



# Neural Networks

## Lecture 3

### Linear auto-associative memories

Igor Farkaš

2019

### Linear NN – analytic solution

Let's consider the train set:  $A_{\text{train}} = \{(\mathbf{x}^{(p)}, \mathbf{y}^{(p)}), p = 1, 2, \dots, N\}$ .

We look for a matrix  $\mathbf{W}$  that satisfies  $\mathbf{y}^{(p)} = \mathbf{W}\mathbf{x}^{(p)}$ , for each  $p$ .

In matrix notation:  $\mathbf{Y} = \mathbf{W}\mathbf{X}$

$$\begin{matrix} [\mathbf{y}^{(1)} & \mathbf{y}^{(2)} & \dots & \mathbf{y}^{(N)}] & = & \mathbf{W} & \times & [\mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(N)}] \\ (m \times N) & & & & & (m \times n) & & & & & (n \times N) \end{matrix}$$

If  $\mathbf{X}$  was regular (i.e., square matrix with  $N = n$ , linearly independent rows)

then  $\mathbf{X}^{-1}$  would exist and  $\mathbf{W} = \mathbf{Y}\mathbf{X}^{-1}$ .

However, in general we cannot assume  $N = n$ , nor linear independence of input vectors ( $\Rightarrow \mathbf{X}^{-1}$  does not exist). Only generalized solutions exists:

$$\mathbf{W} = \mathbf{Y}\mathbf{X}^+$$

$\mathbf{X}^+$  is called (Moore-Penrose) **pseudoinverse matrix** of  $\mathbf{X}$ . Theorem:  $\forall \mathbf{X}, \exists \mathbf{X}^+$  with properties: 1)  $\mathbf{X}\mathbf{X}^+\mathbf{X} = \mathbf{X}$ , 2)  $\mathbf{X}^+\mathbf{X}\mathbf{X}^+ = \mathbf{X}^+$ , 3) symmetric  $\mathbf{X}\mathbf{X}^+$  and  $\mathbf{X}^+\mathbf{X}$ .

a)  $\mathbf{X}^+ = \mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}$ , if  $n < N$  and  $\text{rank}(\mathbf{X}) = n$ .

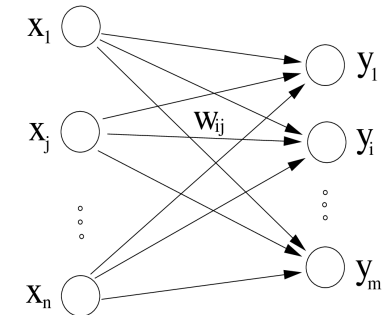
b)  $\mathbf{X}^+ = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ , if  $n > N$  and  $\text{rank}(\mathbf{X}) = N$ .

## Linear NN models

Input vector:  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$

Output vector:  $\mathbf{y} = [y_1, y_2, \dots, y_m]^T$

Weight matrix:  $\mathbf{W} \sim \text{type } [m \times n]$



Linear transformation  $\varphi : \mathcal{R}^n \rightarrow \mathcal{R}^m, \mathbf{y} = \mathbf{W}\mathbf{x}$

- ⊖ ignores saturation property of neurons
- ⊕ allows to find analytic solutions using linear algebra.

(Kohonen, 1970;  
Anderson, 1972;  
Cooper, 1973)

- Adding layers in a linear NN does not appear reasonable (since no complexity is added).
- But: It allows nonlinear learning dynamics in linear deep nets (Saxe, 2015).

### Auto-associator case

Let's consider  $n > N$  (i.e. few examples of large dimension) and the **autoassociative case**:  $\mathbf{y}^{(p)} = \mathbf{x}^{(p)}, m = n$

Model is supposed to remember  $N$  **prototypes**  $[\mathbf{x}^{(1)} \mathbf{x}^{(2)} \dots \mathbf{x}^{(N)}]$ .

*Goal*: train on prototypes and then submit a corrupted version of a prototype. Model should be able to reconstruct it.

Since  $\mathbf{Y} = \mathbf{X}$ , then  $\mathbf{W} = \mathbf{X}\mathbf{X}^+$ . How to interpret  $\mathbf{W}$ ?

In a special case, which is too restrictive ( $N = n$ , linearly independent inputs), we would have a trivial solution  $\mathbf{W} = \mathbf{I}$  (identity).

How about a general case?

## Basics of linear vector spaces

Let's have a linear space  $\mathbb{R}^n$ .

**Linear manifold**  $\mathcal{L} = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} = a_1 \mathbf{x}^{(1)} + a_2 \mathbf{x}^{(2)} + \dots + a_N \mathbf{x}^{(N)}, a_p \neq 0\}$   
 $\mathcal{L} \subset \mathbb{R}^n$

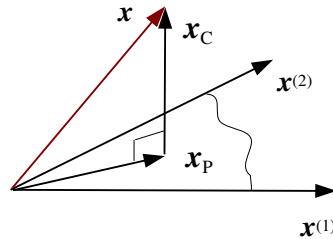
**Orthogonal complement**  $\mathcal{L}^\perp = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} \perp \mathcal{L}\}$

Hence,  $\mathcal{L} \oplus \mathcal{L}^\perp = \mathbb{R}^n$

Each vector  $\mathbf{x} \in \mathbb{R}^n$  can be uniquely decomposed:

$$\mathbf{x} = \mathbf{x}_p + \mathbf{x}_c$$

where  $\mathbf{x}_p \in \mathcal{L}$  and  $\mathbf{x}_c \in \mathcal{L}^\perp$ .



## What does an autoassociative NN do?

Training set  $A_{\text{train}} = \{\mathbf{x}^{(p)}, p = 1, 2, \dots, N\}$  forms the linear manifold  $\mathcal{L}$ .

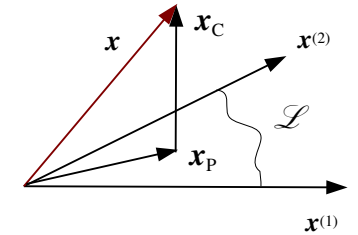
NN considers every departure  $\mathbf{x}$  from  $\mathcal{L}$  as added noise that needs to be filtered out by projecting  $\mathbf{x}$  to  $\mathcal{L}$ :

We need to show that output  $\mathbf{W}\mathbf{x} = \mathbf{X}\mathbf{X}^+\mathbf{x} = \mathbf{x}_p$  (filtered version of  $\mathbf{x}$ ), i.e. that operator  $\mathbf{W} = \mathbf{X}\mathbf{X}^+$  makes an **orthogonal projection** to  $\mathcal{L}$ .

Alternatively, the NN model with operator  $\mathbf{W} = \mathbf{I} - \mathbf{X}\mathbf{X}^+$  is called **novelty detector**, where  $\mathbf{W}\mathbf{x} = \mathbf{x}_c \in \mathcal{L}^\perp$ .

Now assume: you learned  $N$  patterns, and want to add  $(N+1)$ -st pattern.

How to change  $\mathbf{W}$  efficiently?



5

6

## Gram-Schmidt orthogonalization process

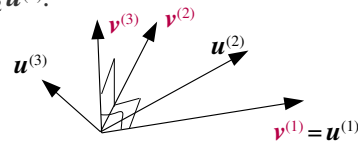
Let's have a base  $\mathbf{u}^{(1)}, \mathbf{u}^{(2)}, \dots, \mathbf{u}^{(k)} \in \mathcal{L}$ , for which we want to create an orthogonal base  $\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(k)} \in \mathcal{L}$ .

Procedure:

Set  $\mathbf{v}^{(1)} = \mathbf{u}^{(1)}$ . In space with base  $\mathbf{v}^{(1)}, \mathbf{u}^{(2)}$  let's find vector  $\mathbf{v}^{(2)}$  such that

$\mathbf{v}^{(2)} \perp \mathbf{v}^{(1)}$ . Hence,  $\mathbf{v}^{(2)} = a_1 \mathbf{v}^{(1)} + a_2 \mathbf{u}^{(2)}$ .

$$\mathbf{v}^{(2)} = \mathbf{u}^{(2)} - \frac{\mathbf{v}^{(1)T} \mathbf{u}^{(2)}}{|\mathbf{v}^{(1)}|^2} \mathbf{v}^{(1)}$$



Recursive formula: we have  $\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(k-1)}$  and compute  $\mathbf{v}^{(k)}$  such that

$\mathbf{v}^{(k)} \perp \mathbf{v}^{(i)}, i = 1, 2, \dots, k-1$

$$\mathbf{v}^{(k)} = \mathbf{u}^{(k)} - \sum_{i=1}^{k-1} \frac{\mathbf{v}^{(i)T} \mathbf{u}^{(k)}}{|\mathbf{v}^{(i)}|^2} \mathbf{v}^{(i)}$$

How to use this recursion for a GI model?

## General Inverse model

We have patterns  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$  and the associated orthogonal base (via Gram-Schmidt process)  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)}$

$\mathbf{W}$  is computed recursively, upon adding a new input  $\mathbf{x}$ .

1. Initialize  $\mathbf{W}^{(0)} = \mathbf{0}, i = 1$

2.  $\mathbf{z}^{(i)} = \mathbf{x}^{(i)} - \mathbf{W}^{(i-1)} \mathbf{x}^{(i)}$

$$\mathbf{W}^{(i)} = \mathbf{W}^{(i-1)} + \frac{\mathbf{z}^{(i)} \mathbf{z}^{(i)T}}{|\mathbf{z}^{(i)}|^2}$$

$i = i+1$

7

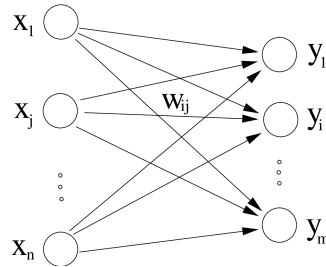
8

# Correlation Matrix Memory

Let's have train set:  $A_{\text{train}} = \{(\mathbf{x}^{(p)}, \mathbf{y}^{(p)}), p = 1, 2, \dots, N\}$ .

and let weights be set using Hebb's (1949) hypothesis:

$$w_{ij} = \sum_{p=1}^N x_j^{(k)} y_i^{(k)} \quad \mathbf{W} = \sum_{p=1}^N \mathbf{y}^{(p)} \mathbf{x}^{(p)T} = \mathbf{YX}^T$$



Let's focus again on auto-associative case:  $\mathbf{W} = \mathbf{XX}^T$ .

$$x_i = \sum_{p=1}^N w_{ij} x_j^{(p)}$$

Recursive computation:

$$\mathbf{W}(N+1) = \mathbf{W}(N) + \mathbf{x}^{(N+1)} \mathbf{x}^{(N+1)T}$$

9

# Cross-talk in CMM

Response to input  $\mathbf{x}^{(p)}, p \in \{1, 2, \dots, N\}$ :

$$\mathbf{Wx}^{(p)} = \mathbf{XX}^T \mathbf{x}^{(p)} = \sum_{k=1}^N \mathbf{x}^{(k)} \mathbf{x}^{(k)T} \mathbf{x}^{(p)} = \mathbf{x}^{(p)} \mathbf{x}^{(p)T} \mathbf{x}^{(p)} + \sum_{k=1; k \neq p}^N \mathbf{x}^{(k)} \mathbf{x}^{(k)T} \mathbf{x}^{(p)}$$

$$\mathbf{Wx}^{(p)} = \mathbf{x}^{(p)} |\mathbf{x}^{(p)}|^2 + C(p) \leftarrow \text{noise vector due to "cross-talk" from other inputs}$$

If the inputs are orthogonal, then  $C(p) = 0$  and  $\mathbf{X}^T = \mathbf{X}^{-1} = \mathbf{X}^+$ , i.e. CMM = GI.

Max. capacity of both GI and CMM equals  $n$ .

For sufficiently dissimilar inputs, CMM is a good faster alternative to GI.

10

# Interpretation of CMM behavior

Let's consider inputs  $\mathbf{x}^{(1)} \mathbf{x}^{(2)} \dots \mathbf{x}^{(N)}$  as realizations of a random variable

$X = (X_1, X_2, \dots, X_n)^T$ , with zero means, i.e.  $\langle X_i \rangle = 0, i = 1, 2, \dots, n$ . Then

$$w_{ij} = \sum_{k=1}^N x_j^{(k)} x_i^{(k)} \propto \text{cov}(X_i, X_j) \quad [\text{unbiased estimate}]$$

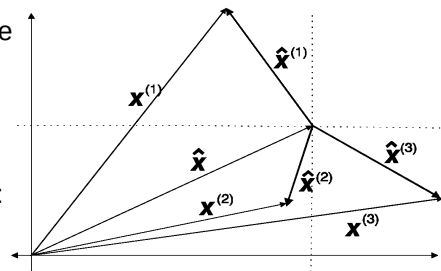
$$\text{cov}(X_i, X_j) = \langle (X_i - \langle X_i \rangle) \cdot (X_j - \langle X_j \rangle) \rangle = \langle X_i X_j \rangle$$

$w_{ij}$  carries information about linear correlations between nodes  $i$  and  $j$ .

Reduced (noisy) input component  $x_i$  can be recovered (bootstrapped) by other active components strongly correlated with  $i$ .

Shifting coordinates decreases "cross-talk":

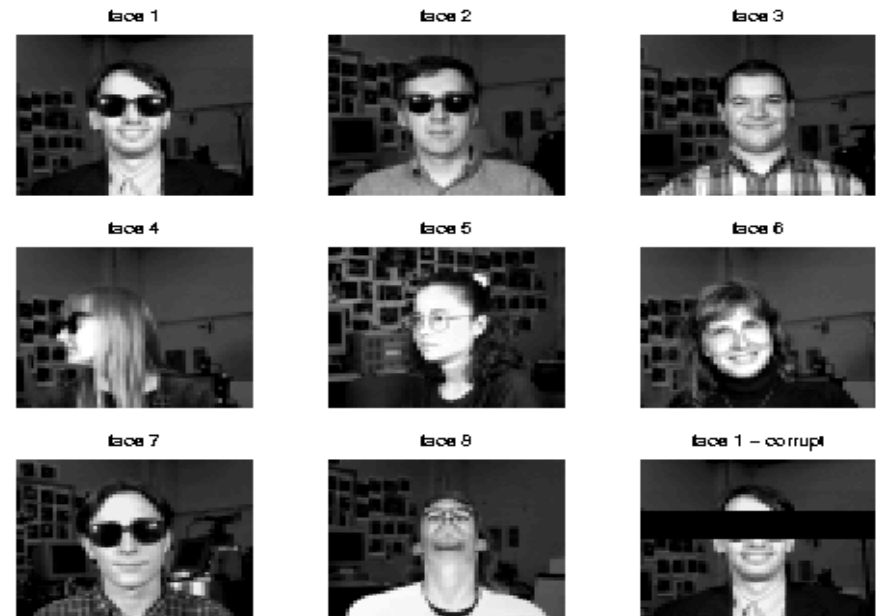
$$\hat{\mathbf{x}}^{(i)} = \mathbf{x}^{(i)} - \frac{1}{N} \sum_{k=1}^N \mathbf{x}^{(k)}$$



(Kvasnička et al, 1997)

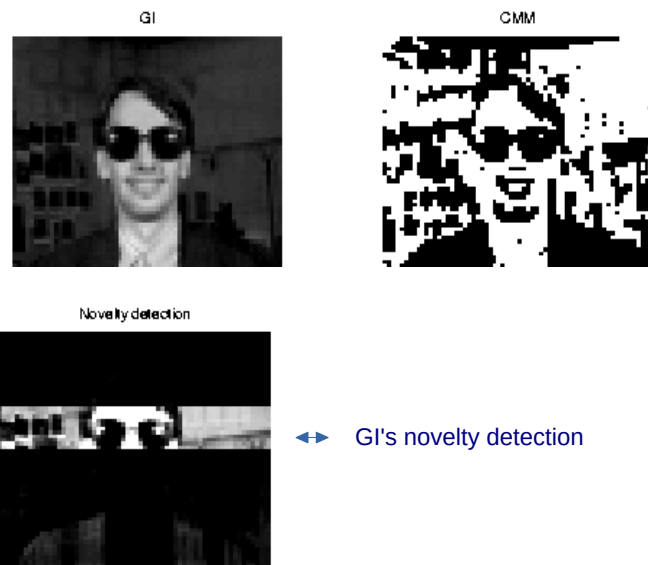
11

# Example: 8 faces from CMU image data repository



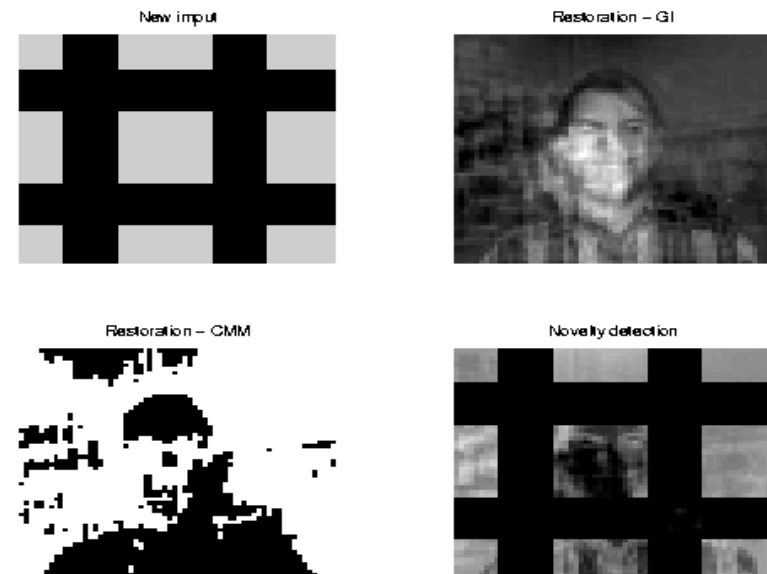
12

## Recall by GI and CMM



13

## Response to a new input



14

## Summary

- Linear models were studied during connectionist depression in the 1970s
- Single layer models as auto-associative memories
- Analytic solutions possible
- General inverse model – noise filtering by projection to linear manifold (of the training data)
- GI – as novelty detector
- Correlation Matrix Memory – Hebbian-based learning, subject to cross-talk
- GI better in general, for sufficiently dissimilar inputs both models are comparable

15