# ENGR-E 533 "Deep Learning Systems" Lecture 04: The First Layer

## Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: http://minjekim.com
Research Group: http://saige.sice.indiana.edu
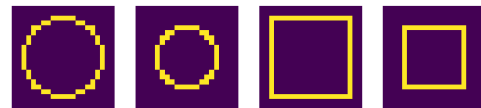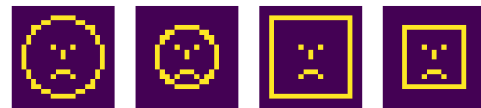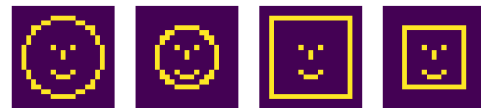Meeting Request: http://doodle.com/minje

INDIANA UNIVERSITY
**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**

# Unsupervised Feature Learning

## - What do we do in the first layer?

- Before deep learning, a common practice was to extract features and then learn a supervised model
  - Source separation
    - Convert the time domain audio signals into matrices by using Short-Time Fourier Transform (STFT) and then learn dictionaries
  - Object recognition
    - Extract a bunch of different features (e.g. HoG, SIFT, etc) and then build a classifier
  - Sentiment analysis
    - Preprocess the text, learn topics, and then build a classifier
  - Speech recognition
    - Extract Mel-Frequency Cepstrum Coefficients (MFCC) and then learn Hidden Markov Models
- YALT
  - In our baby facial expression recognition problem we manually extracted the features first
  - And then built a softmax classifier

# Unsupervised Feature Learning
## - What do we do in the first layer?

o Today we will learn how to systematically learn those features from raw data
  - We call this procedure **unsupervised feature learning**
  - Don't worry, we're not getting away from neural networks

o Why unsupervised?
  - Supervision here means human intervention to solve the problem
    - e.g. When you train a classifier, you need a bunch of pairs: a data sample (facial image) and its label (happy or sad)
    - You let the model know how you think of the problem
    - Mathematically, this can be done by either learning the mapping function between the data sample and its label
      $$\boldsymbol{y} = \mathcal{F}(\boldsymbol{x}) \qquad \hat{\boldsymbol{y}} = \mathcal{G}(\boldsymbol{x}\ ; \mathbb{W})$$
    - Or, finding the parameters that maximize the *a posteriori* probability
      $$P(\boldsymbol{y}|\boldsymbol{x}; \boldsymbol{\Theta})$$
  - Unsupervised learning
    - You don't have the human intervention, or the labels
    - Then, what are we learning?
    - A model that best describes the data

# Unsupervised Feature Learning

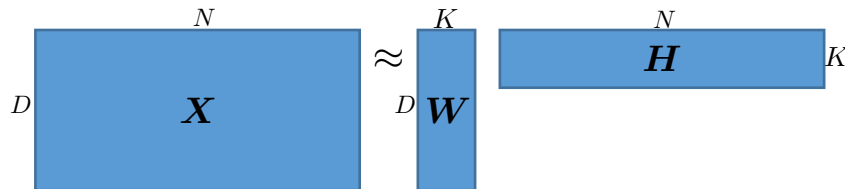## - What do we do in the first layer?

o Modeling the data

- If we use a probabilistic distribution it's to find out the maximum likelihood solution

$$P(\boldsymbol{x}; \boldsymbol{\Theta})$$

- Often, this can be viewed as a latent variable analysis with $K$ latent variables, too

The data matrix $\longrightarrow \boldsymbol{X} \approx \boldsymbol{W}\boldsymbol{H} \longleftarrow$ The activation of the latent variable

The representative vectors of the latent variable



o For example

- If you find the eigenvectors of $\boldsymbol{X}\boldsymbol{X}^\top$, that will correspond to $\boldsymbol{W}$
  - $\boldsymbol{W}^\top\boldsymbol{W} = \boldsymbol{I}$ and, therefore, $\boldsymbol{W}^\top\boldsymbol{X} = \boldsymbol{H}$. Also, rows of $\boldsymbol{H}$ are ordered in their variances
  - Principal Component Analysis
  - $\boldsymbol{H}$ is a set of lower dimensional features that can still describe the original distribution
- If $\boldsymbol{X}_{:,n} \sim \sum_k \boldsymbol{H}_{k,n}\mathcal{N}(\boldsymbol{W}_{:,k}, \boldsymbol{\Sigma}_{:,:,k})$ and $\boldsymbol{H}_{:,n}$ is a one hot vector
  - Gaussian Mixture Model (a.k.a. Vector Quantization)
- If you find $\boldsymbol{W}$ and $\boldsymbol{H}$ that minimize the approximation error, but are nonnegative at the same time
  - Nonnegative Matrix Factorization
  - $\boldsymbol{H}$ gives you the parts-based representation
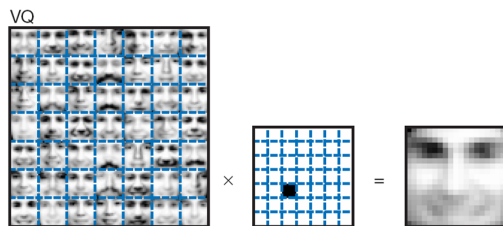
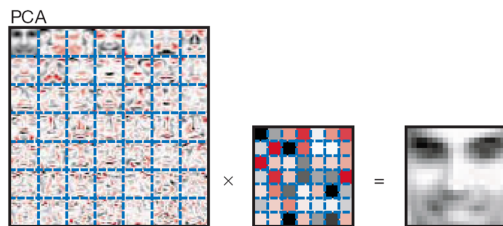# Unsupervised Feature Learning

- What do we do in the first layer?



NMF estimates parts-based representations, something like Lego blocks.
Reconstruction is a linear combination of them, but subtraction is not allowed.

VQ finds a bunch of cluster means.
Reconstruction is choosing the most similar mean.

PCA finds the holistic eigenfaces.
From the important one down to the subtle ones.

INDIANA UNIVERSITY
**SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING**

# Unsupervised Feature Learning

## - What do we do in the first layer?
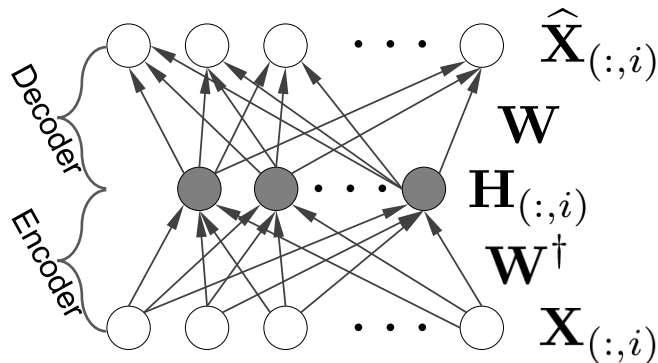
○ There are so many different ways to solve the linear approximation $X \approx WH$

○ So, eventually it's all about how to constrain the problem

$$\arg\min_{W,H} \mathcal{D}(X||WH) + \lambda_W f(W) + \lambda_H g(H)$$

- For clustering, you want a super sparse column vectors in $H$ so that the data sample is one of the cluster means, deviated accordingly
- For PCA, you want the after-projection-samples are with maximal variances
- For NMF, you constrain $W, H$ to be nonnegative (element-wise)
- ...

○ It is convenient to assume that there exists some kind of pseudo inverse of $W$
- For PCA, $W^{\dagger} = W^{\top}$
- For NMF, $W^{\dagger} = W^{\top}$ depending on the choice of the error function $\quad H \leftarrow H \odot \dfrac{W^{\top} X}{W^{\top} W H}$

○ If you assume this pseudo inverse of the basis vectors, you can think of the projection as a way to convert your data into features $\quad W^{\dagger} X \approx W^{\dagger} W H \approx H$

# The Autoencoders

- A unified neural network-based representation of unsupervised learning



$$\underset{\boldsymbol{W},\boldsymbol{W}^\dagger,\boldsymbol{H}}{\arg\min} \ \mathcal{D}(\boldsymbol{X}\|\boldsymbol{W}\boldsymbol{W}^\dagger\boldsymbol{X}) + \lambda_{\boldsymbol{W}}f(\boldsymbol{W}) + \lambda_{\boldsymbol{W}^\dagger}f(\boldsymbol{W}^\dagger) + \lambda_{\boldsymbol{H}}g(\boldsymbol{H})$$

o Once the objective function is setup in this way, you can estimate the parameters using SGD
  • TF or PT will take care of this part using automatic gradient calculation

# The Autoencoders

## - A unified neural network-based representation of unsupervised learning

o Demo
  - Sparse coding
  - NMF

# Nonlinearity in Neural Networks
## - Stacked linear models form another linear model

o So far I haven't said anything about nonlinearity in the lower layers (at least officially)

o Why do we need it?
  • Because stacked linear models form yet another linear model

o Stacked linear models

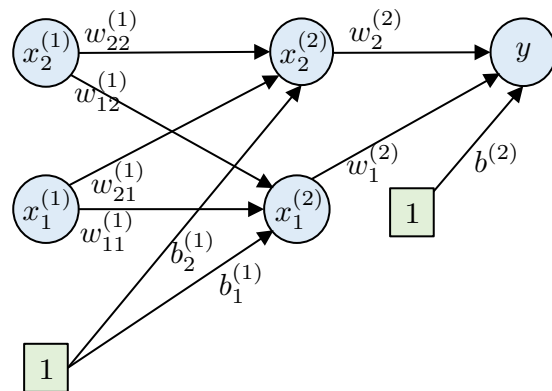$$x_1^{(2)} = w_{11}^{(1)} x_1^{(1)} + w_{12}^{(1)} x_2^{(1)} + b_1^{(1)}$$
$$x_2^{(2)} = w_{21}^{(1)} x_1^{(1)} + w_{22}^{(1)} x_2^{(1)} + b_2^{(1)}$$
$$y = w_1^{(2)} x_1^{(2)} + w_2^{(2)} x_2^{(2)} + b^{(2)}$$

o Form another linear model

$$y = w_1^{(2)} w_{11}^{(1)} x_1^{(1)} + w_1^{(2)} w_{12}^{(1)} x_2^{(1)} + w_1^{(2)} b_1^{(1)}$$
$$+ w_2^{(2)} w_{21}^{(1)} x_1^{(1)} + w_2^{(2)} w_{22}^{(1)} x_2^{(1)} + w_2^{(2)} b_2^{(1)} + b^{(2)}$$
$$= \left( w_1^{(2)} w_{11}^{(1)} + w_2^{(2)} w_{21}^{(1)} \right) x_1^{(1)} + \left( w_1^{(2)} w_{12}^{(1)} + w_2^{(2)} w_{22}^{(1)} \right) x_2^{(1)} + w_1^{(2)} b_1^{(1)} + w_2^{(2)} b_2^{(1)} + b^{(2)}$$
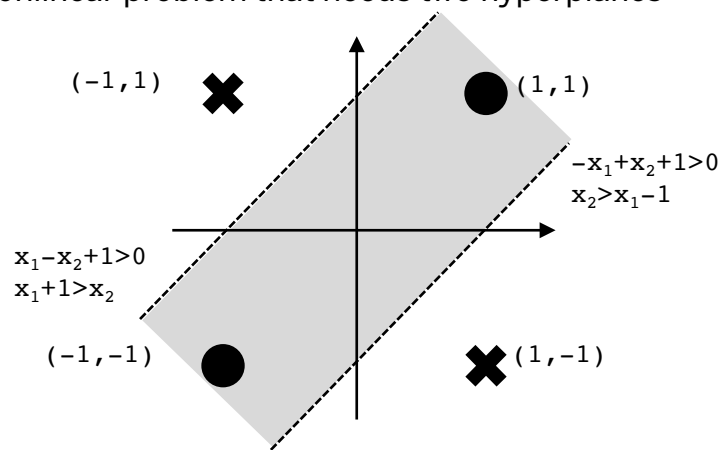
  • Or  $y = \boldsymbol{w}^{(2)} (\boldsymbol{W}^{(1)} \boldsymbol{x} + \boldsymbol{b}^{(1)}) + b^{(2)}$
  $$= \boldsymbol{w}^{(2)} \boldsymbol{W}^{(1)} \boldsymbol{x} + \boldsymbol{w}^{(2)} \boldsymbol{b}^{(1)} + b^{(2)}$$
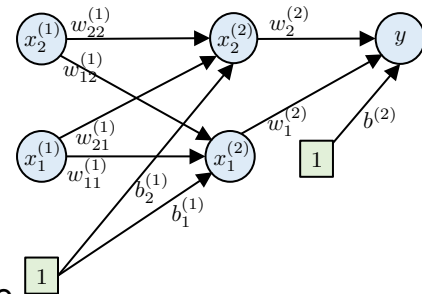
# Nonlinearity in Neural Networks

## - Stacked linear models form another linear model

o What can we do then?

o The XOR problem

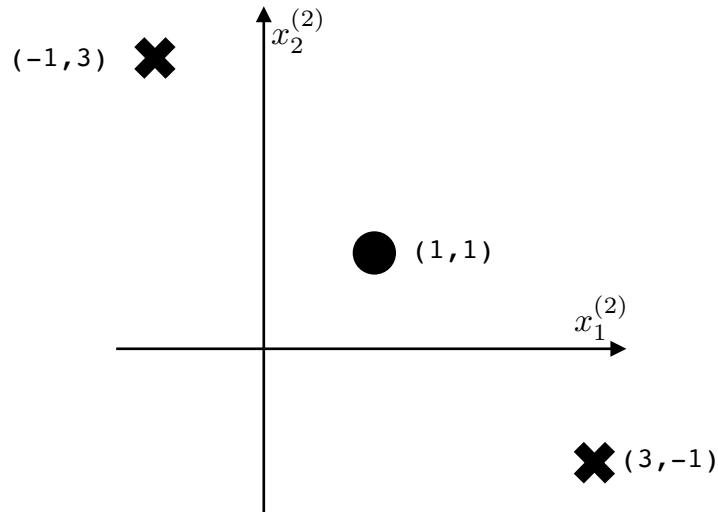• A nonlinear problem that needs two hyperplanes

$$(-1,1) \quad \times \qquad \bullet \quad (1,1)$$

$$-x_1+x_2+1>0$$
$$x_2>x_1-1$$

$$x_1-x_2+1>0$$
$$x_1+1>x_2$$

$$(-1,-1) \quad \bullet \qquad \times \quad (1,-1)$$

o How do you transform the data into features?

• It must have something to do with hyperplanes

$$\boldsymbol{W}^{(1)} = \left[ \begin{array}{ccc} +1 & -1 & 1 \\ -1 & +1 & 1 \end{array} \right]$$

o The new feature space

$$x_2^{(2)}$$

$$(-1,3) \quad \times$$

$$\bullet \quad (1,1)$$

$$x_1^{(2)}$$
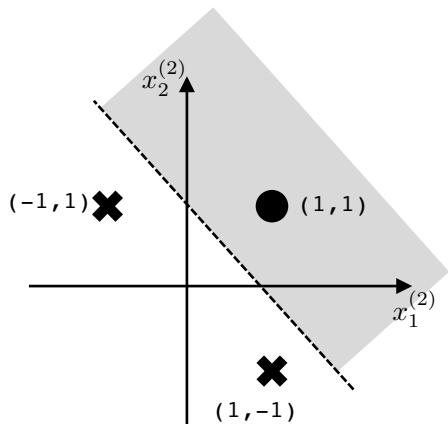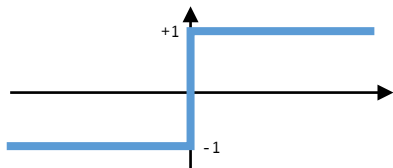
$$\times \quad (3,-1)$$

o What happens in the feature space?

• If it were not for the help from the nonlinearity

• Still NOT linearly solvable!

# Nonlinearity in Neural Networks

## - Adding nonlinearity

○ With a sign function

$$x_i^{(2)} = \text{sign}(\boldsymbol{W}_{i,:}^{(1)} \boldsymbol{x} + b_i^{(1)})$$

○ With a Rectified Linear Unit (ReLU)

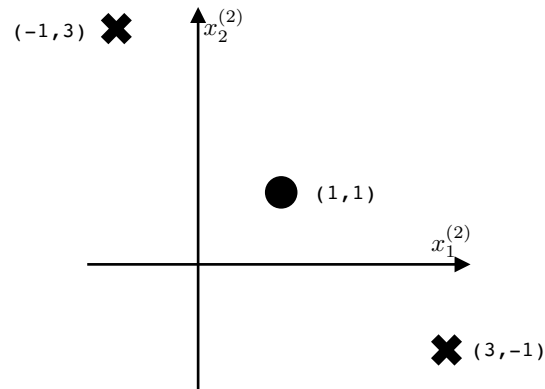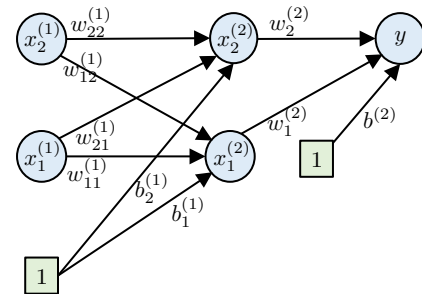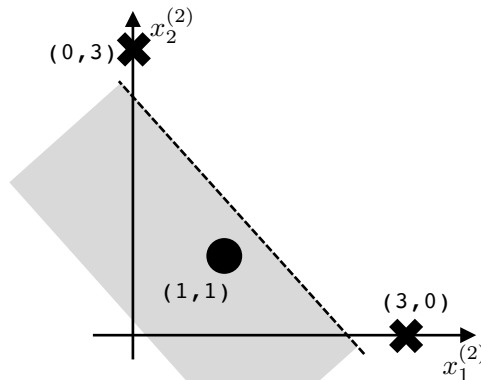$$x_i^{(2)} = \max(\boldsymbol{W}_{i,:}^{(1)} \boldsymbol{x} + b_i^{(1)}, 0)$$

Now they are linearly solvable!

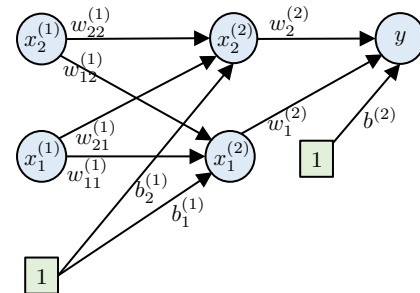New hyperplane corresponds to the last layer weights

# Nonlinearity for AEs?

- What does the XOR example mean for unsupervised learning?



o You can replace the final layer binary output with the input vector
- Classifier turns into an AE

o Suppose you can estimate these weights of the XOR network by using BP

o Can you also estimate these weights with the AE setup?

o You just saw that nonlinearity can help producing better features

o AE might benefit from it to

$$\underset{\boldsymbol{W}, \boldsymbol{W}^{\dagger}, \boldsymbol{H}}{\arg\min} \ \mathcal{D}(\boldsymbol{X}||\boldsymbol{W}\sigma(\boldsymbol{W}^{\dagger}\boldsymbol{X})) + \lambda_{\boldsymbol{W}}f(\boldsymbol{W}) + \lambda_{\boldsymbol{W}^{\dagger}}f(\boldsymbol{W}^{\dagger}) + \lambda_{\boldsymbol{H}}g(\boldsymbol{H})$$

o But it's not clear if the encoder weights are going to be actually create useful features for the following classification

o We will cover this part in the next time

# The Pipeline for Supervised Learning
## - What people have done so far before deep learning

- Feature engineering
  - Tries to come up with a set of features that best describe the data and the problem
  - There can be many different ways (as you've seen)
  - You never know the quality until you actually test out those features for your problem

- Supervised learning (classification/regression)
  - Tries to come up with the best model to best predict the output variable
  - Even if you use a very nice set of features, nonlinearity can be still involved (e.g. kernel methods for SVM)
  - If you don't like the performance, you never know what to blame
    - It could be because of either bad feature engineering or classifier

- A holistic approach
  - In a multilayer perceptron, you have the first layer dedicated to the feature extraction
  - And the last layer as your supervised learning part
  - What if you just optimize the weights of both layers all together?
    - Then, you can skip the feature engineering part
  - Does this work? When does this approach not work? Why?

# Thank You!