# CPSC 340:
# Machine Learning and Data Mining
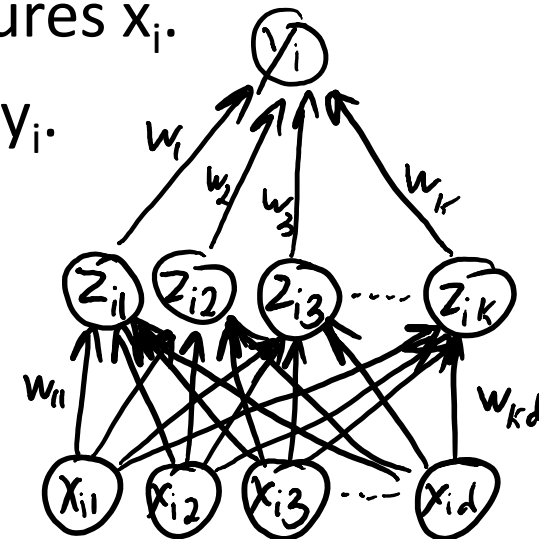
Neural Networks

Fall 2016

# Admin

- **Assignment 4**:
  - 3 late days to hand in Monday.
- **Assignment 5**:
  - Out, due next Friday.
- **Assignment 6**:
  - Out, due last day of class.
- **Final**:
  - December 12 (8:30am – HEBB 100)
  - Covers Assignments 1-6.
  - Final from last year and list of topics will be posted.
  - Closed-book, cheat sheet: 4-pages each double-sided.

# 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 'w' 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 'w' based on change of basis $z_i$ and target $y_i$.
- Part 5: Neural Networks.
  - Jointly learn 'W' and 'w' based on $x_i$ and $y_i$.
  - **Learn basis $z_i$ that is good for supervised learning.**

# Neural Networks: Introducing Non-Linearity

- Natural choice of neural network regression objective would be:

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

- But we saw last time this gives a linear model.

$$w^T W \text{ is a } \underline{vector} \text{ so } \underline{equivalent} \text{ to just having } w^T x_i \text{ for some } 'w'.$$

- Typical fix is to introduce non-linearity 'h':

$$f(w, W) = \frac{1}{2} \sum_{i=1}^{n} (w^T \underbrace{h(W x_i)}_{z_i} - y_i)^2 \quad \text{where 'h' has 'd' inputs and 'k' outputs.}$$

- Most common choice of 'h' is sigmoid applied to elements of $W x_i$.

$$z_{ic} = \frac{1}{1 + \exp(-W_c x_i)}$$

# Notation for Neural Networks

We have our usual supervised learning notation:

$$X = \begin{bmatrix} -\!\!-\!\!- x_1^\top -\!\!-\!\!- \\ -\!\!-\!\!- x_2^\top -\!\!-\!\!- \\ \vdots \\ -\!\!-\!\!- x_n -\!\!-\!\!- \end{bmatrix} \qquad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$n \times d$ 　　　　$n \times 1$

We have our latent features:

$$Z = \begin{bmatrix} -\!\!-\!\!- z_1^\top -\!\!-\!\!- \\ -\!\!-\!\!- z_2^\top -\!\!-\!\!- \\ \vdots \\ -\!\!-\!\!- z_n^\top -\!\!-\!\!- \end{bmatrix}$$

$n \times K$

We have **two** sets of parameters:

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} \qquad W = \begin{bmatrix} -\!\!- W_1 -\!\!- \\ -\!\!- W_2 -\!\!- \\ \vdots \\ -\!\!- W_K -\!\!- \end{bmatrix}$$

$k \times 1$ 　　　　$k \times d$

# Supervised Learning Roadmap



**Hand-engineered features:**

"I think this basis will work"
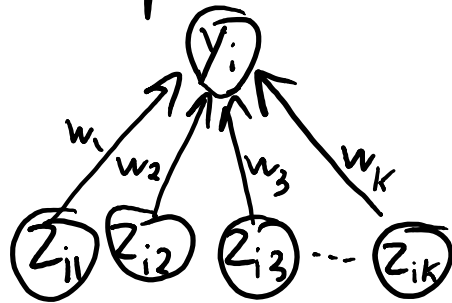
Requires domain knowledge and can be time-consuming
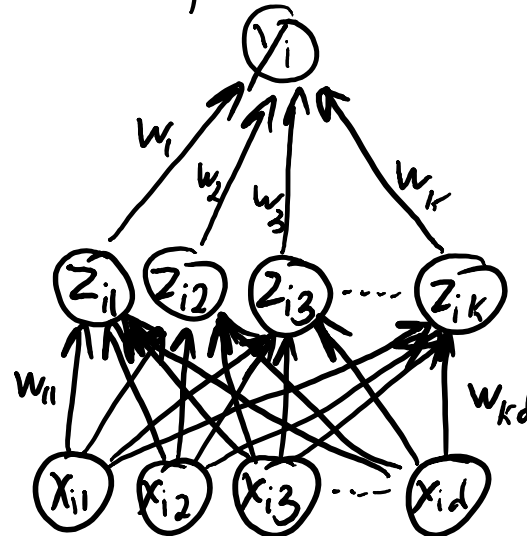
**Learn a latent-factor model:**
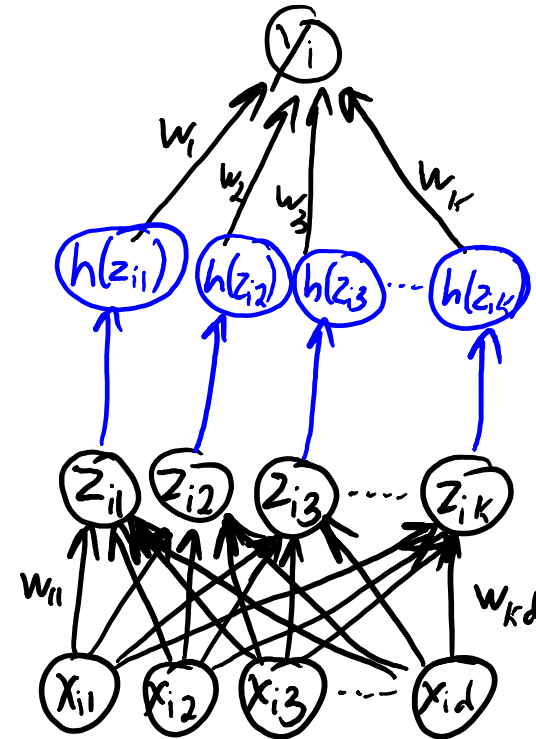
Use latent features in supervised model:

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

**Learn 'w' and 'W' together:**

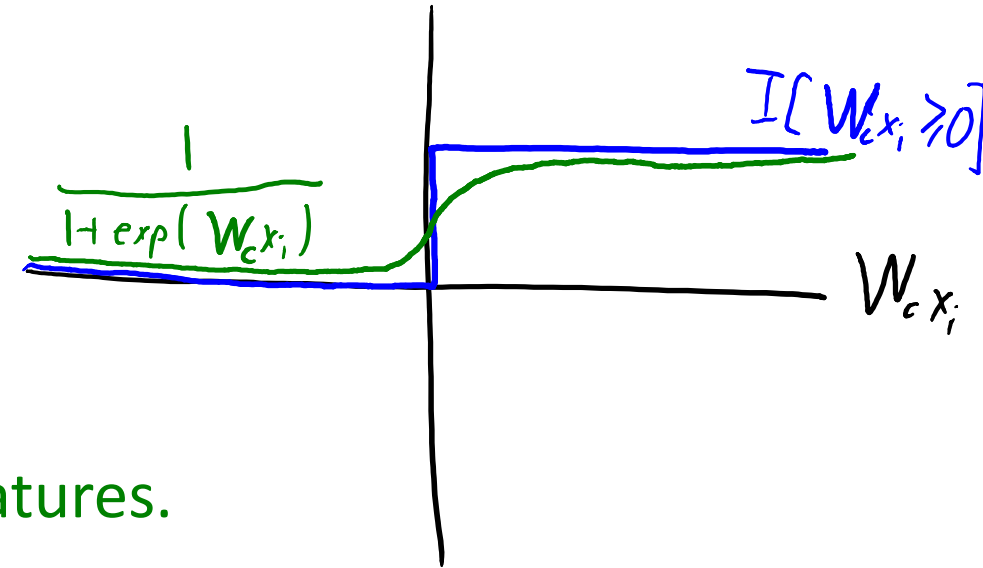But still gives a linear model.

**Neural network:**

Extra non-linear transformation 'h'

# Why Sigmoid?

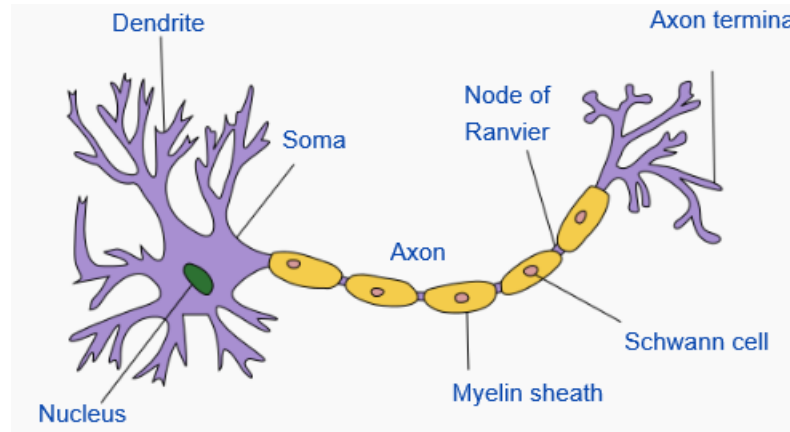- Consider setting 'h' to define binary features $z_i$ using:

$$z_i = I[W_c x_i \geqslant 0]$$

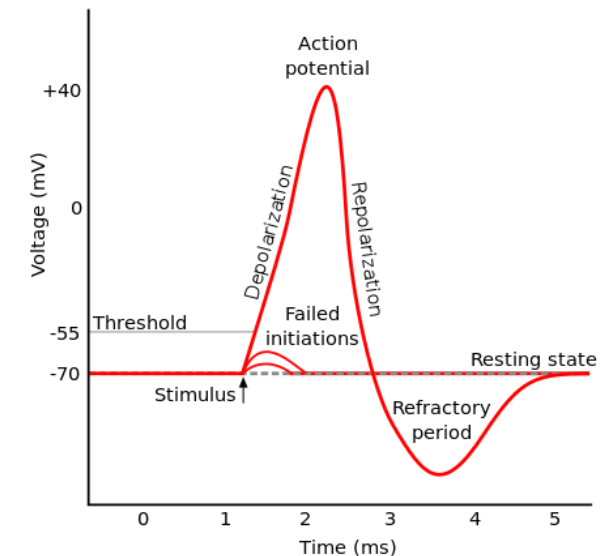$$= \begin{cases} 1 & \text{if } W_c x_i \geqslant 0 \\ 0 & \text{if } W_c x_i < 0 \end{cases}$$



- Vector $z_i = h(Wx_i)$ can be viewed as binary features.
- $z_i$ can take $2^k$ possible values (combinatorial number of "concepts").

- But non-differentiable and discontinuous so hard to optimize.

- Sigmoid is a smooth approximation to these binary features.

# Why "Neural Network"?
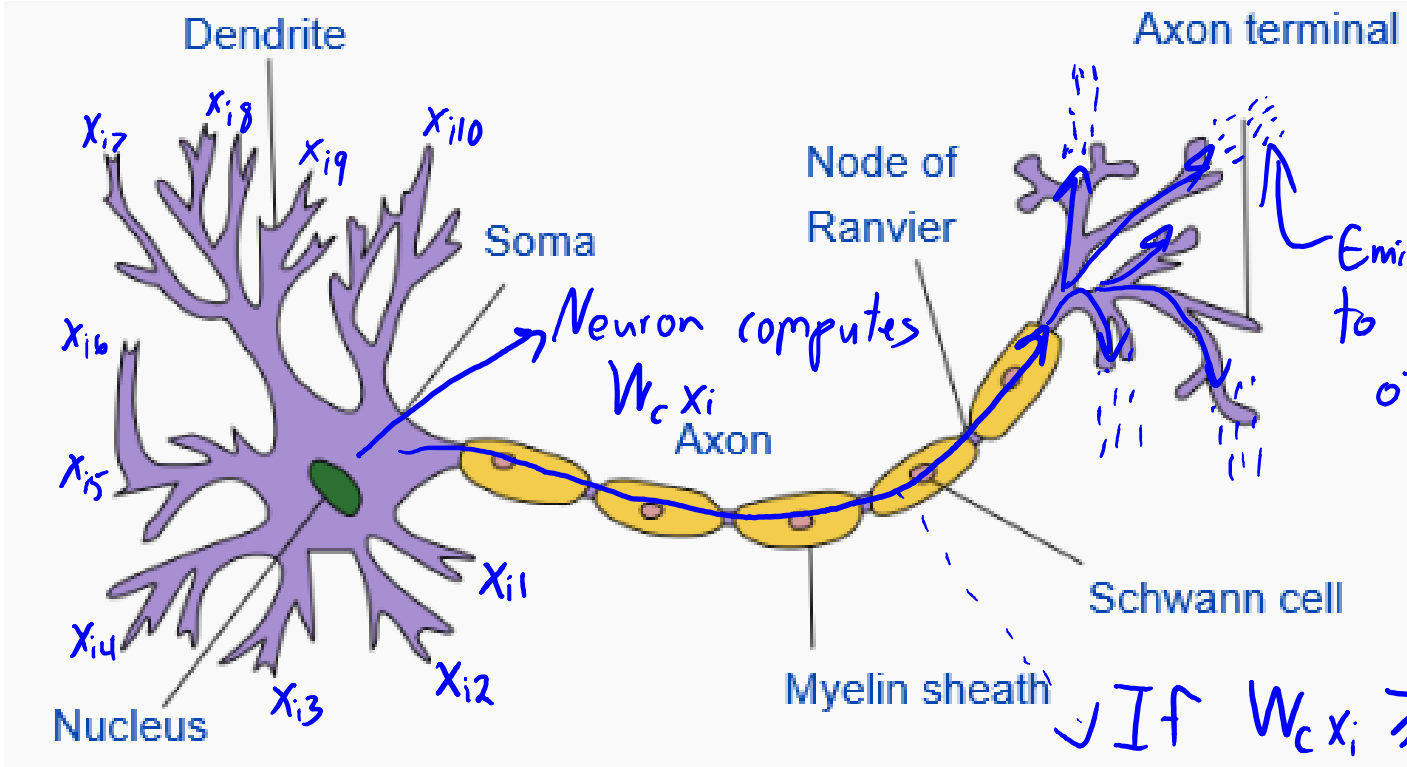
- Cartoon of "typical" neuron:



- 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:
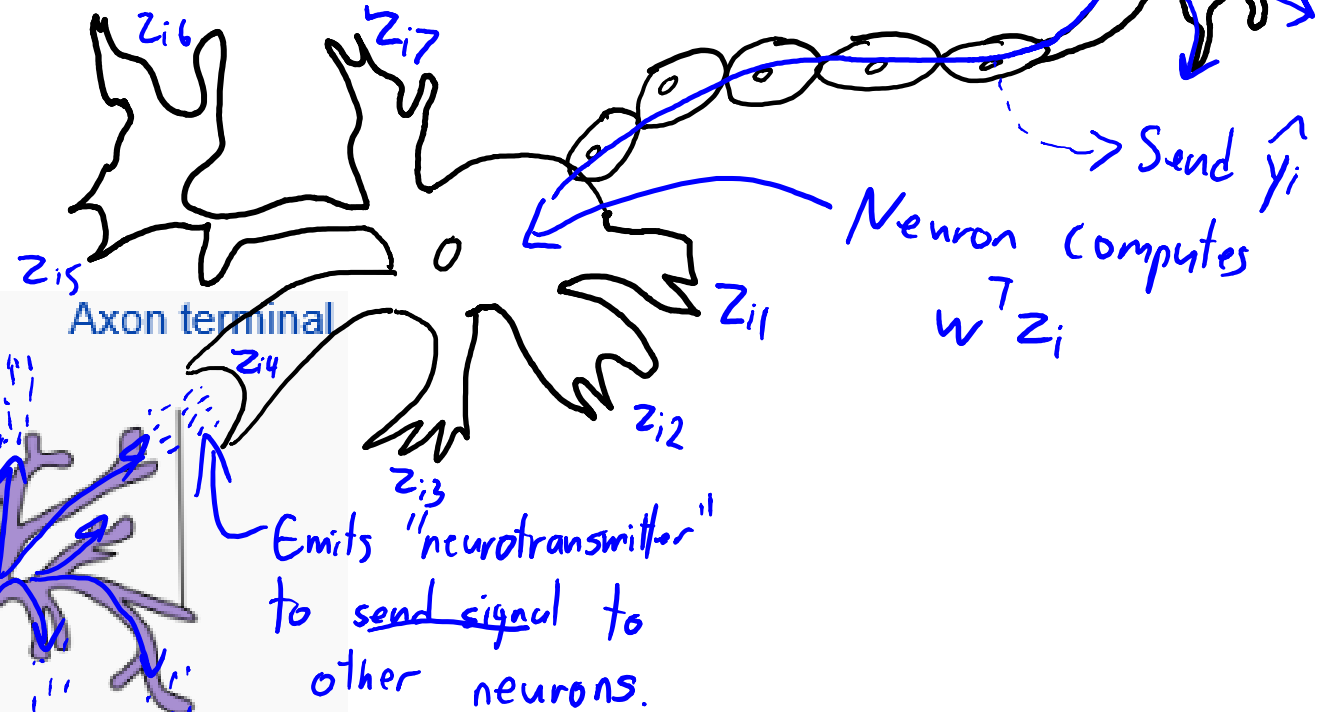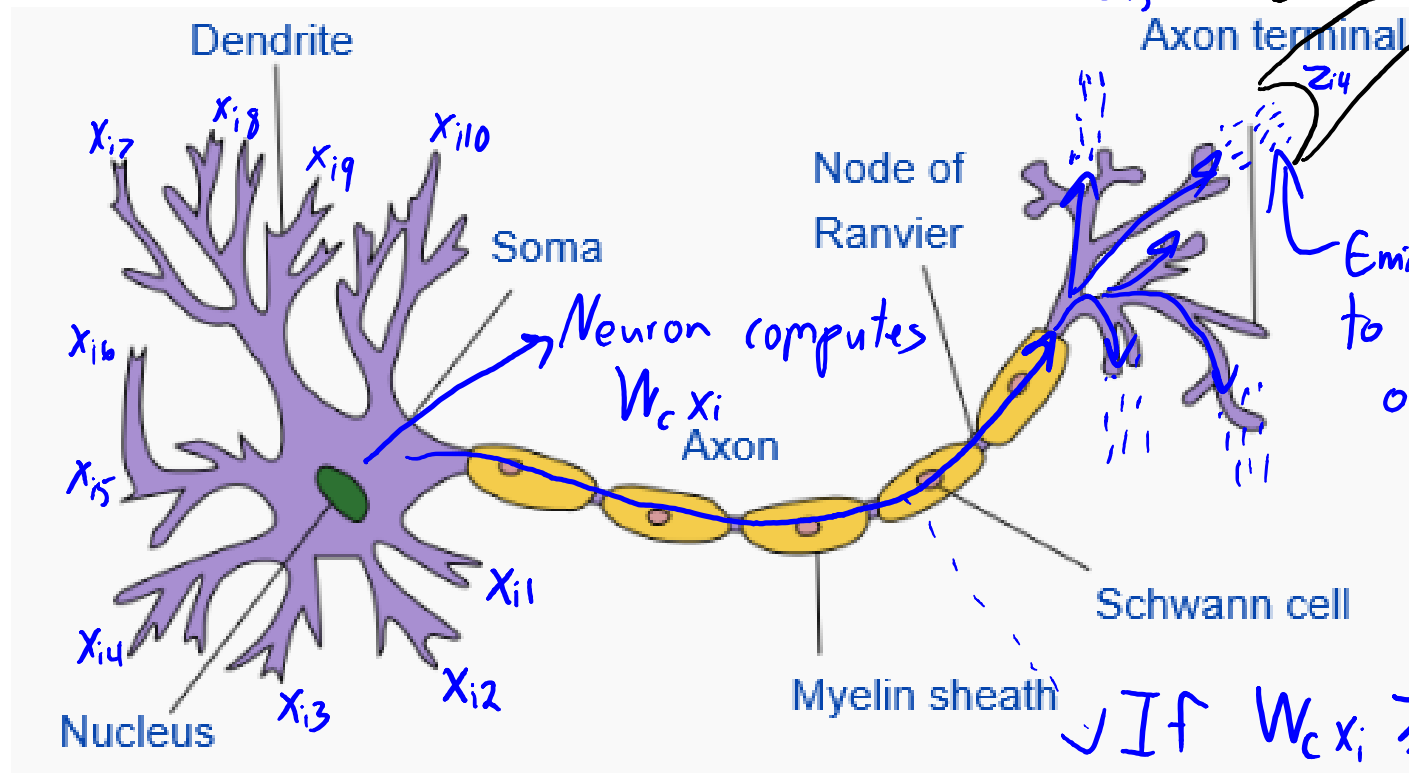  - "Action potential" along axon (like a binary signal):

# Why "Neural Network"?



Dendrite

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

Soma

Node of Ranvier

Axon terminal

$x_{i6}$

$x_{i5}$

Neuron computes $W_c x_i$

Axon

$x_{i1}$

$x_{i4}$

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

Nucleus

Schwann cell

Myelin sheath

Emits "neurotransmitter" to send signal to other neurons.

✓ If $W_c x_i \geq 0$ neuron sends signal along axon.

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

# Why "Neural Network"?



$z_{i6}$ $z_{i7}$

$z_{i5}$

Axon terminal

$z_{i4}$

$z_{i3}$

$z_{i1}$

$z_{i2}$

Send $\hat{y}_i$

Neuron computes $w^T z_i$

Emits "neurotransmitter" to send signal to other neurons.

Dendrite

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

$x_{i6}$

Soma

Node of Ranvier

$x_{i5}$

Neuron computes $W_c x_i$

Axon

$x_{i1}$

$x_{i4}$

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

Nucleus

Schwann cell

Myelin sheath

✓ If $W_c x_i \geqslant 0$ neuron sends signal along axon.

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

# Why "Neural Network"?



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

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

binary signal $h(W_c x_i)$ sent along "axon"

Neuron aggregates signals: $W_c x_i$

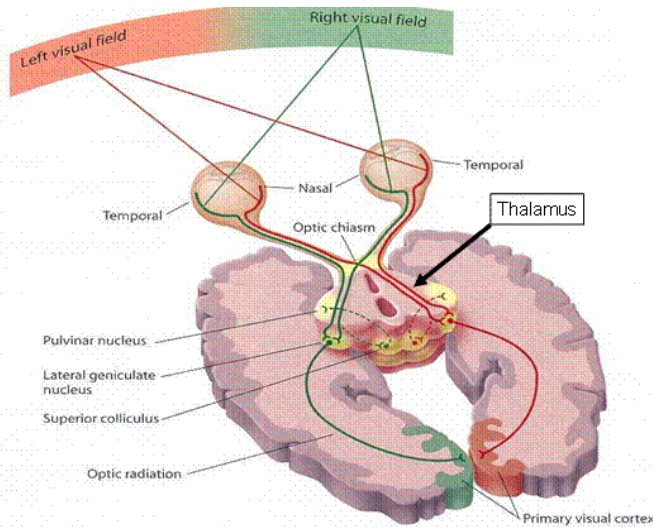"dendrites" for $z_{ik}$ "neuron" are reciving $x_{ij}$ values
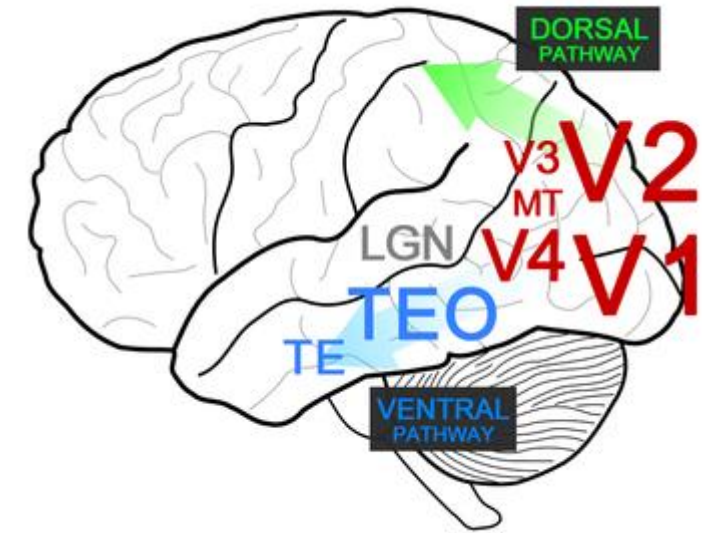
# "Artificial" Neural Nets vs. "Real" Networks Nets



- 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.
  - Neural networks don't reflect 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.
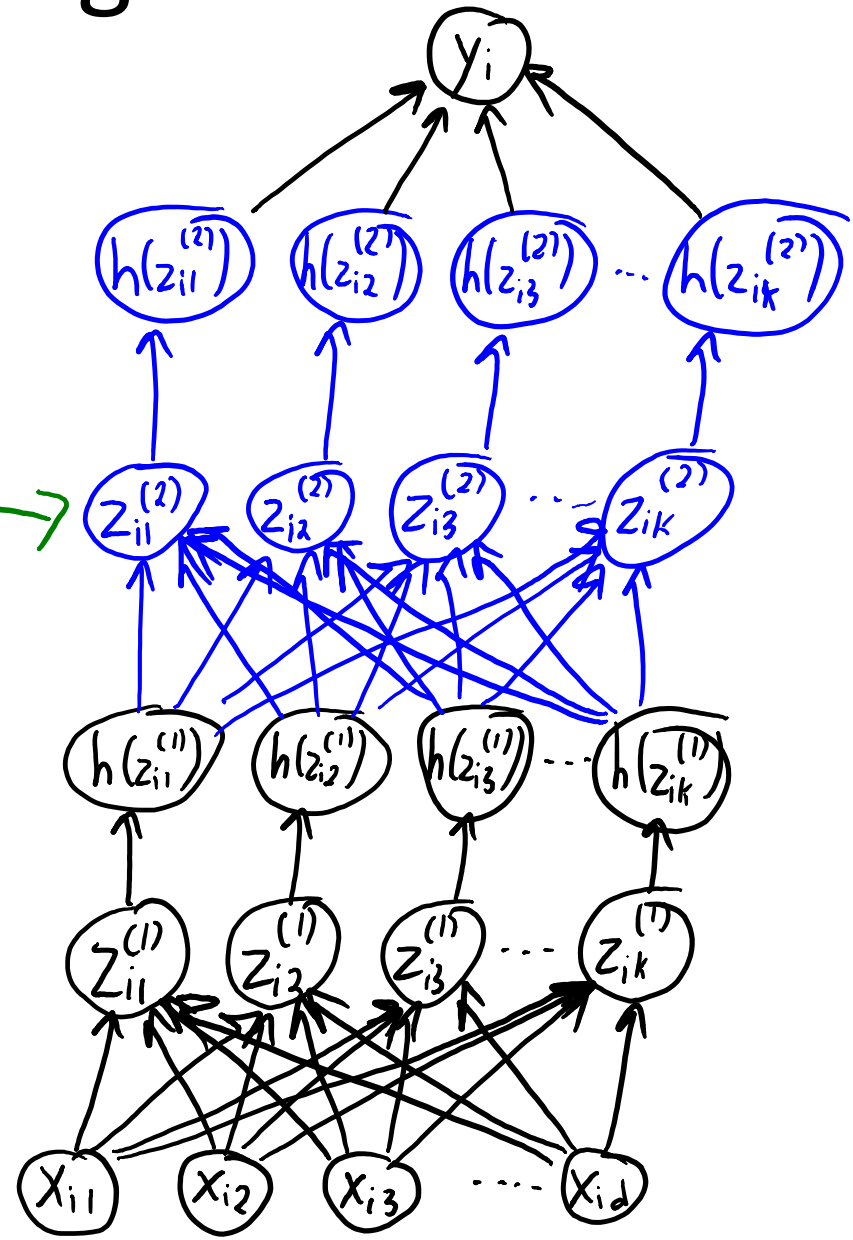
# Deep Hierarchies in the Brain

# Deep Learning



Neural network:

Deep learning

Second "layer" of latent features

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

# Deep Learning

Linear model:

$$y_i = w^T x_i$$

Neural network with 1 hidden layer:

$$y_i = w^T h(\underbrace{Wx_i}_{z_i})$$

Neural network with 2 hidden layers:

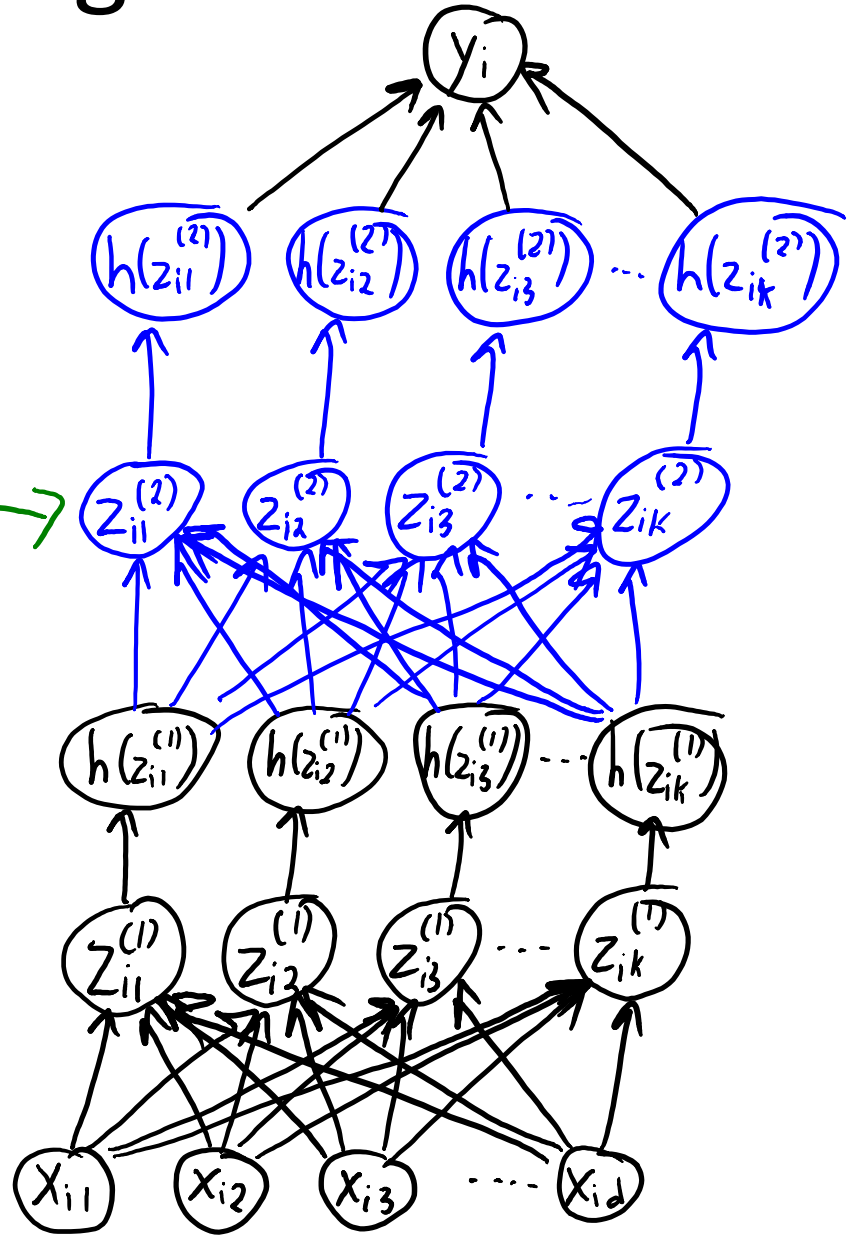$$y_i = w^T h(W^{(2)} \underbrace{h(\underbrace{W^{(1)}x_i}_{z_i^{(1)}})}_{z_i^{(2)}})$$

Neural network with 3 hidden layers

$$y_i = w^T h(W^{(3)} h(W^{(2)} \underbrace{h(\underbrace{\underbrace{W^{(1)}x_i}_{z_i^{(1)}}}_{z_i^{(2)}})}_{z_i^{(3)}}))$$
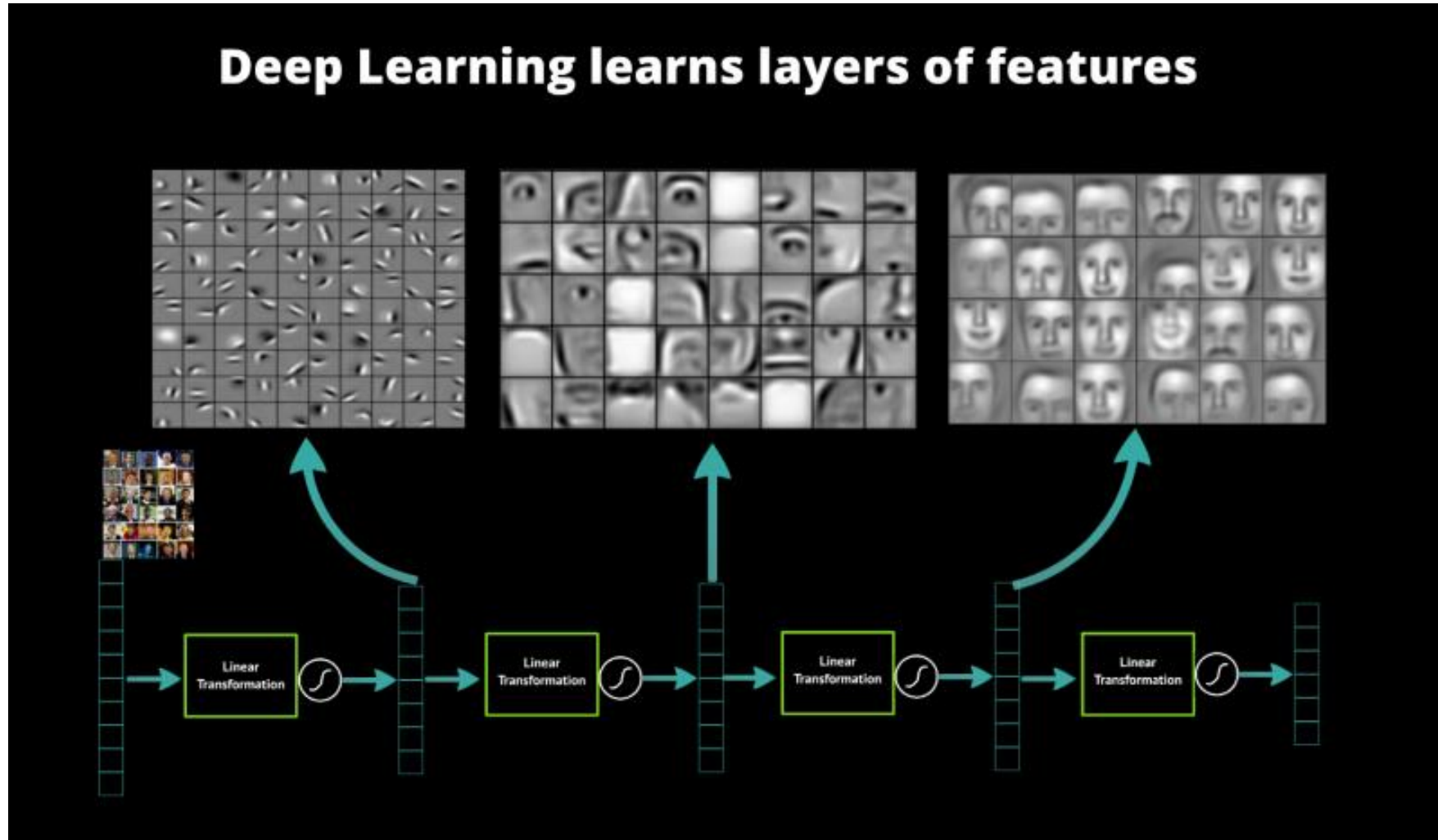
Deep learning:

Second "layer" of latent features
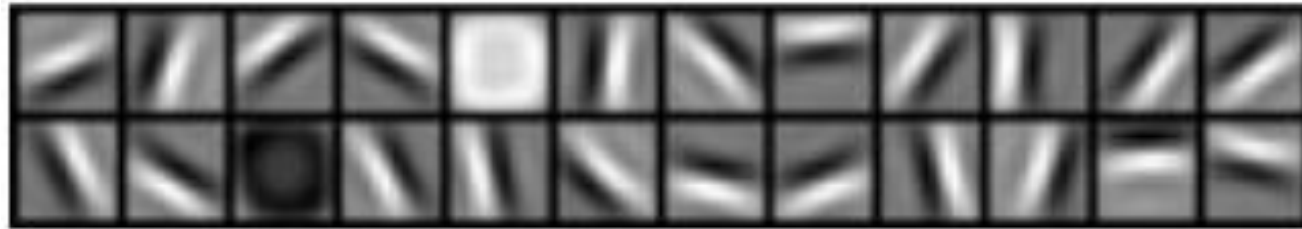
You can add more "layers" to go "deeper"

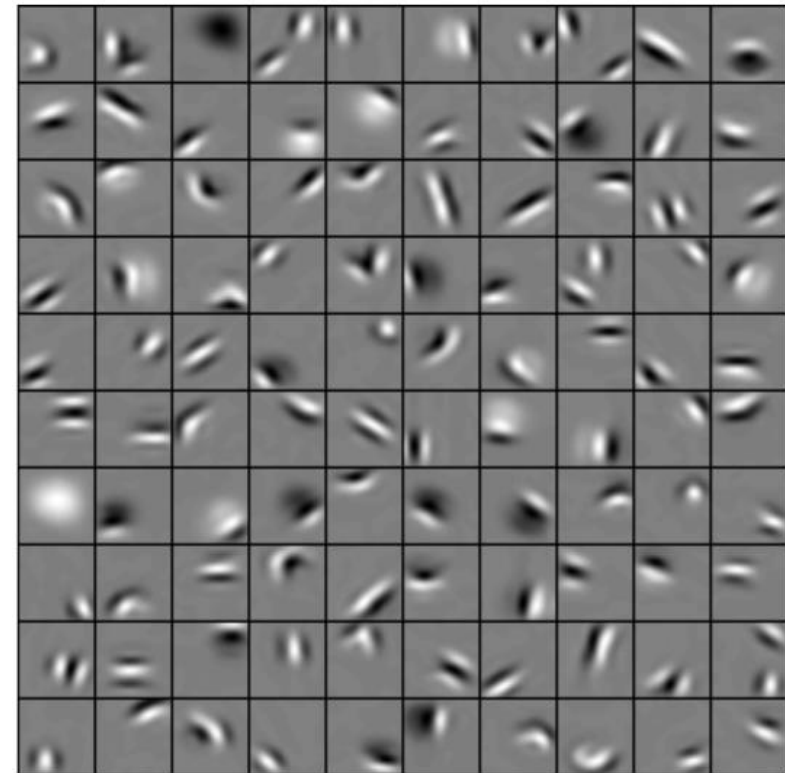# Cool Picture Motivation for Deep Learning

# Cool Picture Motivation for Deep Learning

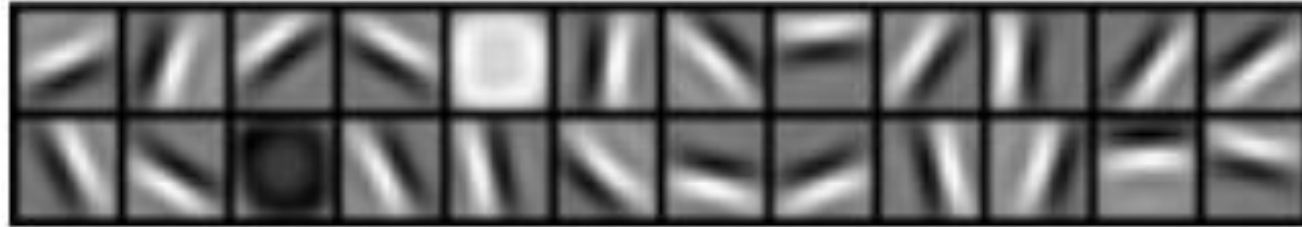- First layer of $z_i$ trained on 10 by 10 image patches:



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

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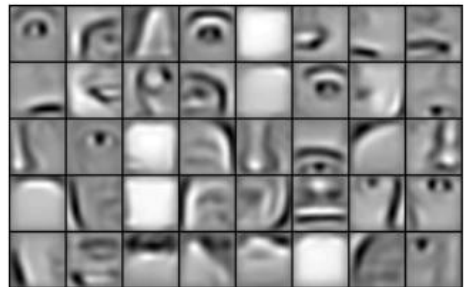  - We'll discuss these next time.

# Cool Picture Motivation for Deep Learning

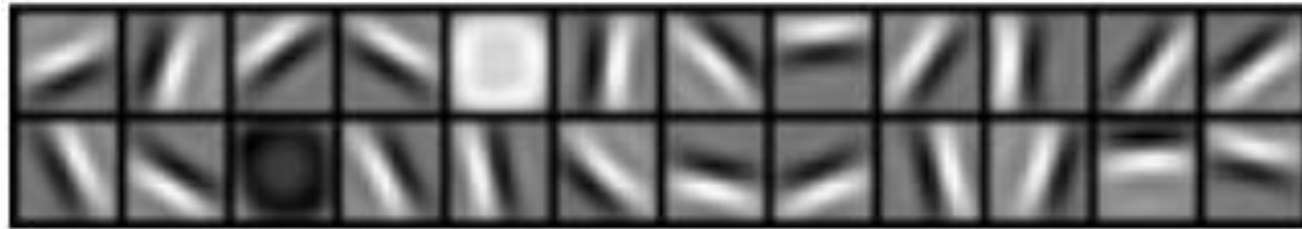- First layer of $z_i$ trained on 10 by 10 image patches:



- 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:



- Visualization of second and third layers trained on specific objects:
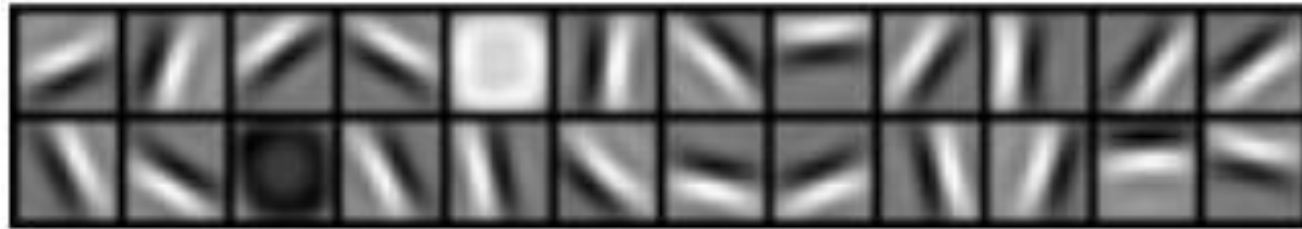
# Cool Picture Motivation for Deep Learning

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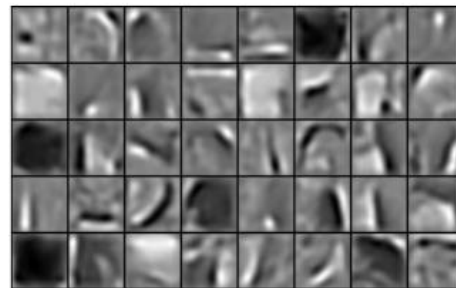


- Visualization of second and third layers trained on specific objects:
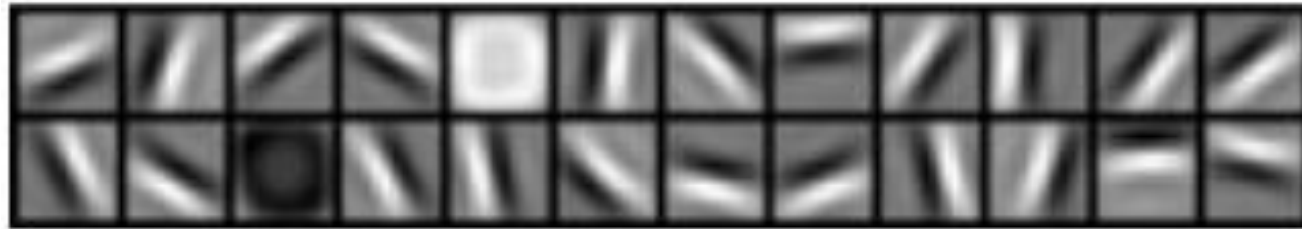


faces      cars      elephants

# Cool Picture Motivation for Deep Learning

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



- 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:
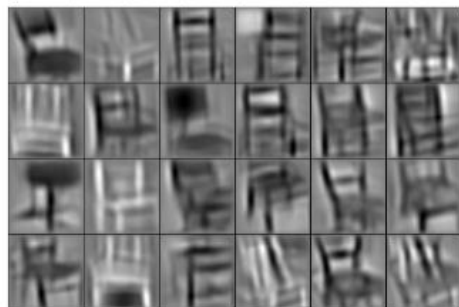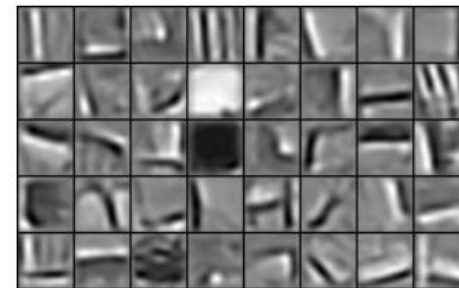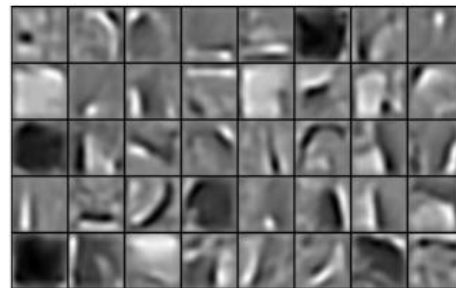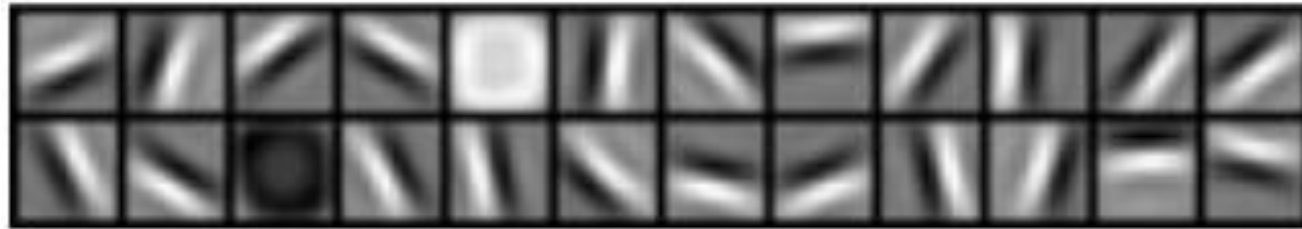


- Visualization of second and third layers trained on specific objects:
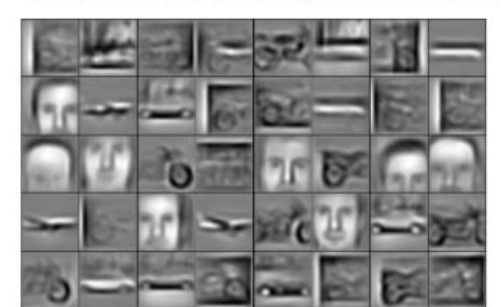


faces        cars        elephants        chairs        faces, cars, airplanes, motorbikes

# 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

- ## Then drop in popularity:

  - Quickly realized limitations of linear models.

# ML and Deep Learning History



- 1970 and 1980s: Connectionism (brain-inspired ML)
  - 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.

  - ML moved closer to other fields (CPSC 540):
    - Numerical optimization.
    - Probabilistic graphical models.
    - Bayesian methods.

# 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.

# 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.

# ImageNet Challenge

- Millions of labeled images, 1000 object classes.



Easy for humans but hard for computers.

# ImageNet Challenge

- ## Object detection task:
  - Single label per image.
  - Humans: ~5% error.

(a) Siberian husky

(b) Eskimo dog

## Image classification

# ImageNet Challenge

- Object detection task:
  - Single label per image.
  - Humans: ~5% error.



(a) Siberian husky



(b) Eskimo dog

## Image classification



3 → usual improvement

# ImageNet Challenge

- Object detection task:
  - Single label per image.
  - Humans: ~5% error.


(a) Siberian husky


(b) Eskimo dog

## Image classification



usual improvement

switch to deep learning (8 layers)

Classification error: 0.3, 0.2, 0.1, 0

2010  2011  2012

# ImageNet Challenge

- ## Object detection task:
  - Single label per image.
  - Humans: ~5% error.

(a) Siberian husky  (b) Eskimo dog

## Image classification



Classification error

- 0.3
- 0.2
- 0.1
- 0

2010  2011  2012  2013

usual improvement

switch to deep learning (8 layers)

tweaks

# ImageNet Challenge

- ## Object detection task:
  - Single label per image.
  - Humans: ~5% error.

(a) Siberian husky

(b) Eskimo dog

## Image classification

*usual improvement*

*switch to deep learning (8 layers)*

*tweaks*

*deeper!*

*GoogLeNet: 6.7% error 22 layers*

# ImageNet Challenge

- Object detection task:
  - Single label per image.
  - Humans: ~5% error.


(a) Siberian husky    (b) Eskimo dog

- 2015 winner: Microsoft
  - 3.6% error.
  - 152 layers.

## Image classification



*usual improvement*

*switch to deep learning (8 layers)*

*tweaks*

*deeper!*

*GoogLeNet: 6.7% error 22 layers*

# Adding a Bias Variable

Remember that in linear models we may nonzero $y$-intercept:

$$y_i = w^T x_i + \beta$$

For neural networks we could have explicit bias:

$$y_i = w^T h(W x_i) + \beta$$

We can just use $y_i = w^T x_i$ if we fix $x_{ij} = 1$ for all '$i$' for some '$j$'.

(constant $x_{ij}$ value)

Or we could set $W_c = 0$ for one ~~row~~ column $W_c$ of 'W':

$$\frac{1}{1 + \exp(-W_c x_i)} = \frac{1}{1 + \exp(0)} = \frac{1}{2}$$

(constant $z_{ic}$ value)

# Artificial Neural Networks

- With squared loss, our objective function is:

$$f(w, W) = \frac{1}{2} \sum_{i=1}^{n} (w^\top h(W x_i) - y_i)^2$$

- Usual training procedure: stochastic gradient.
  - Compute gradient of random example 'i', update both 'w' and 'W'.
  - Highly non-convex and can be difficult to tune.

- Computing the gradient is known as "backpropagation".

# Backpropagation

- Consider the loss for a single example:

$$f(w, W) = \frac{1}{2}\left(\sum_{c=1}^{k} w_c\, h(W_c\, x_i) - y_i\right)^2$$

Element 'c' of 'w' ↙

Row 'c' of $W$ →

- Derivative with respect to '$w_c$':

From squared loss

$$\frac{\partial}{\partial w_c}\left[f(w, W)\right] = \left(\sum_{c=1}^{k} w_c\, h(W_c\, x_i) - y_i\right) h(W_c\, x_i)$$

derivative with respect to $w_c$

- Derivative with respect to '$W_{cj}$'

derivative with respect to $W_{cj}$

$$\frac{\partial}{\partial W_{cj}}\left[f(w, W)\right] = \left(\sum_{c=1}^{k} w_c\, h(W_c\, x_i) - y_i\right) w_c\, h'(W_c\, x_i)\, x_{ij}$$

derivative with respect to $W_{c x_i}$

# Backpropagation

- Notice repeated calculations in gradients:

$$\frac{\partial}{\partial w_c}\left[f(w, W)\right] = \left(\underbrace{\sum_{c=1}^{K} w_c\, h(W_c x_i) - y_i}_{r_i}\right) h(W_c x_i)$$

$$= \underbrace{r_i}\, h(W_c x_i) \qquad \longrightarrow \text{same } r_i \text{ for } \underline{all} \text{ 'c'}$$

$$\frac{\partial}{\partial W_{cj}}\left[f(w, W)\right] = \left(\underbrace{\sum_{c=1}^{K} w_c\, h(W_c x_i) - y_i}_{r_i}\right) \underbrace{w_c\, h'(W_c x_i)}_{v_c} x_{ij}$$

$$= \underline{r_i}\, \underline{v_c}\, x_{ij} \qquad \longrightarrow \text{same } v_c \text{ for } \underline{all} \text{ 'j'}$$
$$\longrightarrow \text{same } r_i \text{ for } \underline{all} \text{ 'c')}$$

# Backpropagation

- Calculation of gradient is split into two phases:

1. "Forward" pass

   (a) Compute $h(W_c x_i)$ for all 'c'

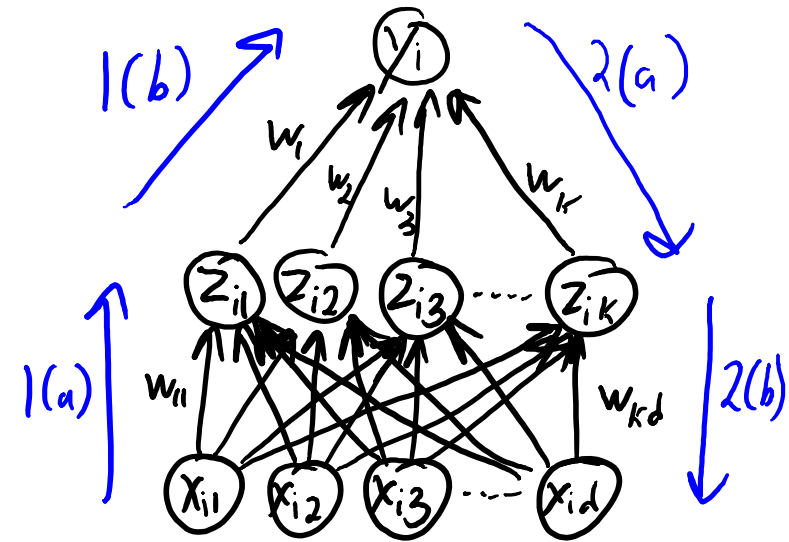   (b) Compute residual $r_i = \left( \sum_{c=1}^{k} w_c h(W_c x_i) - y_i \right)$

2. "Backprogation"

   (a) Compute $\dfrac{\partial f}{\partial w_c} = r_i\, h(W_c x_i)$ for all 'c'

   (b) Compute $v_c = w_c h'(W_c x_i)$ for all 'c'

   (c) Compute $\dfrac{\partial f}{\partial W_{cj}} = r_i\, v_c\, x_{ij}$ for all 'c' and 'j'

# Summary

- Biological motivation for (deep) neural networks.

- Deep learning considers neural networks with many hidden layers.

- Unprecedented performance on difficult pattern recognition tasks.

- Backpropagation computes neural network gradient via chain rule.


- Next time:
  - How deep learners fight the fundamental trade-off.