

Machine Learning with Python
From Linear models
To Deep Learning

MITx 6.86x

Nawar

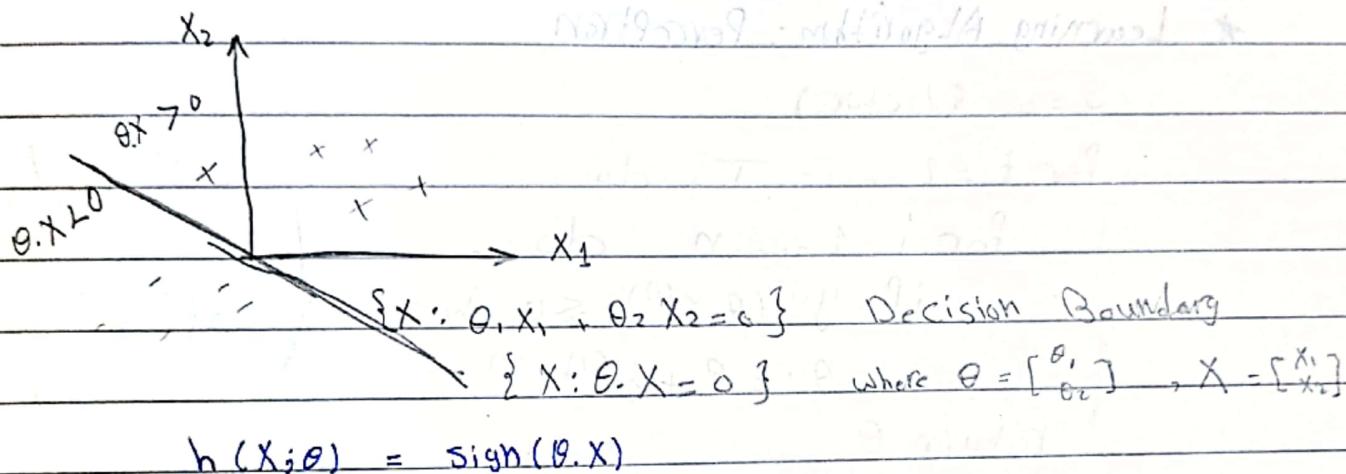
Lecture 2

Linear classifier and Perceptron

3. Linear classifiers Mathematically Revisited

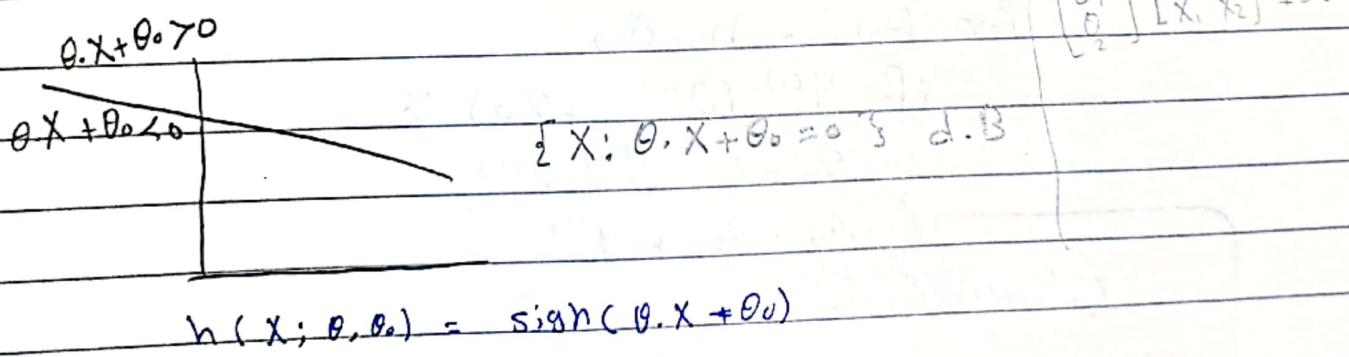
* Linear classifiers through origin

- the dividing line has to go through origin



* Linear classifier

- we can move decision boundary line everywhere



4. * Linear separation: Training examples $S_n = \{(x^{(i)}, y^{(i)}), i=1 \dots n\}$ are linearly separable if there exists a parameter vector $\hat{\theta}$ and offset parameter $\hat{\theta}_0$ such that $y^{(i)}(\hat{\theta} \cdot x^{(i)} + \hat{\theta}_0) > 0$ for all $i=0, 1 \dots n$

5. The Perceptron Algorithm:

* Training error for linear classifier

$$E_n(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n [y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0]$$

if $y=1 \& \text{out}=1 \rightarrow 1 > 0$

if $y=-1 \& \text{out}=-1 \rightarrow -1 > 0$

* Learning Algorithm: Perceptron

$$\theta = 0 \text{ (vector)}$$

for $t = 1 \dots T$ do

 for $i = 1 \dots n$ do

 if $y^{(i)}(\theta \cdot x^{(i)}) \leq 0$ then

$$\theta := \theta + y^{(i)}x^{(i)}$$

return θ

(without
The
offset)

$$\theta = 0 \text{ (vector)} \quad \theta_0 = 0 \text{ (scalar)}$$

for $t = 1 \dots T$ do

 for $i = 1 \dots n$ do

 if $y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0$ then

$$\theta = \theta + y^{(i)}x^{(i)}$$

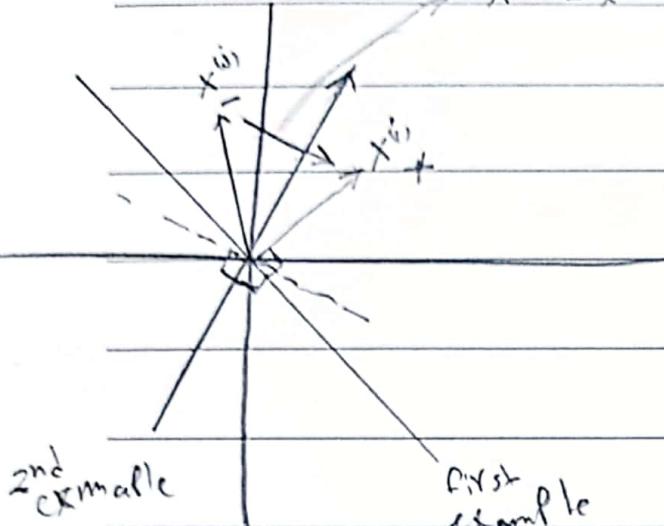
$$\theta_0 = \theta_0 + y^{(i)}$$

return θ, θ_0

$$\begin{bmatrix} \theta \\ \theta_0 \end{bmatrix} = \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix} + y^{(i)} \begin{bmatrix} x^{(i)} \\ 1 \end{bmatrix}$$

example: The Perceptron

$x^{(i)} - x^{(j)}$ ALgorithm



$$\text{let } \theta = 0$$

$$\theta = \theta + y^{(i)} X^{(i)} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{for } x^{(i)}$$

$$\theta = X^{(i)}$$

$$\theta = \theta + y^{(j)} X^{(j)} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{for } x^{(j)}$$

$$= X^{(i)} - X^{(j)}$$

if $y^{(i)} (\theta X^{(i)}) \leq 0 \rightarrow \text{error}$

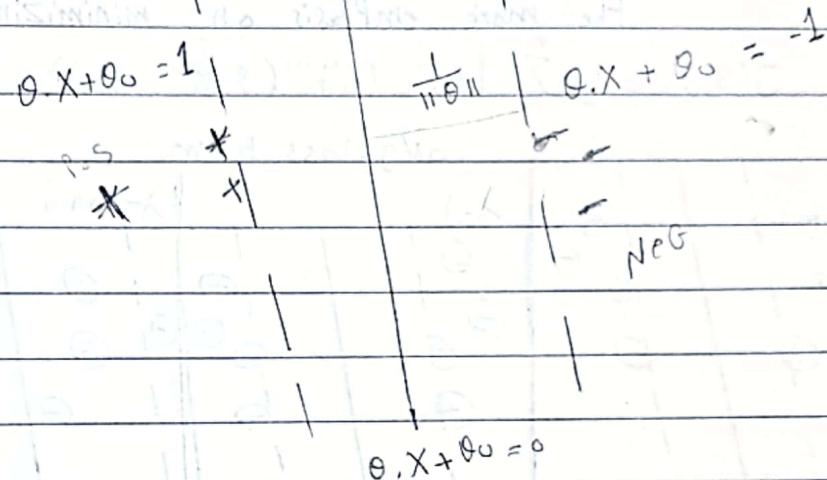
Lecture 3

Hinge Loss, Margin boundaries and Regularization

2. Introduction:

- * Large margin classifier: leaving lots of space on both sides before hitting the training examples.

3. Margin Boundary



4. Hinge Loss and objective function

* Hinge Loss: $\text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) = \begin{cases} 1 - z & \text{if } z \geq 1 \\ 1 & \text{if correct, } -1 \text{ if not correct} \end{cases}$

* Regularization: towards max margin

$$\max \frac{1}{\|\theta\|} \rightarrow \min \frac{1}{2} \|\theta\|^2$$

* The objective:

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \underbrace{\frac{\lambda}{2} \|\theta\|^2}_{\text{Reg}}$$

Avg Loss

Reg. Parameter

Newar

Lecture 4

Linear Classifier and Generalization

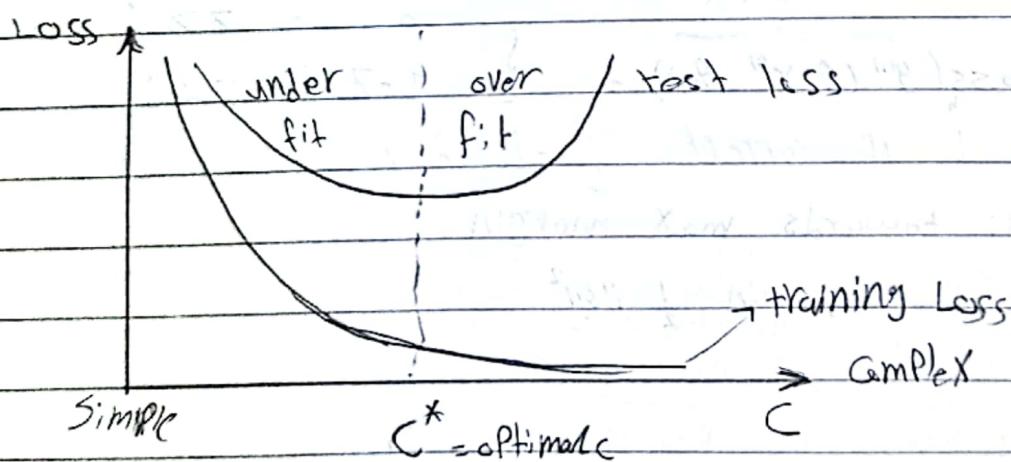
1. Review and the Lambda Parameter

* Lambda: "regularization parameter": changes the balance between the two terms, larger λ is the more we try to push the margin boundary apart, the smaller λ the more emphasis on minimizing the avg loss

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(x^{(i)}\theta + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2$$



3. Regularization and Generalization:



$$\frac{1}{\lambda} J(\theta, \theta_0) = \frac{1}{\lambda n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(x^{(i)}\theta + \theta_0)) + \frac{1}{2} \|\theta\|^2$$

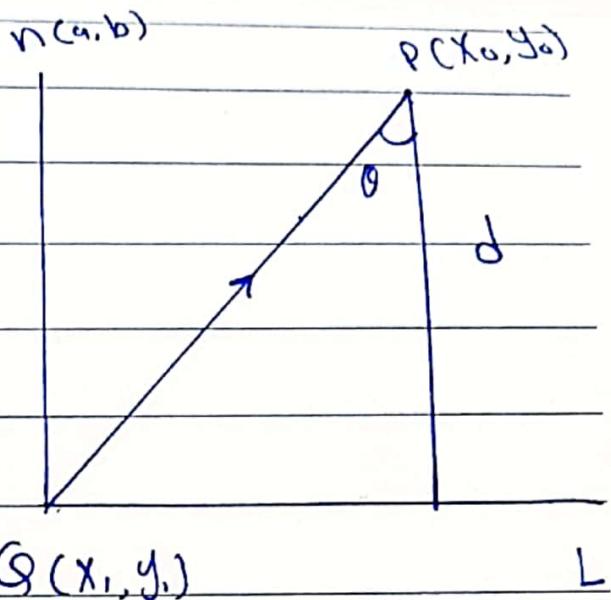
Line $L: ax + by + c = 0$

P: (x_0, y_0)

Let \vec{n} be a vector normal to the line

The distance d is the orthogonal projection of the vector \vec{QP}

$$\rightarrow d = \|\vec{QP}\| \cos \theta \quad \text{①}$$



$$d = \|\vec{QP}\| \cos \theta \cdot \frac{\|\vec{n}\|}{\|\vec{n}\|} \quad \text{②}$$

$$\therefore \vec{QP} \cdot \vec{n} = \|\vec{QP}\| \|\vec{n}\| \cos \theta \quad \therefore d = \frac{\vec{QP} \cdot \vec{n}}{\|\vec{n}\|} \quad \text{③}$$

$$\therefore \vec{QP} = (x_0 - x_1, y_0 - y_1), \vec{QP} \cdot \vec{n} = (x_0 - x_1, y_0 - y_1) \cdot (a, b)$$

$$\|\vec{n}\| = \sqrt{a^2 + b^2} \quad \therefore d = \frac{|a(x_0 - x_1) + b(y_0 - y_1)|}{\sqrt{a^2 + b^2}}$$

$$\therefore C = -ax_1 - by_1 \quad \& \quad d = \frac{|ax_0 + by_0 - c|}{\sqrt{a^2 + b^2}}$$

$$\therefore d = \frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$$

4. Gradient Descent

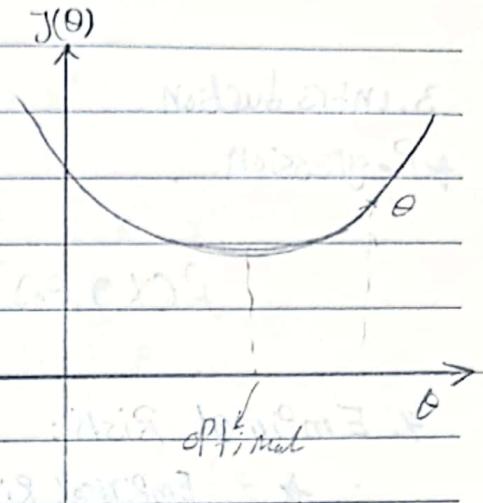
$$\theta \leftarrow \theta - \gamma \frac{\partial J}{\partial \theta}$$

step size \downarrow

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial J}{\partial \theta_j}$$

$$\theta \leftarrow \theta - \gamma \nabla J(\theta)$$

$$J(\theta)$$



5. Stochastic Gradient Descent

* It's more efficient way to do optimization.

* Perform gradient descent update with respect to the selected sampled term

$$\theta \leftarrow \theta - \gamma \nabla J_i(\theta)$$

→ Random sampling is often used in stochastic gradient descent.

$$\theta \leftarrow \theta - \gamma \nabla_{\theta} [\text{Loss}_i(y^{(i)} \cdot \theta \cdot x^{(i)}) + \frac{\lambda}{2} \|\theta\|^2]$$

$$\theta \leftarrow \theta - \gamma \begin{cases} 0 & \text{if Loss} = 0 \\ -y^{(i)}x^{(i)} & \text{if Loss} > 0 \end{cases} + \lambda \theta$$

6. The Realizable Case Quadratic Program

* Support vector machine finds the maximum margin linear separator by solving the quadratic program that corresponds to $J(\theta, \theta_0)$

* In the realizable case, if we disallow any margin violations, the quadratic program we have to solve is:

→ Find θ, θ_0 that minimize $\frac{1}{2} \|\theta\|^2$ subject to

$$y^{(i)} (\theta \cdot x^{(i)} + \theta_0) \geq 1$$

UNIT 2Lecture 5
Linear Regression

3. Introduction:

* Regression :

$$f(X, \theta, \theta_0) = \sum_{i=0}^d \theta_i X_i + \theta_0 = \theta X + \theta_0$$

4. Empirical Risk:

* Empirical Risk "The objective"

$$R_n(\theta) = \frac{1}{n} \sum_{t=1}^n (y^{(t)} - \theta X^{(t)})^2 / 2, \text{ assume } \theta_0 = 0$$

* There are 2 types of mistakes

→ structural : when the mapping between vectors and y 's is non linear

→ Estimation : when the mapping is Linear but you have a very limited training data, so you can't estimate it correctly

5. Gradient-Based Approach

$$\nabla_{\theta} \frac{(y^{(t)} - \theta X^{(t)})^2}{2} = (y^{(t)} - \theta X^{(t)}) \nabla_{\theta} (y^{(t)} - \theta X^{(t)})$$

$$= -(y^{(t)} - \theta X^{(t)}) \cdot X^{(t)}$$

→ initialize $\theta = 0$ then randomly pick $t = \{1, \dots, n\}$

$$\rightarrow \text{update } \theta := \theta + \gamma (y^{(t)} - \theta X^{(t)}) \cdot X^{(t)}$$

6. Closed Form Solution

$$\begin{aligned}\nabla_{\theta} R_n(\theta)_{|\theta=\hat{\theta}} &= -\frac{1}{n} \sum_{t=1}^n (y^{(t)} - \hat{\theta} x^{(t)}) x^{(t)} \\ &= -\frac{1}{n} \underbrace{\sum_{t=1}^n y^{(t)} x^{(t)}}_b + \frac{1}{n} \sum_{t=1}^n \hat{\theta} x^{(t)} x^{(t)} \\ &= -b + \underbrace{\frac{1}{n} \sum_{t=1}^n x^{(t)} (x^{(t)})^T}_{A} \hat{\theta} \\ &= -b + A \hat{\theta} = 0, \quad \hat{\theta} = A^{-1} b\end{aligned}$$

- Caveats : - the Matrix should be reversible, it'll be reversible if $(x^{(1)}, x^{(2)}, \dots, x^{(n)})$ span \mathbb{R}^d
- $n \gg d$, $d \rightarrow$ dimensions of the feature vector
- Cost, the operation is $O(d^3)$
- Matrix is non invertable which $\Rightarrow \det(X) = 0$

8. Regularization:

$$J_{\lambda, \theta}(\theta) = \frac{\lambda}{2} \|\theta\|^2 + R_n(\theta)$$

Keep all the θ_s very close to zero, \leftarrow entry: how much loss we are incurring from training set, so because it pulls the θ_s to be in some area, it'll not jumping for every noise in any training example. Find θ_s as possible match training examples

$$\begin{aligned}\nabla_{\theta} J_{\lambda, \theta}(\theta) &= \lambda \theta - (y^{(t)} - \theta \cdot x^{(t)}) x^{(t)} \\ &\rightarrow \theta = \theta - \gamma (\lambda \theta - (y^{(t)} - \theta \cdot x^{(t)}) x^{(t)}) = (1-\gamma \lambda) \theta + \gamma (y^{(t)} - \theta \cdot x^{(t)}) x^{(t)}\end{aligned}$$

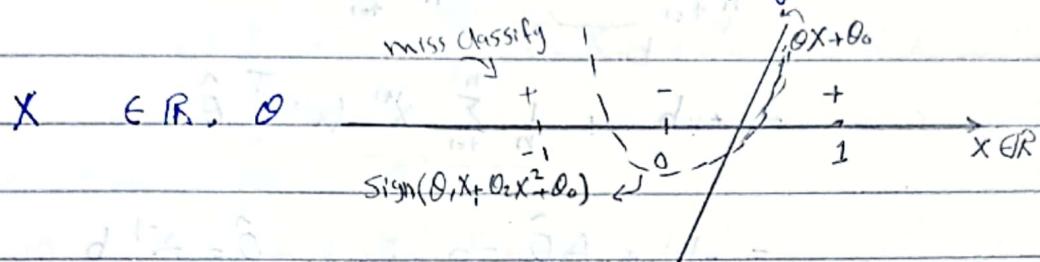
This achieves 2 goals: - pushes θ_s in the right direction to minimize the loss, also pushes θ_s down

Lecture 6

Non-Linear Classification

2. Higher order Feature Vector

Linear classification on the real line: $h(x_i; \theta, \theta_0) = \text{sign}(\theta x + \theta_0)$



Feature space

$$h(x_i; \theta, \theta_0) = \text{sign}(\theta_1 x_1 + \theta_2 x_1^2 + \theta_0)$$

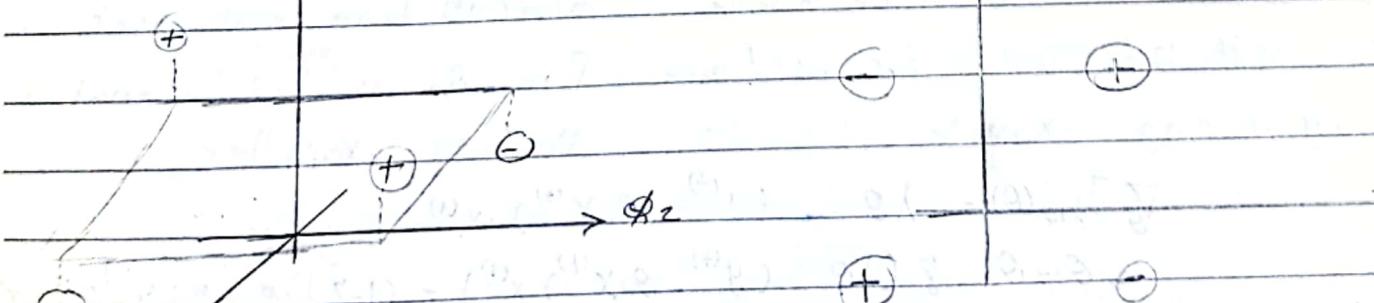
$$\phi(x) = \begin{bmatrix} x_1 \\ x_1^2 \end{bmatrix} \in \mathbb{R}^2, \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$



→ Linear classifier in the new feature coordinates implies a nonlinear classifier in the original x space

2-D example

$$\phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \end{bmatrix}, \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_0 \end{bmatrix}, \theta_0$$



→ based on x_1, x_2 +ve labeled examples have coordinates that are all +ve for ϕ_3 , and negatively labeled examples have negative values for that coordinate $\rightarrow \theta_3 = [0], \theta_0 = 0$

3. Introduction to nonlinear classification:

* Polynomial Features:

→ We can add more Polynomial terms $x \in \mathbb{R}$, $\phi(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}$

→ Means lots of features in high dimensions

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2, \quad \phi(x) = [x_1 \ x_2 \ x_1^2 \ x_2^2 \ \sqrt{x_1 x_2}]^T \in \mathbb{R}^5$$

→ We can do the same thing with Linear Regression

* Why Not Feature Vector?

→ The vectors can be quite high dimensional.

→ so we will use Kernels

4. Motivation for Kernels: Computational Efficiency

* Computing inner products between 2 feature vectors can be cheap even if the vectors are high dimensional

$$\phi(x) = [x_1, x_2, x_1^2, \sqrt{x_1 x_2}, x_2^2]^T$$

$$\phi(x') = [x'_1, x'_2, x'_1^2, \sqrt{x'_1 x'_2}, x'_2^2]^T$$

$$K(x, x') = \phi(x) \cdot \phi(x') = (x \cdot x') + (x \cdot x')^2 = x_1 x'_1 + x_2 x'_2 + x_1^2 x'_1^2 + 2x_1 x'_1 x_2 x'_2 - (x_1 x'_1 + x_2 x'_2)^2$$

→ For some feature maps we can evaluate the inner products very efficiently $K(x, x') = \phi(x) \cdot \phi(x') = (1 + x \cdot x')^p$, $p=2$

→ In these cases it's advantageous to express the linear classifier in terms of kernels rather than explicitly constructing feature vector

→ we need an equation to solve the classification problem using kernel trick, so we recall Perceptron with some difference.

Nawar

5. The Kernel Perceptron Algorithm

$$\theta = 0$$

run through $i = 1, \dots, n$

if $y^{(i)} \theta \cdot \phi(x^{(i)}) \leq 0$

$$\rightarrow \theta \leftarrow \theta + y^{(i)} \phi(x^{(i)})$$

from $\theta = \sum_{j=1}^n \alpha_j y^{(j)} \phi(x^{(j)})$ represented the parameter vector
 that lives in high-dim feature space
 \downarrow # of mistakes in times of number of mistakes

For prediction $\rightarrow \theta \phi(x^{(i)})$ per input j

$$\theta \phi(x^{(i)}) = \sum_{j=1}^n \alpha_j y^{(j)} \phi(x^{(j)}) \phi(x^{(i)}) = \sum_{j=1}^n \alpha_j y^{(j)} K(x^{(j)}, x^{(i)})$$

How similar the j -th example is

to the i -th example

The Kernel Perceptron will be:

$$\alpha_1 = \dots = \alpha_n = 0$$

run through $i = 1, \dots, n$

if $y^{(i)} \sum_{j=1}^n \alpha_j y^{(j)} K(x^{(j)}, x^{(i)}) \leq 0$

$$\alpha_i \leftarrow \alpha_i + 1$$

6. Kernel Composition Rules:

* Instead of directly constructing feature vectors by adding coordinates and then taking it in the product and seeing how it collapses into a kernel we can construct kernels directly from simpler kernels.

* Composition Rules:

1) $K(x, x') = 1$ is a kernel function if $\phi(x) = 1$

2) If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a scalar function, $K(x, x')$ is a kernel.

Then: $\tilde{K}(x, x') = f(x) K(x, x') f(x') \rightarrow \tilde{\phi}(x) = f(x) \phi(x)$

3) if $K_1(x, x')$ and $K_2(x, x')$ are kernels, then

$K(x, x') = K_1(x, x') + K_2(x, x')$ is a kernel $\phi(x) = \phi^{(1)}(x) + \phi^{(2)}(x)$

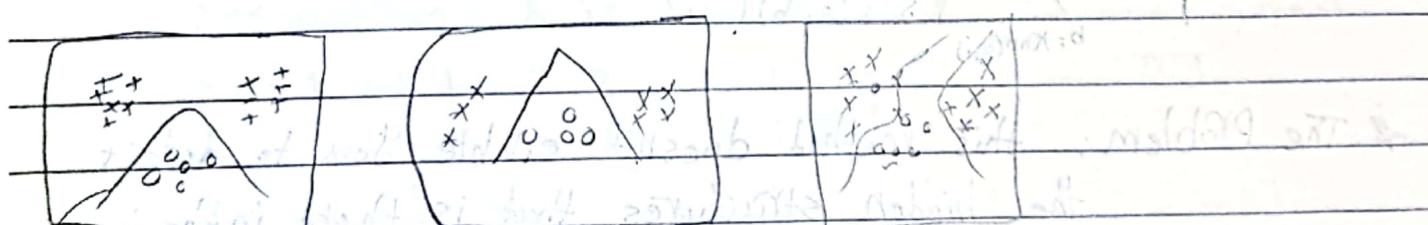
4) If $K_1(x, x')$ and $K_2(x, x')$ are kernels then

$K(x, x') = K_1(x, x') K_2(x, x')$ is a kernel

7. The Radial basis Kernel

* The feature vector can be infinite dimensional... this means that they have unlimited expressive power

$$K(x, x') = \exp(-\frac{1}{2} \|x - x'\|^2)$$



The DB now mapped into the original coordinates, its linear boundary in infinite dim space. But Mapped into an original space as (nonlinear)

$\{x_i: \sum_{j=1}^n \alpha_j y^{(j)} K(x^{(j)}, x_i) = 0\}$ & we can also use Random Forest to create non-linear Classification.

Naive

Lecture 7 Recommender systems

3. K-Nearest Neighbors

- * How many neighbors you want to look at and this can be one of the hyperparameters
- & we will find K number of users who are similar to the desired user, and aggregate them all and take the average

$$\hat{Y}_{ai} = \frac{\sum_{b \in \text{Knn}(a, i)} Y_{bi}}{K}$$

	User <i>i</i>			
User <i>a</i>	?	5	4	1
Users	1	4	4	2
1	5	4	5	5
2	5	5	5	5

- * we can use any technique to calculate the similarity like: Cos-similarity, Euclidean

- * weighting how similar ~~this~~ the user is to me?

$$\hat{Y}_{ai} = \sum_{b \in \text{Knn}(a, i)} \text{sim}(a, b) Y_{bi}$$

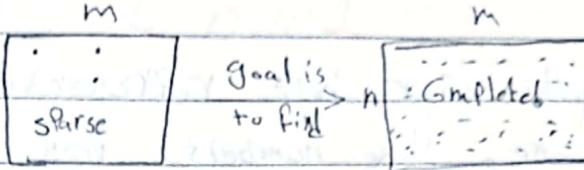
$$\sum_{b \in \text{Knn}(a, i)} |\text{sim}(a, b)|$$

- * The problem: this method doesn't enable you to detect the hidden structures that is there in the data "you may be similar to some pool of users in one dim, but similar to some other users in different dim"

* 4. Collaborative Filtering: The Naive approach.

$$J(X) = \sum_{(a,i) \in D} \frac{1}{2} (Y_{ai} - X_{ai})^2 + \frac{\lambda}{2} \sum_{(a,i)} X_{ai} \|X\|_F^2$$

$$D = \{(a,i) | Y_{ai} \text{ is given}\}$$



Given

Output

* assumption: ① $(a,i) \in D \rightarrow \frac{\partial J(X)}{\partial X_{ai}} = \frac{\partial}{\partial X_{ai}} (Y_{ai} - X_{ai})^2/2 + \frac{\lambda}{2} X_{ai}^2 = 0$

where no solution from diff. $X_{ai} = \frac{Y_{ai}}{1+\lambda}$

② $(a,i) \notin D \rightarrow X_{ai} = 0$

* this solution is worse than before

* there something here: the derivative was taken for each entry independently

5. Collaborative Filtering with matrix Factorization:

* The problem with the naive approach was that there is no connection between for all different entries of X

* we want to decrease number of parameters so we will use factorization.

* Assumption: X is Low Rank

$$\text{rank } 1: \begin{bmatrix} 1 & 2 & 3 \\ 5 & 10 & 15 \end{bmatrix} = \begin{bmatrix} 1 \\ 5 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} 10 \\ 50 \end{bmatrix} \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix}$$

rank 2: $D(n \times m)$

rank 3: $D(n \times m)$

rank 4: $D(n \times m)$

rank 5: $D(n \times m)$

rank 6: $D(n \times m)$

rank 7: $D(n \times m)$

rank 8: $D(n \times m)$

rank 9: $D(n \times m)$

rank 10: $D(n \times m)$

rank 11: $D(n \times m)$

rank 12: $D(n \times m)$

rank 13: $D(n \times m)$

rank 14: $D(n \times m)$

rank 15: $D(n \times m)$

rank 16: $D(n \times m)$

rank 17: $D(n \times m)$

rank 18: $D(n \times m)$

rank 19: $D(n \times m)$

rank 20: $D(n \times m)$

rank 21: $D(n \times m)$

rank 22: $D(n \times m)$

rank 23: $D(n \times m)$

rank 24: $D(n \times m)$

rank 25: $D(n \times m)$

rank 26: $D(n \times m)$

rank 27: $D(n \times m)$

rank 28: $D(n \times m)$

rank 29: $D(n \times m)$

rank 30: $D(n \times m)$

rank 31: $D(n \times m)$

rank 32: $D(n \times m)$

rank 33: $D(n \times m)$

rank 34: $D(n \times m)$

rank 35: $D(n \times m)$

rank 36: $D(n \times m)$

rank 37: $D(n \times m)$

rank 38: $D(n \times m)$

rank 39: $D(n \times m)$

rank 40: $D(n \times m)$

rank 41: $D(n \times m)$

rank 42: $D(n \times m)$

rank 43: $D(n \times m)$

rank 44: $D(n \times m)$

rank 45: $D(n \times m)$

rank 46: $D(n \times m)$

rank 47: $D(n \times m)$

rank 48: $D(n \times m)$

rank 49: $D(n \times m)$

rank 50: $D(n \times m)$

rank 51: $D(n \times m)$

rank 52: $D(n \times m)$

rank 53: $D(n \times m)$

rank 54: $D(n \times m)$

rank 55: $D(n \times m)$

rank 56: $D(n \times m)$

rank 57: $D(n \times m)$

rank 58: $D(n \times m)$

rank 59: $D(n \times m)$

rank 60: $D(n \times m)$

rank 61: $D(n \times m)$

rank 62: $D(n \times m)$

rank 63: $D(n \times m)$

rank 64: $D(n \times m)$

rank 65: $D(n \times m)$

rank 66: $D(n \times m)$

rank 67: $D(n \times m)$

rank 68: $D(n \times m)$

rank 69: $D(n \times m)$

rank 70: $D(n \times m)$

rank 71: $D(n \times m)$

rank 72: $D(n \times m)$

rank 73: $D(n \times m)$

rank 74: $D(n \times m)$

rank 75: $D(n \times m)$

rank 76: $D(n \times m)$

rank 77: $D(n \times m)$

rank 78: $D(n \times m)$

rank 79: $D(n \times m)$

rank 80: $D(n \times m)$

rank 81: $D(n \times m)$

rank 82: $D(n \times m)$

rank 83: $D(n \times m)$

rank 84: $D(n \times m)$

rank 85: $D(n \times m)$

rank 86: $D(n \times m)$

rank 87: $D(n \times m)$

rank 88: $D(n \times m)$

rank 89: $D(n \times m)$

rank 90: $D(n \times m)$

rank 91: $D(n \times m)$

rank 92: $D(n \times m)$

rank 93: $D(n \times m)$

rank 94: $D(n \times m)$

rank 95: $D(n \times m)$

rank 96: $D(n \times m)$

rank 97: $D(n \times m)$

rank 98: $D(n \times m)$

rank 99: $D(n \times m)$

rank 100: $D(n \times m)$

rank 101: $D(n \times m)$

rank 102: $D(n \times m)$

rank 103: $D(n \times m)$

rank 104: $D(n \times m)$

rank 105: $D(n \times m)$

rank 106: $D(n \times m)$

rank 107: $D(n \times m)$

rank 108: $D(n \times m)$

rank 109: $D(n \times m)$

rank 110: $D(n \times m)$

rank 111: $D(n \times m)$

rank 112: $D(n \times m)$

rank 113: $D(n \times m)$

rank 114: $D(n \times m)$

rank 115: $D(n \times m)$

rank 116: $D(n \times m)$

rank 117: $D(n \times m)$

rank 118: $D(n \times m)$

rank 119: $D(n \times m)$

rank 120: $D(n \times m)$

rank 121: $D(n \times m)$

rank 122: $D(n \times m)$

rank 123: $D(n \times m)$

rank 124: $D(n \times m)$

rank 125: $D(n \times m)$

rank 126: $D(n \times m)$

rank 127: $D(n \times m)$

rank 128: $D(n \times m)$

rank 129: $D(n \times m)$

rank 130: $D(n \times m)$

rank 131: $D(n \times m)$

rank 132: $D(n \times m)$

rank 133: $D(n \times m)$

rank 134: $D(n \times m)$

rank 135: $D(n \times m)$

rank 136: $D(n \times m)$

rank 137: $D(n \times m)$

rank 138: $D(n \times m)$

rank 139: $D(n \times m)$

rank 140: $D(n \times m)$

rank 141: $D(n \times m)$

rank 142: $D(n \times m)$

rank 143: $D(n \times m)$

rank 144: $D(n \times m)$

rank 145: $D(n \times m)$

rank 146: $D(n \times m)$

rank 147: $D(n \times m)$

rank 148: $D(n \times m)$

rank 149: $D(n \times m)$

rank 150: $D(n \times m)$

rank 151: $D(n \times m)$

rank 152: $D(n \times m)$

rank 153: $D(n \times m)$

rank 154: $D(n \times m)$

rank 155: $D(n \times m)$

rank 156: $D(n \times m)$

rank 157: $D(n \times m)$

rank 158: $D(n \times m)$

rank 159: $D(n \times m)$

rank 160: $D(n \times m)$

rank 161: $D(n \times m)$

rank 162: $D(n \times m)$

rank 163: $D(n \times m)$

rank 164: $D(n \times m)$

rank 165: $D(n \times m)$

rank 166: $D(n \times m)$

rank 167: $D(n \times m)$

rank 168: $D(n \times m)$

rank 169: $D(n \times m)$

rank 170: $D(n \times m)$

rank 171: $D(n \times m)$

rank 172: $D(n \times m)$

rank 173: $D(n \times m)$

rank 174: $D(n \times m)$

rank 175: $D(n \times m)$

rank 176: $D(n \times m)$

rank 177: $D(n \times m)$

rank 178: $D(n \times m)$

rank 179: $D(n \times m)$

rank 180: $D(n \times m)$

rank 181: $D(n \times m)$

rank 182: $D(n \times m)$

rank 183: $D(n \times m)$

rank 184: $D(n \times m)$

rank 185: $D(n \times m)$

rank 186: $D(n \times m)$

rank 187: $D(n \times m)$

rank 188: $D(n \times m)$

rank 189: $D(n \times m)$

rank 190: $D(n \times m)$

rank 191: $D(n \times m)$

rank 192: $D(n \times m)$

rank 193: $D(n \times m)$

rank 194: $D(n \times m)$

rank 195: $D(n \times m)$

rank 196: $D(n \times m)$

rank 197: $D(n \times m)$

rank 198: $D(n \times m)$

rank 199: $D(n \times m)$

rank 2

Rank 2 :

$$\begin{matrix} u \\ \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \end{matrix} \quad \begin{matrix} v \\ \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

i. the user here represented by 2 numbers instead of one, these numbers may represent their opinion about the first class & second class of movies

v. now we have the movies represented in 2 numbers
so we will have 2 groups / kind of aspects

* selecting "K" the rank by try different values on the validation set

6. Alternating Minimization

$$x_{ai} = u_a v_i : \text{Rank 1}$$

$$Y(u, v) = \sum_{(a, i) \in D} (y_{ai} - u_a v_i)^2 + \frac{\lambda}{2} \sum_{a=1}^n u_a^2 + \frac{\gamma}{2} \sum_{i=1}^m v_i^2$$

$$\text{example: } Y = \begin{bmatrix} 5 & ? & 7 \\ 1 & 2 & ? \end{bmatrix}, u_1 = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, v^T = [v_1, v_2, v_3]$$

$$\text{Initialize } v = \begin{bmatrix} 2 \\ 7 \\ 8 \end{bmatrix} \rightarrow uv^T = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} [v_1, v_2, v_3] = \begin{bmatrix} 2u_1 & 7u_1 & 8u_1 \\ 2u_2 & 7u_2 & 8u_2 \end{bmatrix}$$

→ trying to find u_1 & u_2 to get value as \bigcirc

$$\therefore u_1 = \frac{(5-2u_1)^2}{2} + \frac{(7-8u_1)^2}{2} + \frac{1}{2} u_1^2$$

$$\frac{\partial Y}{\partial u_1} = -66 + (68+1)u_1 = 0 \quad \rightarrow u_1 = \frac{66}{1+68} = \frac{66}{69} \quad \text{assume } \lambda=1$$

$$u_2 = \frac{16}{1+53} = \frac{16}{54}$$

PAGE _____
DATE _____

- * Then we take the New u_1 & u_2 and Compute V_1 & $V_2 \& V_3$
- * Repeat until Convergence
- * Local optimum لـ ادنا بـ نوصل لـ optimal لـ مـ حـ مـ حـ مـ حـ

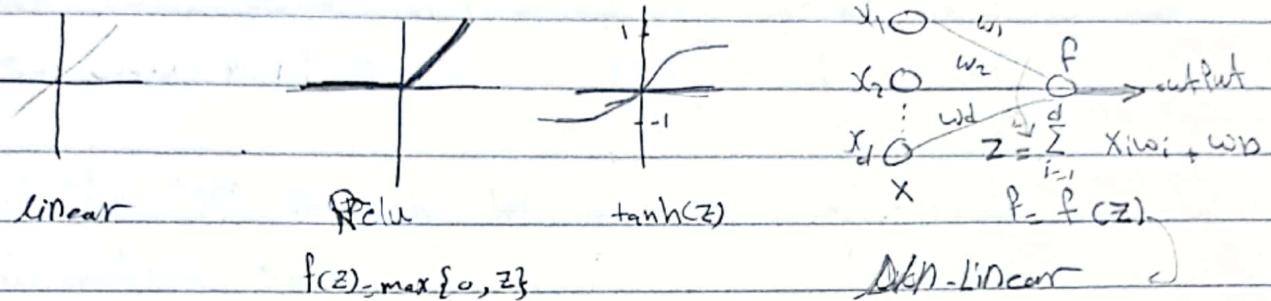
Lecture 8

Intro to FF Neural Networks

3. Motivation:

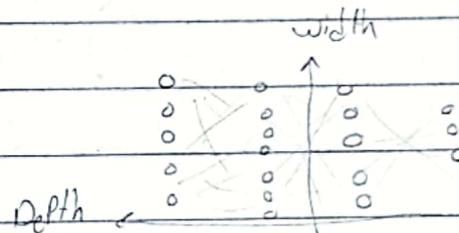
→ Neural Nets tries to optimize the feature representation for the task that we are trying to solve.

4.



5. Deep Neural Nets

Deep: many layers

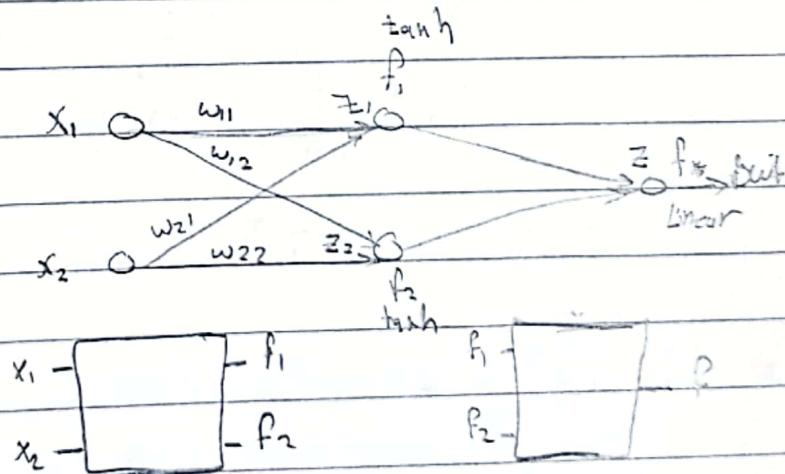


6. Hidden Layer Models

$$z_1 = \sum_{j=1}^2 x_j w_{j1} + w_0$$

$$f_1 = f(z_1) = \tanh(z_1)$$

$$z_2 = f_1$$

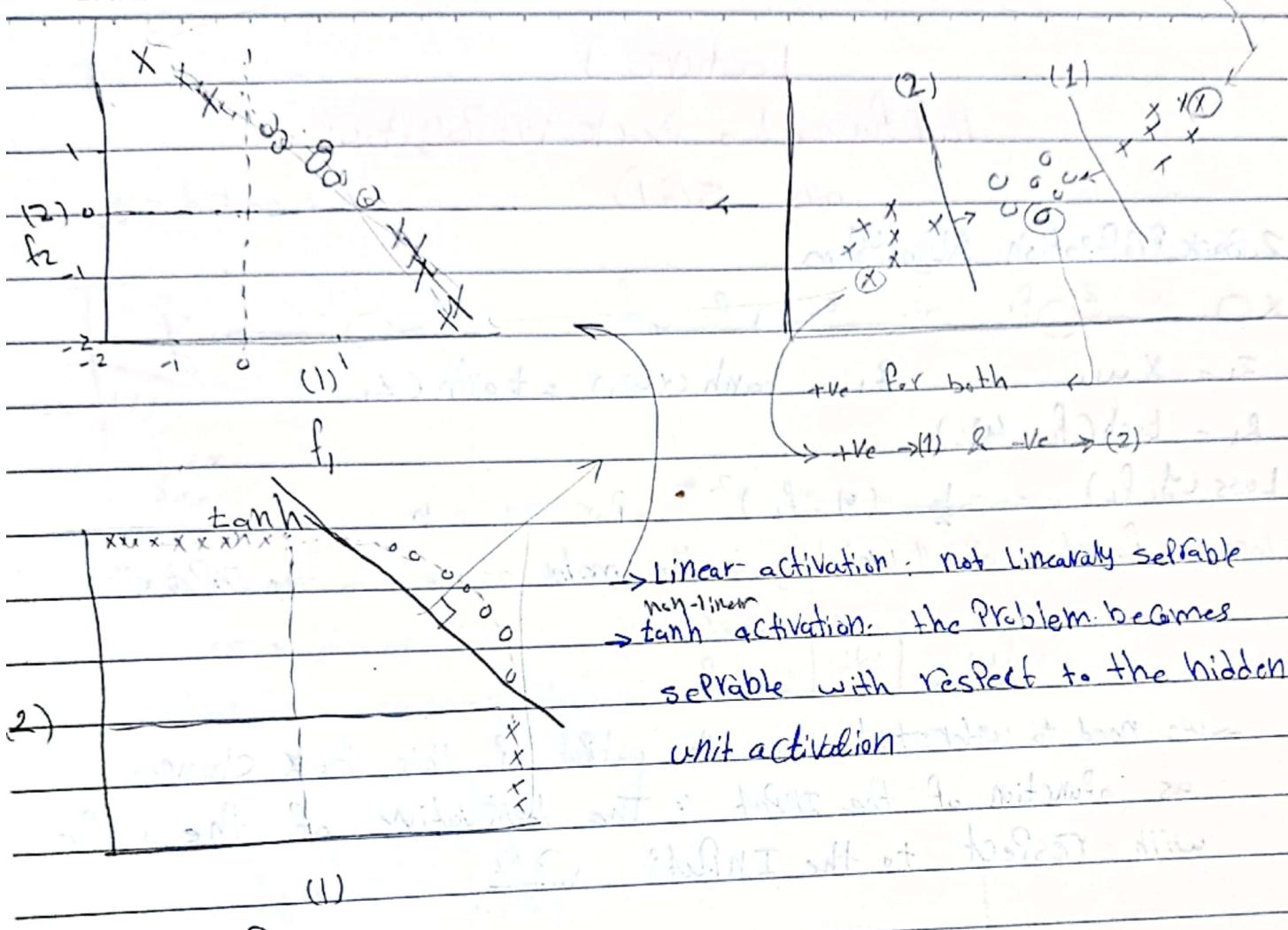


→ we can understand the first hidden unit as a linear classifier

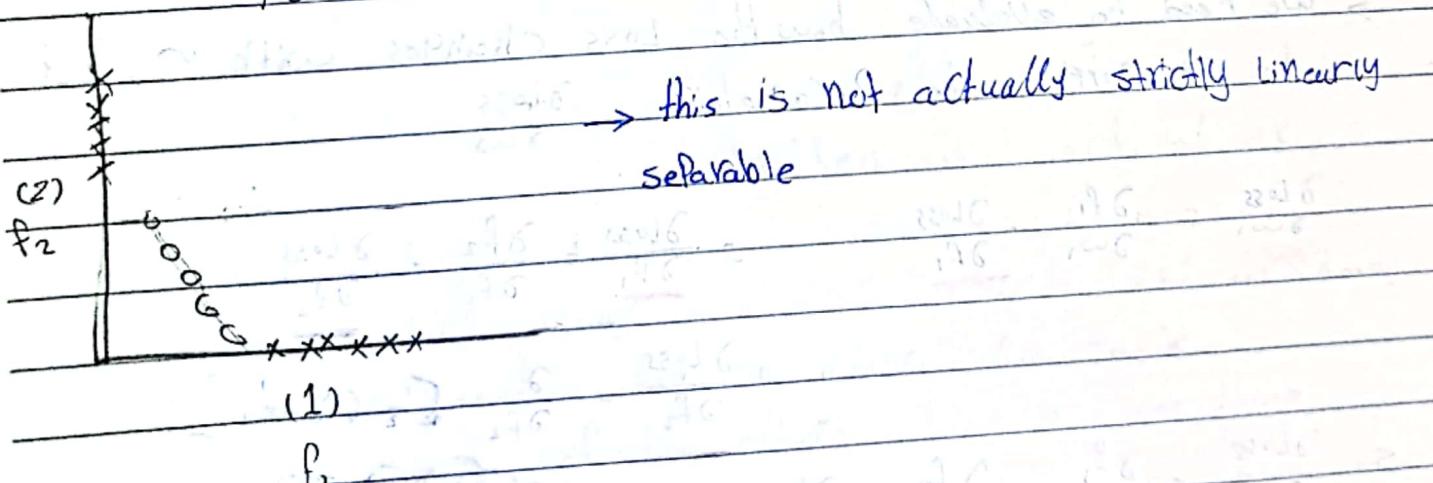
$$z_1 = \vec{w}_1 \cdot \vec{x} + w_0$$

$$\begin{bmatrix} w_{11} \\ w_{21} \end{bmatrix}$$

$$x_1 \quad x_2$$

v_c value for (1) & v_o for (2)

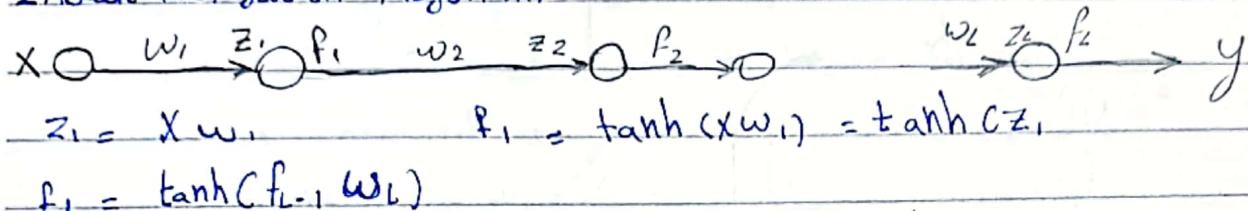
ReLU



Lecture 9

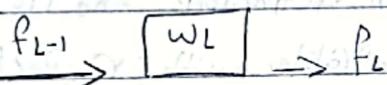
Feed Forward, Back Propagation
and SGD

2. Back Propagation Algorithm



$$\text{Loss}(y, f_L) = \frac{1}{2}(y - f_L)^2 \rightarrow \text{For regression}$$

Loss is function of all weights in the model as well as the Input example



→ we need to understand how the output of this box changes as a function of the input "The derivative of the output with respect to the input": $\frac{\partial f_L}{\partial f_{L-1}}$

→ we need to evaluate how the Loss changes with respect to changing the parameters $\frac{\partial \text{Loss}}{\partial w}$

$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{\partial f_1}{\partial w_1} \cdot \frac{\partial \text{Loss}}{\partial f_1}$$

$$\frac{\partial \text{Loss}}{\partial f_1} = \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial \text{Loss}}{\partial f_2}$$

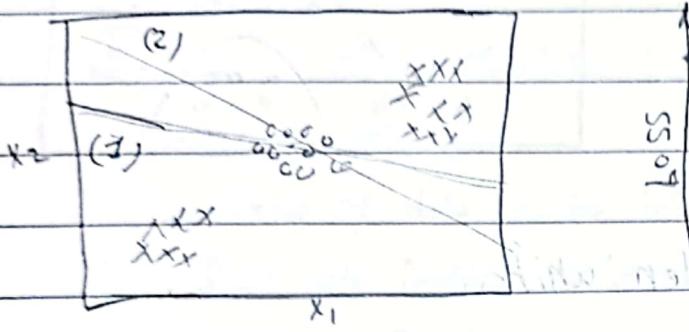
$$\frac{\partial \text{Loss}}{\partial f_2} = \frac{\partial}{\partial f_2} \left[\frac{1}{2}(y - f_2)^2 \right]$$

$$\text{So } \frac{\partial \text{Loss}}{\partial w_1} = \frac{\partial f_1}{\partial w_1} \cdot \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial \text{Loss}}{\partial f_2}$$

Chain Rule
etc

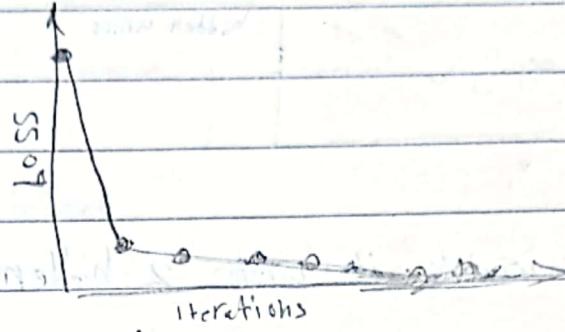
3. Training Models with 1 hidden layer:

* 2 hidden units



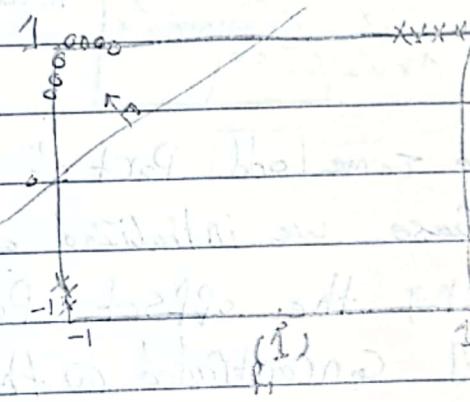
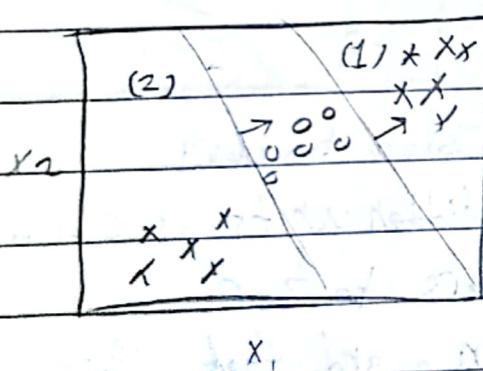
initial network

(2 hidden units)



hinge loss

→ after ~10 passes through the data



hidden units act various

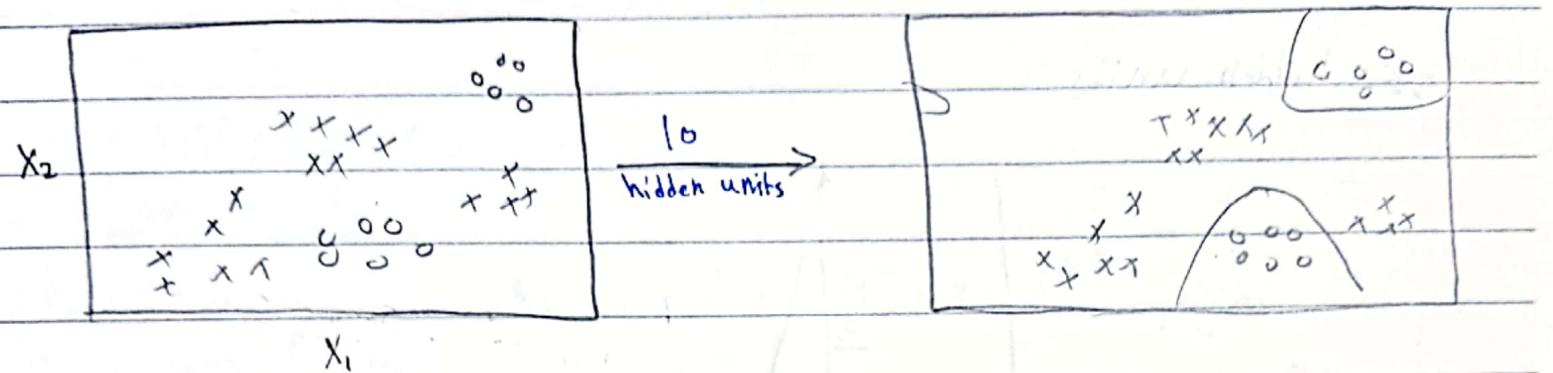
* 10 hidden units

→ after training many of those hidden units end up doing something useful for the classification task

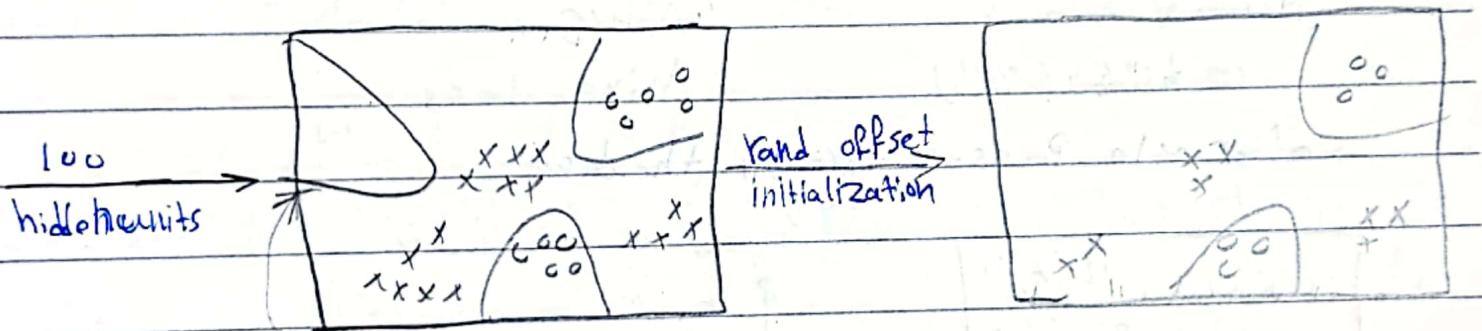
But some of the others doing random things

→ this issue called "OverCapacity"

* 2 hidden units with harder task



→ Can't solve it with 2 hidden units



→ You can see some [odd] part of the decision boundary
→ this is because we initialized all the hidden units as classifiers
by setting the offset parameters to zero
so they concentrated @ the origin, and that effect
persists in the solution because there is nothing
to change that

* Size Vs optimization

→ RELU is (cheap to evaluate, sparsity)
→ easier to learn as Large Models.

Lecture 10 Recurrent Neural Networks

2. Introduction:

* How to cast (Temporal / sequence) Problems as a supervised learning problem?

→ Historical data can be broken down into feature vectors and target values (sliding window)

$$x(t) = \begin{bmatrix} 0.82 \\ 0.8 \\ 0.73 \\ 0.72 \end{bmatrix} \quad y(t) = 0.89$$

* Language modeling: what comes next?

"This course has been a tremendous..."

$$\text{tremendous} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad y(t) = ?$$

$$a = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

(a(t))

→ The problem here is that we don't know exactly the length of the history that we should pay attention to.

→ The more history we include the longer the feature vectors and the more complex the learning and prediction.

→ Another problem is that sometimes the history is variable length, There may be words at the beginning of the sentence that are quite relevant for predicting what happens towards the end.

→ we need more a more flexible way of turning sequences into feature vector.

3. why we need RNNs?

* Learning to encode/decode

→ language modeling

"This course has been al?"



Success

→ sentiment classification

"I have seen better lectures"



1

→ Machine translation

"I have seen better lectures"



"Josic ciklara nisijid"

Encoding

Decoding

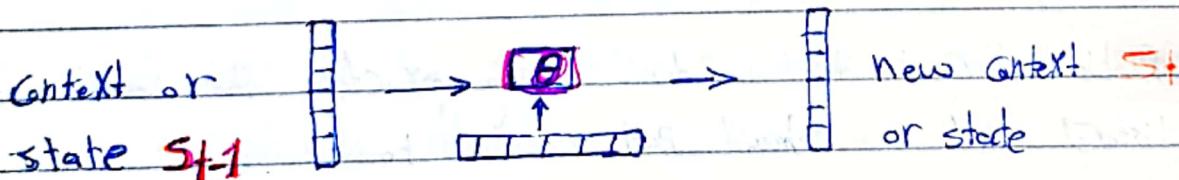
* Encoding is mapping a sequence to a vector.

* Decoding is mapping a vector to a sequence.

* we can encode everything as vectors (images, sentences, events, etc...), so we can turn the text into image and vice versa.

4. Encoding with RNN

* Encoding sentences:



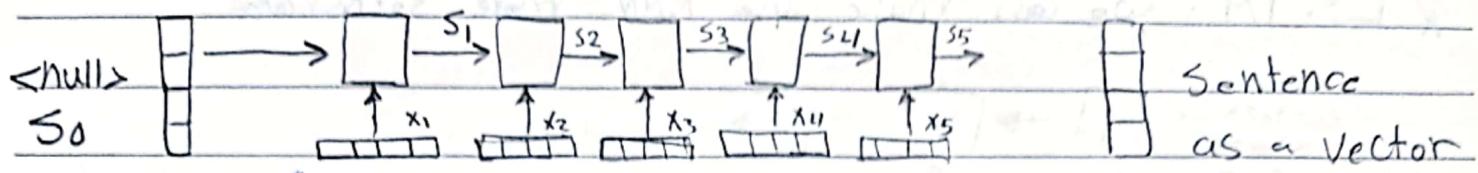
(*) New Information θ

$$s_t = \tanh(W^{s,s} \cdot s_{t-1} + W^{s,x} \cdot x_t)$$

$m \times 1$ $m \times m$ $m \times 1$ $m \times d$ $d \times 1$

→ This is simple transformation that takes previous state and new information, has some parameters represented by these two ~~matrices~~ matrices ($W^{s,s}, W^{s,x}$) which are adjustable

→ we can adjust these parameters (θ) so that this box and this transformation behaves as we want it to



"Efforts and Garage are not"

→ the sentence vector representation depends on theta, as well as the sequence and we can adjust those parameter to make this representation appropriate for our task

5. Gating and LSTM

- * The Problem with the previous is that we forget the information that we had, we override it completely to get the new information
- * we need to control over what we write and retain in the summaries or the state is typically done via another network called a gating network

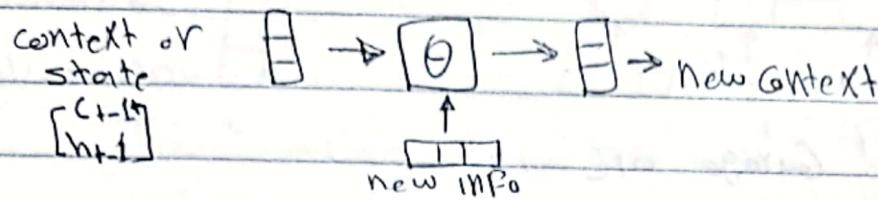
$$g_t = \text{Sigmoid}(W_{s,t}^{g,s} + W_{x,t}^{g,x} \cdot x_t)$$

$$s_t = (1 - g_t) \odot s_{t-1} + g_t \odot \tanh(W_{s,t}^{s,s} s_{t-1} + W_{x,t}^{s,x} \cdot x_t)$$

→ g_t control what we should write and retain into the state

→ when (g_{t-1}) is zero → that means that we don't retain any information of the previous state

* LSTM : we can make the RNN more sophisticated



$$f_t = \text{Sigmoid}(W^{f,h} \cdot h_{t-1} + W^{f,x} \cdot x_t) \quad \text{Forget gate}$$

$$i_t = \text{Sigmoid}(W^{i,h} \cdot h_{t-1} + W^{i,x} \cdot x_t) \quad \text{Input gate}$$

$$o_t = \text{Sigmoid}(W^{o,h} \cdot h_{t-1} + W^{o,x} \cdot x_t) \quad \text{Output gate}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W^{c,h} \cdot h_{t-1} + W^{c,x} \cdot x_t) \quad \text{Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \quad \text{Visible state}$$

- Forget, input and output gates control how to read information into the memory cell, how to forget information that we've had previously, and how to output information from the memory cell to the visible form

→ The previous state now in 2 parts: c_{t-1} & h_{t-1}

→ Input gate: Controls to what degree of these coordinates of this suggested signal should enter the memory cell

→ Forget gate: Controls which of the coordinates from the previous memory cell we should erase information from.

Lecture 11 Recurrent Neural Networks 2

2. Markov Models:

* Next word in a sentence depends on Previous symbols

already written (history = one, two, or more words)

* we assume that limited Vocabulary V , the words that are not in our V are introducing an UNK symbol

→ <beg> symbol specifying the start of the sentence

→ <end> symbol specifying the end of the sentence

<beg> The Lecture leaves me UNK <end>

w₀ w₁ w₂ w₃ w₄ w₅ w₆

* In first order Markov model (bigram) the next symbol only depends on the previous one

$$P(w_1, \dots, w_6) = P(w_1 | w_0) \cdot \dots \cdot P(w_6 | w_5)$$

<beg>

* A first order Markov model:

→ Each symbol (except <beg>) is predicted using the same conditional probability table until all <end> symbol is seen

	M1	Course	is	UNK	end
<beg>	0.7	0.1	0.1	0.1	0
M1	0.1	(0.5) <small>(= 0.5)</small>	0.2	0.1	0.1
Course	0	0	0.7	0.1	0.2
is	0.1	0.3	0	0.6	0
UNK	0.1	0.2	0.2	0.3	0.2

→ Row sums to 1 and defines the probability of the next word

<beg> Course is ^{UNK} get "end"

The probability $P(\text{course}|\text{beg}) \cdot P(\text{is}|\text{course}) \cdot P(\text{UNK}|is) \cdot P(\text{end}|\text{UNK})$

of this sentence 0.1 . . . 0.7 . . . 0.6 . . . 0.2

Never

* Simple example of how would we use the markov model:

~~the first word~~

→ zeroth symbol is always <beg>

→ sample from first row "Probability of distribution over the first word" which comes after the <beg> symbol by throwing a weighted five-way dice

→ we most likely get the word 'ME' but we might also get the word 'Course'

→ Now move to the "Course" row, we might get the most likely outcome "is" and so on.

→

* Maximum likelihood estimation

→ The goal is to maximize the probability that the model can generate all the observed sentences (corpus S)

$$S \in S, s = \{w_1^s, w_2^s, w_3^s, \dots, w_{|s|}^s\}$$

$$\log \left\{ \prod_{s \in S} \left[\prod_{i=1}^{1st} P(w_i^s | w_{i-1}^s) \right] \right\} \rightarrow \text{The Max Likelihood estimation for a single sentence}$$

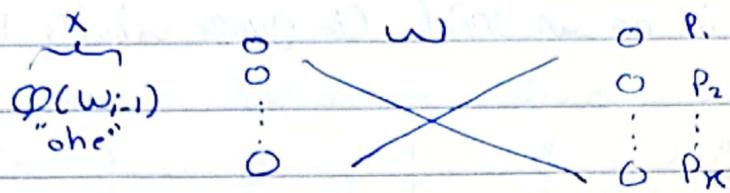
* Count (w, w')

$$\hat{P}(w, w') = \frac{\text{Count}(w, w')}{\sum_{\tilde{w}} \text{Count}(w, \tilde{w})}$$

The Max Likelihood estimate of its row in the normalized probability table.

3. Markov Model to Feedforward Neural Nets

→ we can represent the Markov model as a FFNN



$$p_k = P(w_{t,k} | w_{t-1}) \quad p_k > 0 \quad \sum p_k = 1$$

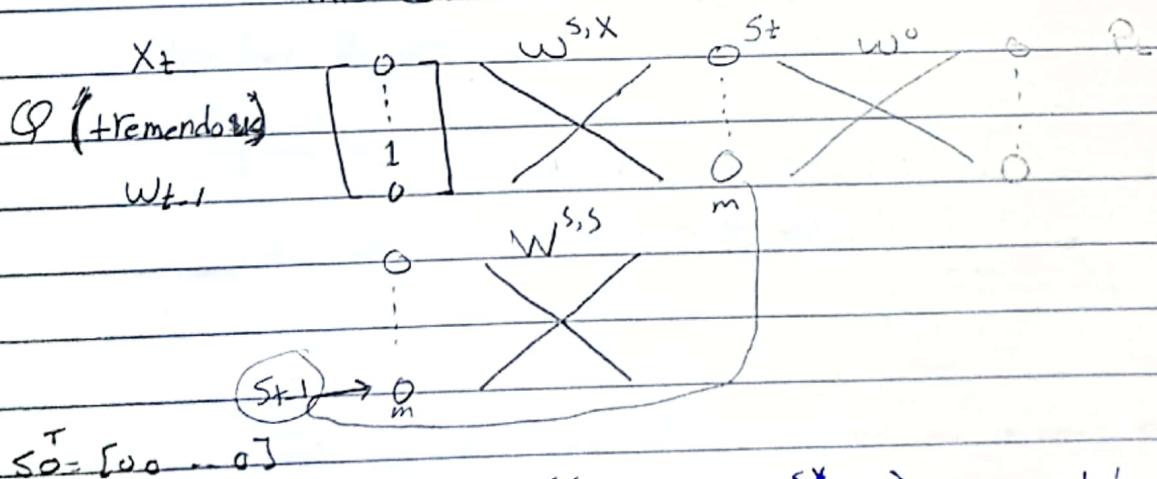
$$z_k = \sum_j x_j w_{jk} + w_0 \in \mathbb{R}$$

$$p_k = \frac{e^{z_k}}{\sum_i e^{z_i}} \rightarrow \text{softmax}$$

* with FFNN we can extend the Input by adding $\phi(w_{t-2})$ and add more hidden units

4. RNN Decoder Dive

← This course has been a tremendous...?

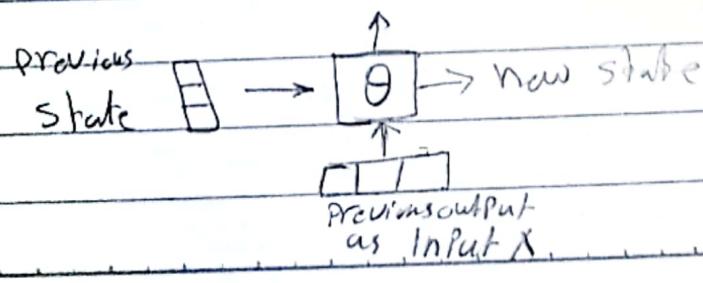


$$s_0^T = [0 \dots 0]$$

$$s_t = \tanh(w^{s,s} s_{t-1} + w^{s,x} x_t) \quad \text{state}$$

$$p_t = \text{softmax}(w^o s_t)$$

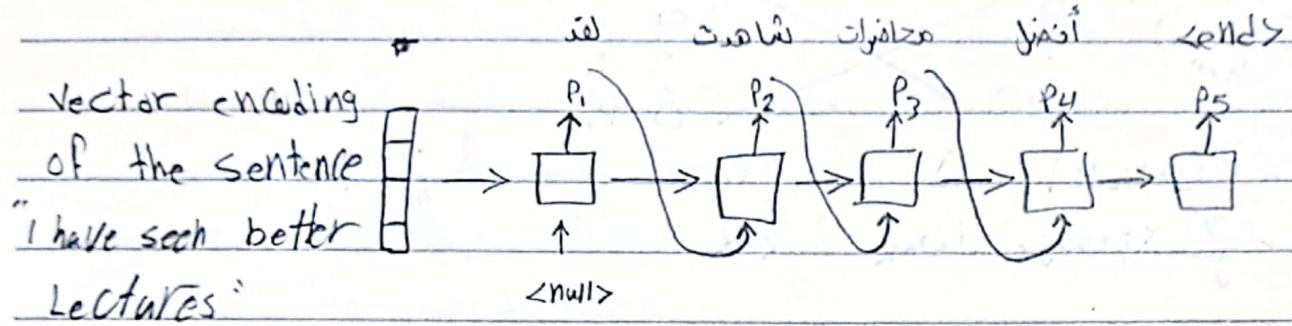
p_t : output distribution



Nawar

5. RNN Decoding :

* The output is fed in as an input (to gauge what's left)

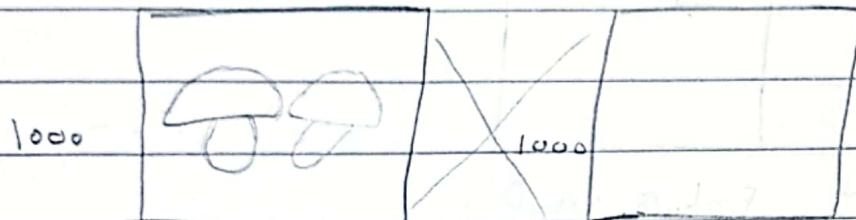


Lecture 12

Convolution Neural Networks

1. Convolution Neural Networks

* Feed-Forward Networks



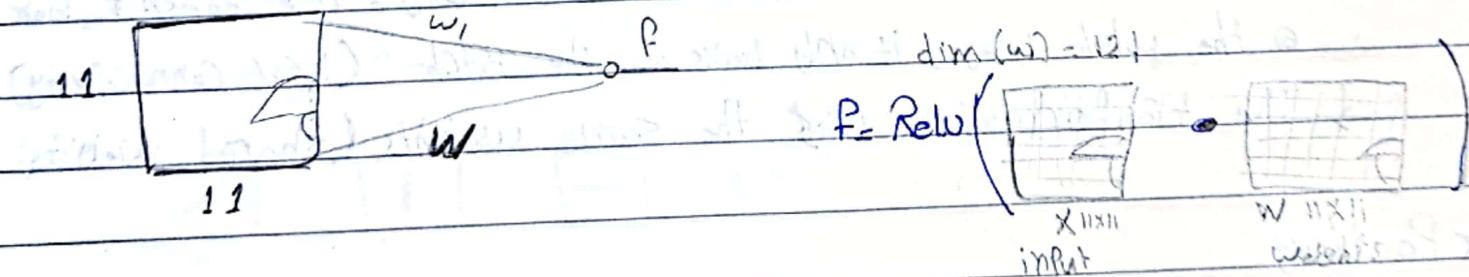
1000
Input

W
 $10^6 \times 10^6$

1000
layer 1

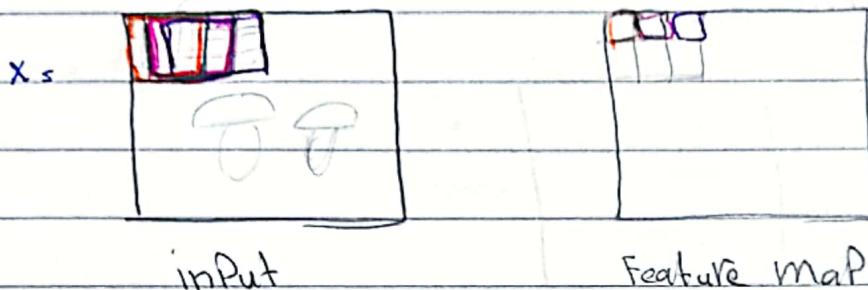
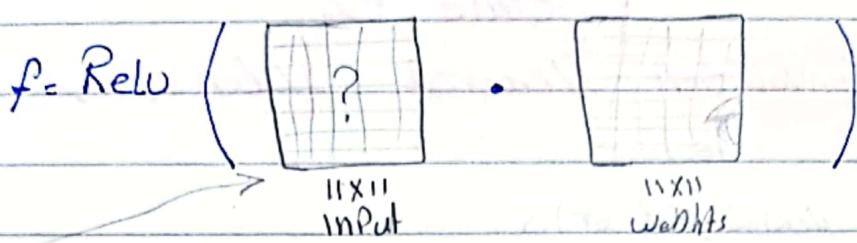
→ Problems: - There is a lot of weights are shared $10^6 \times 10^6$
Same object that appears in slightly different locations would have
to ~~learn~~ be learned again with weights.

* Patch Classifier/Filter



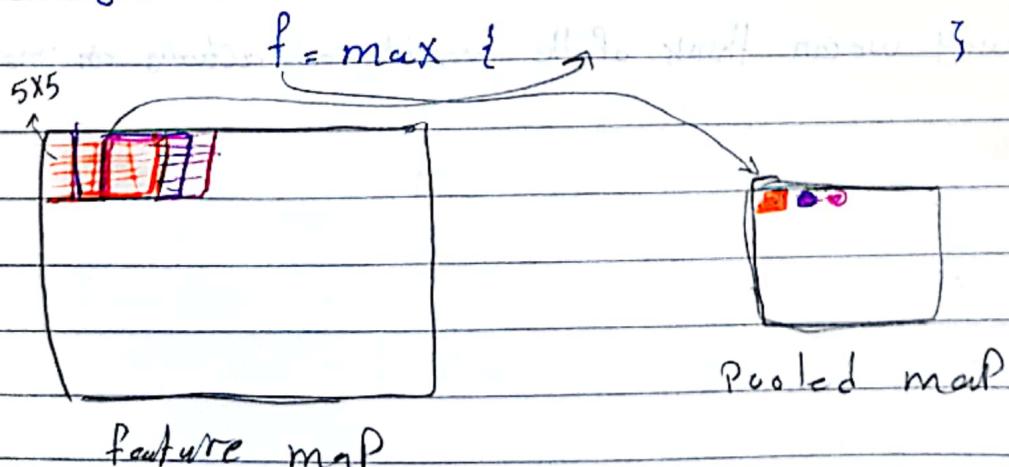
→ At this way we can think of the weights as actually an image

* Convolution:



- If we have looking for edges and we pass the patch ~~here~~ on Hashrum we might get a higher value in locations where such features appears (Hashrums)
- The transformation from the input to a feature map involves only 121 weights
- Each location on feature map "activation of units" tells us how much that feature present in the corresponding location
- each unit on feature map connects only locally, it doesn't look @ the whole image, it only looks at the patch (Local Connectivity)
- The transformation uses the same weights (shared weights)

* Pooling

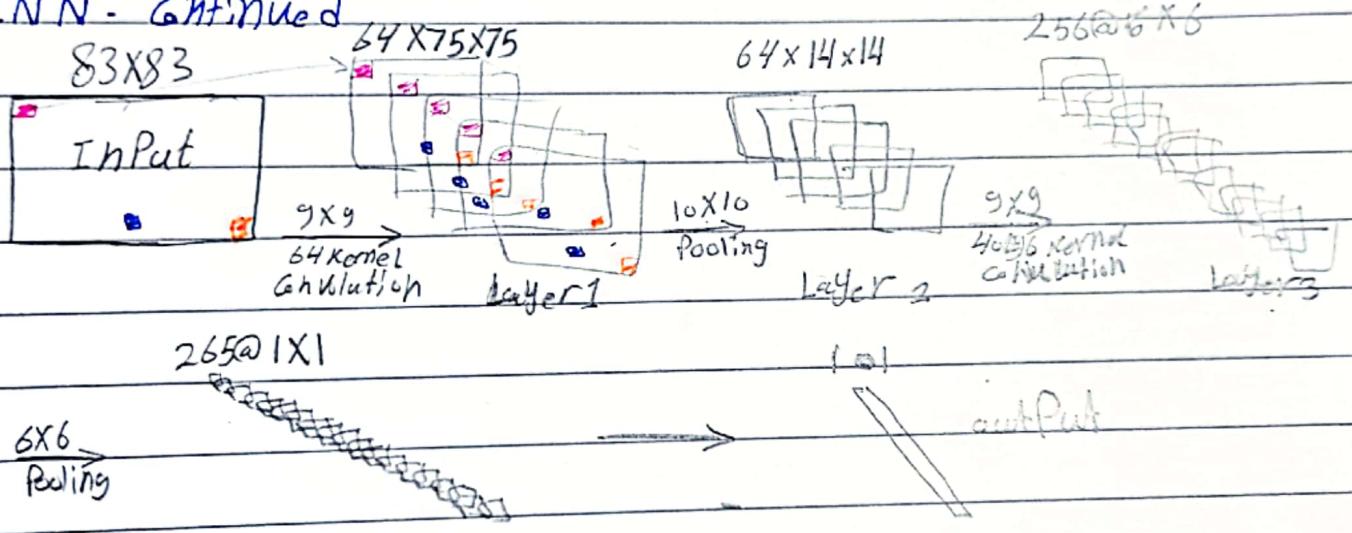


→ we wish to know whether the feature was there, but not exactly where it was.

→ The image we will get has 4 times less pixels

→ Each location in the new image corresponds to the maximum of the 25 pixels that it is looking at in the original image

3. CNN - Continued



→ each convolution corresponds to some weight matrix, it corresponds to looking for a particular type of features like edges below:



$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

So, the goal is find the weights of each kernel.