# CPSC 340:
# Machine Learning and Data Mining

Deep Learning
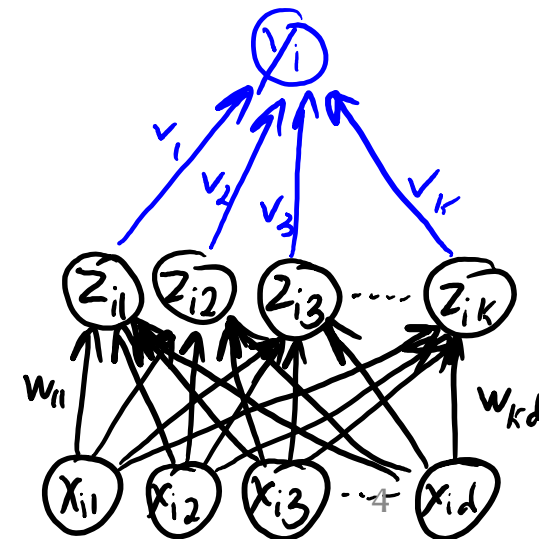
Summer 2021

# In This Lecture

1. Encoder-Predictor Learning
   - aka deep learning
2. Artificial Neural Networks
3. "Biological" Motivations for Deep Learning
4. History of Deep Learning

# Part 5: Deep Learning

# Supervised Learning Roadmap

- Part 1: "Direct" Supervised Learning.
  - We learned parameters 'w' based on the original features $x_i$ and target $y_i$.
- Part 3: Change of Basis.
  - We learned parameters 'v' based on a change of basis $z_i$ and target $y_i$.
- Part 4: Latent-Factor Models.
  - We learned parameters 'W' for basis $z_i$ based on only on features $x_i$.
  - You can then learn 'v' based on change of basis $z_i$ and target $y_i$.
- Part 5: Neural Networks.
  - Jointly learn 'W' and 'v' based on $x_i$ and $y_i$.
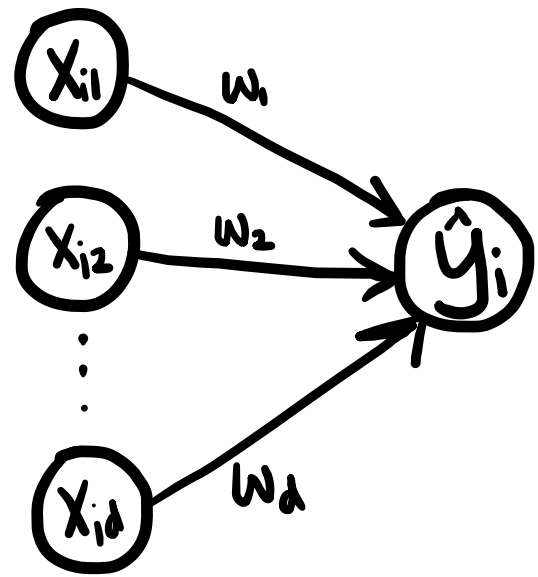  - **Learn features $z_i$ that is good for supervised learning.**

Coming Up Next

# ENCODER-PREDICTOR LEARNING

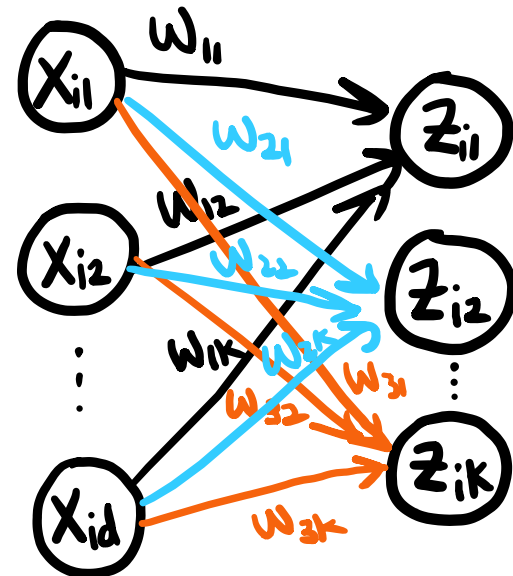# "Graph" View of Matrix Multiplication

$$\hat{y}_i = W^T x_i$$

- "target" node
- edge weights
- "input" nodes



$$\hat{y}_i = [w_1 \; w_2 \cdots w_d] \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}$$

$$z_i = W x_i$$

- "output" nodes
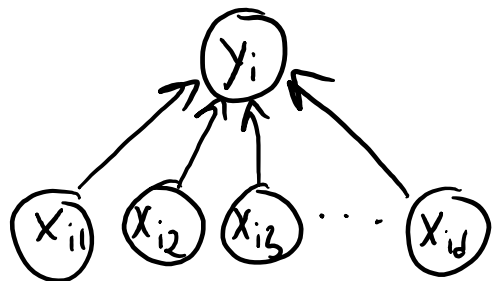- edge weights
- "input" nodes



$$\begin{bmatrix} z_{i1} \\ z_{i2} \\ \vdots \\ z_{ik} \end{bmatrix}_k = \begin{bmatrix} -w_1- \\ -w_2- \\ \vdots \\ -w_k- \end{bmatrix}_d \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}$$
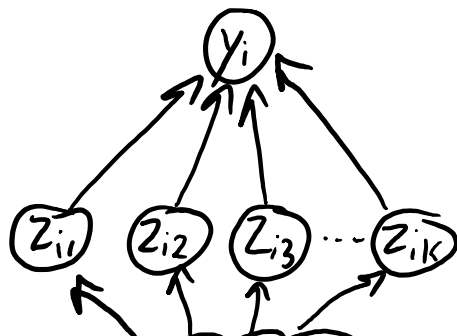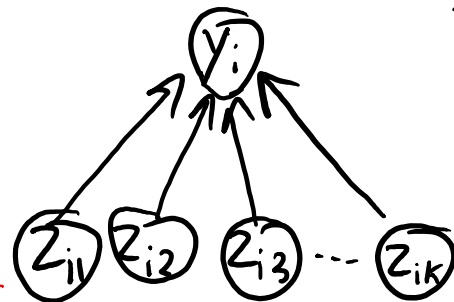
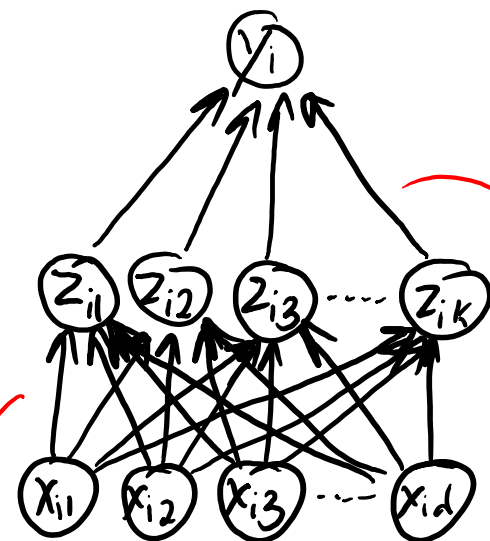# A Graphical Summary of CPSC 340 Parts 1-5



Part 1: "I have features $x_i$"

Part 2: "What is the group of $x_i$?"

Part 3: change of basis

"I think this basis will work"

Part 4: basis from latent-factor model

"PCA will give me good features"

"What are the 'parts' of $x_i$?"

Part 5: Neural networks

Trained separately

Learn features and classifier at the same time.

# Recall: Encoder Learning

$$(X, \require{cancel}\cancel{\phantom{Y}}) \rightarrow E$$

Learned encoder

$$X_i \rightarrow E \rightarrow Z_i$$

$d \times 1$

$k \times 1$

Example i in
learned feature space

# "Encoder-Predictor Learning"

- **Encoder-Predictor learning** problem:
  - Input:  Labeled examples
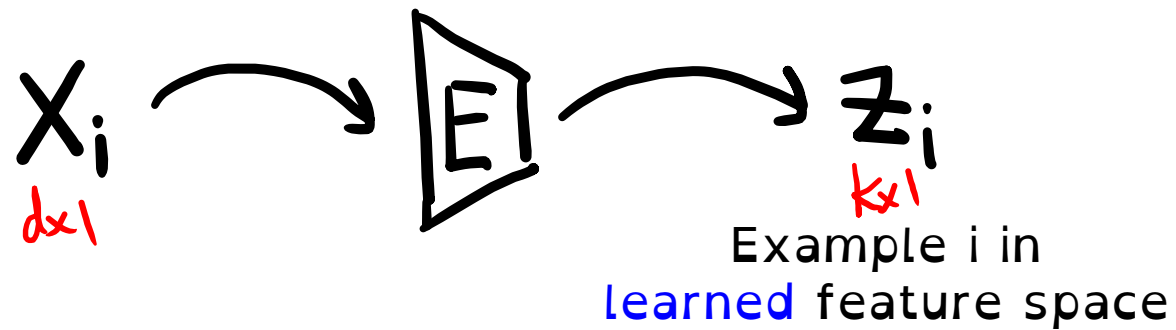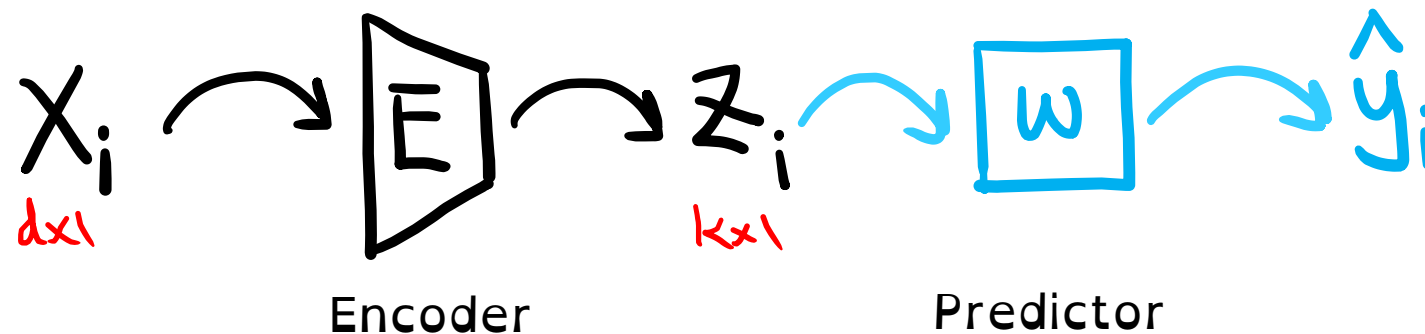  - Output:  Encoder E and predictor w

$$(X, y) \rightarrow \boxed{E} \quad \boxed{w}$$

Encoder  Predictor

- Using learned encoder and predictor:

$$X_i \rightsquigarrow \boxed{E} \rightsquigarrow Z_i \rightsquigarrow \boxed{w} \rightsquigarrow \hat{y}_i$$

$d \times 1$ $\qquad k \times 1$

Encoder  Predictor

# "Artificial Neural Networks"



$$X_i \rightarrow \boxed{\text{Mat Mul}} \rightarrow Z_i \rightarrow \boxed{\text{Linear Model}} \rightarrow \hat{y}_i$$

- "Artificial neural network" := encoder-predictor model using matrix multiplication for encoder
  - Must use non-linear activations (soon)
  - Usually use linear model as predictor
- "Deep neural network" := artificial neural network that uses more than one matrix multiplication

10

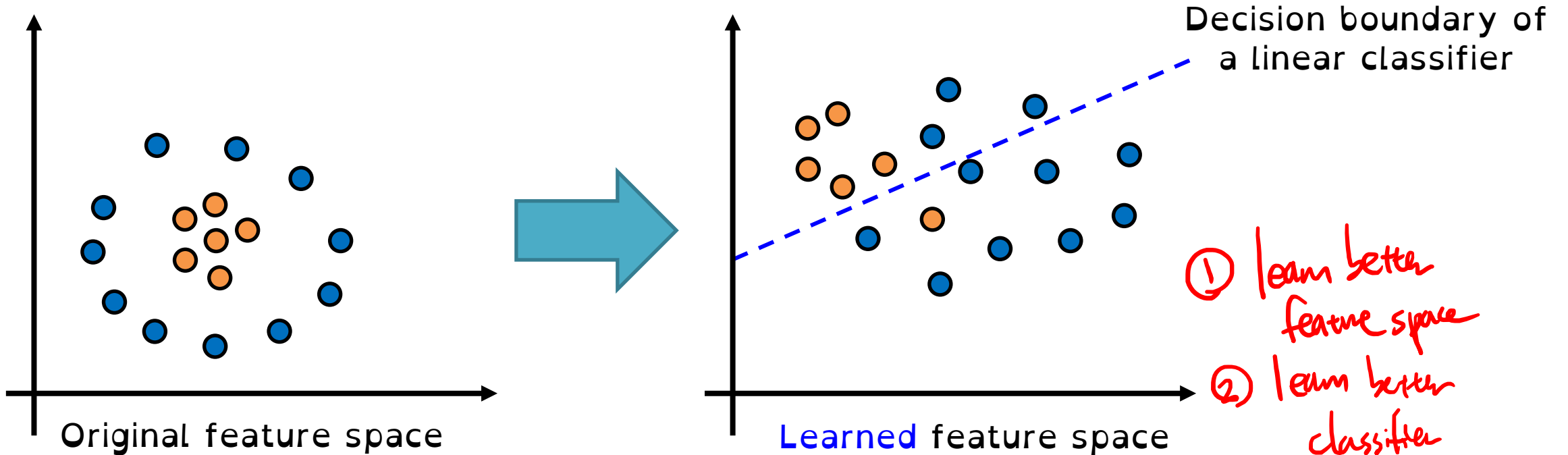# Visualizing Encoder-Predictor



Original feature space

Decision boundary of a linear classifier

Learned feature space

① learn better feature space

② learn better classifier

Q: How can we make
our model better than this?

# Visualizing Encoder-Predictor



Decision boundary of a linear classifier

Original feature space

Learned feature space

- Intuition: supervised latent factor model
  - Loss function based on labels
  - Encourage encoder to produce more linearly separable results

Coming Up Next

# MORE FORMAL DETAILS ON NEURAL NETWORKS

# Notation for Neural Networks (MEMORIZE)

We have our usual supervised learning notation:

$$X = \begin{bmatrix} \underline{\quad} & x_1^T & \underline{\quad} \\ \underline{\quad} & x_2^T & \underline{\quad} \\ & \vdots & \\ \underline{\quad} & x_n & \underline{\quad} \end{bmatrix} \qquad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$n \times d$ $\qquad$ $n \times 1$
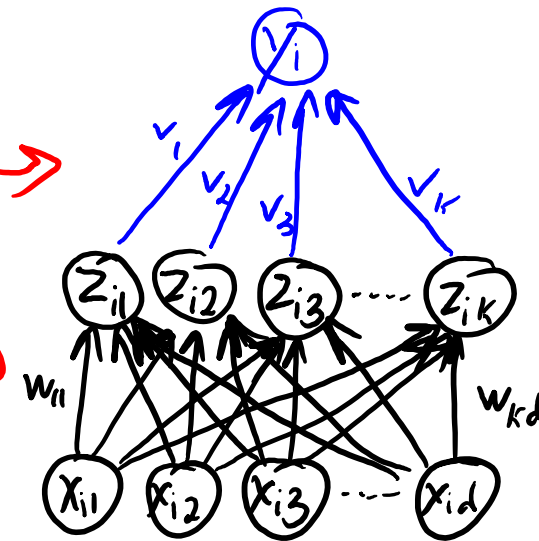
We have our latent features: $\quad$ We have two sets of parameters:

$$Z = \begin{bmatrix} \underline{\quad} & z_1^T & \underline{\quad} \\ \underline{\quad} & z_2^T & \underline{\quad} \\ & \vdots & \\ \underline{\quad} & z_n^T & \underline{\quad} \end{bmatrix}$$

$n \times K$

$$V = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_K \end{bmatrix} \qquad W = \begin{bmatrix} \underline{\quad} & w_1 & \underline{\quad} \\ \underline{\quad} & w_2 & \underline{\quad} \\ & \vdots & \\ \underline{\quad} & w_K & \underline{\quad} \end{bmatrix}$$

$k \times 1$ $\qquad\qquad k \times d$

# Linear-Linear Neural Net

- Obvious choice: linear latent-factor encoder with linear regression predictor

Use features from latent-factor model: $z_i = W x_i$

Make predictions using a linear model: $y_i = v^T z_i$

- We want to train 'W' and 'v' jointly, so we could minimize:

$$f(W, v) = \frac{1}{2} \sum_{i=1}^{n} (v^T z_i - y_i)^2 = \frac{1}{2} \sum_{i=1}^{n} (v^T (W x_i) - y_i)^2$$

linear regression with $z_i$ as features

$z_i$ come from latent-factor model

Q: What can go wrong with this?

15

# Linear-Linear Neural Net

- Obvious choice: linear latent-factor encoder with linear regression predictor

Use features from latent-factor model: $z_i = W x_i$

Make predictions using a linear model: $y_i = v^T z_i$

- We want to train 'W' and 'v' jointly, so we could minimize:

$$f(W, v) = \frac{1}{2} \sum_{i=1}^{n} \underbrace{(v^T z_i - y_i)^2}_{\substack{\text{linear regression} \\ \text{with } z_i \text{ as features}}} = \frac{1}{2} \sum_{i=1}^{n} \left( v^T \underbrace{(W x_i)}_{\substack{z_i \text{ come from} \\ \text{latent-factor model}}} - y_i \right)^2$$

- This is just a linear model:

$$\hat{y}_i = v^T z_i = v^T (W x_i) = \underbrace{\overbrace{(v^T W)}^{1 \times d}}_{\text{some vector 'w'}} x_i = w^T x_i$$

# Introducing Non-Linearity

- To increase flexibility, something needs to be non-linear.
- Typical choice: transform $z_i$ by non-linear function 'h'.

$$z_i = W_{x_i} \qquad y_i = v^T h(z_i)$$

  - Here the function 'h' transforms 'k' inputs to 'k' outputs.
- Common choice for 'h': applying sigmoid function element-wise:
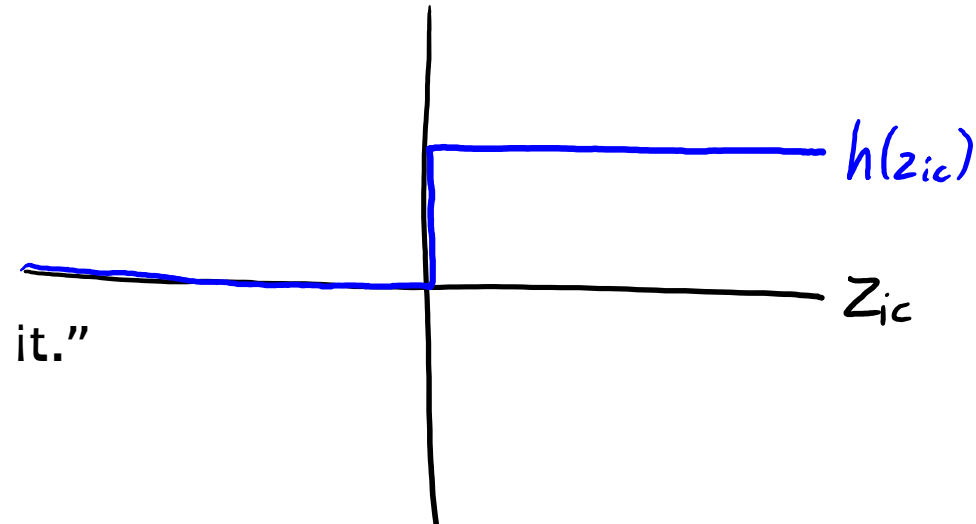
$$h(z_{ic}) = \frac{1}{1 + exp(-z_{ic})}$$

- So this takes the $z_{ic}$ in $(-\infty, \infty)$ and maps it to $(0, 1)$ .

# Why Sigmoid?

- Consider setting 'h' to define binary features $z_i$ using:

$$h(z_{ic}) = \begin{cases} 1 & \text{if } z_{ic} \geq 0 \\ 0 & \text{if } z_{ic} < 0 \end{cases}$$



- Each $h(z_i)$ can be viewed as binary feature.
  - "You either have this 'part' or you don't have it."
- We can make $2^k$ objects by all the possible "part combinations".



Motivation: Pixels vs. Parts

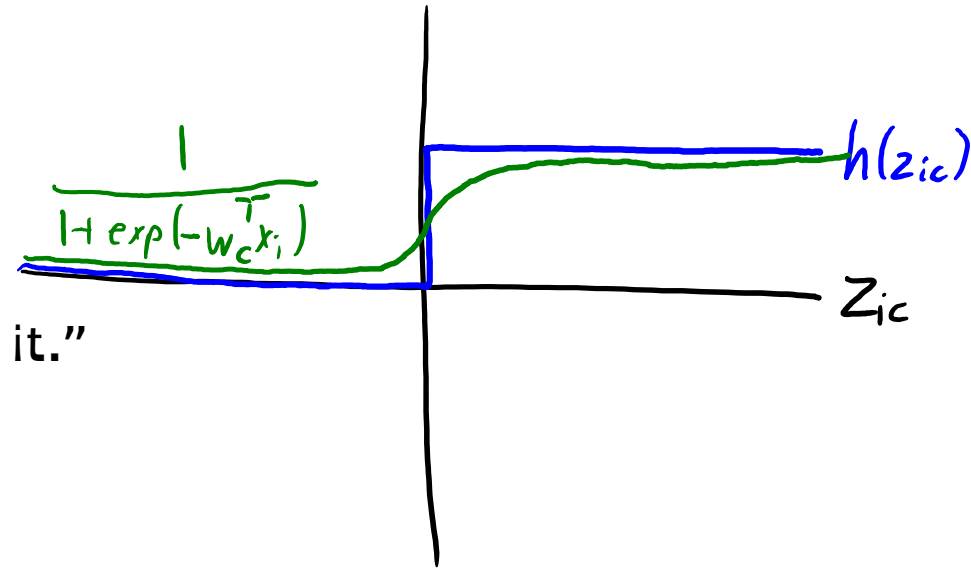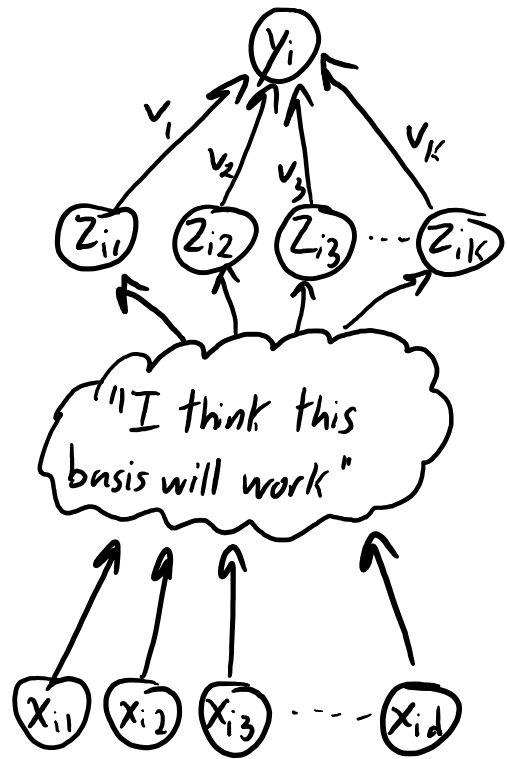- We could represent other digits as different combinations of "parts":

# Why Sigmoid?

- Consider setting 'h' to define binary features $z_i$ using:

$$h(z_{ic}) = \begin{cases} 1 & \text{if } z_{ic} \geq 0 \\ 0 & \text{if } z_{ic} < 0 \end{cases}$$



$$\frac{1}{1 + exp(-w_c^T x_i)}$$

$h(z_{ic})$

$z_{ic}$

$\mathbb{Z}$

$n \times k$

  - Each $h(z_i)$ can be viewed as binary feature.
    - "You either have this 'part' or you don't have it."

- But this is hard to optimize (non-differentiable, discontinuous).
- Sigmoid is a smooth approximation to these binary features.
  - Non-parametric version is a universal approximator:
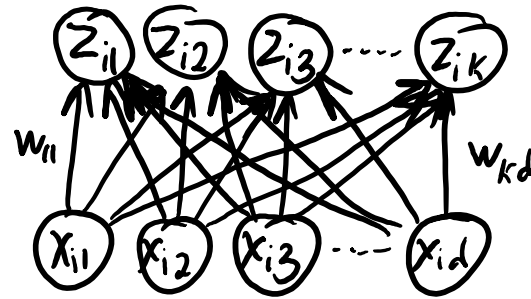    - If 'k' grows appropriately with 'n', can model any continuous function.

# Supervised Learning Roadmap



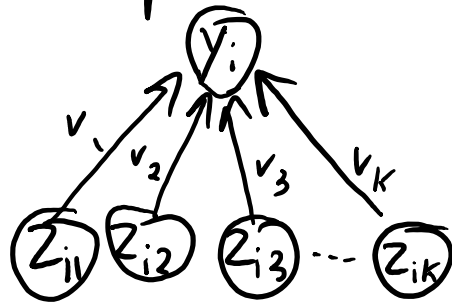Hand-engineered features:

"I think this basis will work"

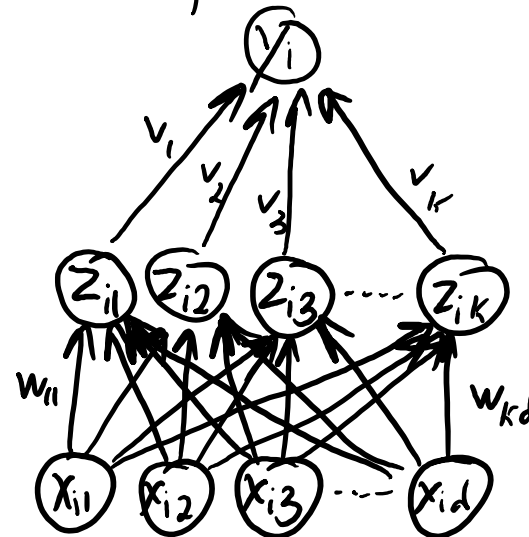Requires domain knowledge and can be time-consuming

Learn a latent-factor model:

Use latent features in supervized model:

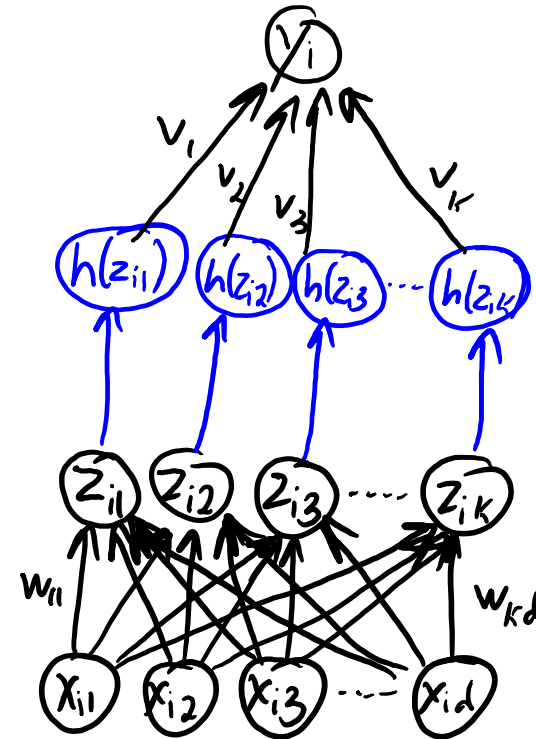Good representation of $x_i$ might be bad for predicting $y_i$

Learn 'w' and 'W' together:

But still gives a linear model.
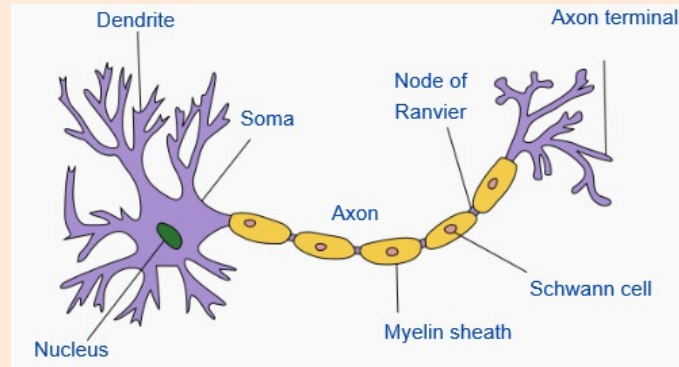
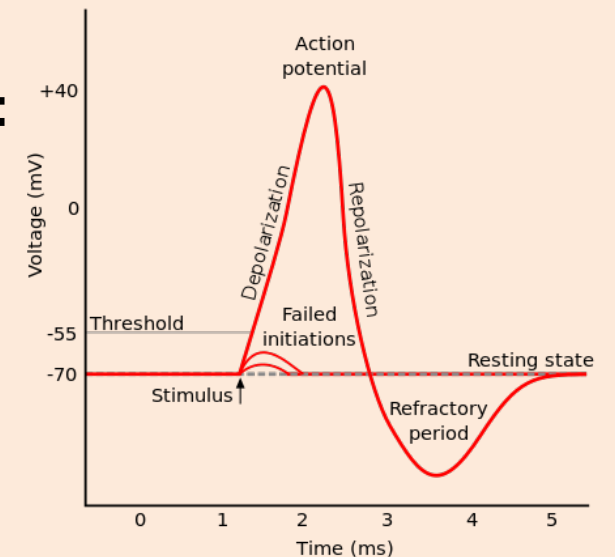Neural network:

Extra non-linear transformation 'h'

20

# (SUPPOSEDLY) BIOLOGICAL MOTIVATION FOR ARTIFICIAL NEURAL NETWORKS
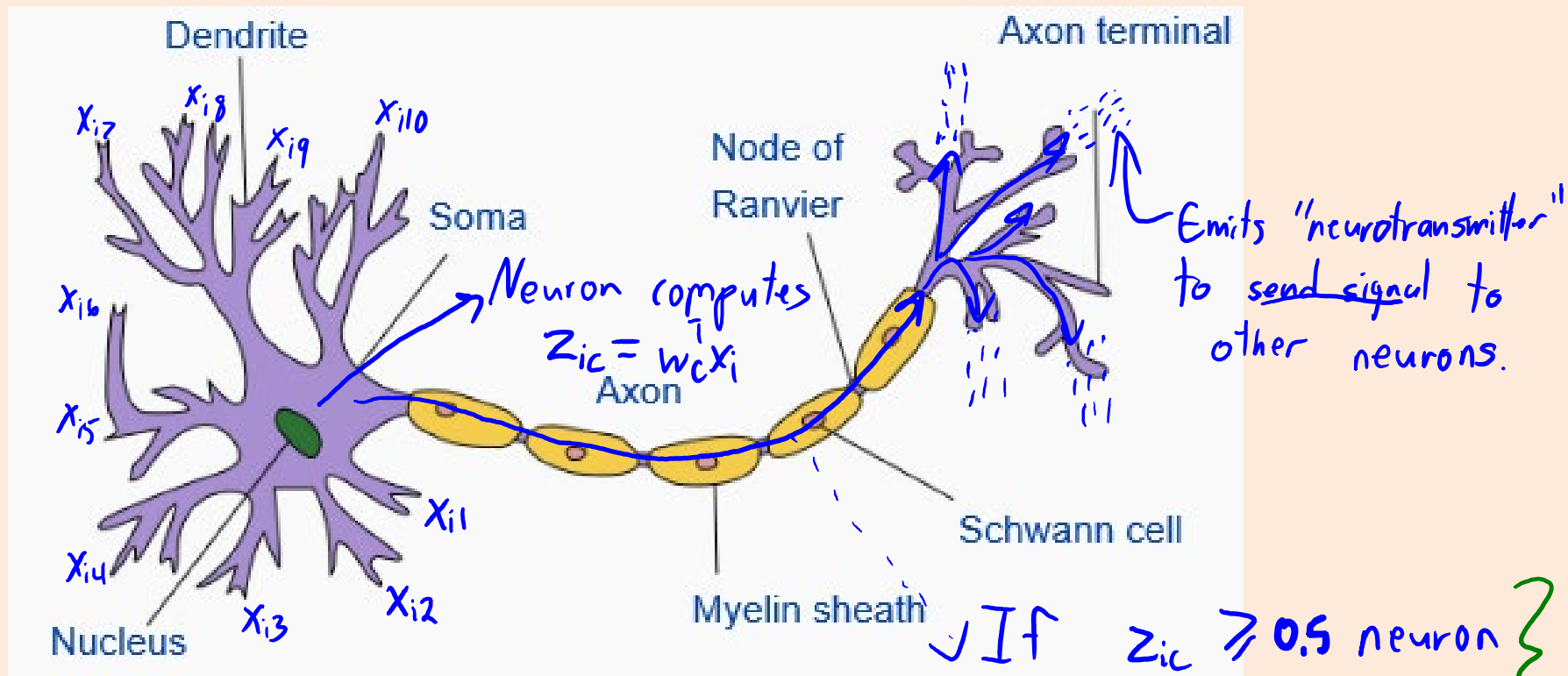
# Why "Neural Network"?

- Cartoon of "typical" neuron:



- Neuron has many "dendrites", which take an input signal.
- Neuron has a single "axon", which sends an output signal.
- With the right input to dendrites:
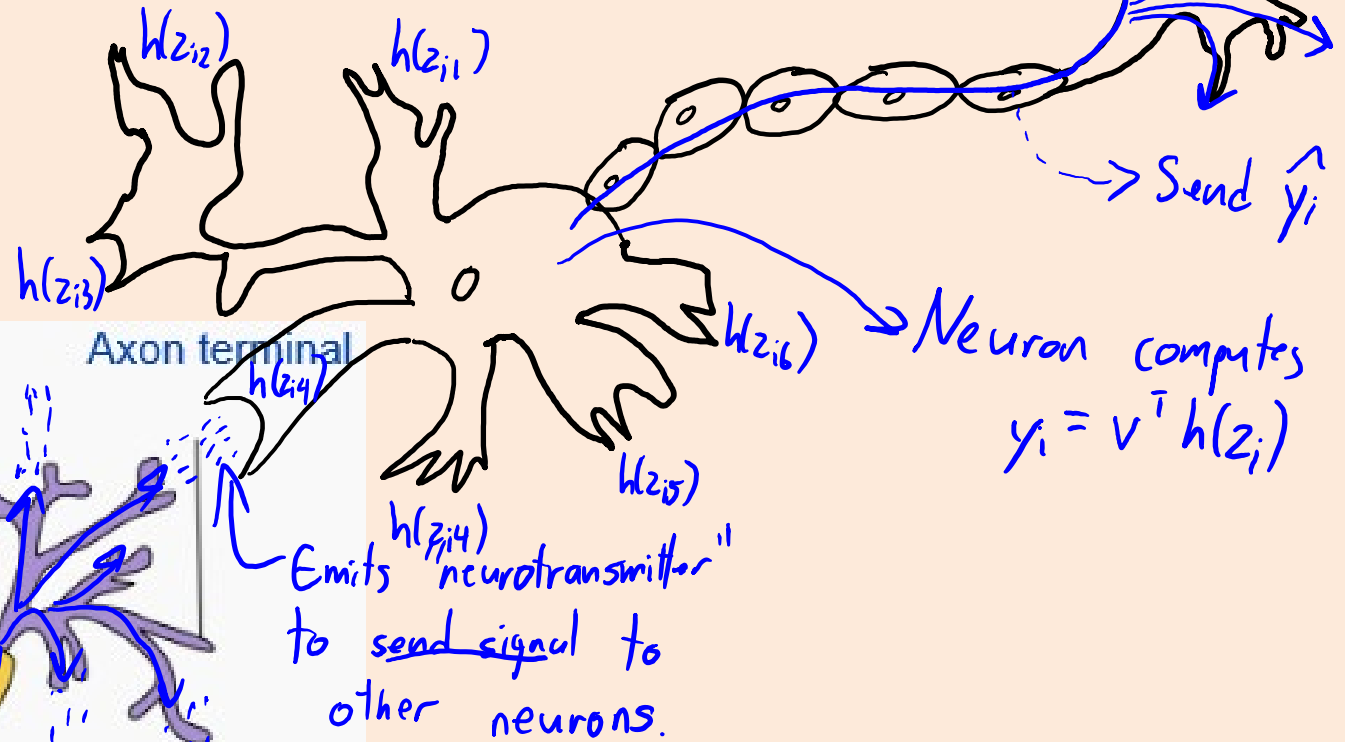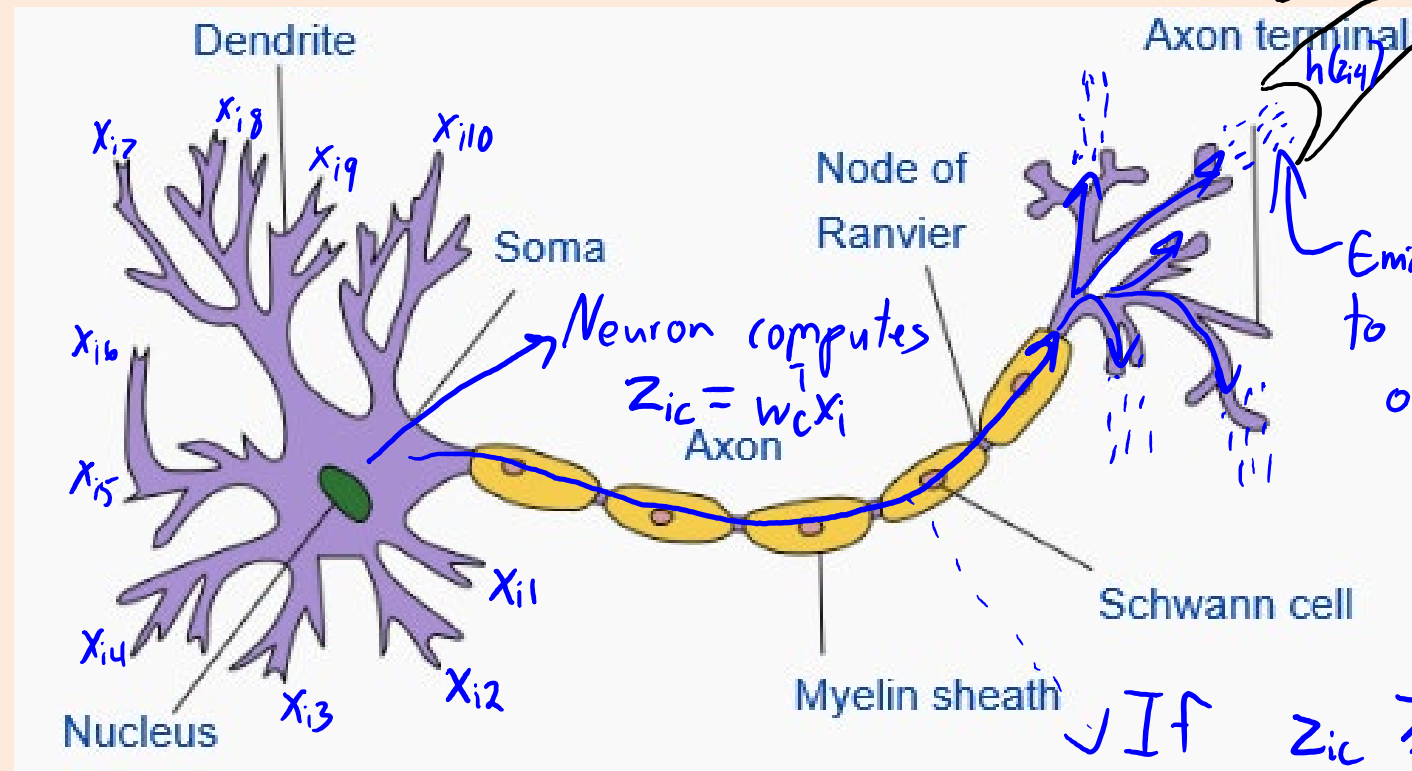  - "Action potential" along axon (like a binary signal):

# Why "Neural Network"?



Neuron computes $z_{ic} = w_c x_i$

Emits "neurotransmitter" to send signal to other neurons.

✓ If $z_{ic} \geq 0.5$ neuron sends signal along axon.

We approximate binary signal with $\dfrac{1}{1 + \exp(-z_{ic})}$

23

# Why "Neural Network"?



$h(z_{i2})$

$h(z_{i1})$

$h(z_{i3})$

$\rightarrow$ Send $\hat{y}_i$

$h(z_{i6})$

$\rightarrow$ Neuron computes
$y_i = v^T h(z_i)$

$h(z_{i5})$

$h(z_{i4})$

Emits "neurotransmitter"
to send signal to
other neurons.

Dendrite

$x_{i7}$  $x_{i8}$  $x_{i9}$  $x_{i10}$

Soma

Node of
Ranvier

Axon terminal

$x_{i6}$

Neuron computes
$z_{ic} = w_c^T x_i$

Axon

$x_{i5}$

$x_{i1}$

$x_{i4}$

Schwann cell

$x_{i3}$  $x_{i2}$

Nucleus

Myelin sheath

If $z_{ic} \geq 0.5$ neuron
Sends signal along axon.

We approximate binary
signal with $\dfrac{1}{1 + \exp(-z_{ic})}$

24

# Why "Neural Network"?
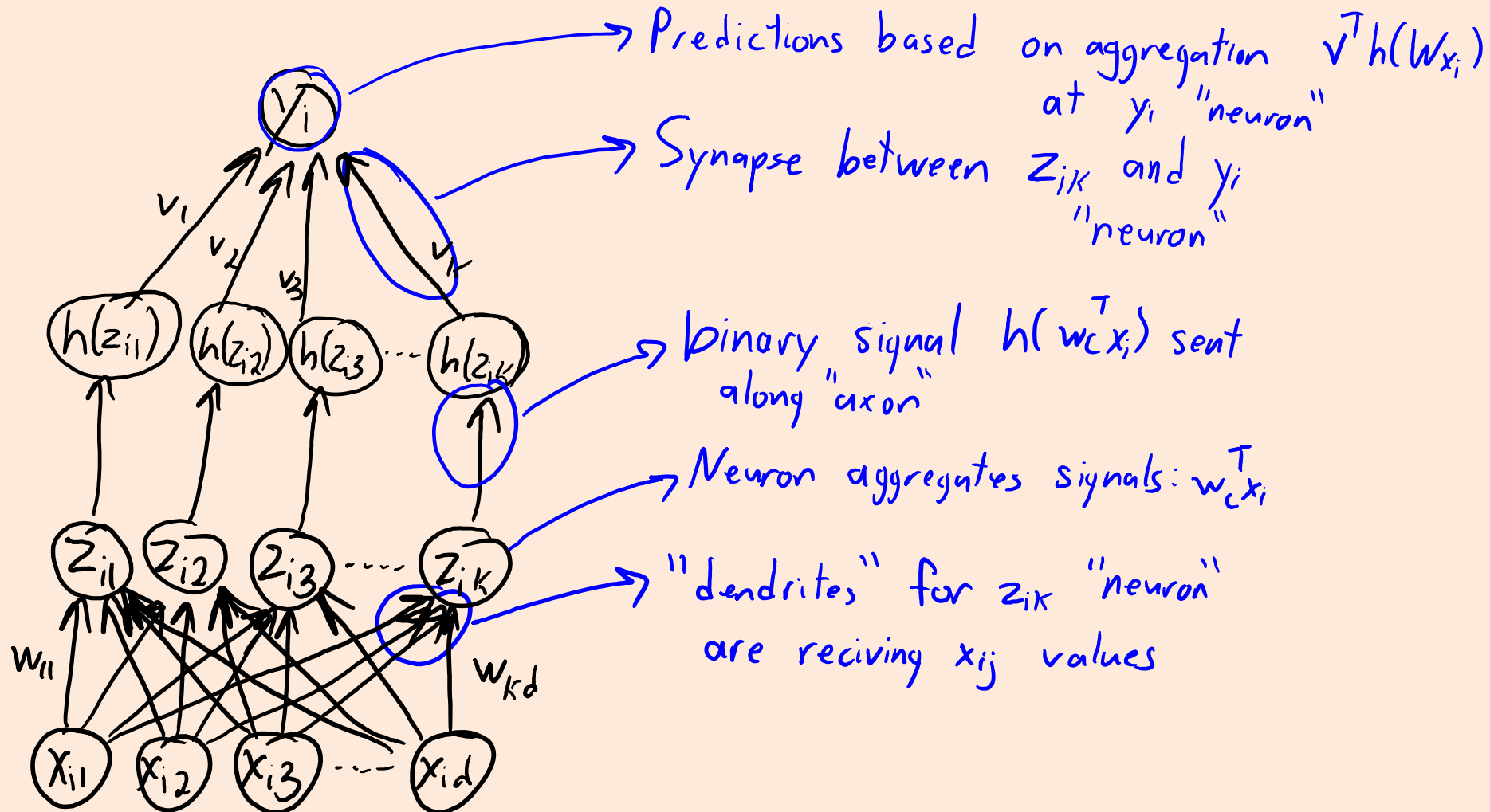


Predictions based on aggregation $v^T h(W_{x_i})$ at $y_i$ "neuron"

Synapse between $z_{jk}$ and $y_i$ "neuron"

binary signal $h(w_c^T x_i)$ sent along "axon"

Neuron aggregates signals: $w_c^T x_i$
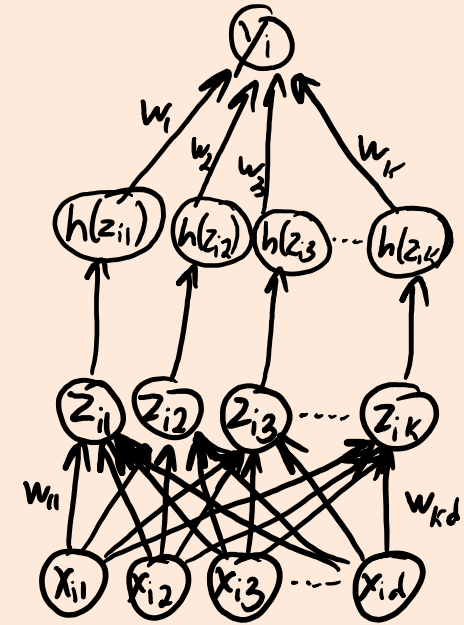
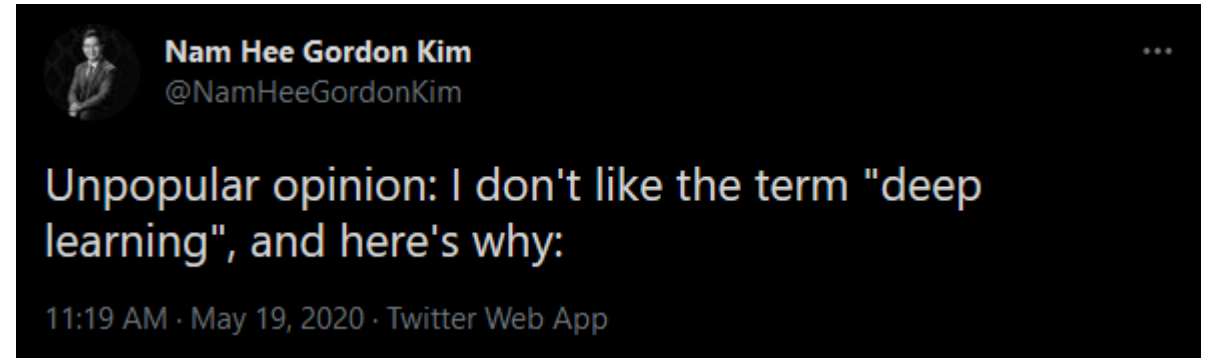"dendrites" for $z_{ik}$ "neuron" are reciving $x_{ij}$ values

# "Artificial" Neural Nets vs. "Real" Networks Nets



- Artificial neural network:
  - $x_i$ is measurement of the world.
  - $z_i$ is internal representation of world.
  - $y_i$ is output of neuron for classification/regression.

- Real neural networks are more complicated:
  - Timing of action potentials seems to be important.
    - "Rate coding": frequency of action potentials simulates continuous output.
  - Sparsity of action potentials.
  - How much computation is done inside neuron?
  - Brain is highly organized (e.g., substructures and cortical columns).
  - Connection structure changes.
  - Different types of neurotransmitters.

Nam Hee Gordon Kim
@NamHeeGordonKim

Unpopular opinion: I don't like the term "deep learning", and here's why:

11:19 AM · May 19, 2020 · Twitter Web App

Coming Up Next
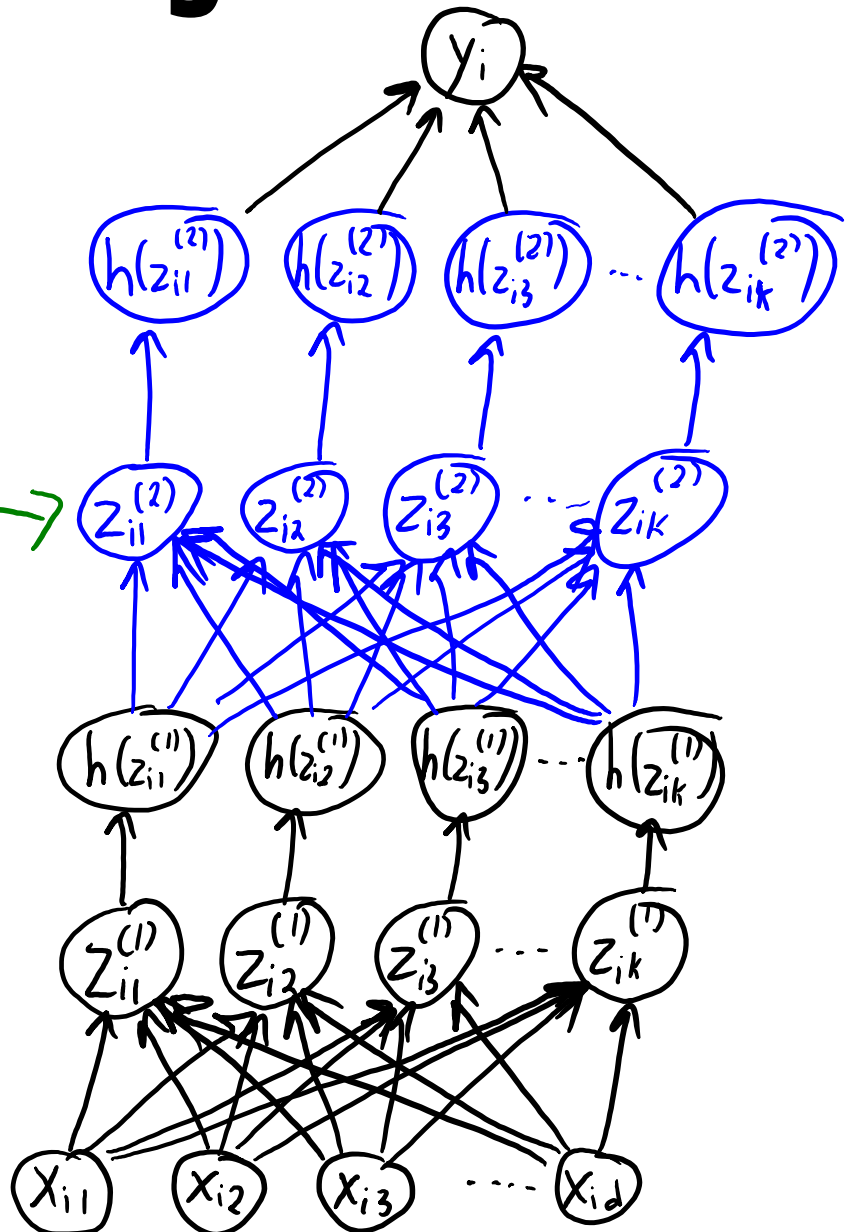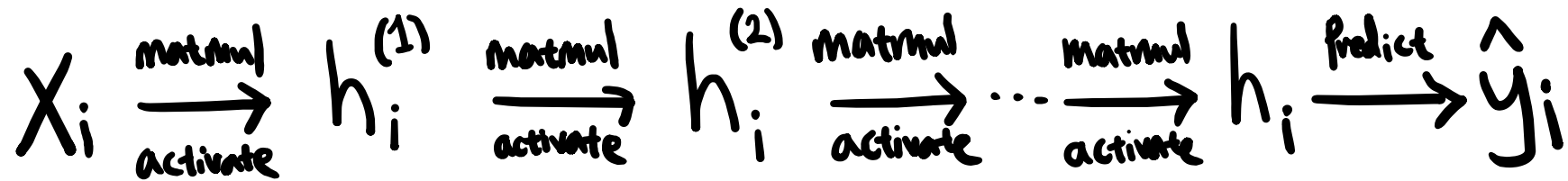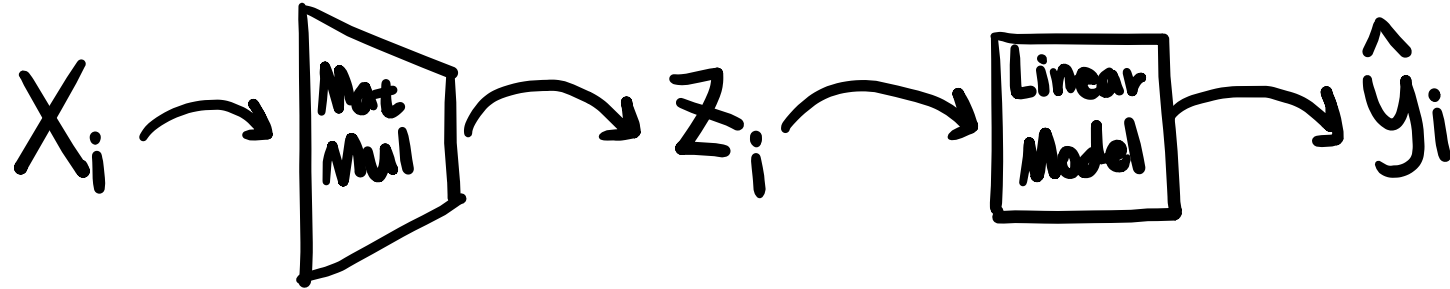
# WHAT IS DEEP LEARNING?

# Deep Learning



Neural network:

Deep learning:

Second "layer" of latent features

You can add more "layers" to go "deeper"

# Encoder-Predictor View of Deep Learning

$$X_i \rightsquigarrow \boxed{\substack{\text{Mat} \\ \text{Mul}}} \rightsquigarrow Z_i \rightsquigarrow \boxed{\substack{\text{Linear} \\ \text{Model}}} \rightsquigarrow \hat{y}_i$$

$$X_i \xrightarrow[\text{activate}]{\text{matmul}} h_i^{(1)} \xrightarrow[\text{activate}]{\text{matmul}} h_i^{(2)} \xrightarrow[\text{activate}]{\text{matmul}} \cdots \xrightarrow[\text{activate}]{\text{matmul}} h_i \xrightarrow{\text{predict}} \hat{y}_i$$
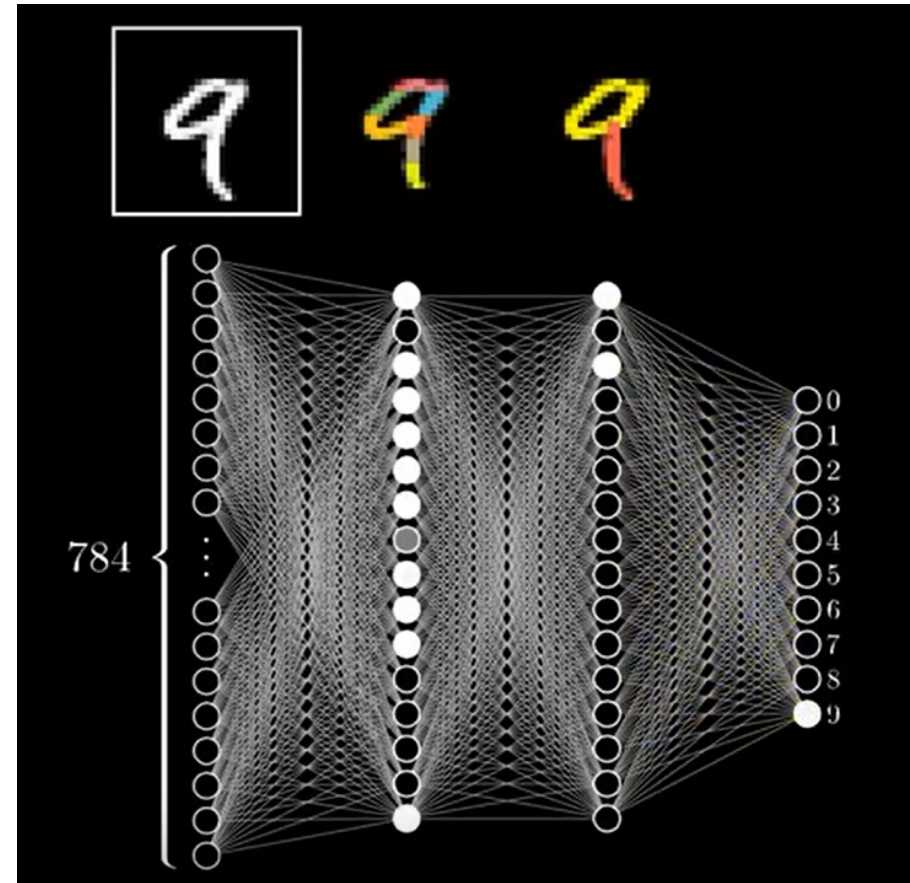
- Compose multiple non-linear encoders
- Overall idea is still the same:
  - Train encoder and predictor at the same time
    - (we have a "bigger" encoder now)

# "Hierarchies of Parts" Motivation for Deep Learning

- Each "neuron" might recognize a "part" of digit.
  - "Deeper" neurons might recognize combinations of parts.
  - Represent complex objects as hierarchical combinations of re-useable parts (a simple "grammar").

- Watch the full video here:
  - https://www.youtube.com/watch?v=aircAruvnKk

- Theory:
  - 1 big-enough hidden layer already gives universal approximation.
  - But some functions require exponentially-fewer parameters to approximate with more layers (can fight curse of dimensionality).

# Deep Learning

Linear model:
$$\hat{y}_i = w^\top x_i$$

Neural network with 1 hidden layer:
$$\hat{y}_i = v^\top h(\underbrace{W x_i}_{z_i})$$

Neural network with 2 hidden layers:
$$\hat{y}_i = v^\top h(\underbrace{W^{(2)} \underbrace{h(\underbrace{W^{(1)} x_i}_{z_i^{(1)}})}}_{z_i^{(2)}})$$

Neural network with 3 hidden layers
$$\hat{y}_i = v^\top h(\underbrace{W^{(3)} h(\underbrace{W^{(2)} \underbrace{h(\underbrace{W^{(1)} x_i}_{z_i^{(1)}})}}_{z_i^{(2)}})}_{z_i^{(3)}})$$

Deep learning:

Second "layer" of latent features

You can add more "layers" to go "deeper"

# Deep Learning
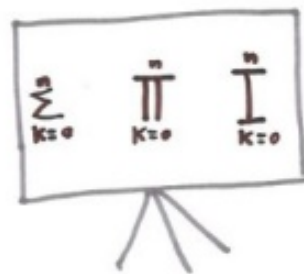
- For 4 layers, we could write the prediction as:

$$\hat{y}_i = v^T h(W^{(4)} h(W^{(3)} h(W^{(2)} h(W^{(1)} x_i))))$$

- For 'm' layers, we could use

$$\hat{y}_i = v^T \left( \underset{\ell=1}{\overset{m}{\text{I}}} h(W^{(\ell)} x_i) \right)$$

Symbol: $\underset{k=0}{\overset{n}{\text{I}}} f_k(t)$

Meaning: $f_n \circ f_{n-1} \circ f_{n-2} \circ \ldots \circ f_2 \circ f_1 \circ f_0(t)$

$\underset{k=0}{\overset{n}{\Sigma}} \quad \underset{k=0}{\overset{n}{\Pi}} \quad \underset{k=0}{\overset{n}{\text{I}}}$

— I knew something was missing!

Coming Up Next

# HISTORY OF DEEP LEARNING

# 'Godfathers of AI' honored with Turing Award, the Nobel Prize of computing

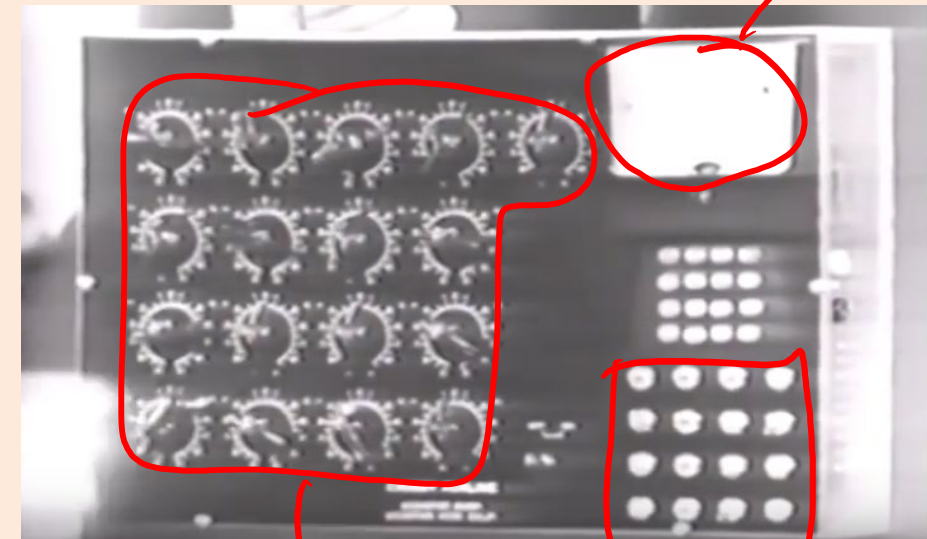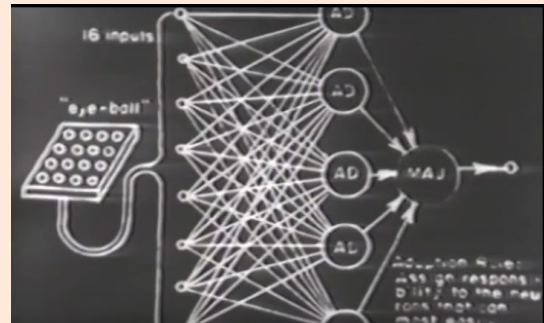*Yoshua Bengio, Geoffrey Hinton, and Yann LeCun laid the foundations for modern AI*

By James Vincent | Mar 27, 2019, 6:02am EDT



From left to right: Yann LeCun | Photo: Facebook; Geoffrey Hinton | Photo: Google; Yoshua Bengio | Photo: Botler AI
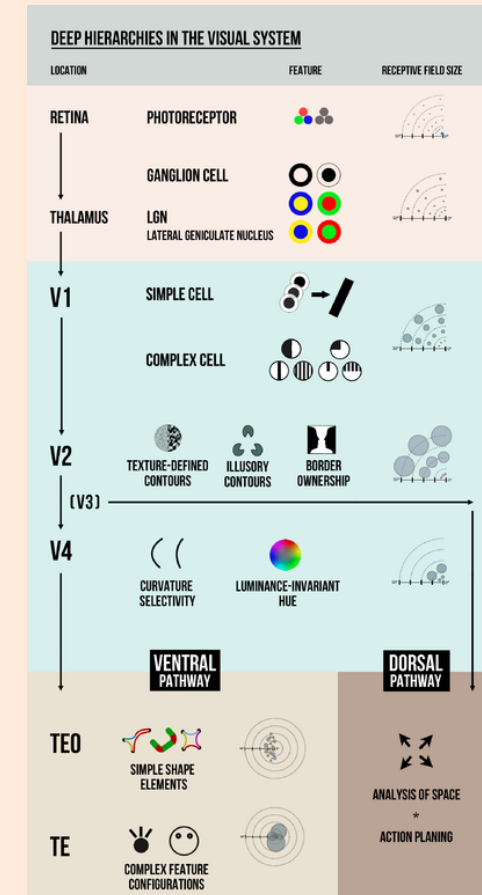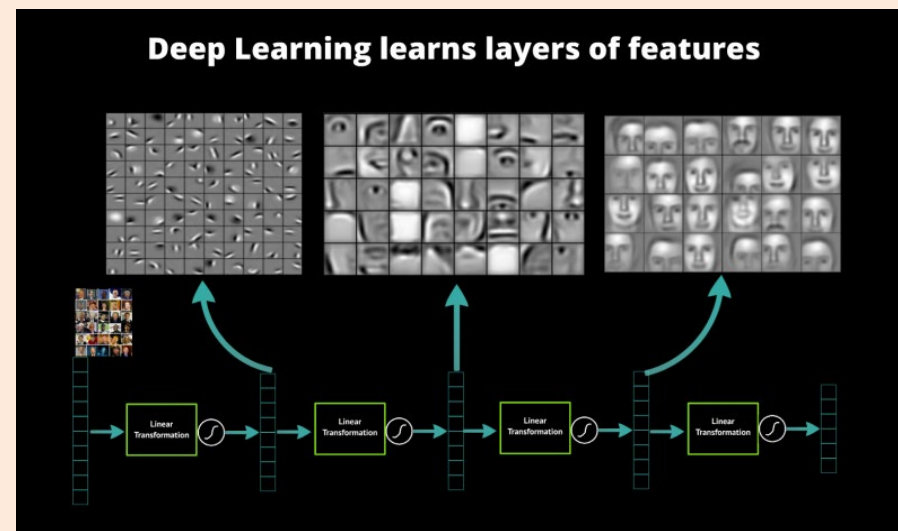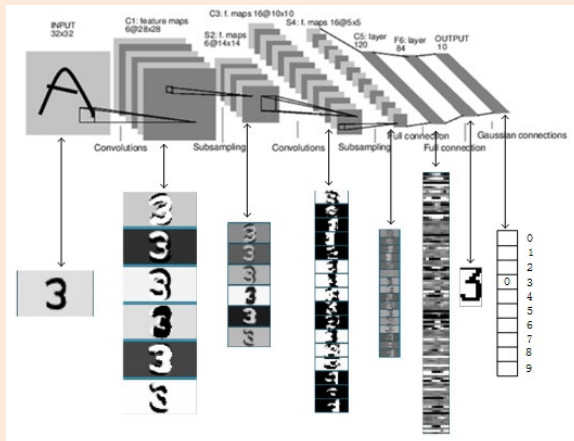
# ML and Deep Learning History

- ## 1950 and 1960s: Initial excitement.
  - Perceptron: linear classifier and stochastic gradient (roughly).
  - "the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence." New York Times (1958).
    - https://www.youtube.com/watch?v=IEFRtz68m-8
  - Object recognition assigned to students as a summer project

- ## Then drop in popularity:
  - Quickly realized limitations of linear models.

# ML and Deep Learning History

- **1970 and 1980s: Connectionism (brain-inspired ML)**
  - Want "connected networks of simple units".
    - Use parallel computation and distributed representations.
  - Adding hidden layers $z_i$ increases expressive power.
    - With 1 layer and enough sigmoid units, a universal approximator.
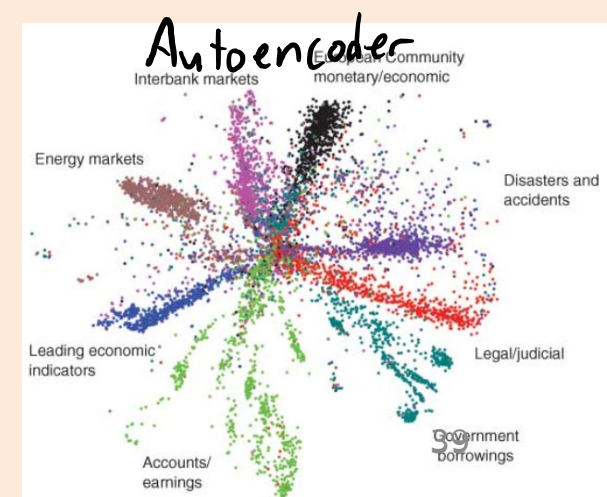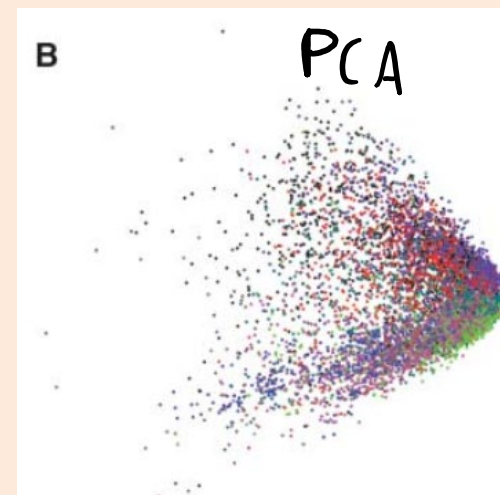  - Success in optical character recognition.

# ML and Deep Learning History

- **1990s and early-2000s: drop in popularity.**
  - It proved really difficult to get multi-layer models working robustly.

  - We obtained similar performance with simpler models:
    - Rise in popularity of logistic regression and SVMs with regularization and kernels.

  - Lots of internet successes (spam filtering, web search, recommendation).

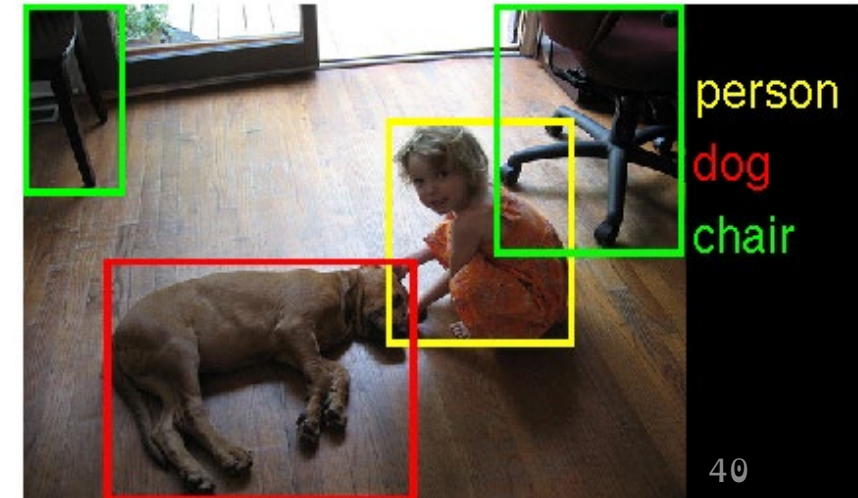  - ML moved closer to other fields like numerical optimization and statistics.

# ML and Deep Learning History

- ## Late 2000s: push to revive connectionism as "deep learning".
  - Canadian Institute For Advanced Research (CIFAR) NCAP program:
    - "Neural Computation and Adaptive Perception".
    - Led by Geoff Hinton, Yann LeCun, and Yoshua Bengio ("Canadian mafia").
  - Unsupervised successes: "deep belief networks" and "autoencoders".
    - Could be used to initialize deep neural networks.
    - https://www.youtube.com/watch?v=KuPai0ogiHk

# 2010s: DEEP LEARNING!!!

- Bigger datasets, bigger models, parallel computing (GPUs/clusters).
  - And some tweaks to the models from the 1980s.
- Huge improvements in automatic speech recognition (2009).
  - All phones now have deep learning.
- Huge improvements in computer vision (2012).
  - Changed computer vision field almost instantly.
  - This is now finding its way into products.



person
dog
chair

http://www.image-net.org/challenges/LSVRC/2014/

# 2010s: DEEP LEARNING!!!

- **Media hype:**
  - "How many computers to identify a cat? 16,000"

    New York Times (2012).
  - "Why Facebook is teaching its machines to think like humans"

    Wired (2013).
  - "What is 'deep learning' and why should businesses care?"

    Forbes (2013).
  - "Computer eyesight gets a lot more accurate"

    New York Times (2014).

- **2015:** huge improvement in language understanding.

# Cut-off for Final Exam

(Final exam will have materials from everything before this slide)

# Summary

- **Neural networks** learn features $z_i$ for supervised learning.
- **Sigmoid function** avoids degeneracy by introducing non-linearity.
  - Universal approximator with large-enough 'k'.
- **Biological motivation** for (deep) neural networks.
- **Deep learning** considers neural networks with many hidden layers.
  - Can more-efficiently represent some functions.
- **Unprecedented performance** on difficult pattern recognition tasks.

- Next time:
  - Training deep networks.

# Please Do Course Evaluation!

# Review Questions

- Q1: What is the problem with using a linear encoder and a linear predictor for a neural network?

- Q2: What is the motivation for using multiple layers of encoders?

- Q3: Does it make sense to use neural networks for classifying linearly separable data?

# Why $z_i = Wx_i$?

- In PCA we had that the optimal $Z = XW^T(WW^T)^{-1}$.
- If $W$ had normalized+orthogonal rows, $Z = XW^T$ (since $WW^T = I$).
  - So $z_i = Wx_i$ in this normalized+orthogonal case.

- Why we would use $z_i = Wx_i$ in neural networks?
  - We didn't enforce normalization or orthogonality.

- Well, the value $W^T(WW^T)^{-1}$ is just "some matrix".
  - You can think of neural networks as just directly learning this matrix.

# Cool Picture Motivation for Deep Learning

- Faces might be composed of different "parts":

# Cool Picture Motivation for Deep Learning

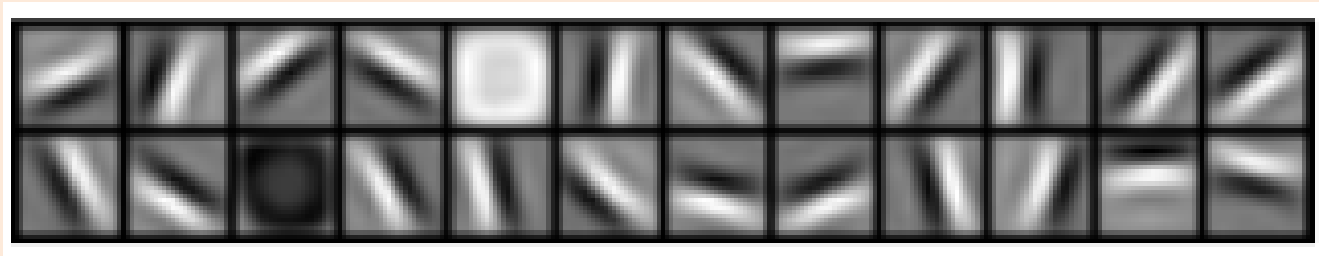- First layer of $z_i$ trained on 10 by 10 image patches:



"Gabor filters"

- **Attempt to visualize second layer:**
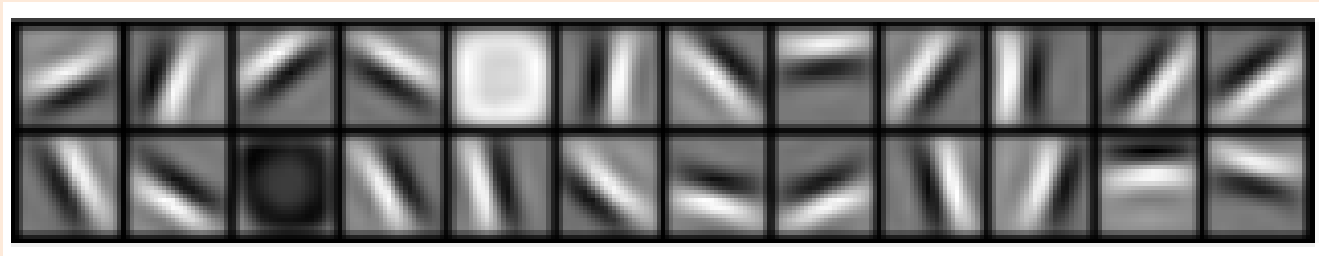  - Corners, angles, surface boundaries?

- **Models require many tricks to work.**
  - We'll discuss these next time.

# Cool Picture Motivation for Deep Learning

- First layer of $z_i$ trained on 10 by 10 image patches:



"Gabor filters"

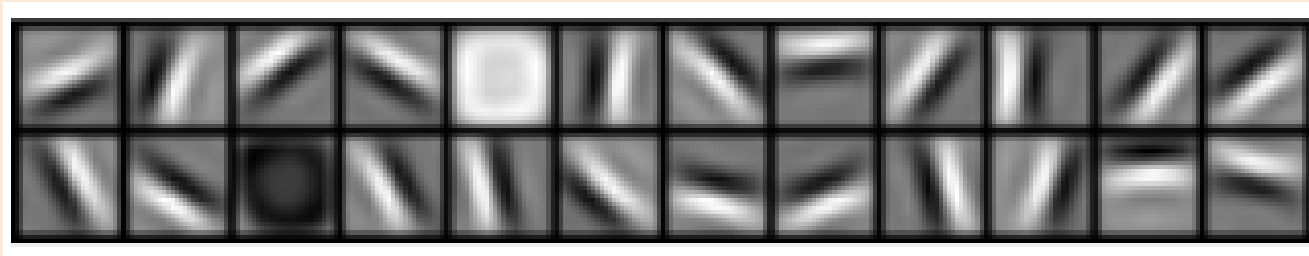- Visualization of second and third layers trained on specific objects:

faces

# Cool Picture Motivation for Deep Learning
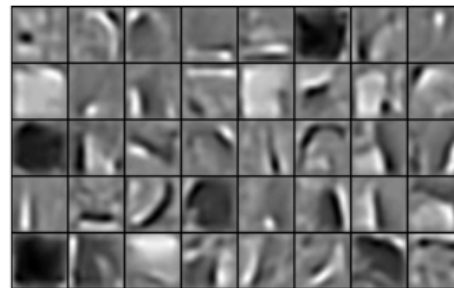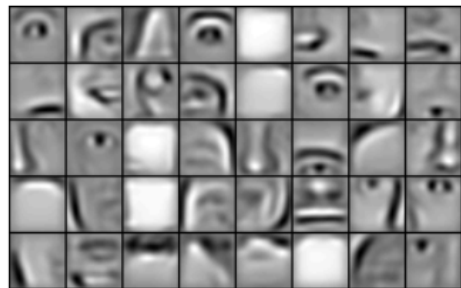
- First layer of $z_i$ trained on 10 by 10 image patches:



{ "Gabor filters"

- Visualization of second and third layers trained on specific objects:
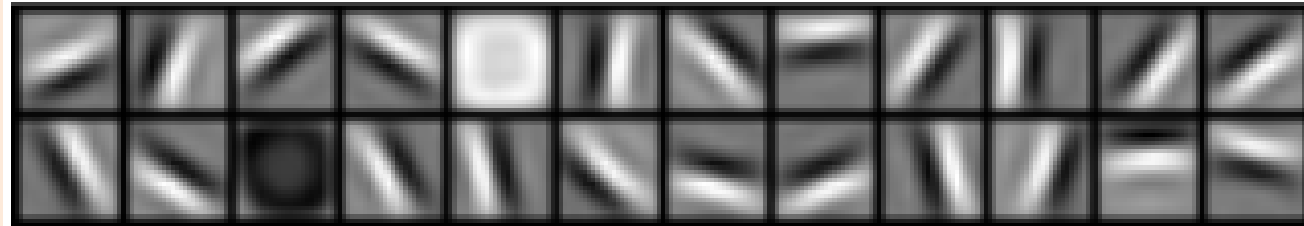


faces

cars

http://www.cs.toronto.edu/~rgrosse/icml09-cdbn.pdf

# Cool Picture Motivation for Deep Learning

- First layer of $z_i$ trained on 10 by 10 image patches:
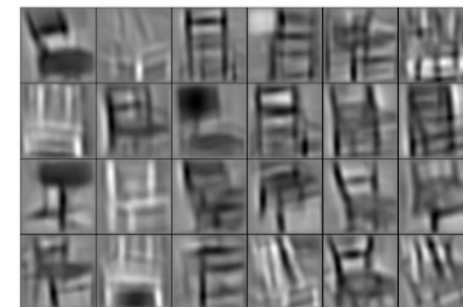


"Gabor filters"

- Visualization of second and third layers trained on specific objects:



faces    cars    elephants

# Cool Picture Motivation for Deep Learning

- First layer of $z_i$ trained on 10 by 10 image patches:
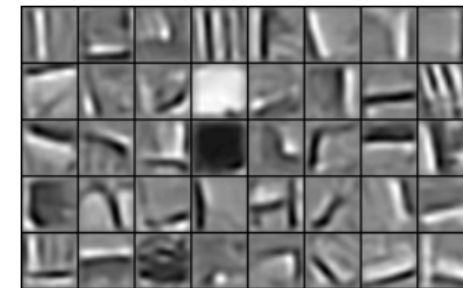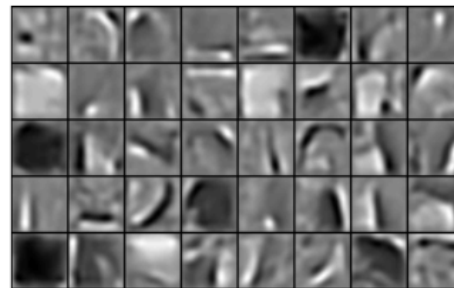


"Gabor filters"

- Visualization of second and third layers trained on specific objects:
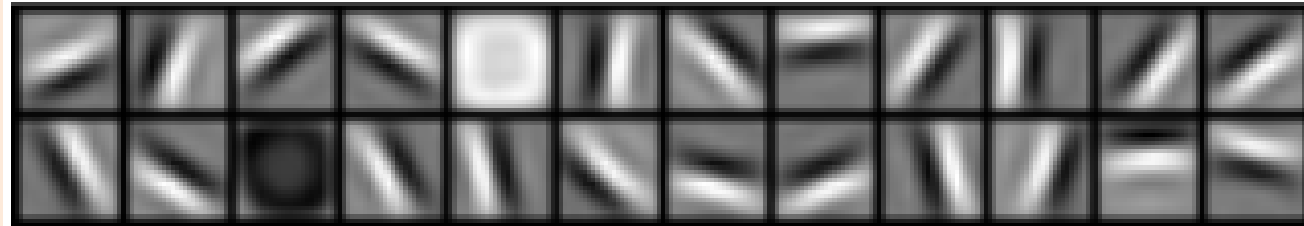


faces    cars    elephants    chairs

# Cool Picture Motivation for Deep Learning

- First layer of $z_i$ trained on 10 by 10 image patches:



"Gabor filters"

- Visualization of second and third layers trained on specific objects:



faces     cars     elephants     chairs     faces, cars, airplanes, motorbikes