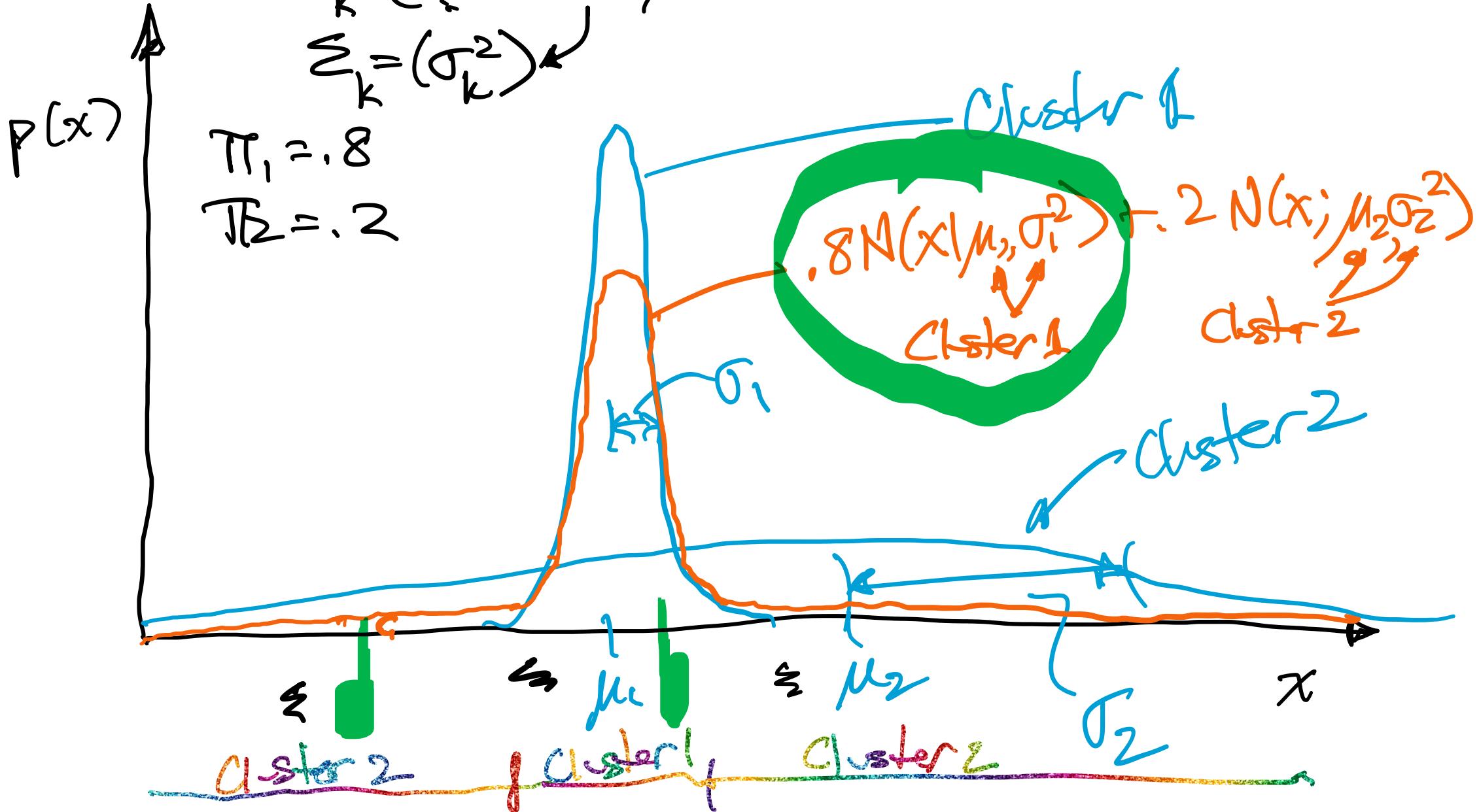
Scalar example (1-D)

Sunday, September 29, 2024

10:06 PM

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) = \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(x - \mu_k)^2}{2\sigma_k^2}}$$

$\Sigma_k = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \dots & \sigma_{1D}^2 \\ \sigma_{21}^2 & \dots & \dots & \end{pmatrix}$
 $\Sigma_k = (\sigma_k^2)$



Example: Mixture of 3 Gaussians

Monday, September 30, 2024

10:43 AM

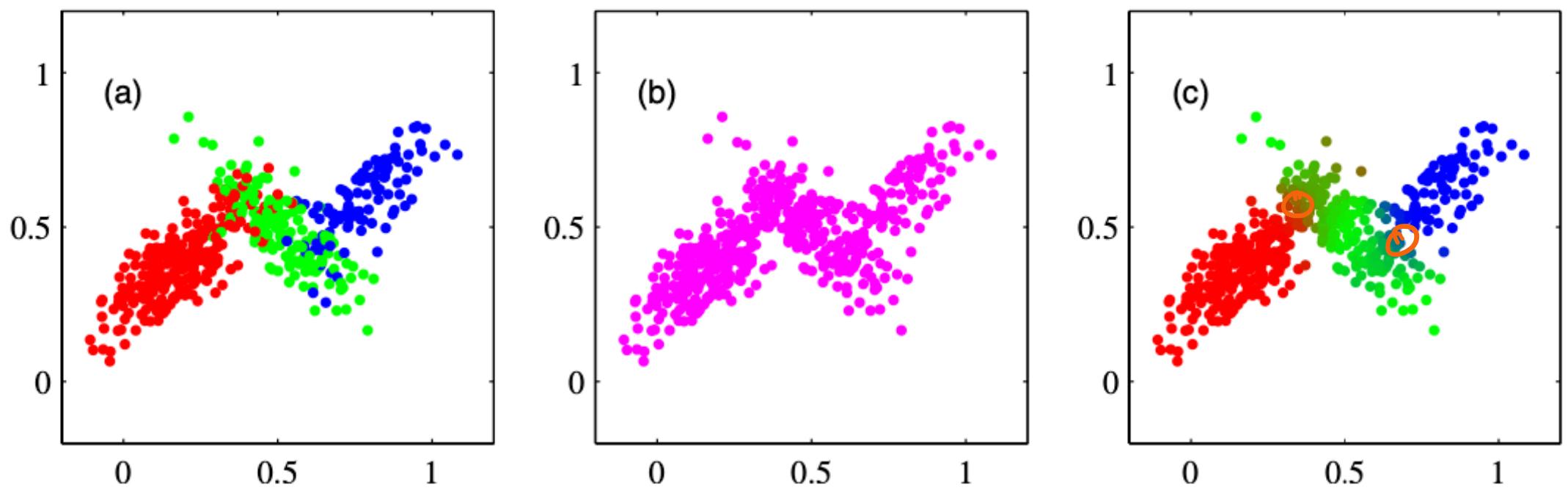


Figure 9.5 Example of 500 points drawn from the mixture of 3 Gaussians shown in Figure 2.23. (a) Samples from the joint distribution $p(z)p(x|z)$ in which the three states of z , corresponding to the three components of the mixture, are depicted in red, green, and blue, and (b) the corresponding samples from the marginal distribution $p(x)$, which is obtained by simply ignoring the values of z and just plotting the x values. The data set in (a) is said to be *complete*, whereas that in (b) is *incomplete*. (c) The same samples in which the colours represent the value of the responsibilities $\gamma(z_{nk})$ associated with data point x_n , obtained by plotting the corresponding point using proportions of red, blue, and green ink given by $\gamma(z_{nk})$ for $k = 1, 2, 3$, respectively

Likelihood of Data - MoG

Sunday, September 29, 2024

10:14 PM

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

variable params

Data likelihood

$$\mathbf{X} = \underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$$

$$p(\mathbf{X}) = \prod_{n=1}^N p(\underline{x}_n) \quad (\text{IID assumption})$$

Data log-likelihood

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

Notation: $f(x|\mu)$, we say "density of x given μ "

Relating k-means & MoG

Sunday, September 29, 2024

10:12 PM

$$\text{k-means error: } E[X] = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

$$\text{MoG log-likelihood: } \ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right\}$$

Assume that for \underline{x}_n almost all probab.: likely is accounted for by one component, denoted by r_{nk}

$$\begin{aligned} \log P(X) &= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log(\pi_k N(x_n | \mu_k, \Sigma_k)) \\ &= \sum_{n=1}^N \sum_{k=1}^K r_{nk} [\log \pi_k + \log N(x_n | \mu_k, \Sigma_k)] \\ &= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \left[\log \pi_k - \frac{1}{2} \log(2\pi \sigma_k^2) - \frac{1}{2} \frac{(x_n - \mu_k)^2}{\sigma_k^2} \right] \end{aligned}$$

One-hot encoding

Scalar case -

Wrap Up

Sunday, September 29, 2024

10:22 PM

- ▢ k-means clustering
- ▢ Mixture of Gaussians (MoG)
- ▢ k-means as narrow case of MoG

Next class

- ▢ EM Algorithm for learning the MoG

Announcements

- ▢ Midterm, Oct 21, 9am - 11am
- ▢ Free days during midterm week
- ▷ Assignment 2 coming soon

Learning a MoG

Tuesday, October 1, 2024 9:56 PM

Recall: Prob of \underline{x} under cluster k : $N(\underline{x} | \underline{\mu}_k, \Sigma_k)$

Total prob of \underline{x} : $p(\underline{x}) = \sum_{k=1}^K \pi_k N(\underline{x} | \underline{\mu}_k, \Sigma_k)$

logP(\underline{x}): $\ln p(\underline{x} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\underline{x}_n | \underline{\mu}_k, \Sigma_k) \right\}$

MoG: Hard assignment Alg

Tuesday, October 1, 2024

10:15 PM

Hard assignment algorithm:

Like k-means, but assign r's using $N(\cdot)$ instead of $\|\cdot\|^2$

Initialize: $\pi_k = \frac{1}{K}$, $\Sigma_k = \text{COV DATA}$, $\mu_k = \text{random } \underline{x}'s$.

Iterate until convergence

For $n=1 \dots N$ set $r_{nk} = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{argmax}} \pi_j^T N(\underline{x}) \mu_j, \Sigma_j \\ 0 & \text{otherwise} \end{cases}$

For $k=1 \dots K$, set

$$\pi_k = \frac{1}{N} \sum_{n=1}^N r_{nk} \quad \mu_k = \frac{\sum_{n=1}^N r_{nk} \underline{x}_n}{\sum_{n=1}^N r_{nk}}$$

$$\Sigma_k = \frac{1}{\sum_{n=1}^N r_{nk}} \sum_{n=1}^N r_{nk} (\underline{x}_n - \mu_k) (\underline{x}_n - \mu_k)^T$$

Vector
Transpose

Sample covariance

"Soft assignment" algorithm

Initialize $\pi_k = \frac{1}{K}$, $\Sigma_k = \text{COV DATA}$, $\mu_k = \text{random } \underline{x}'s$.

Iterate until convergence:

E-Step For $n=1 \dots N$, for $k=1 \dots K$ set $r_{nk} = \frac{\pi_k N(\underline{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(\underline{x}_n | \mu_j, \Sigma_j)}$

Note: $\sum_k r_{nk} = 1$. Responsibilities

M-Step

For $k=1 \dots K$, set

$$\pi_k = \frac{1}{N} \sum_{n=1}^N r_{nk} \quad \mu_k = \frac{\sum_{n=1}^N r_{nk} \underline{x}_n}{\sum_{n=1}^N r_{nk}}$$

$$\Sigma_k = \frac{1}{\sum_{n=1}^N r_{nk}} \sum_{n=1}^N r_{nk} (\underline{x}_n - \mu_k) (\underline{x}_n - \mu_k)^T$$

Vector Transpose

Sample covariance

Transforming Data

Wednesday, October 2, 2024 9:55 AM



Example of EM

Tuesday, October 1, 2024

10:29 PM

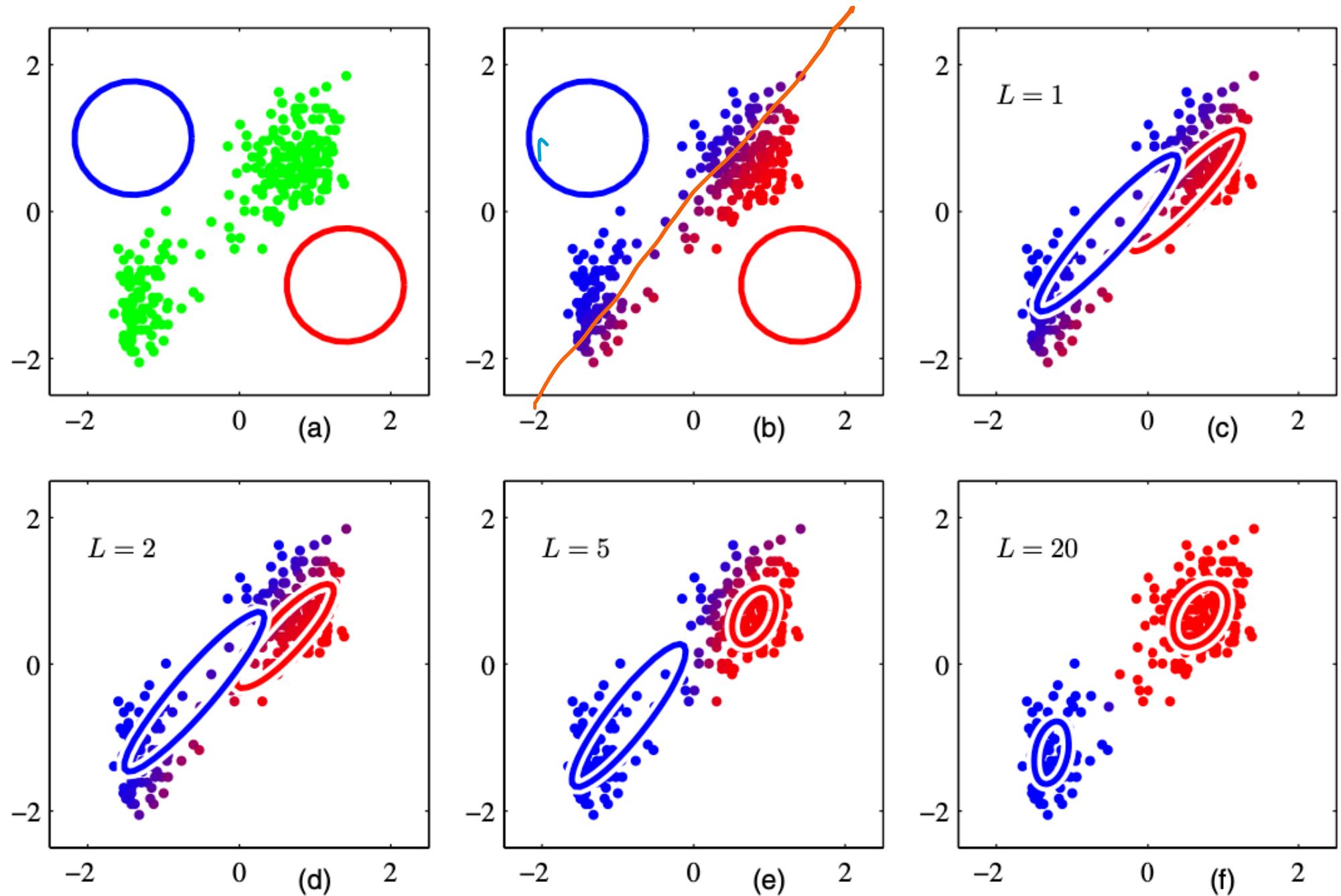
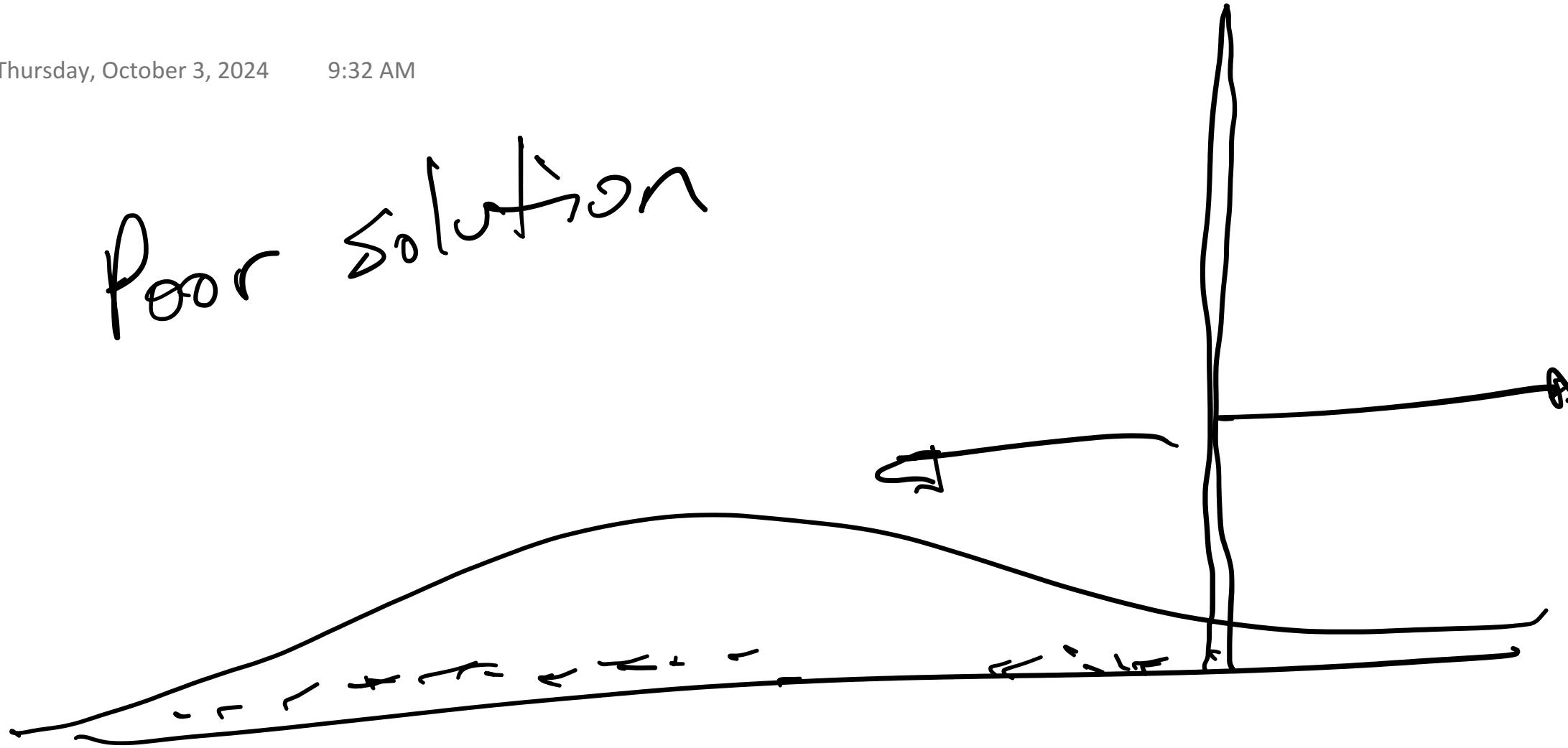


Figure 9.8 Illustration of the EM algorithm using the Old Faithful set as used for the illustration of the K -means algorithm in Figure 9.1. See the text for details.

Poor Solution



within the Alg:

If $\sum_{n=1}^N r_{nk} < 5$ (must have at least 5 dp)

Re-initialize Gaussian k.

Wrap Up

Tuesday, October 1, 2024

10:42 PM

- ▢ k-means clustering
- ▢ Mixture of Gaussians (MoG)
- ▢ k-means as narrow case of MoG
- ▢ EM Algorithm for learning the MoG

Announcements

- ▢ Midterm, Oct 21, 9am - 11am
- ▢ Free days during mid-term week
- ▢ Assignment 2 coming soon

In today's news!!!

Wednesday, October 9, 2024

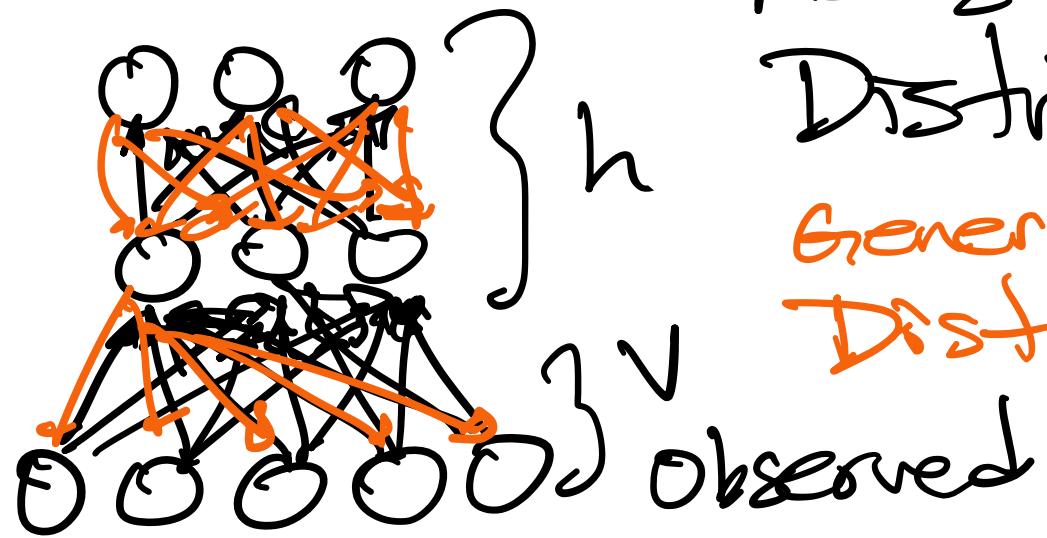
2:53 PM

- Geoff Hinton & John Hopfield win Nobel Prize in Physics for their work on neural networks
 - Very relevant to this week's lectures!
 - Big deal for UofT!
- David Baker, Demis Hassabis & John Jumper win Nobel Prize in Chemistry for AlphaFold, the deep learning system for protein structure prediction

My work with Hinton

Wednesday, October 9, 2024 9:18 AM

The Helmholtz Machine



Recognition net:

Distribution over $h \sim Q(h)$

Generative net

Distribution over $h, v \sim P(h, v)$

Physics, free energy, $F = \langle \text{Energy} \rangle - \text{Entropy}$

Energy = $-\log P(h, v)$, Entropy = $-\sum_h Q(h) \log Q(h)$

$$F = -\sum_h Q(h) \log P(h, v) + \sum_h Q(h) \log Q(h)$$

$$= \sum_h Q(h) \log \frac{Q(h)}{P(h, v)}$$

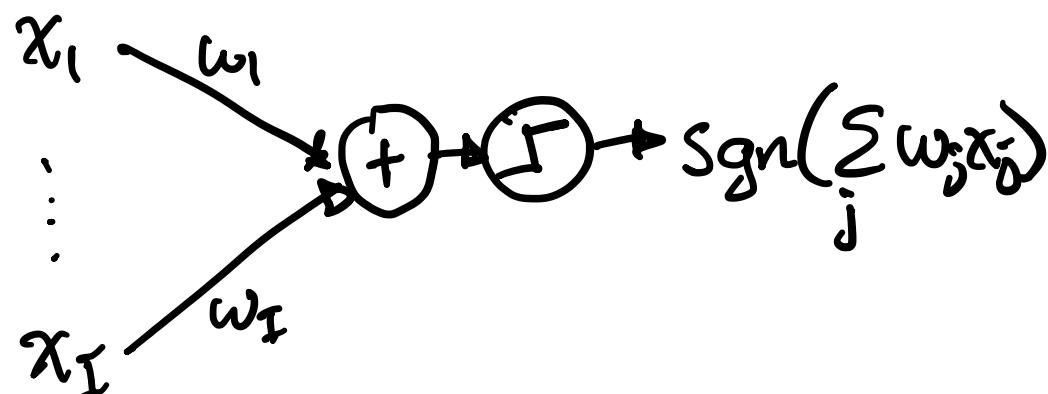
$$\therefore F = \text{KL}\{Q, P\} - \log P(v)$$

$$F \geq -\log P(v)$$

Deep Learning

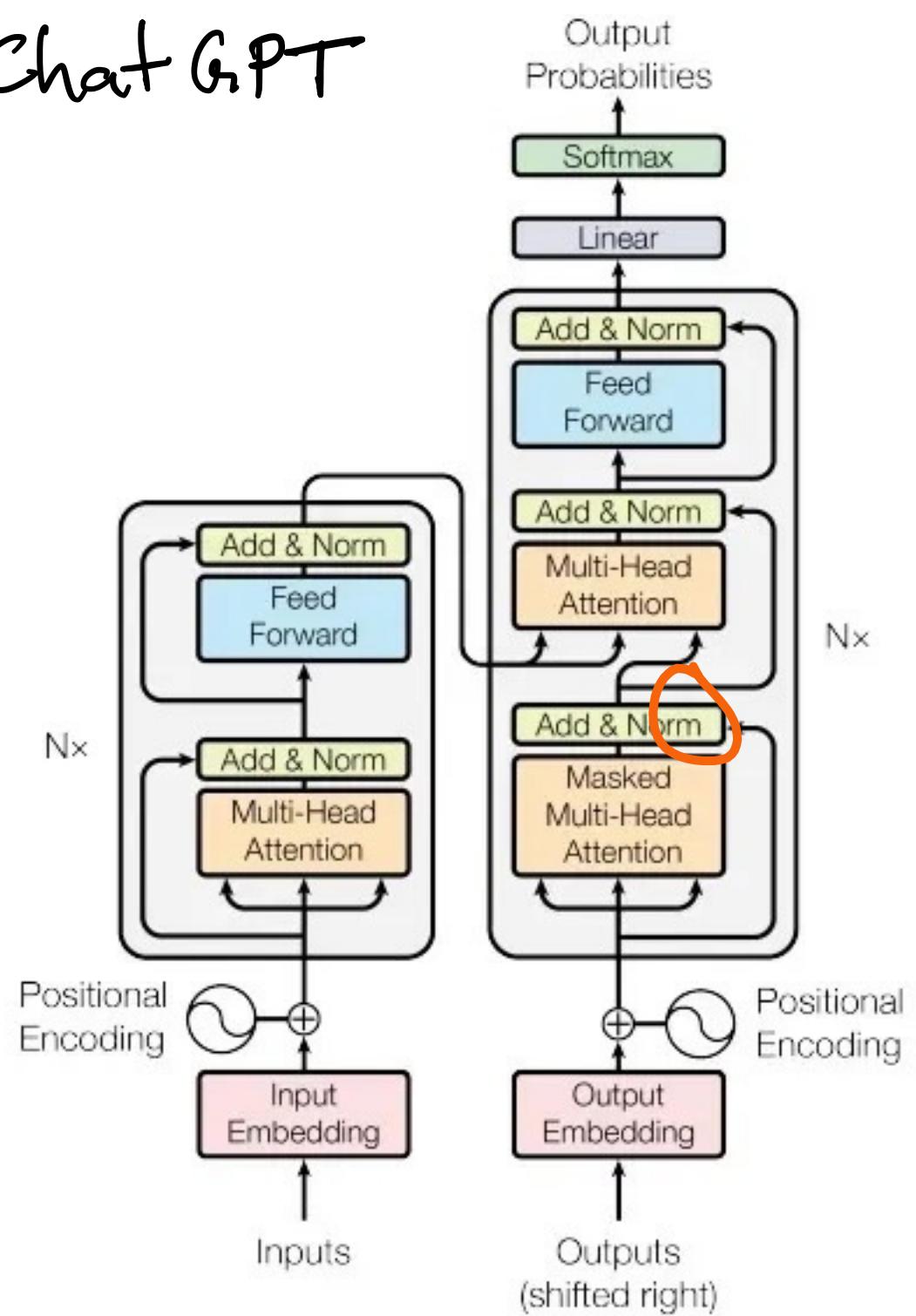
Saturday, October 5, 2024 1:52 PM

A) Perception



How do we go
from A to B?

B) Chat GPT



Recall Perception

Saturday, October 5, 2024

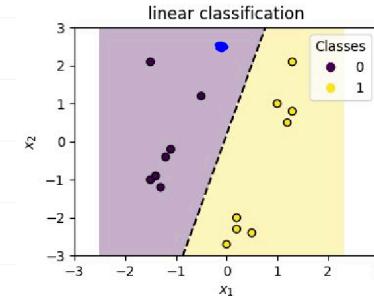
1:59 PM

Determine a classification rule

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

to minimize classification error

We will see Perceptron Learning algorithm,
logistic regression and Gradient descent



Perceptron Learning Algorithm (PLA)

Input: training set D that is linearly separable

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

Initialization: choose arbitrary \underline{w} , e.g., $\underline{w} = \underline{0}$

Step 1: check if $E_{in}(\underline{w}) = 0$. If yes, stop and return \underline{w} .

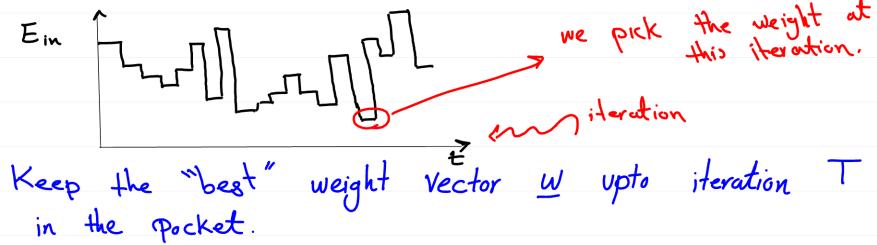
Step 2: Let (\underline{x}_n, y_n) be a miss-classified point,
i.e., $y_n \neq \hat{y}_n$ (including the points on the boundary)

$$\text{If } y_n = +1, \underline{w} \leftarrow \underline{w} + \underline{x}_n \quad \text{If } y_n = -1, \underline{w} \leftarrow \underline{w} - \underline{x}_n$$

Go to Step 1.

Pocket Algorithm

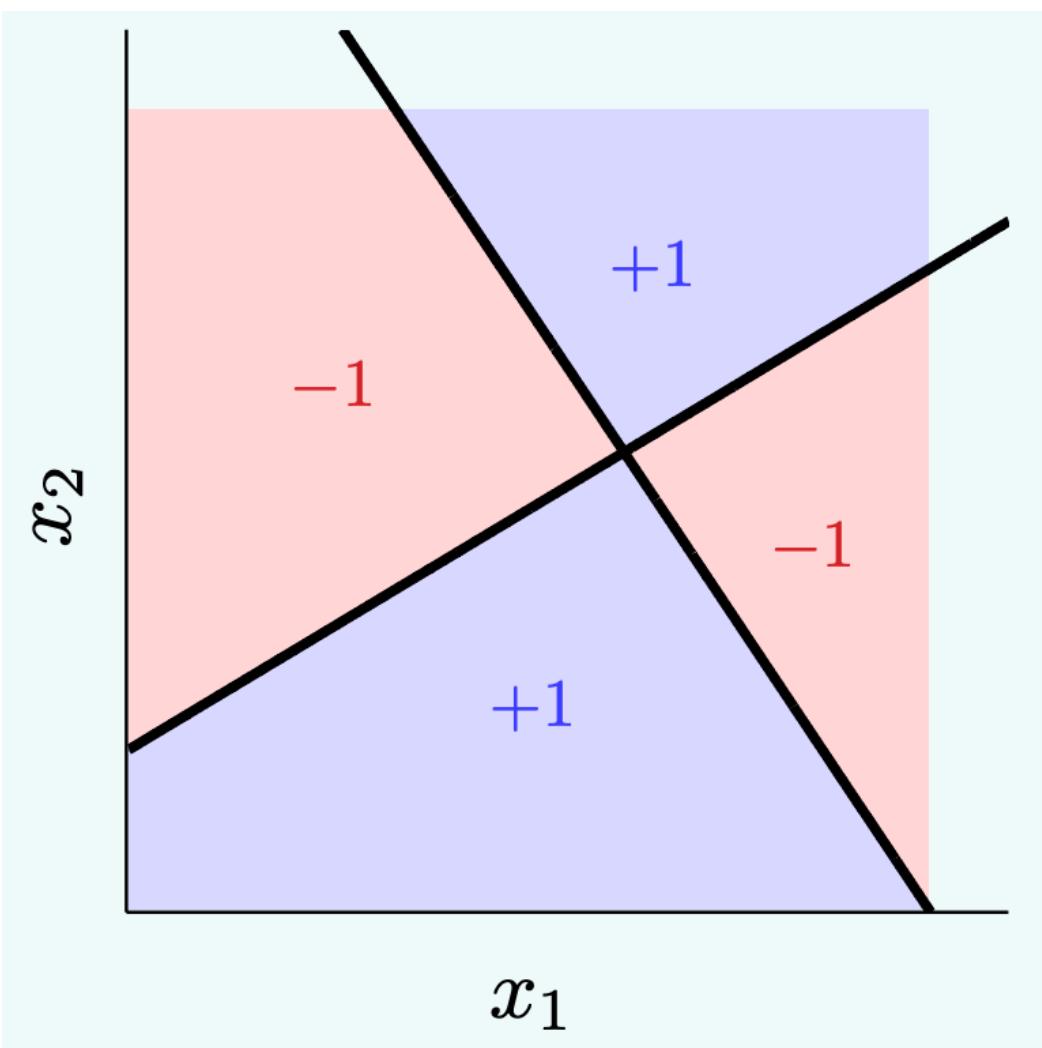
■ Pocket Algorithm extends PLA for dataset that are not linearly separable.



Limitation of Perceptron

Saturday, October 5, 2024

2:07 PM

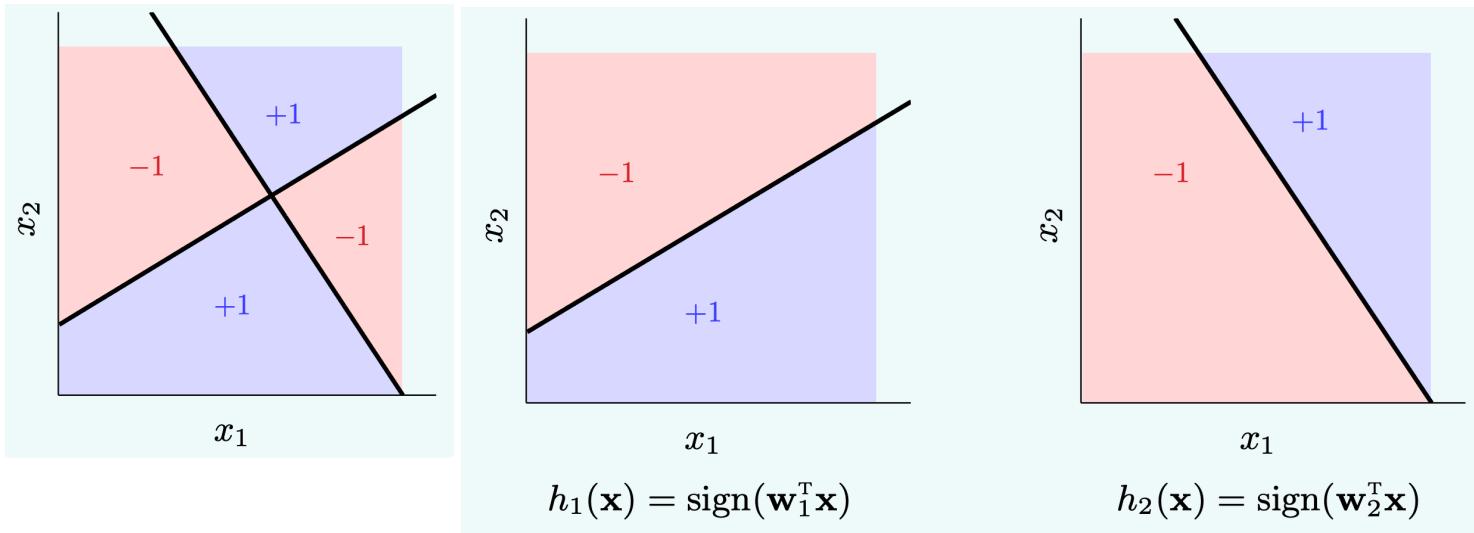


Problem!

Decomposing a complex problem

Saturday, October 5, 2024

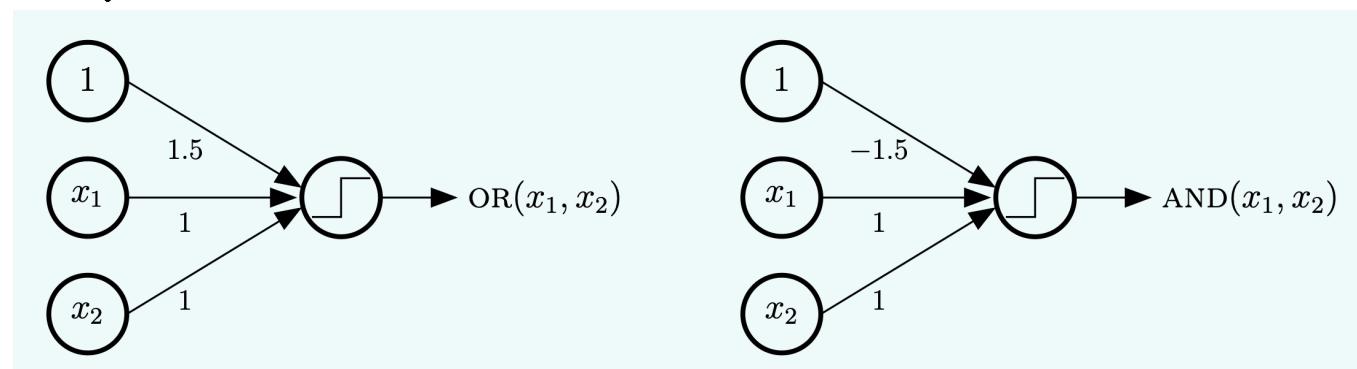
2:11 PM



$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$

AND NOT OR

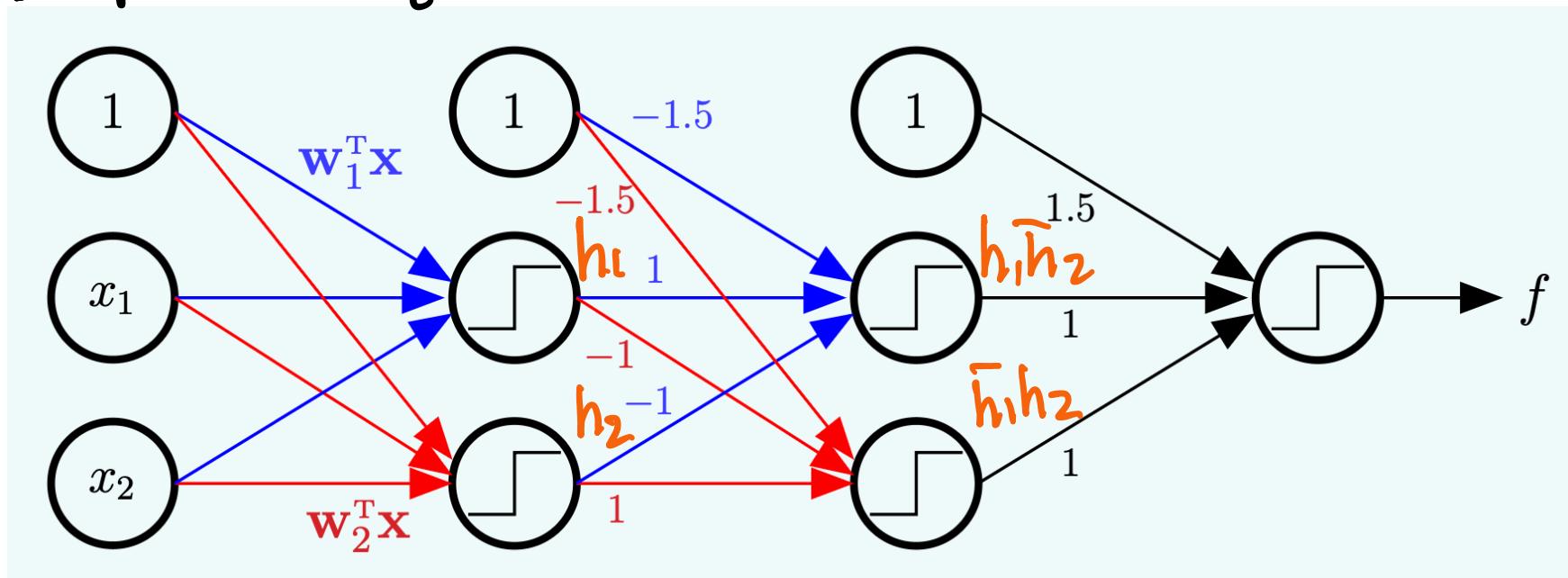
Implementing OR, AND



$x_1 \in \{-1, 1\}$
 $x_2 \in \{-1, 1\}$

False
True

Implementing $f = h_1 \overline{h_2} + \overline{h_1} h_2$ "Multilayer Perceptron"



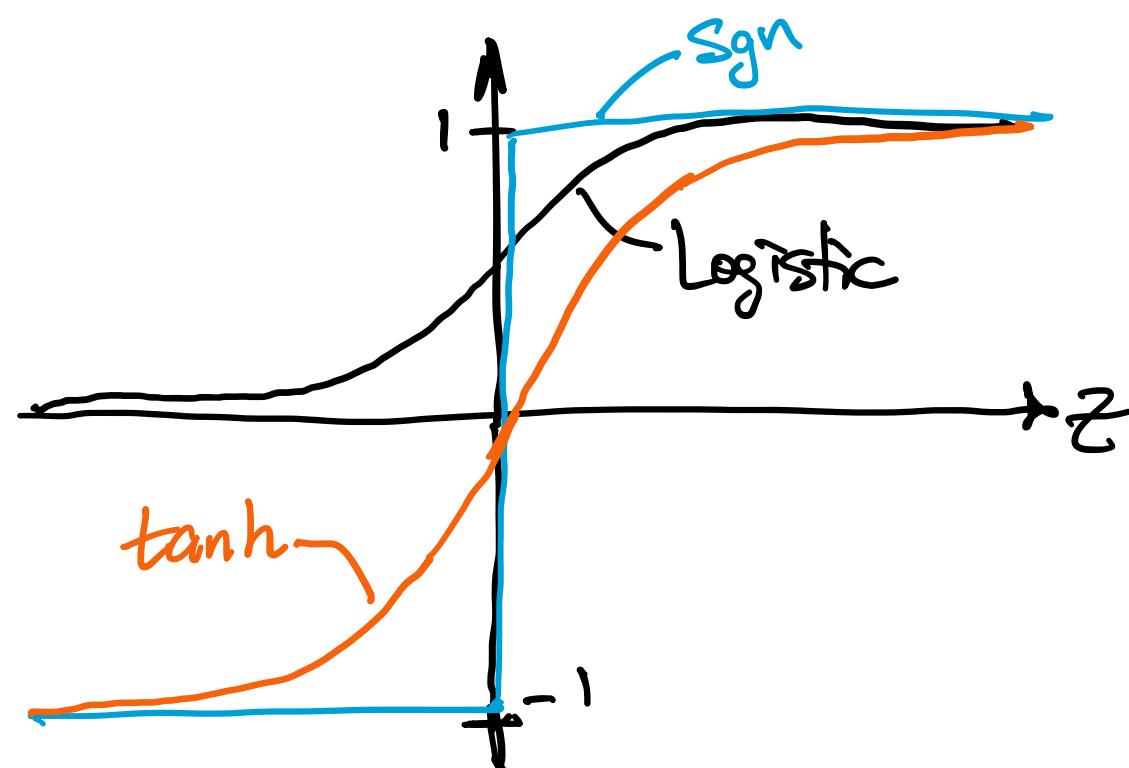
Differentiable activation

Saturday, October 5, 2024 3:22 PM

- Recall that for linear/logistic regression, k-means clustering, MoG, etc, we differentiated the objective function to derive a learning algorithm.
- $\text{sgn}(\underline{w}^T \underline{x})$ is not differentiable
- Introduce a continuous activation function

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

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

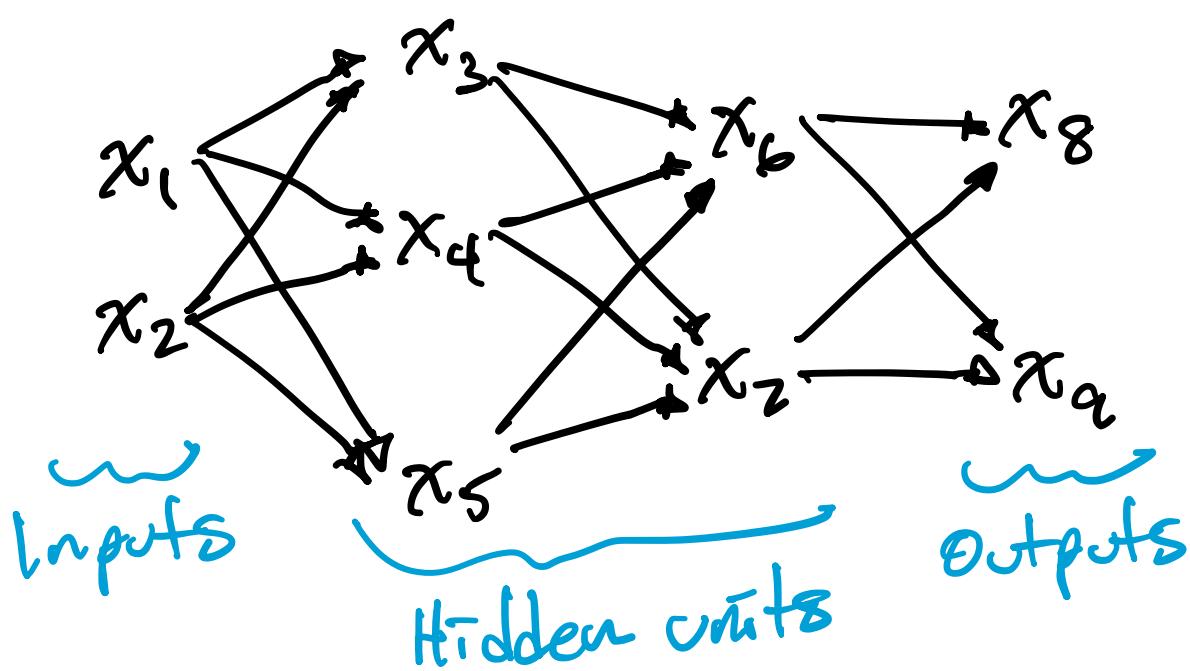


Neural network

Saturday, October 5, 2024

4:32 PM

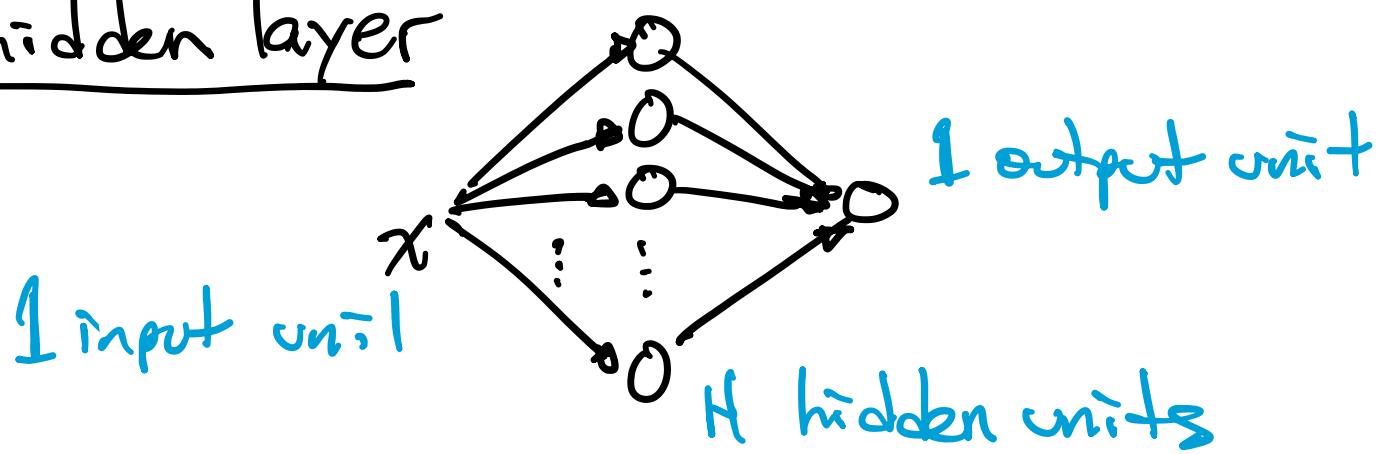
Example



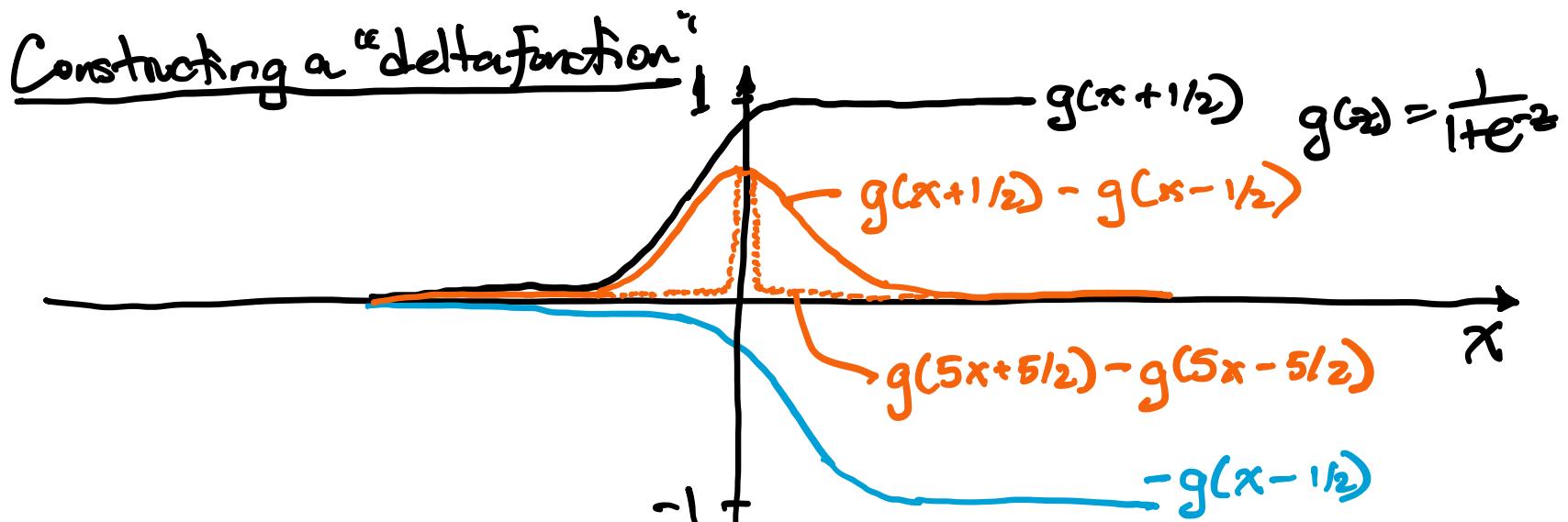
Arbitrary 1-D fn approximation

Saturday, October 5, 2024 3:40 PM

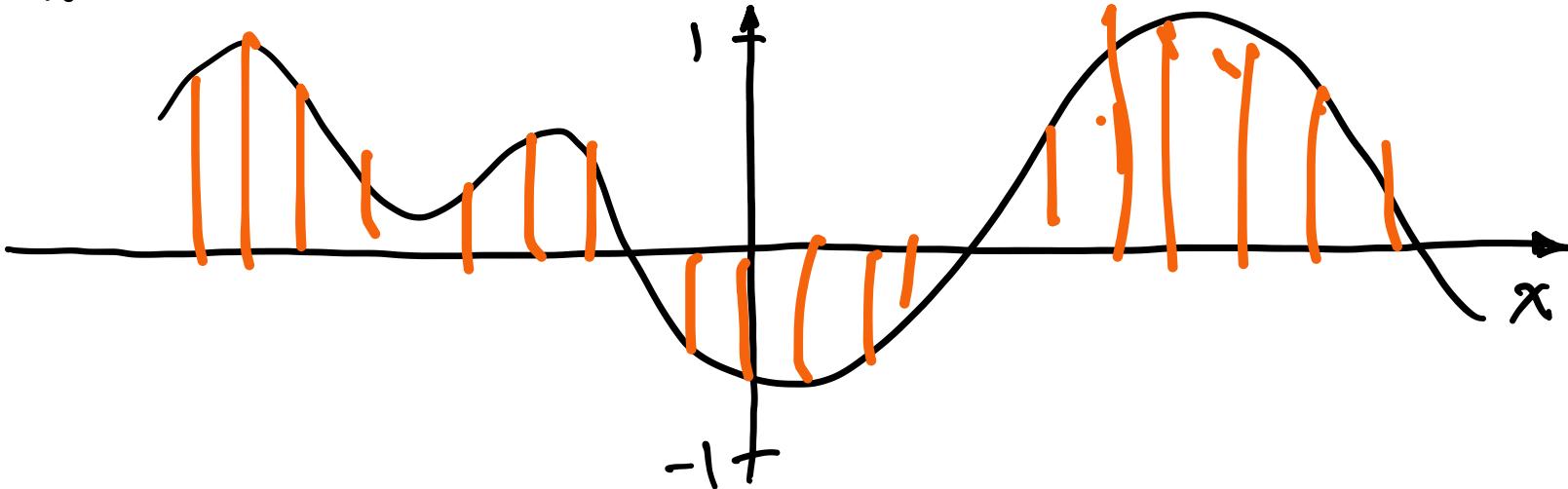
One hidden layer



1 output unit



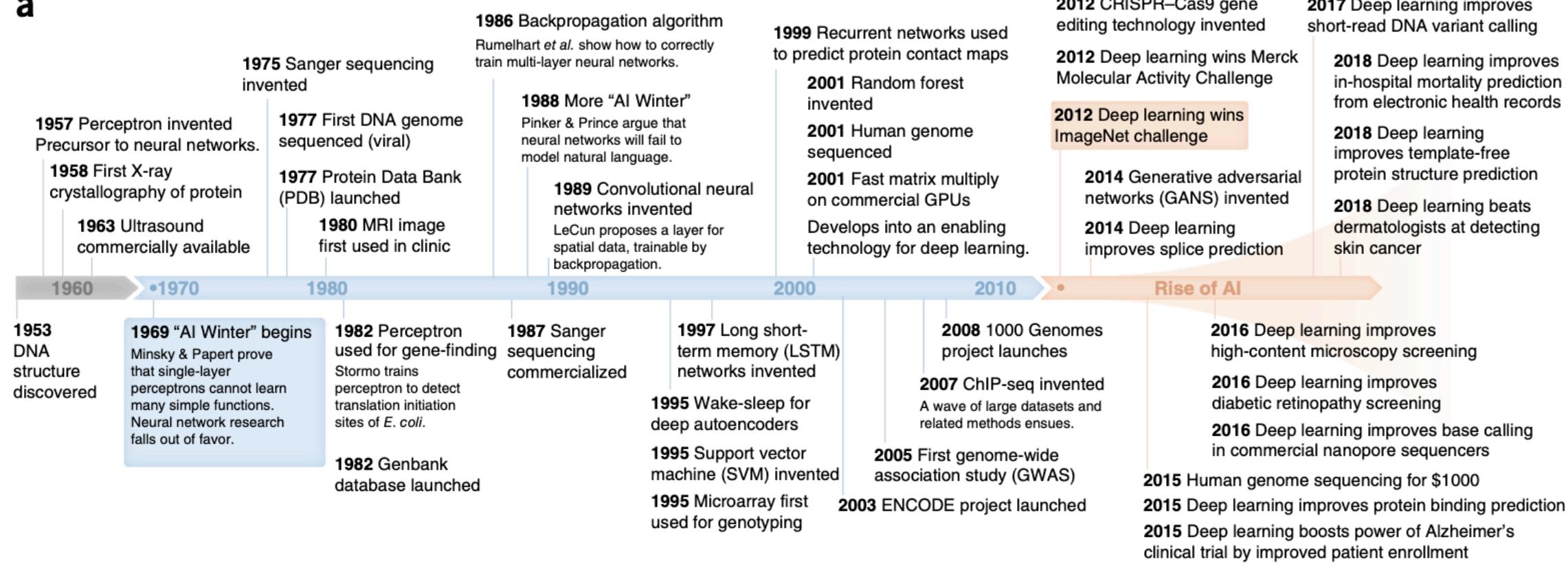
Approximating any fn



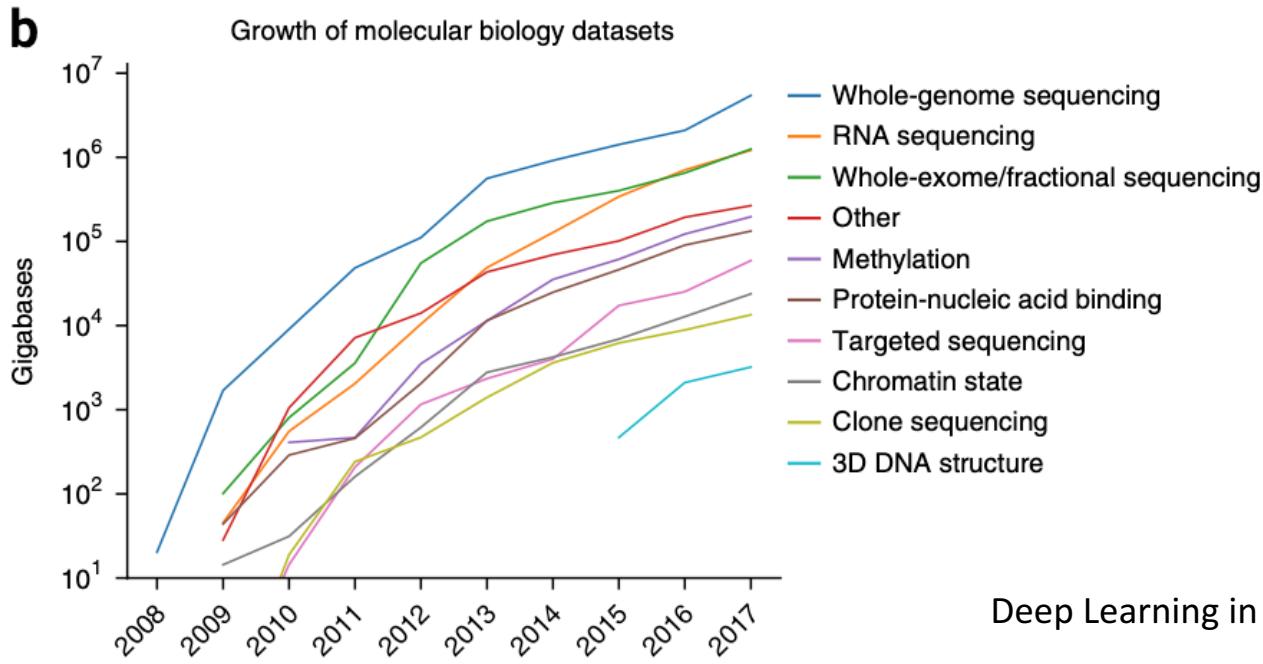
Deep Learning in genomics for medicine

Sunday, October 6, 2024 9:59 PM

a



b



Deep Learning in Biomedicine, Nature Biotechnology, 2018

Detecting protein binding sites in DNA

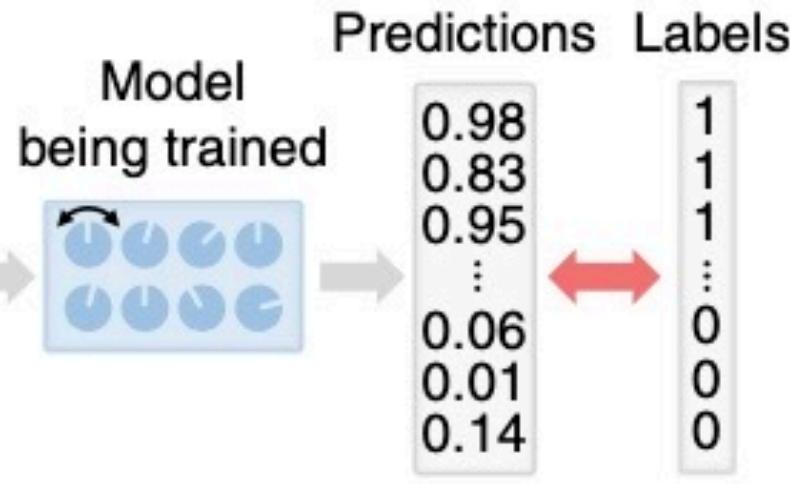
Sunday, October 6, 2024

10:13 PM

a

Training sequences

```
ACGCCCTACACGTGCTGCAATT  
GGCACGTGGTTCTTTCACGGAA  
CCAGTCATGCACACATGTGATGT  
⋮  
GTACCGGATCAGTCCGTTAGGT  
GAAATTAAATGGATAGGGGGATG  
CGATCGTTGTTACGATTAGAAA
```



Tune weights to minimize errors

c

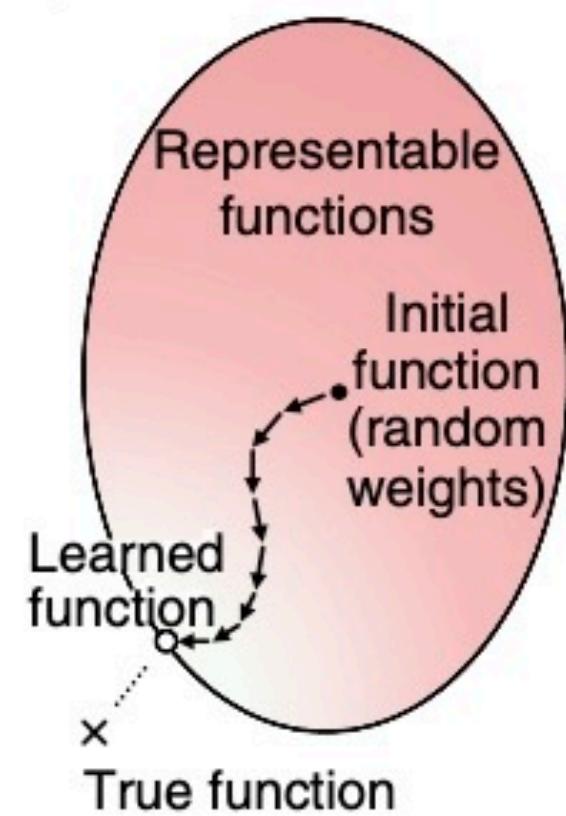
Novel sequences

```
ACGACTACACGTGCTGCAATT  
GGCACGTGGTTCTCACGGAA  
CCAGTCATGCACACTGGATGT
```

Trained model



b



Deep Learning in Biomedicine, Nature Biotechnology, 2018

Seq2Seq DNA model

Sunday, October 6, 2024

10:11 PM

d

Convolution on a biological sequence

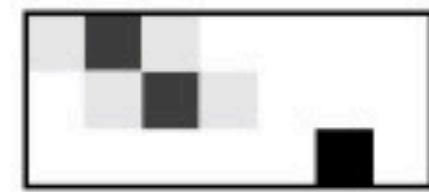
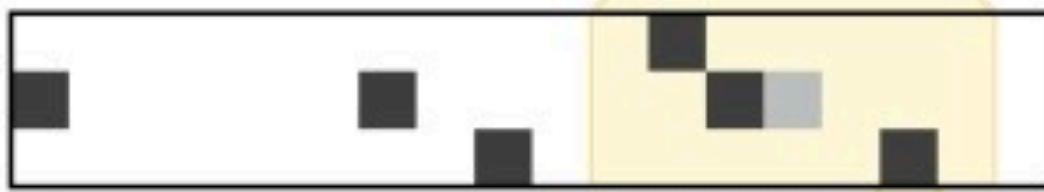
Output track



Pattern detected

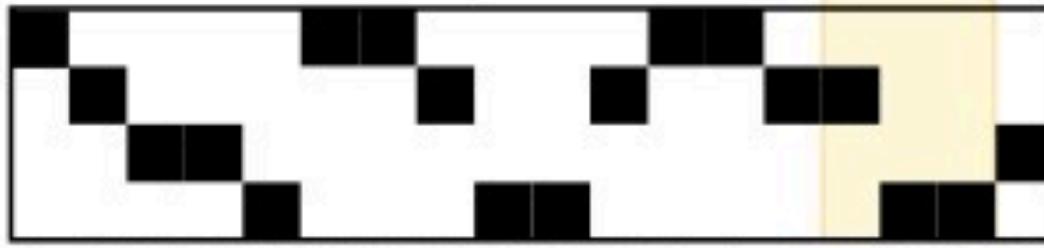
W^2

2nd layer weights
(1 filter of size 7×3)

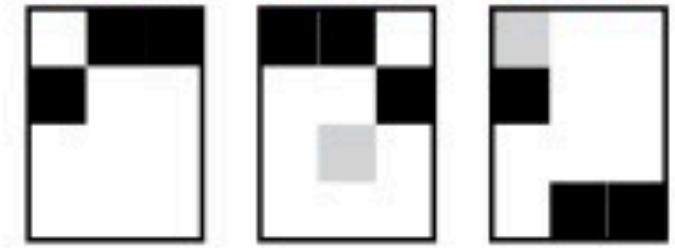


Feature map

W^1



CAA AAC CTT



Input sequence

A C G G T A A C T T C A A C C T T G

1st layer weights
(3 filters of size 3×4)

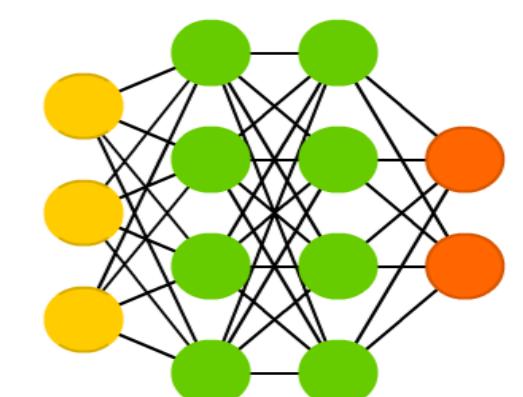
A mostly complete chart of

Neural Networks

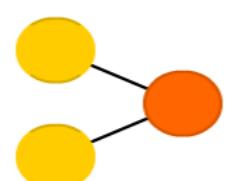
©2016 Fjodor van Veen - asimovinstitute.org

- (○) Backfed Input Cell
- (○) Input Cell
- (△) Noisy Input Cell
- (●) Hidden Cell
- (○) Probabilistic Hidden Cell
- (△) Spiking Hidden Cell
- (●) Output Cell
- (○) Match Input Output Cell
- (●) Recurrent Cell
- (○) Memory Cell
- (△) Different Memory Cell
- (●) Kernel
- (○) Convolution or Pool

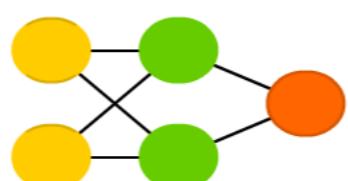
Deep Feed Forward (DFF)



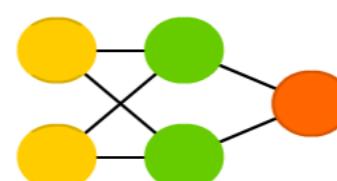
Perceptron (P)



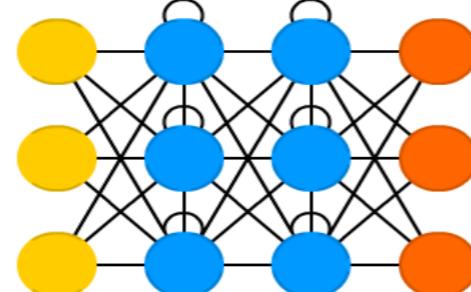
Feed Forward (FF)



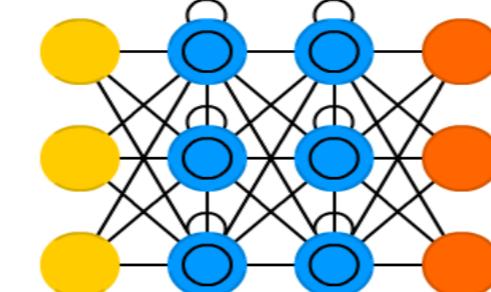
Radial Basis Network (RBF)



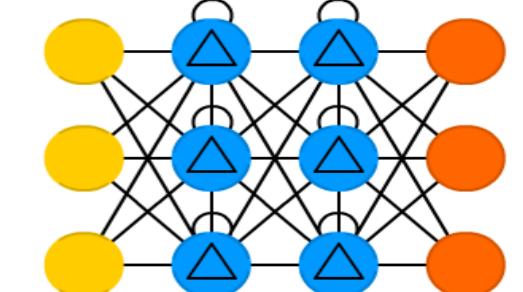
Recurrent Neural Network (RNN)



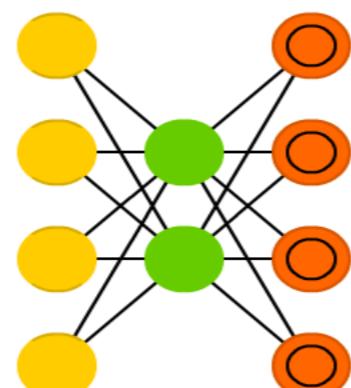
Long / Short Term Memory (LSTM)



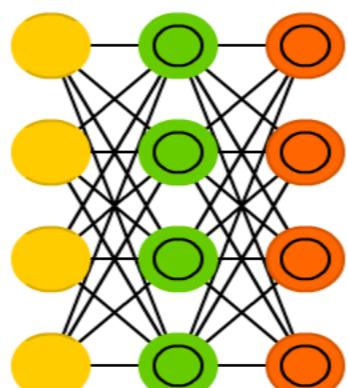
Gated Recurrent Unit (GRU)



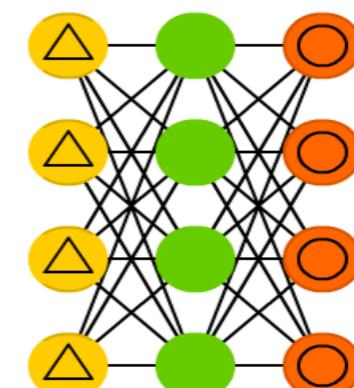
Auto Encoder (AE)



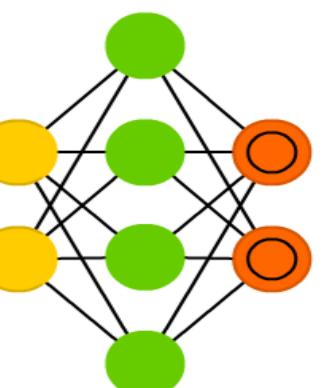
Variational AE (VAE)



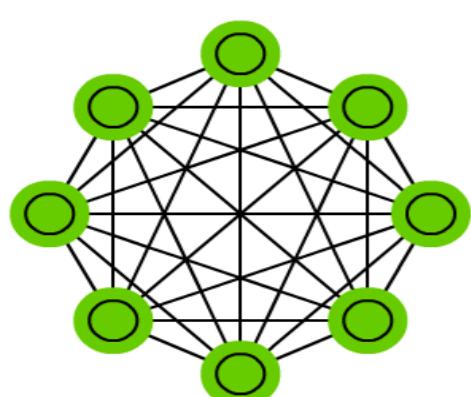
Denoising AE (DAE)



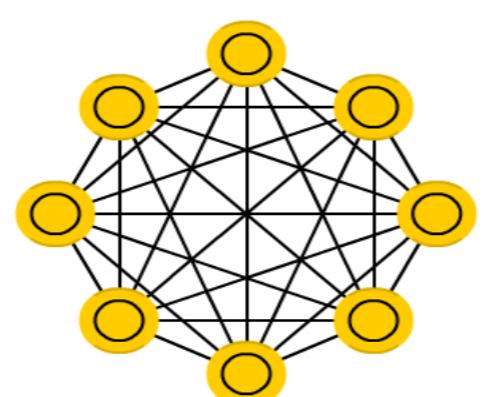
Sparse AE (SAE)



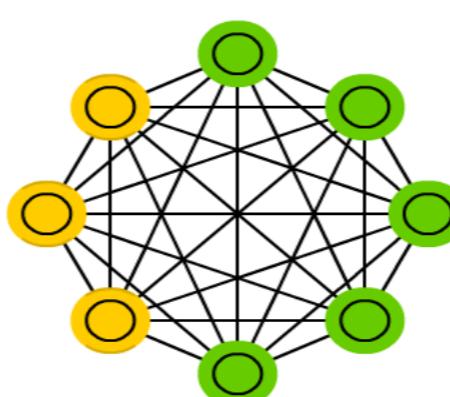
Markov Chain (MC)



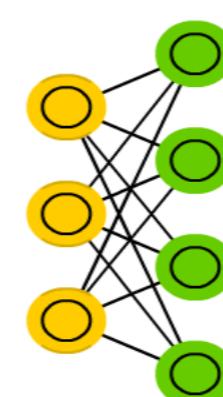
Hopfield Network (HN)



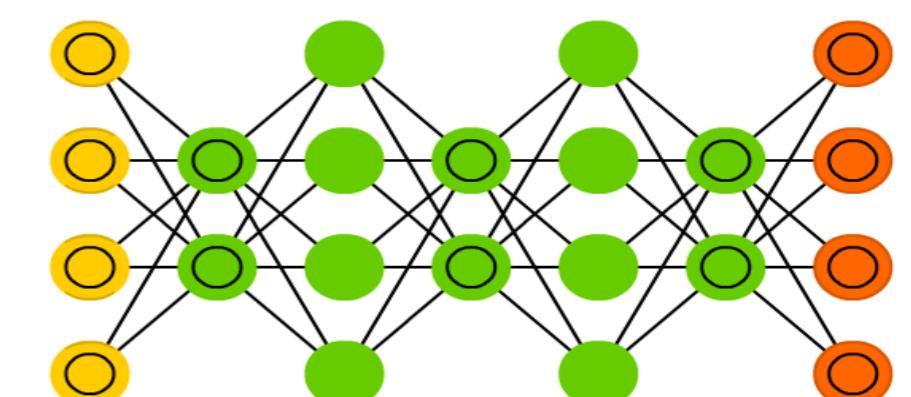
Boltzmann Machine (BM)



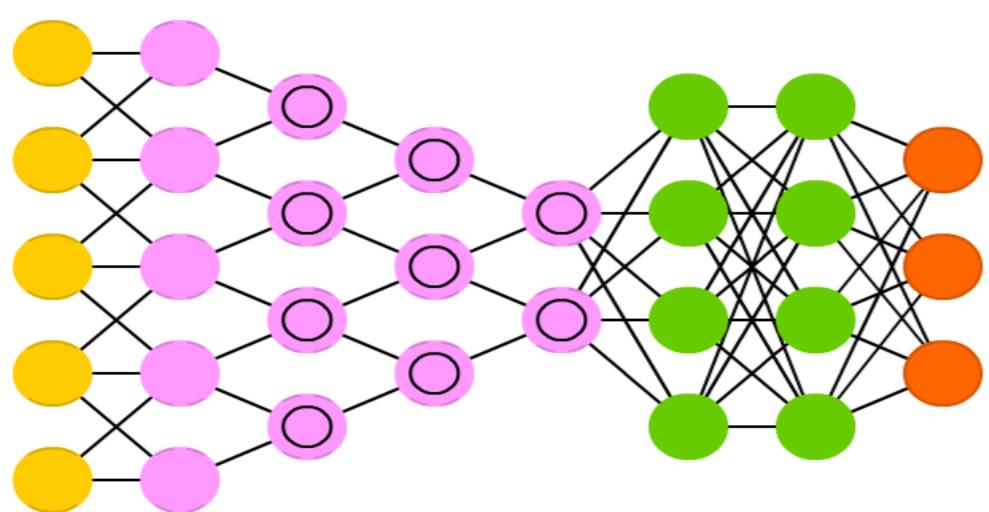
Restricted BM (RBM)



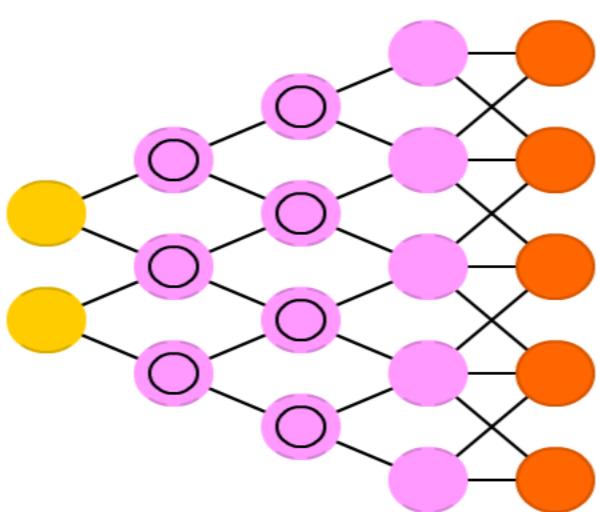
Deep Belief Network (DBN)



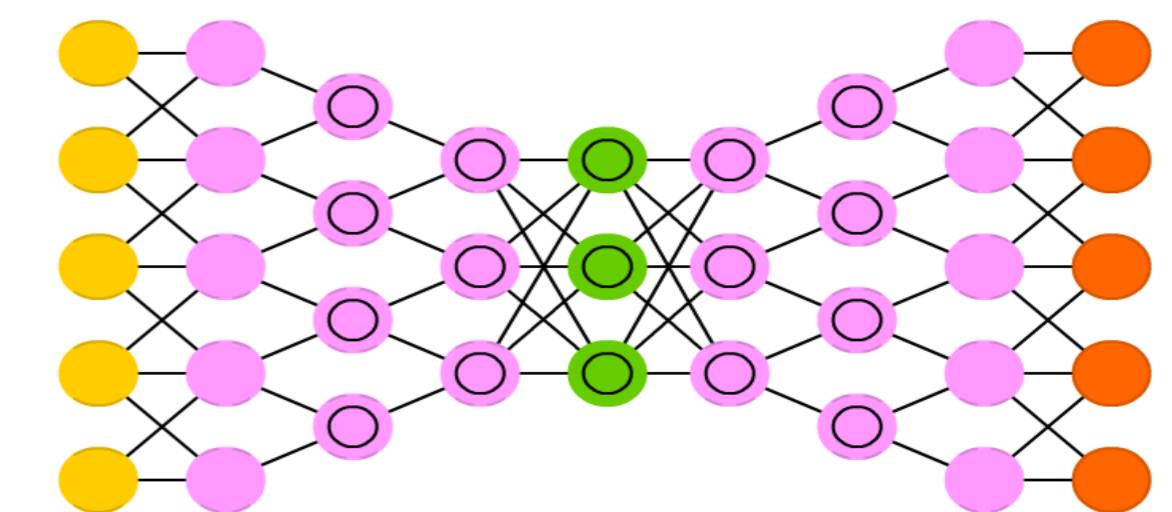
Deep Convolutional Network (DCN)



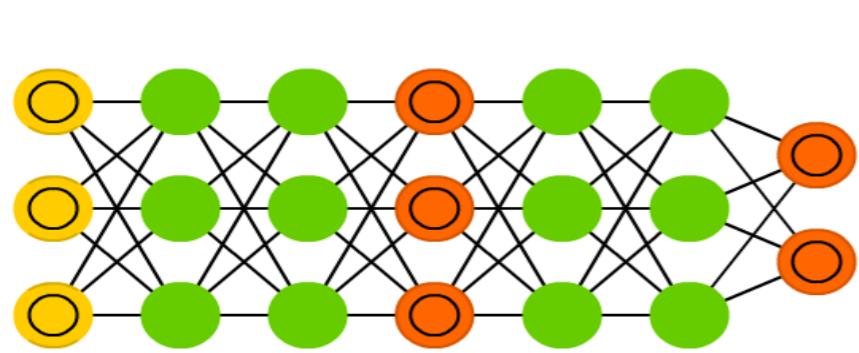
Deconvolutional Network (DN)



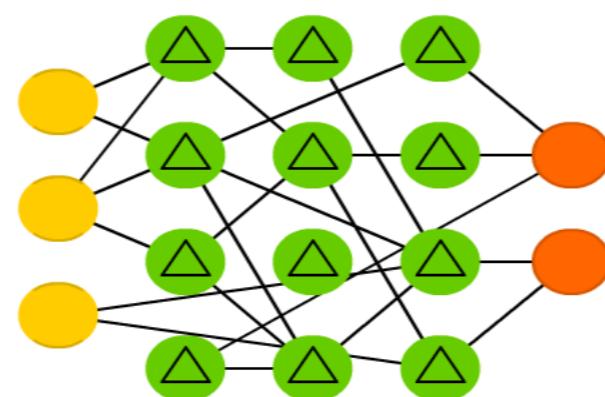
Deep Convolutional Inverse Graphics Network (DCIGN)



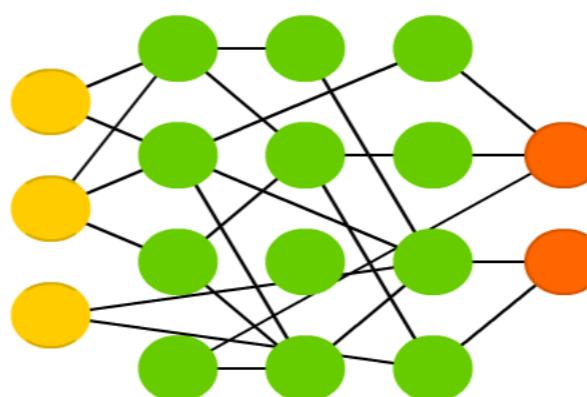
Generative Adversarial Network (GAN)



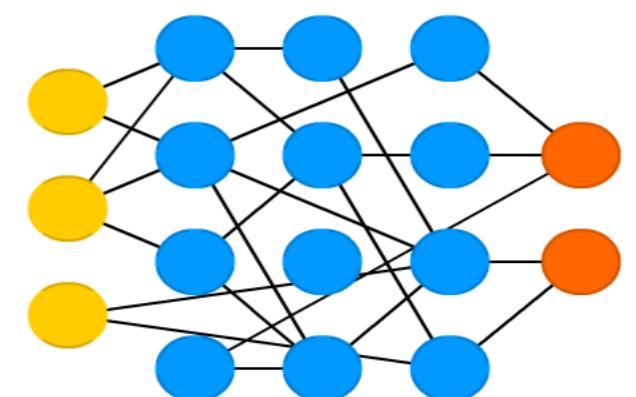
Liquid State Machine (LSM)



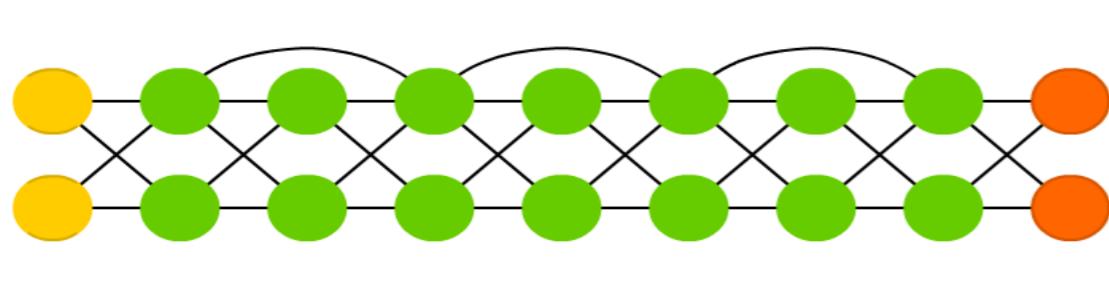
Extreme Learning Machine (ELM)



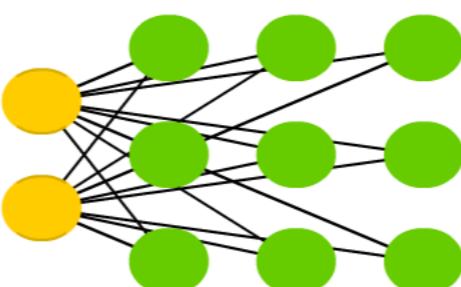
Echo State Network (ESN)



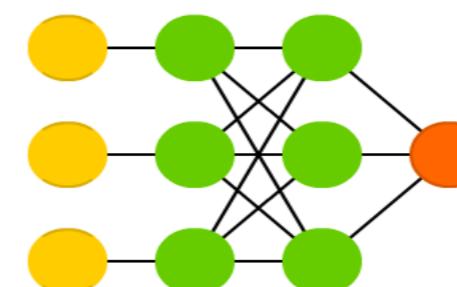
Deep Residual Network (DRN)



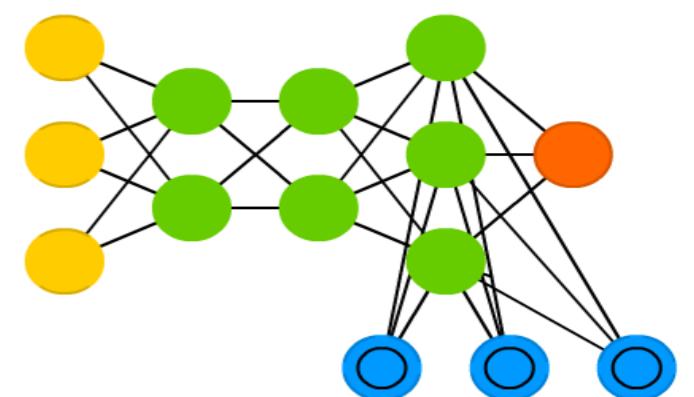
Kohonen Network (KN)



Support Vector Machine (SVM)



Neural Turing Machine (NTM)



Forward propagation

Saturday, October 5, 2024

4:21 PM

- I inputs, O outputs, H hidden units
- Total units $M = I + H + O$
- Index units so that connection $i \rightarrow j$ satisfies $i < j$
- Ordered activations: Standard form
 $x_1 \dots x_I, x_{I+1} \dots x_{I+H}, x_{M-O+1} \dots x_M$
- It is convenient to let x_i be the pre-activation value for unit i
- Combining the above, we have

$$\text{for } j > I, \quad x_j = \sum_{i=1}^{j-1} w_{ij} g_i(x_i)$$

If i is not connected to j , set $w_{ij} = 0$

Activation function for x_i

Forward propagation: Notes

Sunday, October 6, 2024

10:20 PM

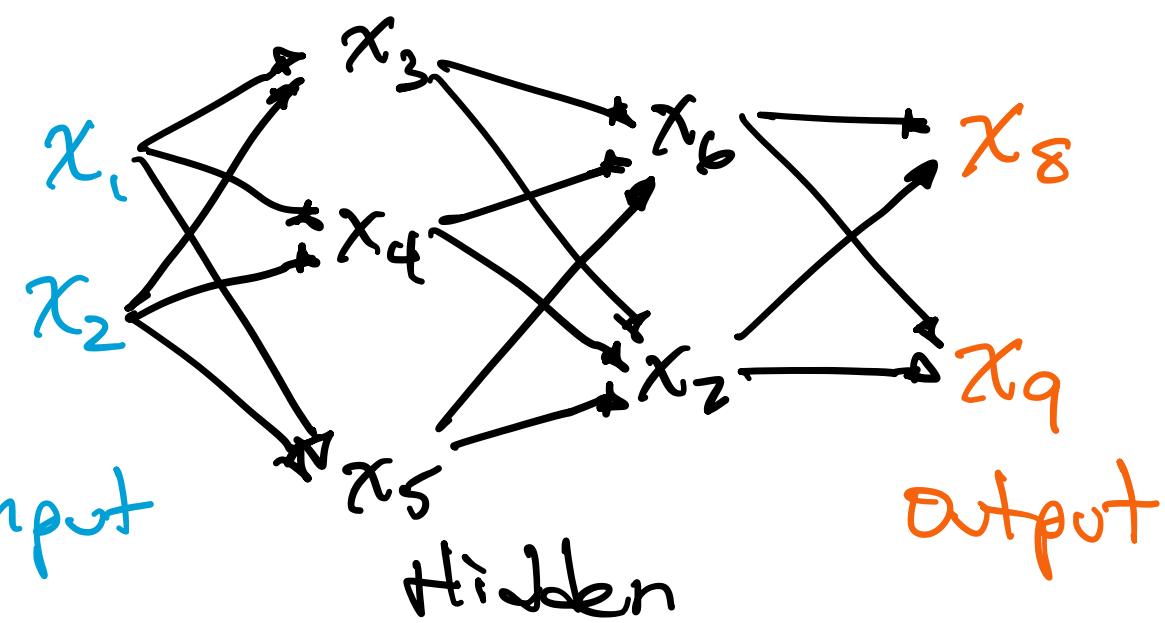
$$\text{for } j > I, \quad x_j = \sum_{i=1}^{j-1} w_{ij} g_i(x_i)$$

if i is not connected to j , set $w_{ij} = 0$

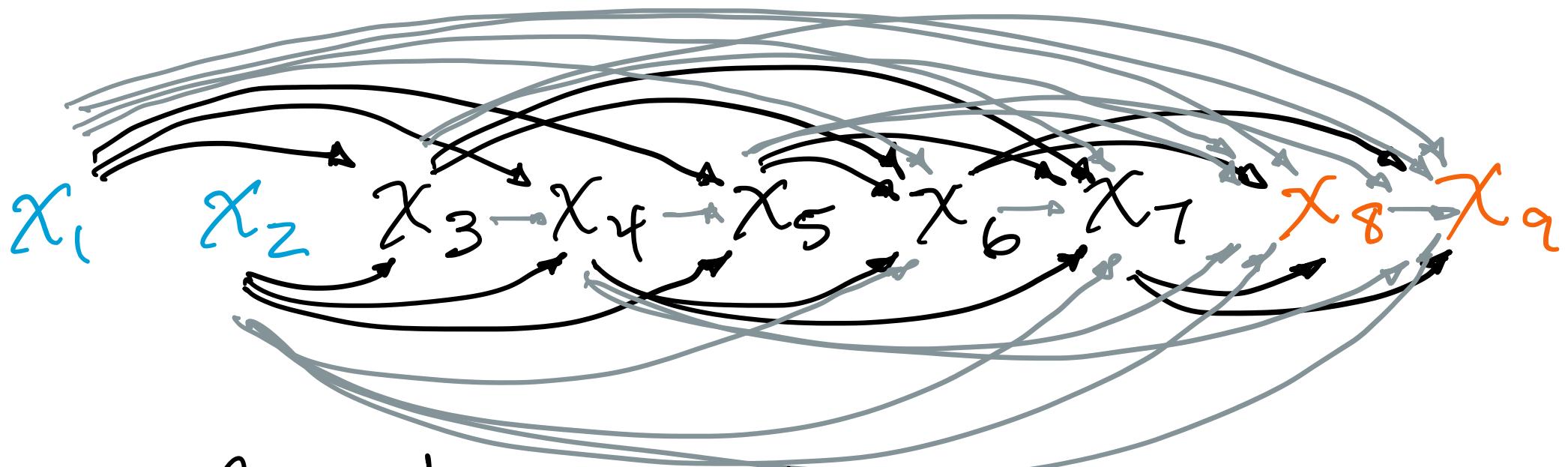
Activation function for x_i

- If input unit $x_i, i \in \{1 \dots I\}$ is not transformed, set $g_i(x_i) = x_i$, ie, $g_i()$ is identity,
- The outputs of the neural network are $g_{n-0+1}(x_{n-0+1}) \dots g_n(x_n)$

layered net



Standard form (fully connected DAG)



→ Connections with learnable w

→ Connections for which w clamped to 0

Summary

Sunday, October 6, 2024

9:56 PM

- Limitations of the perceptron
- Decomposing a complex problem
- Differentiable activation functions ($\text{sgn} \rightarrow \text{logistic, tanh}$)
- Neural Networks
 - Applications to genomics
 - Forward propagation
- Next: Backpropagation learning

Learning Neural Networks

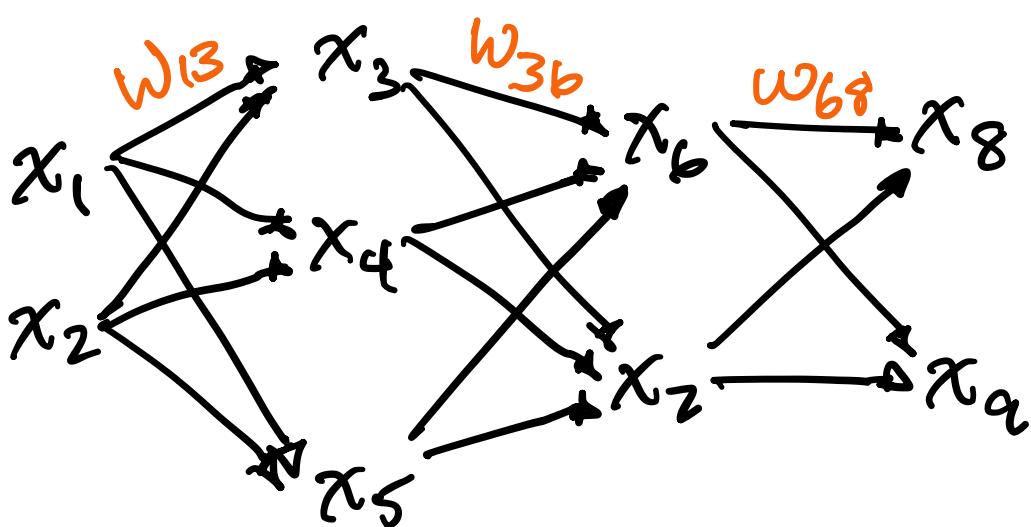
Saturday, October 5, 2024

4:48 PM

+1 or -1

Recall for perceptron: $\underline{w} \leftarrow \underline{w} + y \underline{x_n}$

How do we learn the weights $\{w_{ij}\}$?



Deriving $\frac{\partial E}{\partial w_{lm}}$

Tuesday, October 8, 2024

4:55 PM

Forward recursion

$$\text{We know that for } j = l+1 \dots M, x_j = \sum_{i=1}^{j-1} w_{ij} g_i(x_i) \quad \text{Eq1}$$

$\frac{\partial E}{\partial w_{lm}}$ for one training case: w_{lm} influences E only

$$\text{via } x_m, \text{ so } \frac{\partial E}{\partial w_{lm}} = \frac{\partial E}{\partial x_m} \frac{\partial x_m}{\partial w_{lm}}, \quad \frac{\partial E}{\partial w_{lm}} = \frac{\partial E}{\partial x_m} g_l(x_l) \quad \text{Apply Eq1}$$

If x_m is not an output unit, then x_m influences E via $x_{m+1} \dots x_M$, and we obtain a recursion:

$$\frac{\partial E}{\partial x_m} = \sum_{k=m+1}^M \frac{\partial E}{\partial x_k} \frac{\partial x_k}{\partial x_m}, \quad \frac{\partial E}{\partial x_m} = \sum_{k=m+1}^M \frac{\partial E}{\partial x_k} w_{mk} g'_m(x_m) \quad \text{Apply Eq1}$$

If x_m is an output unit, then

$$\frac{\partial E}{\partial x_m} = \frac{\partial (y_m - g_m(x_m))^2}{\partial x_m}, \quad \frac{\partial E}{\partial x_m} = -2(y_m - g_m(x_m))g'_m(x_m)$$

Backpropagation

Tuesday, October 8, 2024 5:32 PM

forward propagation: for $j = I+1 \dots M$, $x_j = \sum_{i=1}^{j-1} w_{ij} g_i(x_i)$

Output unit error derivatives:

$$\text{For } m = M \dots M-O+1 : \frac{\partial E}{\partial x_m} = -2(y_m - g_m(x_m))g'_m(x_m)$$

Back propagation:

$$\text{For } m = M-O \dots M-O-H : \frac{\partial E}{\partial x_m} = \sum_{k=m+1}^M \frac{\partial E}{\partial x_k} w_{mk} g'_k(x_m)$$

Weight derivatives:

$$\text{For } l = 1 \dots M-O+1, \text{ for } m = l+1 \dots M$$

learning rate

$$\frac{\partial E}{\partial w_{lm}} = \frac{\partial E}{\partial x_m} g_l(x_l), w_{lm} \leftarrow w_{lm} - \gamma \frac{\partial E}{\partial w_{lm}}$$

Relationship to perceptron learning

Wednesday, October 9, 2024

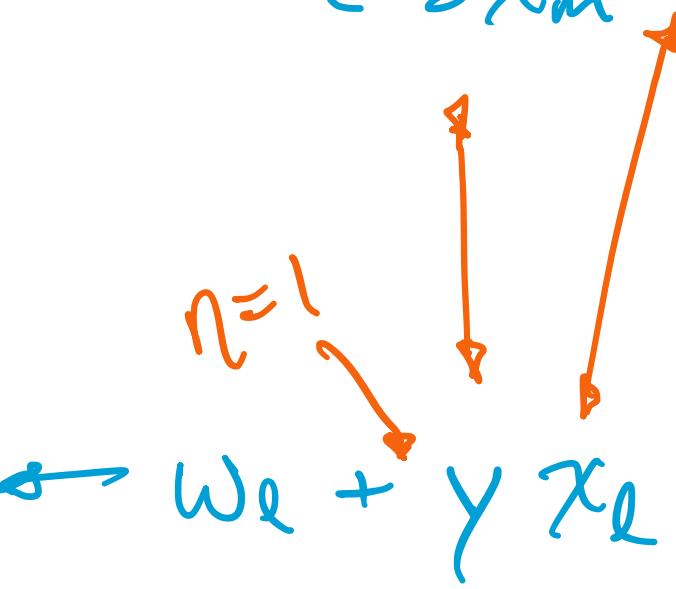
3:13 PM

Neural net:

$$\begin{array}{ccc} \vec{x}_l & \xrightarrow{w_{lm}} & \vec{x}_m \\ \text{Say } g(x_l) = x_l & & \end{array}$$
$$w_{lm} \leftarrow w_{lm} - \eta \frac{\partial E}{\partial x_m} x_l$$

Perceptron:

$$x_l \xrightarrow{w_e} \hat{y}$$

$$w_e \leftarrow w_e + \gamma x_l$$


Homework exercise

Wednesday, October 9, 2024

3:07 PM

forward propagation: for $j = I+1 \dots M$, $x_j = \sum_{i=1}^{j-1} w_{ij} g_i(x_i)$

Output unit error derivatives:

$$\text{For } m = M \dots M-D+1: \frac{\partial E}{\partial x_m} = -2(y_m - g_m(x_m))g'_m(x_m)$$

Backpropagation:

$$\text{For } m = M-D \dots M-D-H: \frac{\partial E}{\partial x_m} = \sum_{k=m+1}^M \frac{\partial E}{\partial x_k} w_{mk} g'_m(x_m)$$

Weight derivatives:

$$\text{For } l = I \dots M-D-H, \text{ for } m = l+1 \dots M$$
$$\frac{\partial E}{\partial w_{lm}} = \frac{\partial E}{\partial x_m} g_l(x_l), w_{lm} \leftarrow w_{lm} + \text{learning rate} \cdot \frac{\partial E}{\partial w_{lm}}$$

Rewrite these formulas in terms of the post activation values,
 $\hat{x}_i \triangleq g_i(x_i)$
 $i = 1 \dots M$.

Stochastic gradient descent

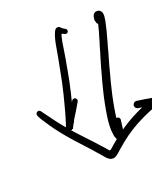
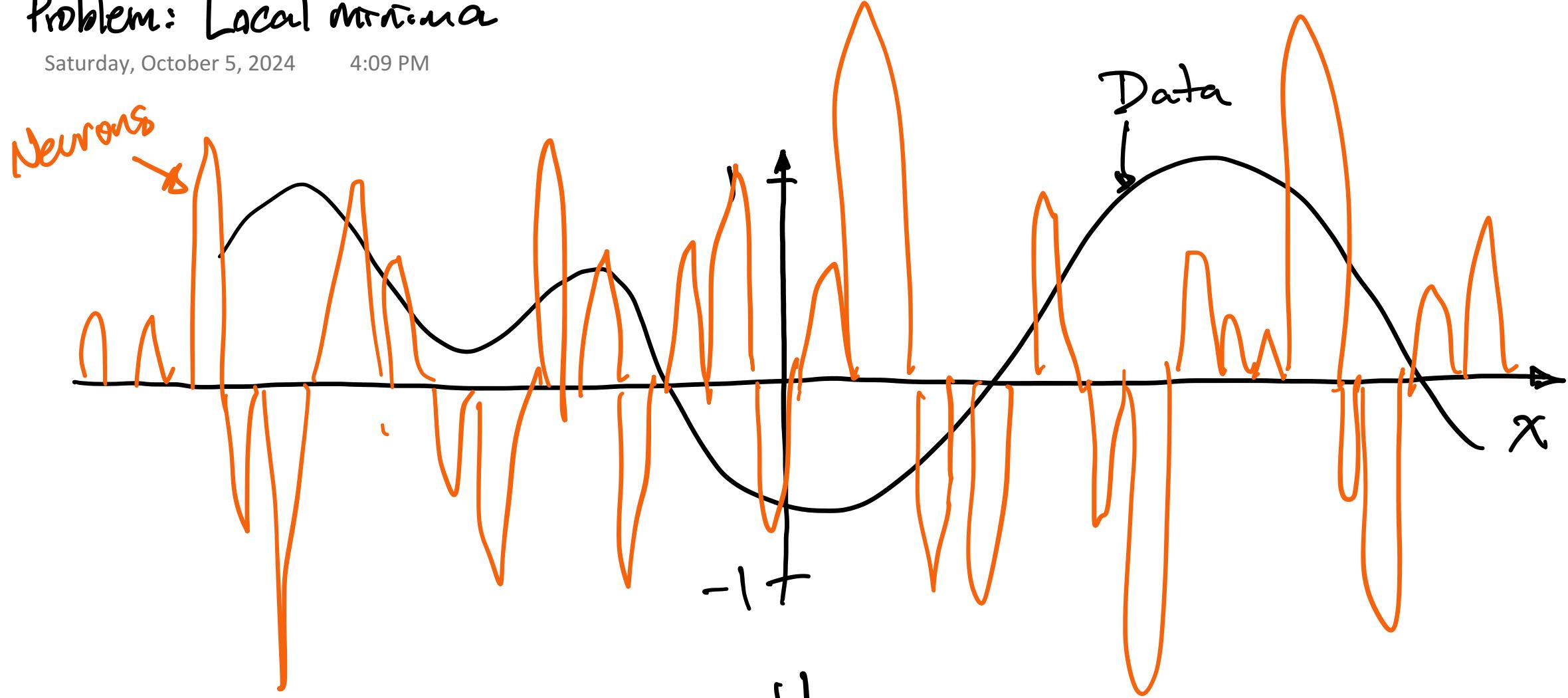
Saturday, October 5, 2024 4:48 PM

- Initialize weights, e.g., to small random values
- Iterate until stopping criteria met:
 - Select a training case
 - Apply forward propagation to get x 's
 - Compute output unit error derivatives
 - Apply backpropagation to get $\frac{\partial E}{\partial x}$'s
 - Compute weight derivatives $\frac{\partial E}{\partial w_m}$ and update weights, $w_m \leftarrow w_m - \eta \frac{\partial E}{\partial w_m}$

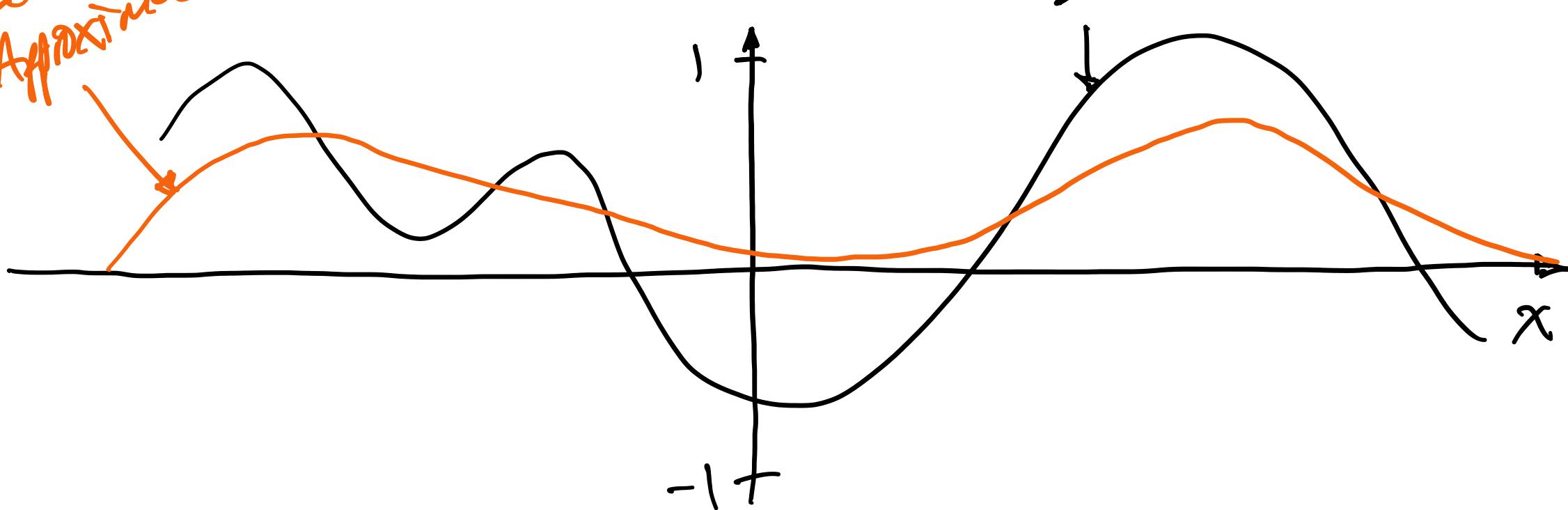
Problem: Local minima

Saturday, October 5, 2024

4:09 PM



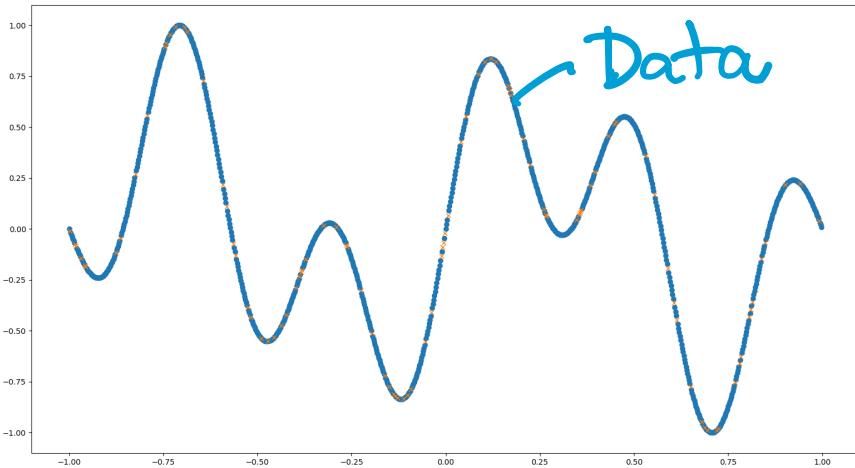
Learned Approximation



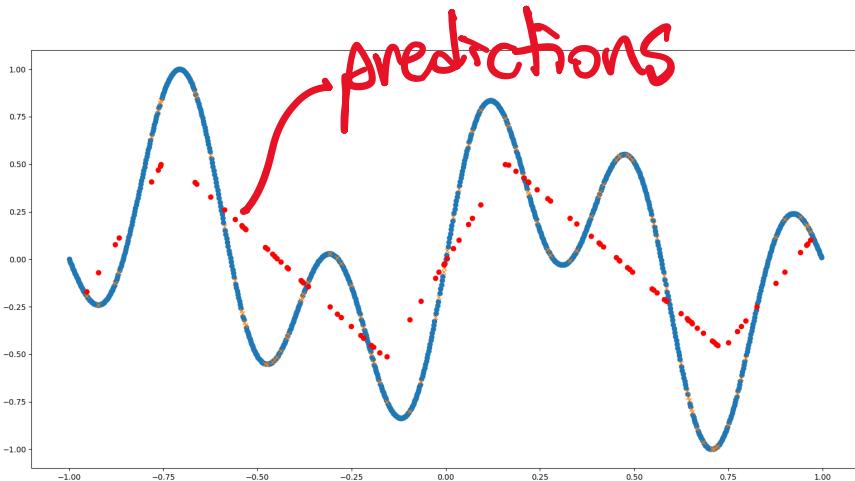
Example

Saturday, October 5, 2024

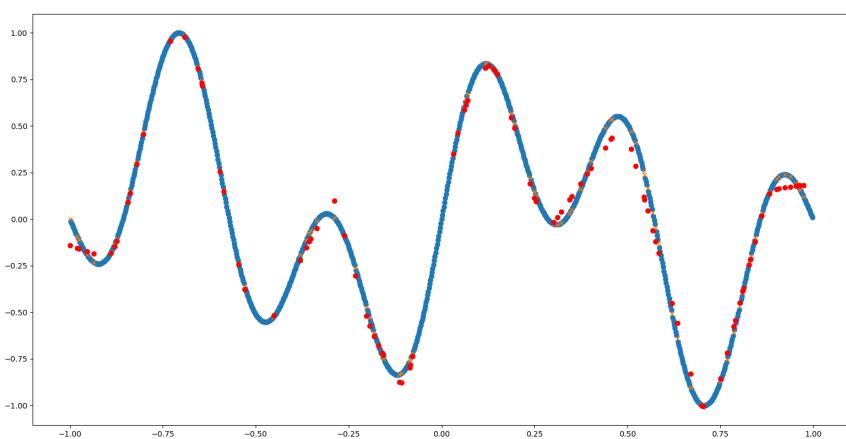
4:01 PM



Training data



Neural network predictions
1 hidden layer of 30 units
after training to convergence



Neural network predictions
4 hidden layers of 20, 40,
50 and 30 units
after training to convergence