

CPSC 340: Machine Learning and Data Mining

Deep Learning

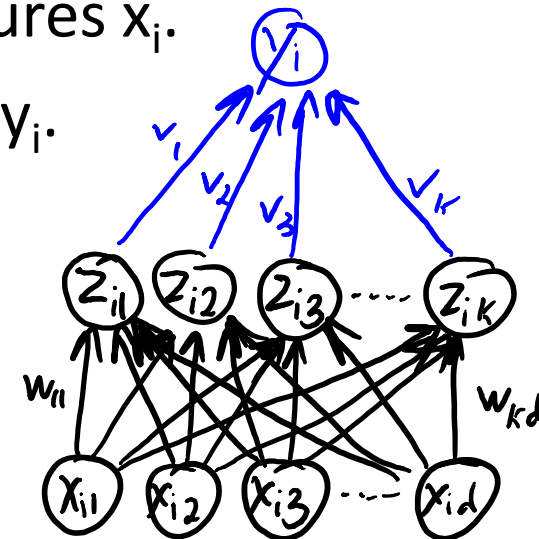
Fall 2017

Admin

- **Assignment 4:**
 - 2 late days for tonight Wednesday.
- **Assignment 5:**
 - Due Monday.
- **Final:**
 - Details and previous exams posted on Piazza.
- **Extra office hours:**
 - 3:00 next Thursday.

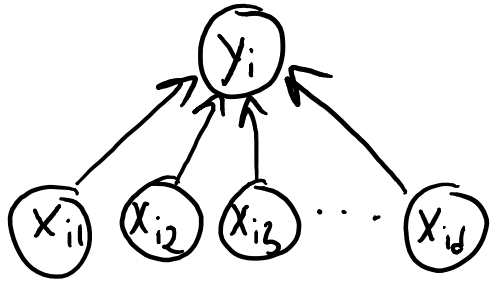
Supervised Learning Roadmap

- 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.**

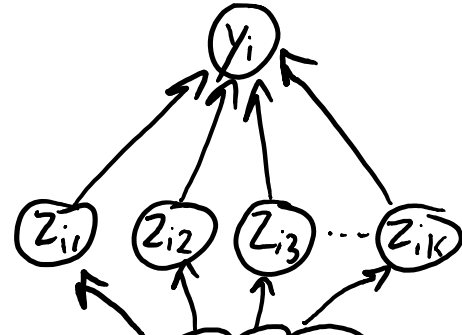


A Graphical Summary of CPSC 340 Parts 1-5

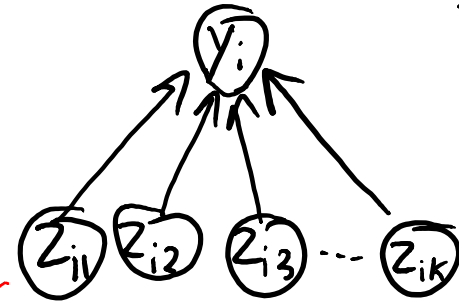
Part 1: "I have features x_i "



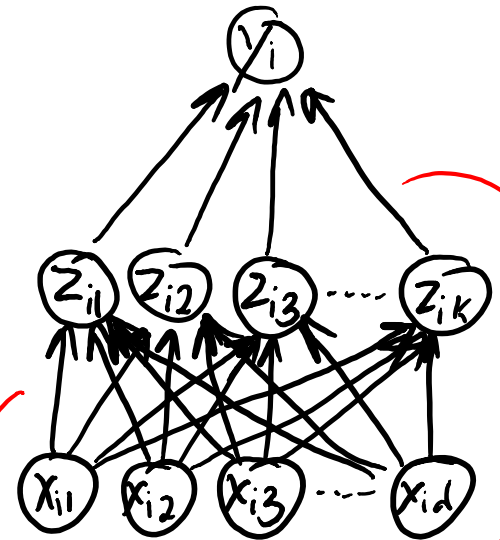
Part 3: change of basis



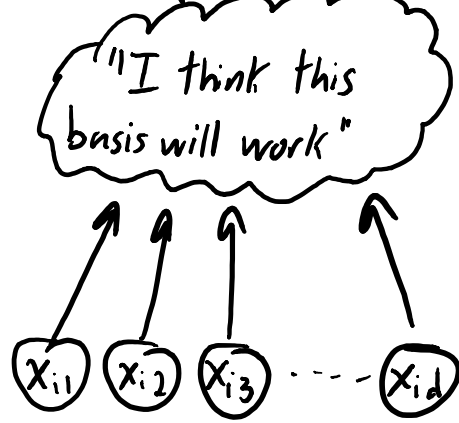
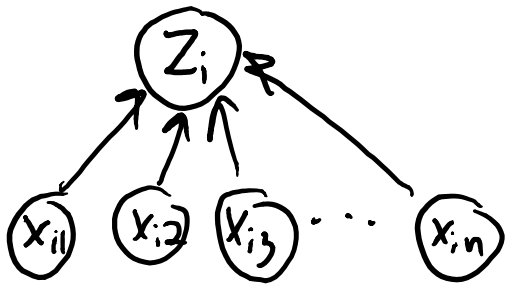
Part 4: basis from latent-factor model



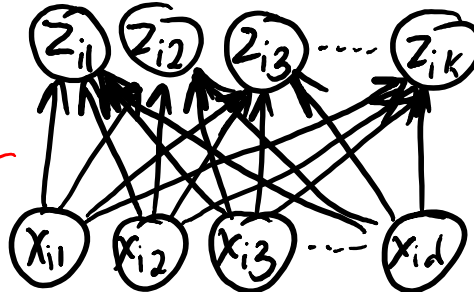
Part 5: Neural networks



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



"PCA will give me good features"



"What are the 'parts' of x_i ?"

Trained separately

Learn features and classifier at the same time.

Notation for Neural Networks

We have our usual supervised learning notation:

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

$n \times d$ $n \times 1$

We have our latent features:

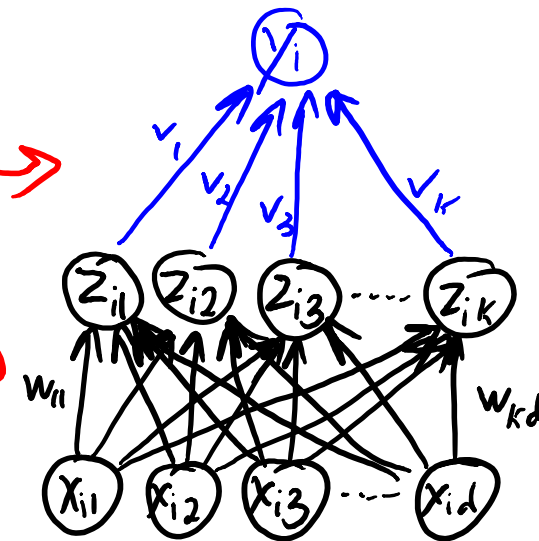
$$Z = \begin{bmatrix} \text{---} z_1^T \text{---} \\ \text{---} z_2^T \text{---} \\ \vdots \\ \text{---} z_n^T \text{---} \end{bmatrix}$$

$n \times k$

We have two sets of parameters:

$$V = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_k \end{bmatrix} \quad W = \begin{bmatrix} \text{---} w_1 \text{---} \\ \text{---} w_2 \text{---} \\ \vdots \\ \text{---} w_k \text{---} \end{bmatrix}$$

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



Linear-Linear Model

- Obvious choice: linear latent-factor model with linear regression.

Use features from latent-factor model: $z_i = Wx_i$

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

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

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

- But this is just a linear model:

$$y_i = v^T z_i = v^T (Wx_i) = \overbrace{(v^T W)}^{1 \times d} x_i = \underbrace{\text{some vector } 'w'} x_i = w^T x_i$$

Introducing Non-Linearity

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

$$z_i = Wx_i \quad y_i = v^T h(z_i)$$

– Here the function 'h' transforms 'k' inputs to 'k' outputs.

- Common choice for 'h': applying **sigmoid** function element-wise:

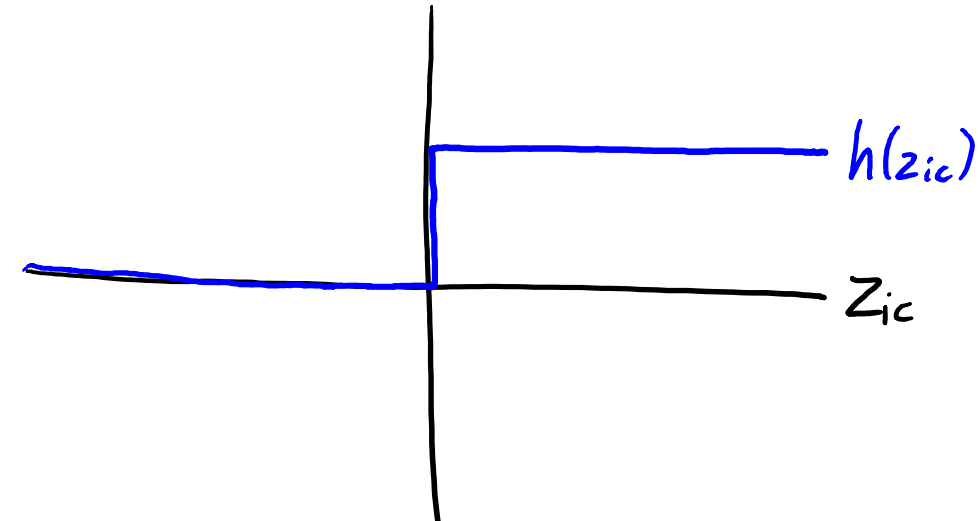
$$h(z_{ic}) = \frac{1}{1 + \exp(-z_{ic})}$$

- So this takes the z_{ic} in $(-\infty, \infty)$ and maps it to $(0, 1)$.
- This is called a “multi-layer perceptron” or a “**neural network**”.

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”:

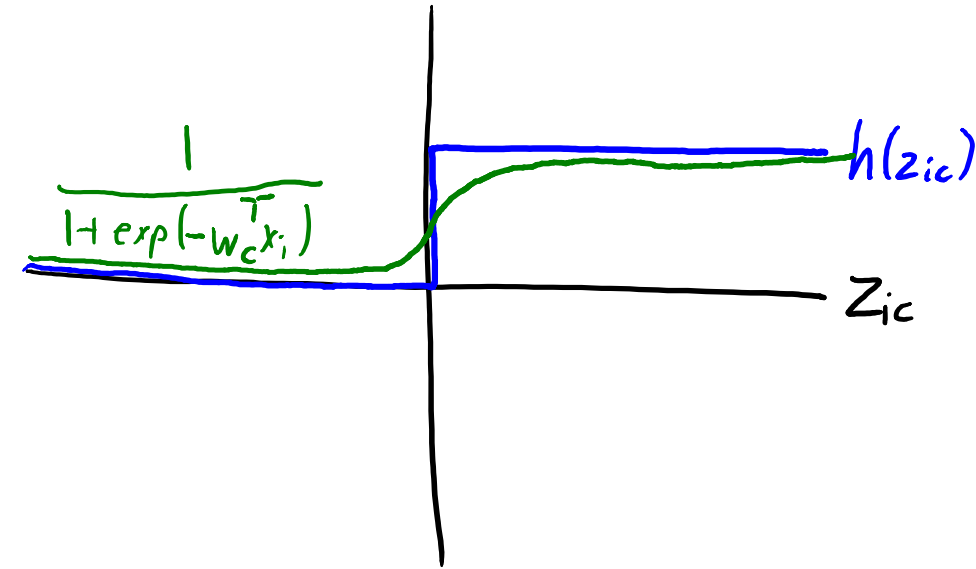
3	=	1	+	1	+	1	+	1	+	1	+	0	+	0
5	=	1	+	0	+	1	+	1	+	1	+	0	+	1
8	=	1	+	1	+	1	+	1	+	1	+	1	+	1

The figure shows three rows of digit representations. Each row starts with a digit (3, 5, or 8) followed by an equals sign and a sequence of parts. The parts are small boxes containing segments of the digit. The parts are combined with plus signs (+) and some are preceded by a coefficient (1 or 0). The coefficients are color-coded: green for 1 and red for 0.

Why Sigmoid?

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

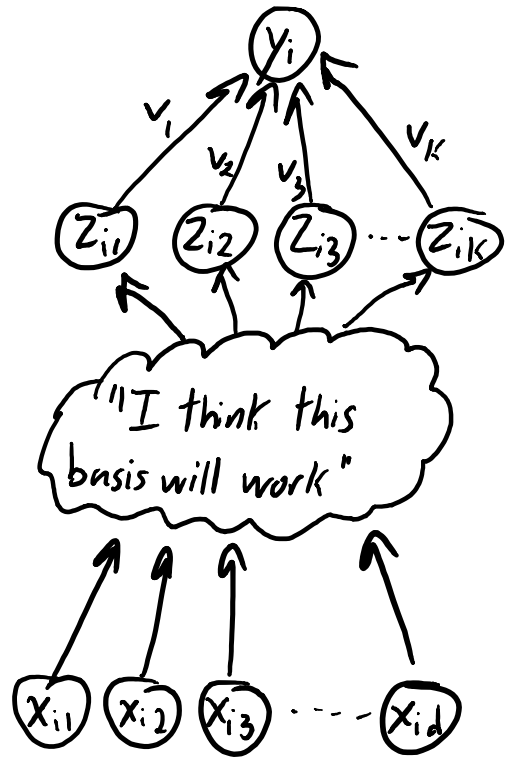
$$h(z_{ic}) = \begin{cases} 1 & \text{if } z_{ic} \geq 0 \\ 0 & \text{if } z_{ic} < 0 \end{cases}$$



- 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”.
- But this is hard to optimize (**non-differentiable/discontinuous**).
- Sigmoid is a smooth approximation to these binary features.

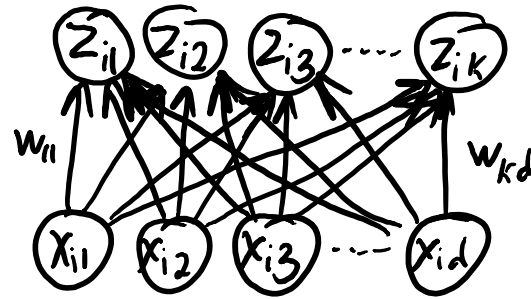
Supervised Learning Roadmap

Hand-engineered features:

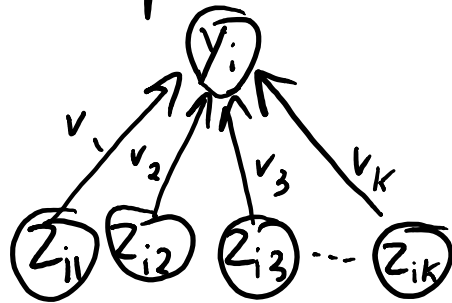


Requires domain knowledge and can be time-consuming

Learn a latent-factor model:

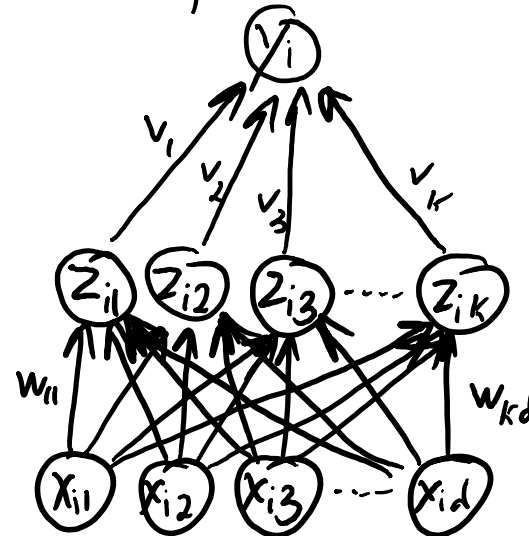


Use latent features in supervised model:



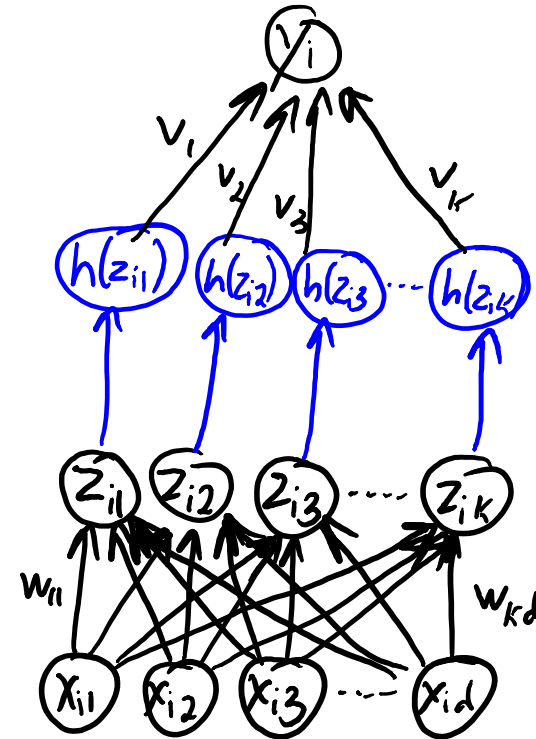
Good representation of x_i might be bad for predicting y_i

Learn ' v ' and ' W ' together:



But still gives a linear model.

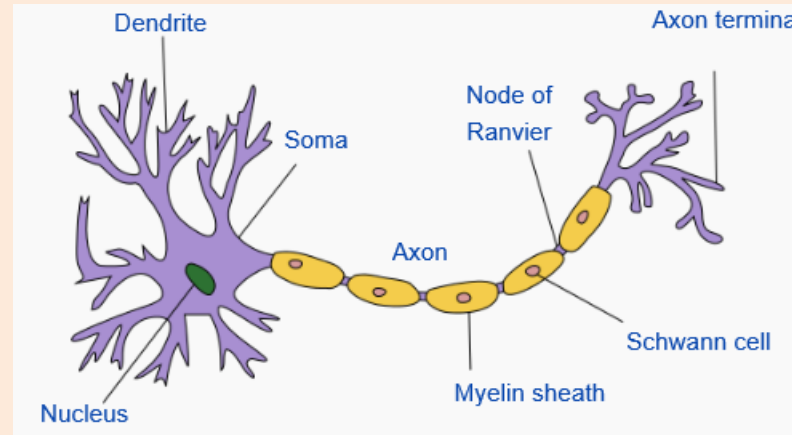
Neural network:



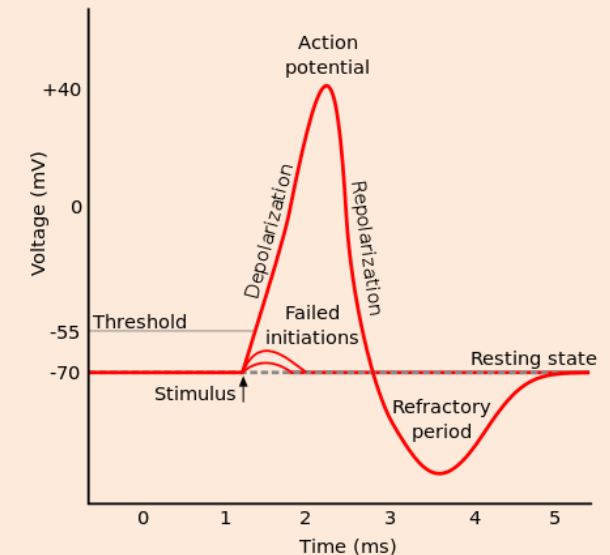
Extra non-linear transformation ' h '

Why “Neural Network”?

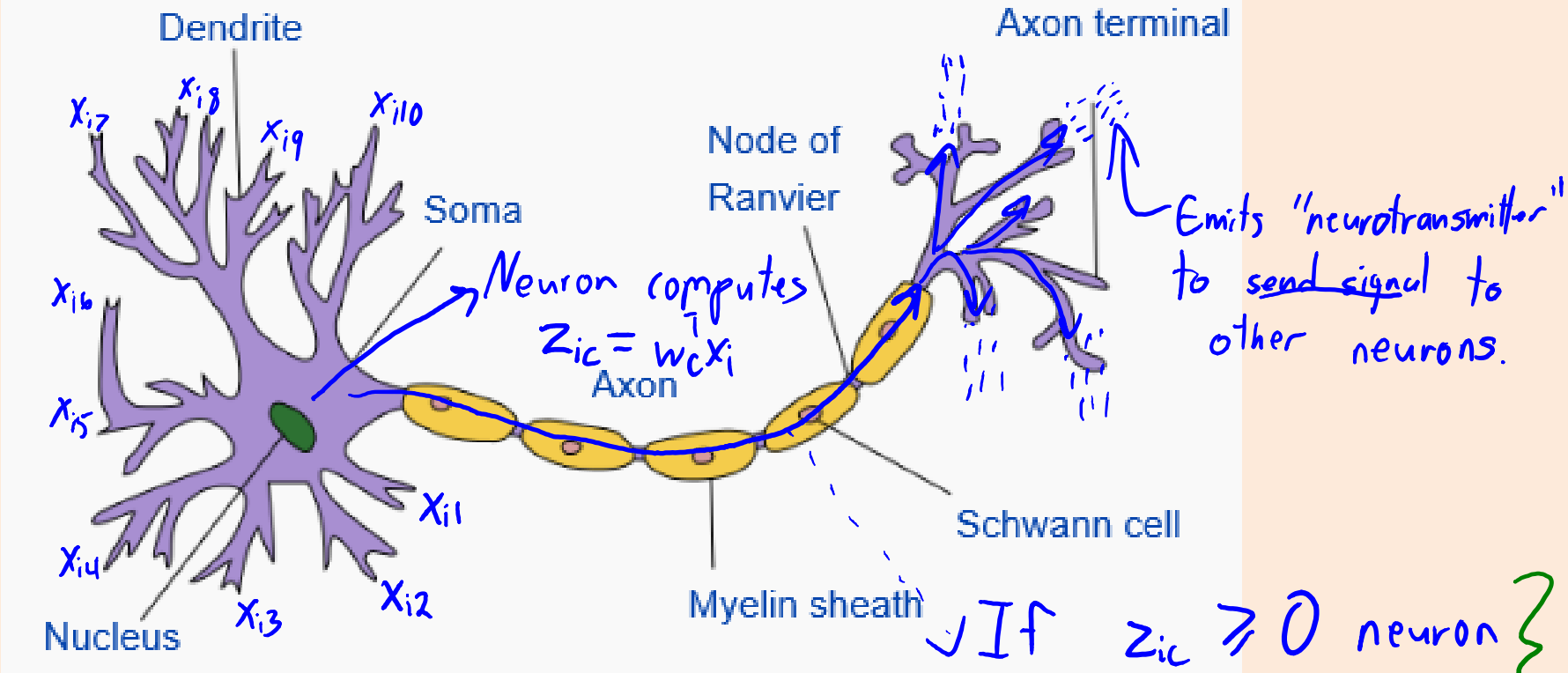
- Cartoon of “typical” neuron:



- 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):

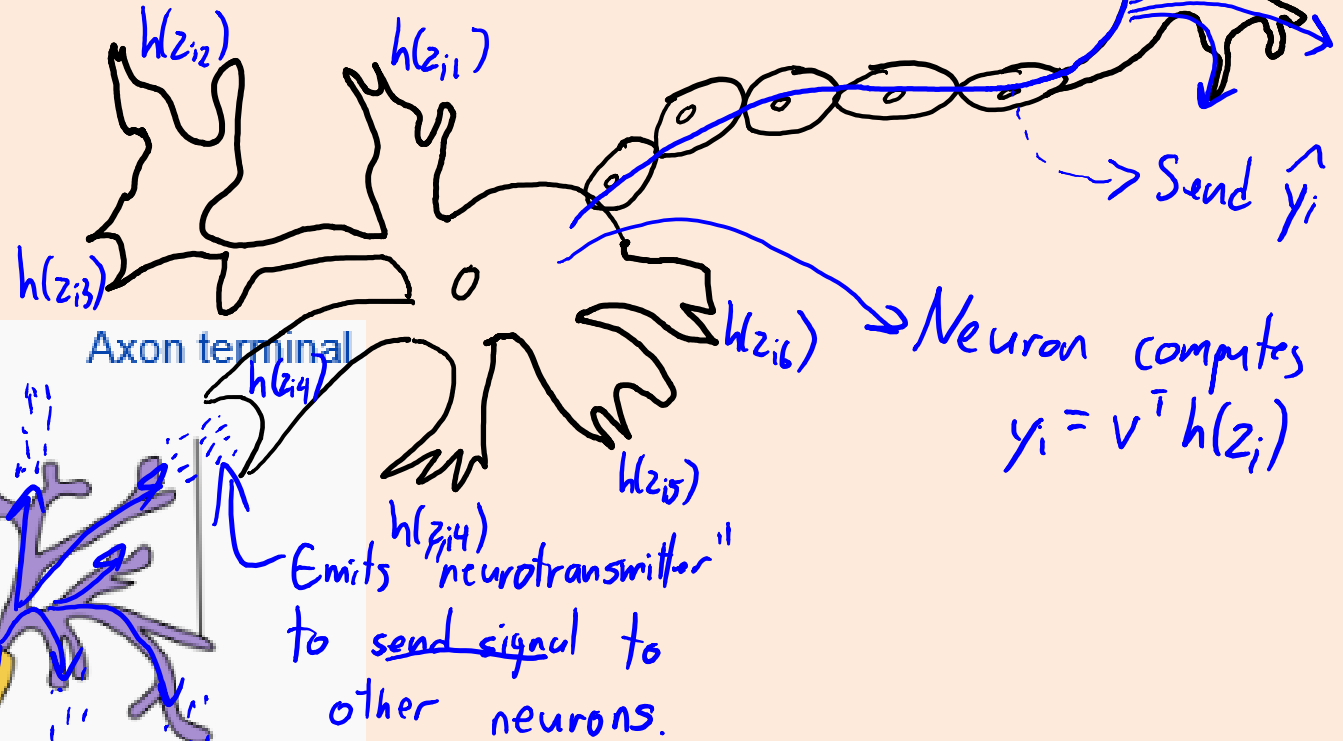
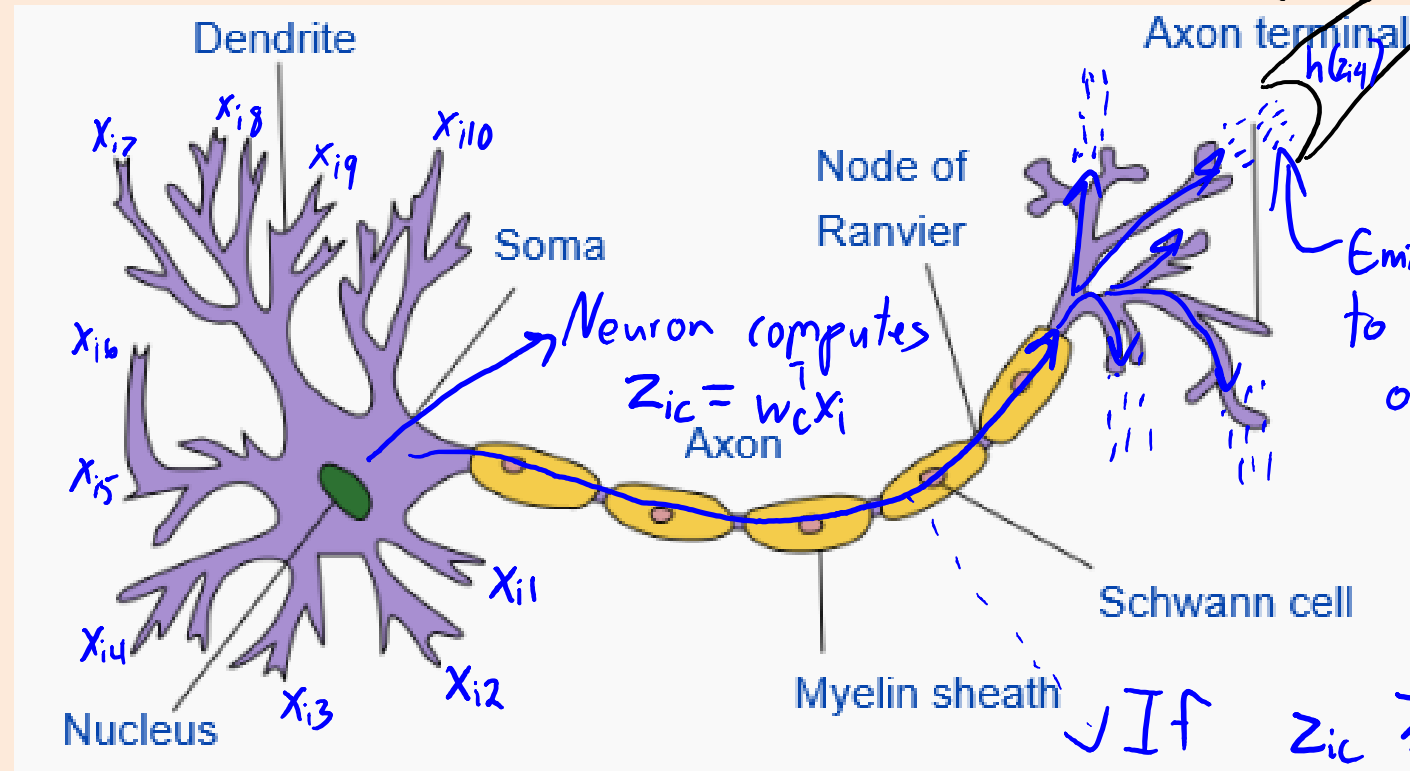


Why "Neural Network"?



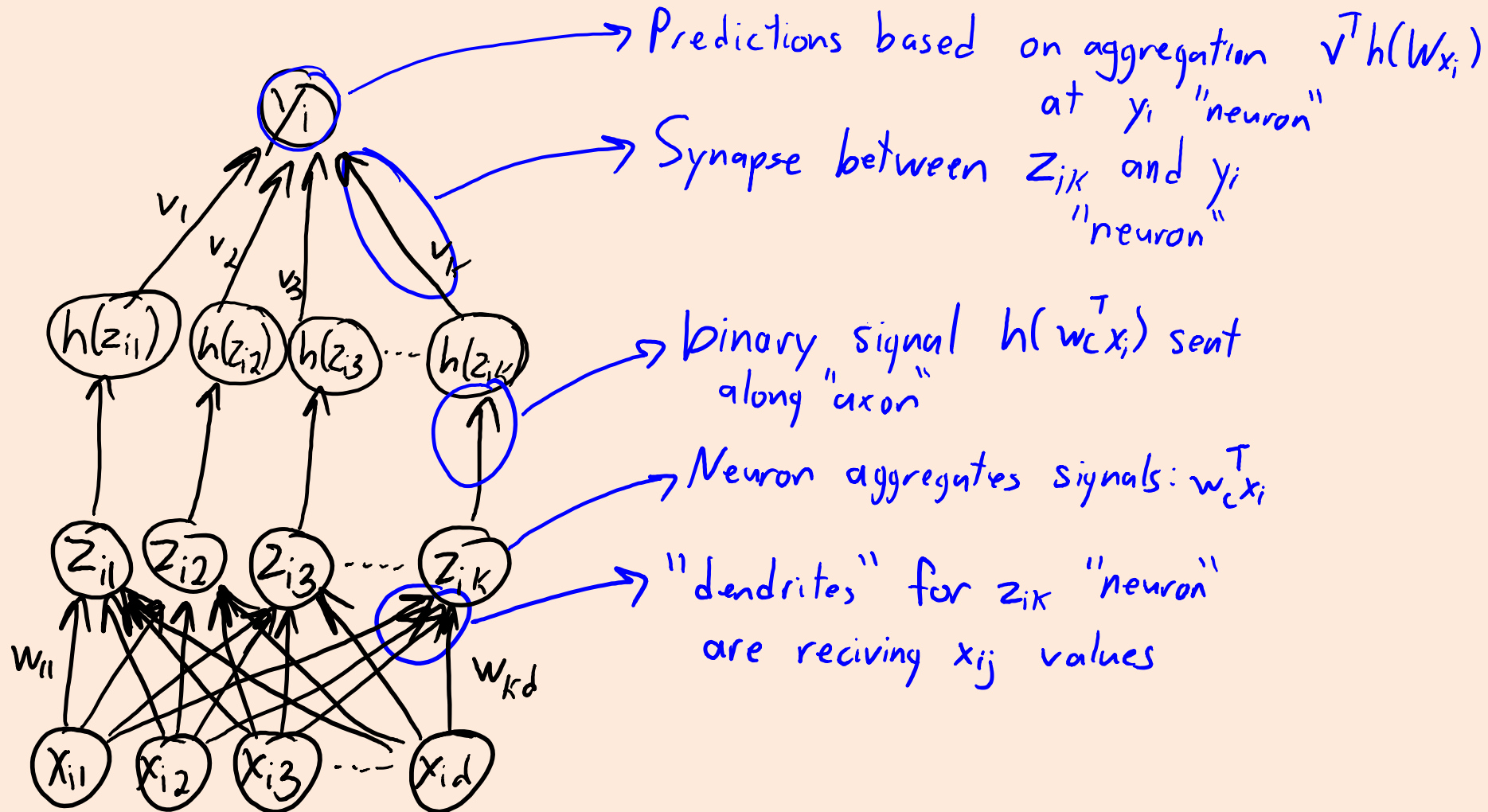
If $z_{ic} \geq 0$ neuron } We approximate binary
 sends signal along axon. } signal with $\frac{1}{1 + \exp(-z_{ic})}$

Why "Neural Network"?

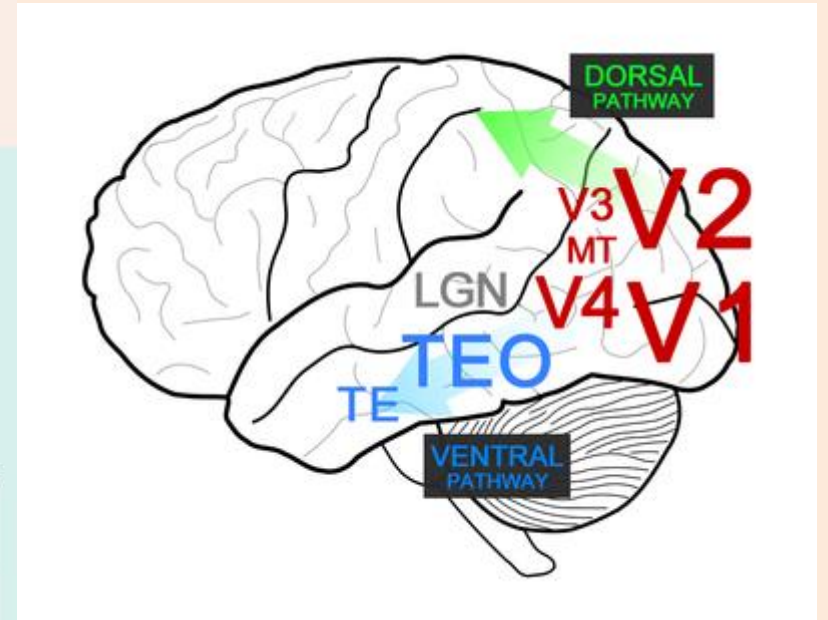
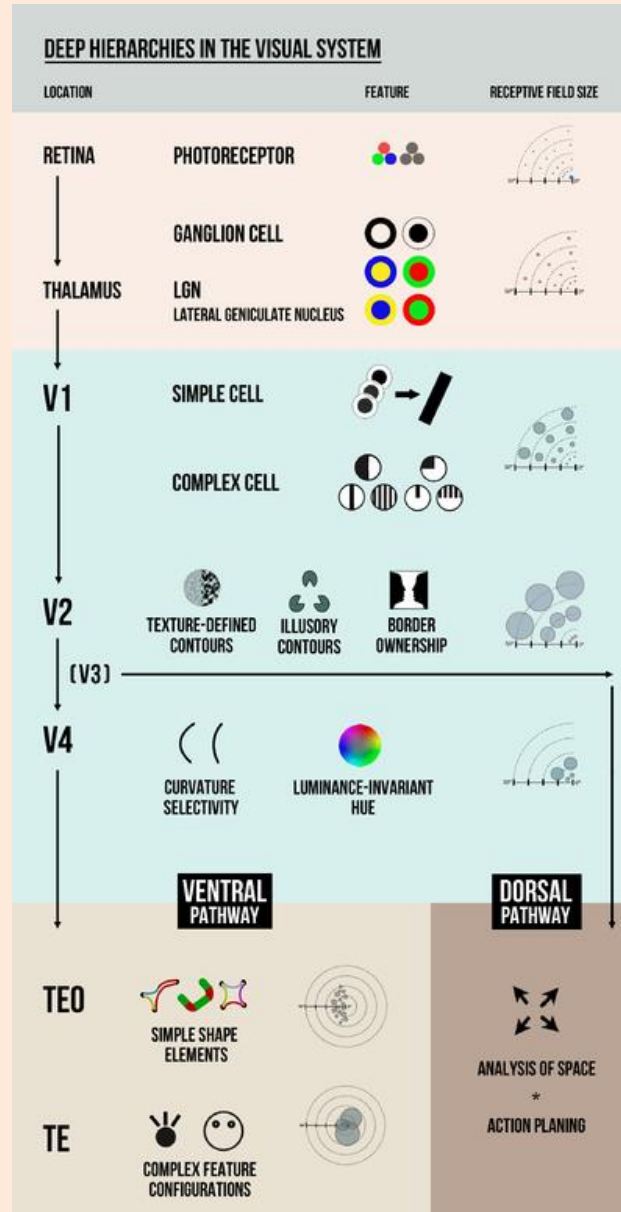
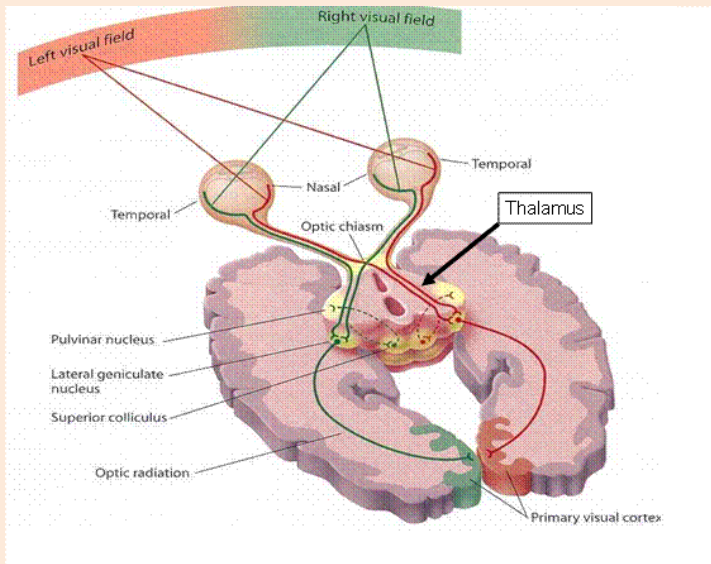


\checkmark If $z_{ic} \geq 0$ neuron } We approximate binary
Sends signal along axon. } signal with $\frac{1}{1 + \exp(-z_{ic})}$

Why "Neural Network"?

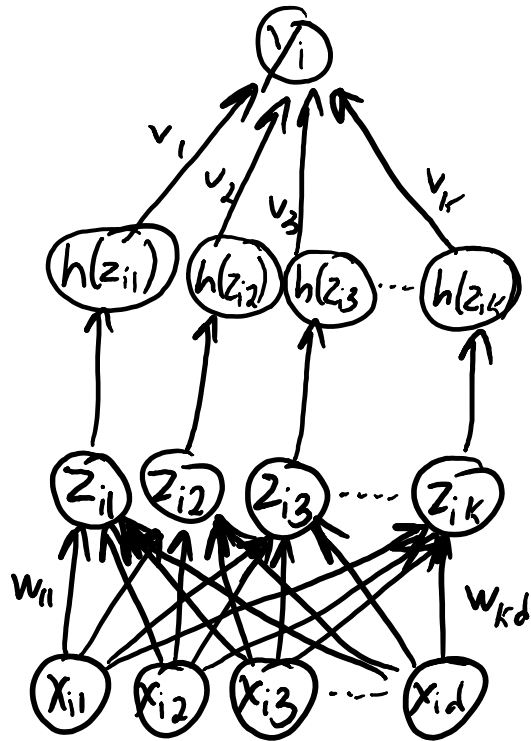


Deep Hierarchies in the Brain



Deep Learning

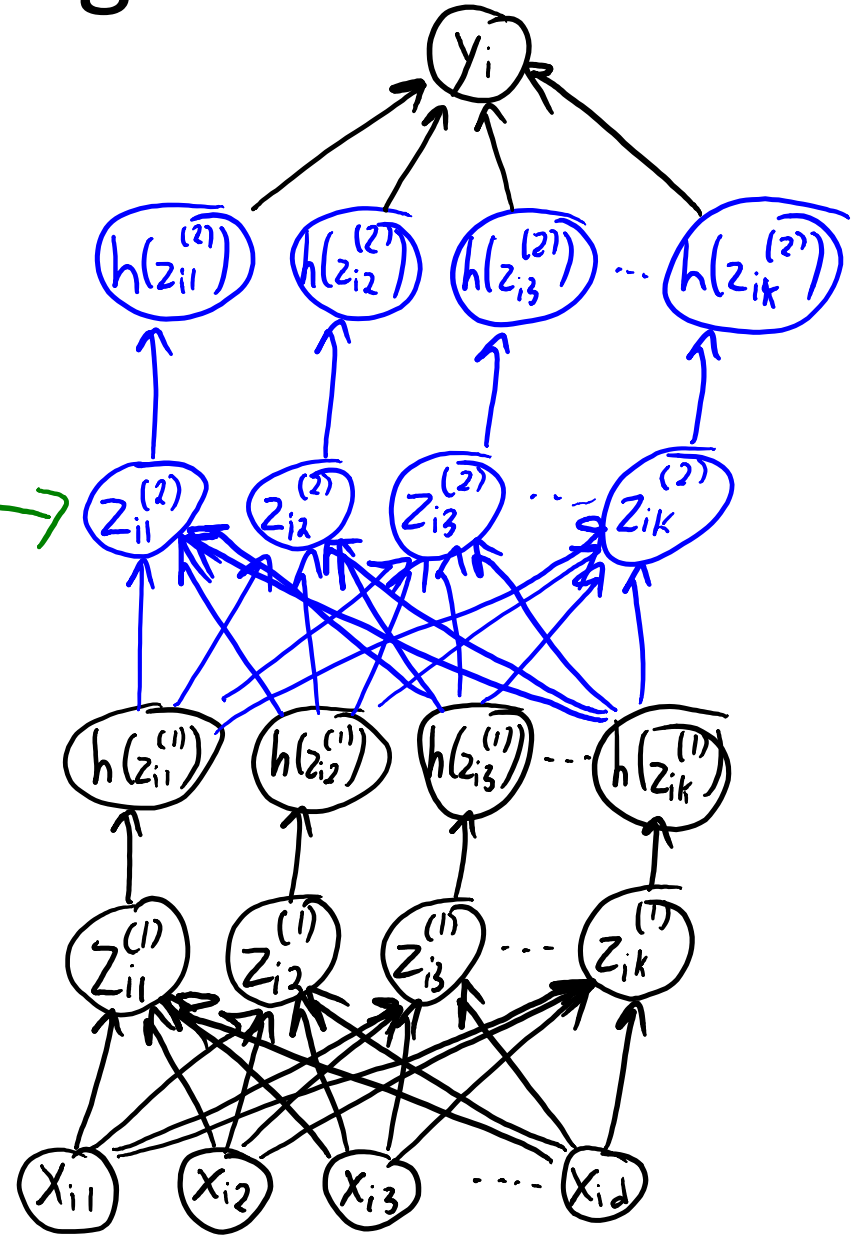
Neural network:



Deep learning:

Second "layer" of latent features

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

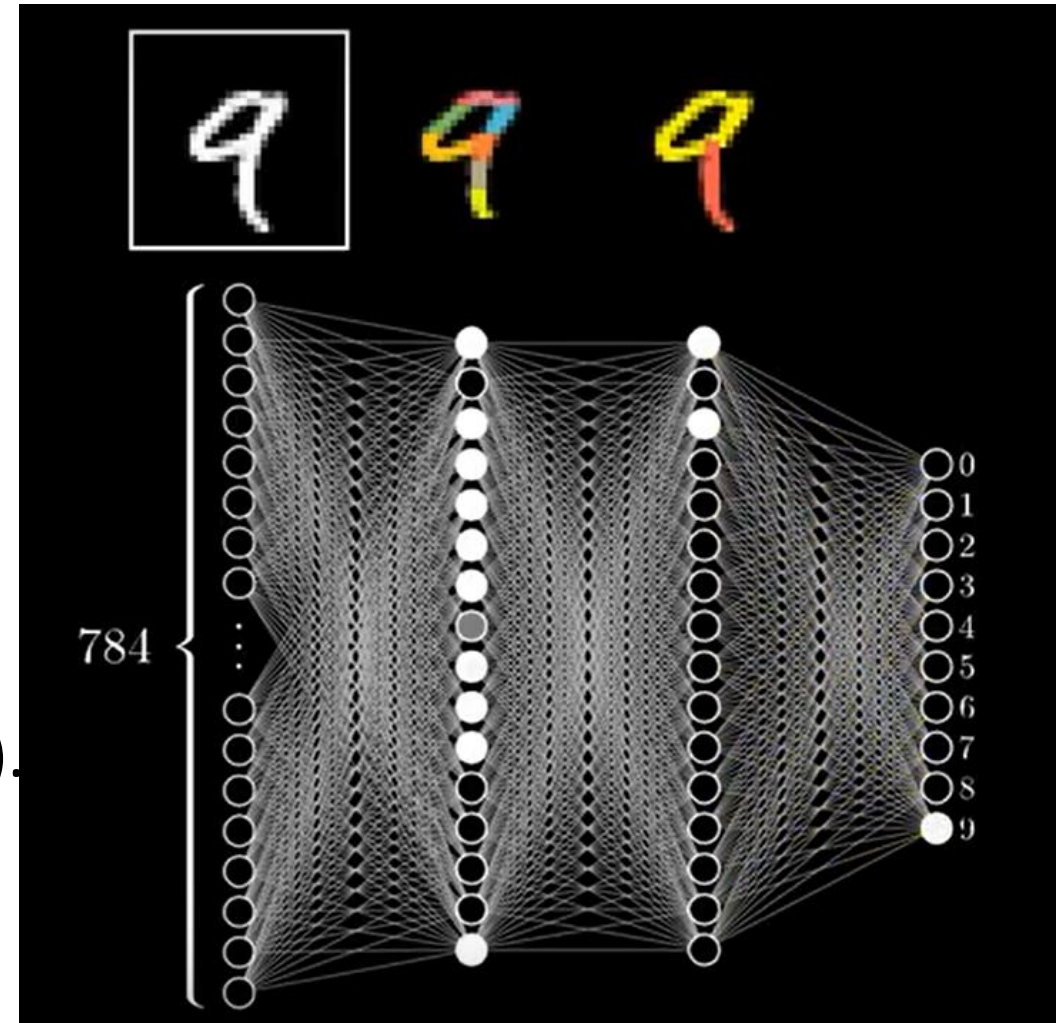


“Hierarchies of Parts” Motivation for Deep Learning

- Each “neuron” might recognize a “part” of a 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>



Deep Learning

Linear model:

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

Neural network with 1 hidden layer:

$$\hat{y}_i = v^T h(Wx_i)$$

z_i

Neural network with 2 hidden layers:

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

$z_i^{(2)}$ $z_i^{(1)}$

Neural network with 3 hidden layers:

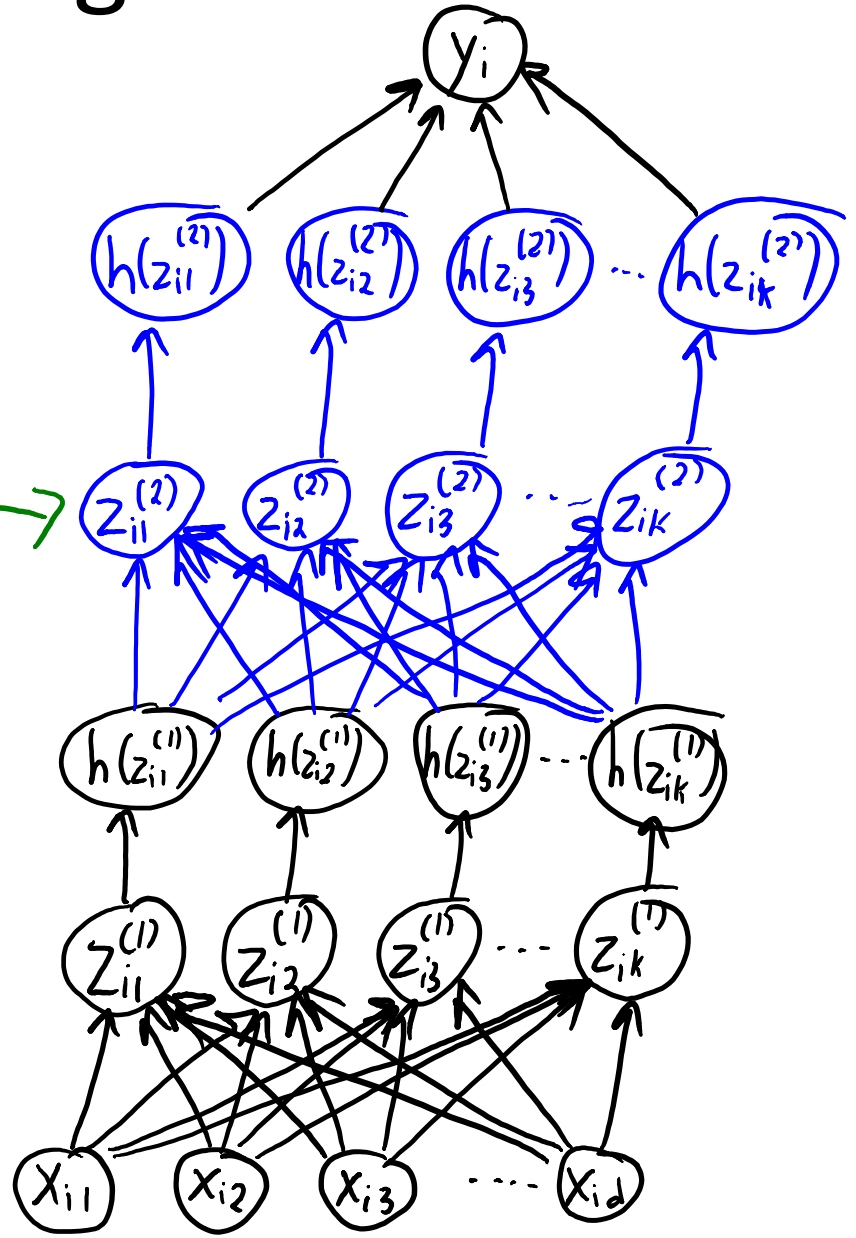
$$\hat{y}_i = v^T h(W^{(3)} h(W^{(2)} h(W^{(1)} x_i)))$$

$z_i^{(3)}$ $z_i^{(2)}$ $z_i^{(1)}$

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:

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

Symbol:

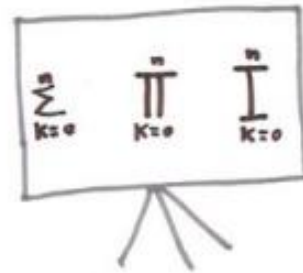
$$\prod_{k=0}^n f_k(+)$$

Meaning:

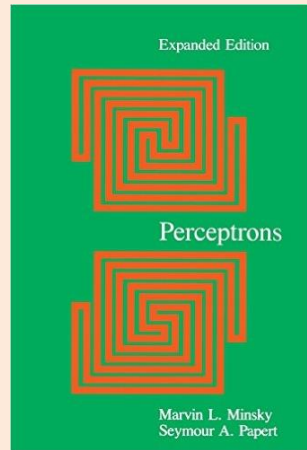
$$f_n \circ f_{n-1} \circ f_{n-2} \circ \dots \circ f_2 \circ f_1 \circ f_0(+)$$

- For 'm' layers, we could use:

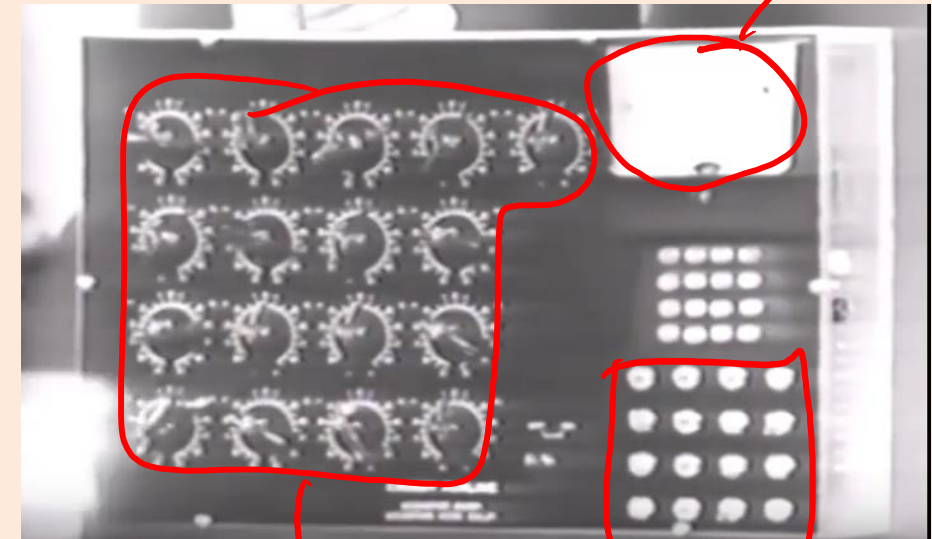
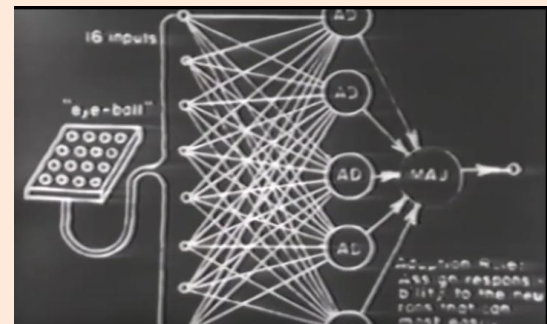
$$\hat{y}_i = w^T \left(\prod_{l=1}^m h(W^{(l)} x_i) \right)$$



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>
 - Marvin Minsky assigns object recognition to his students as a summer project
- Then drop in popularity:
 - Quickly realized **limitations of linear models.**



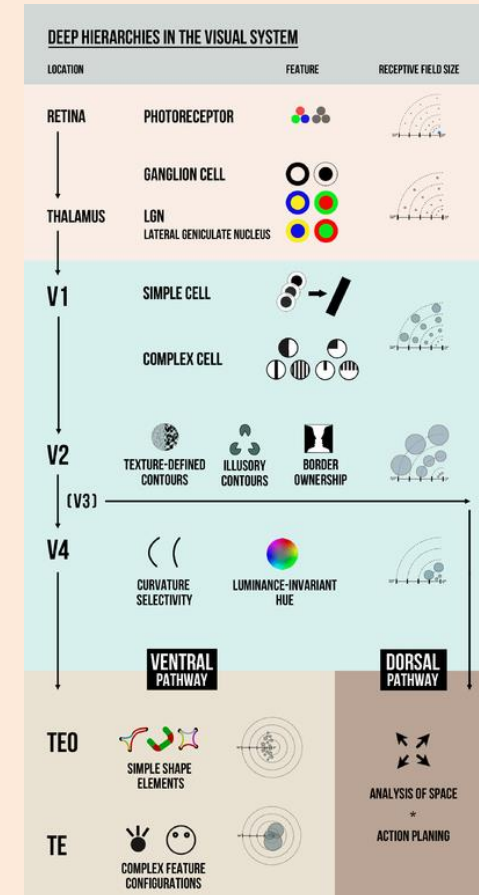
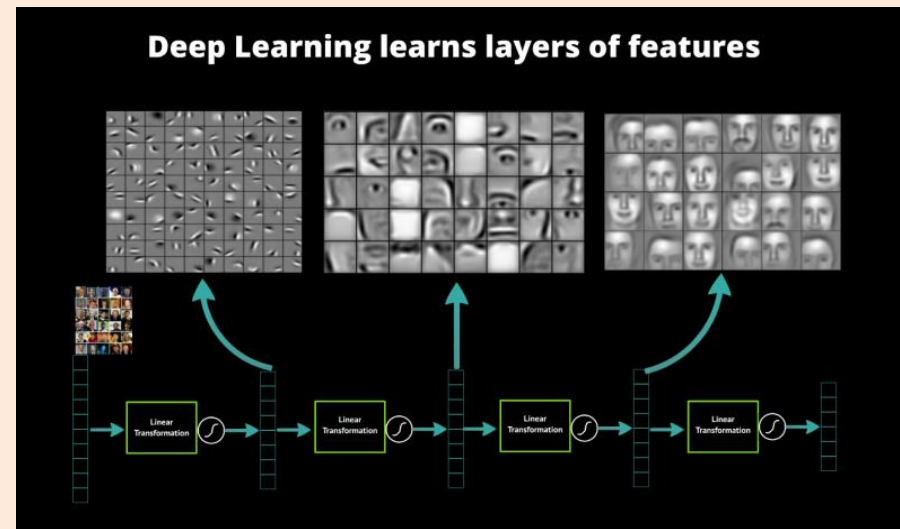
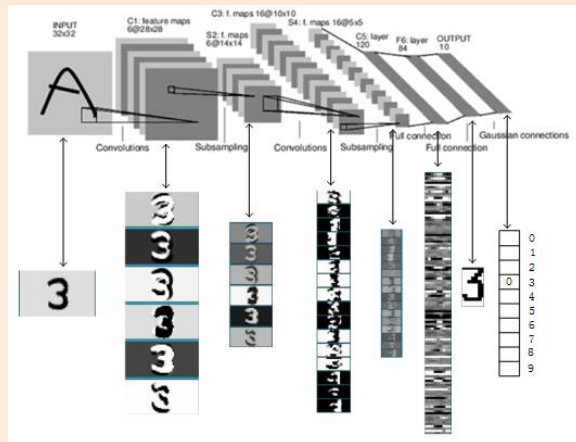
w

x_i

$w^T x_i$

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.

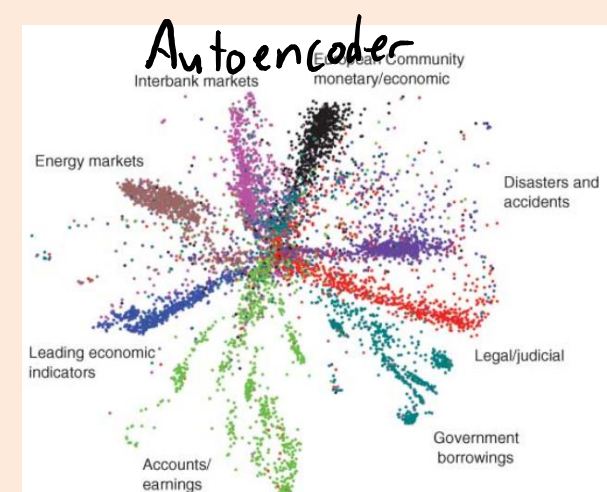
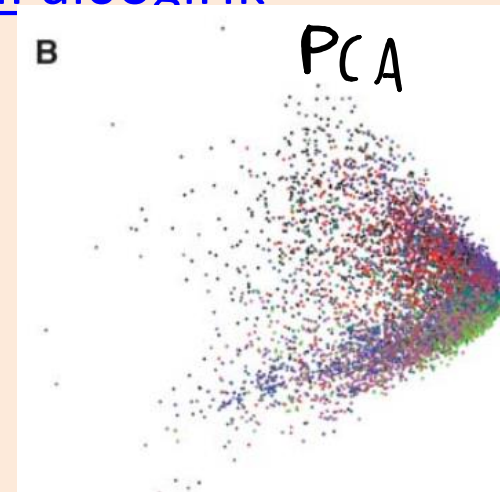


ML and Deep Learning History

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

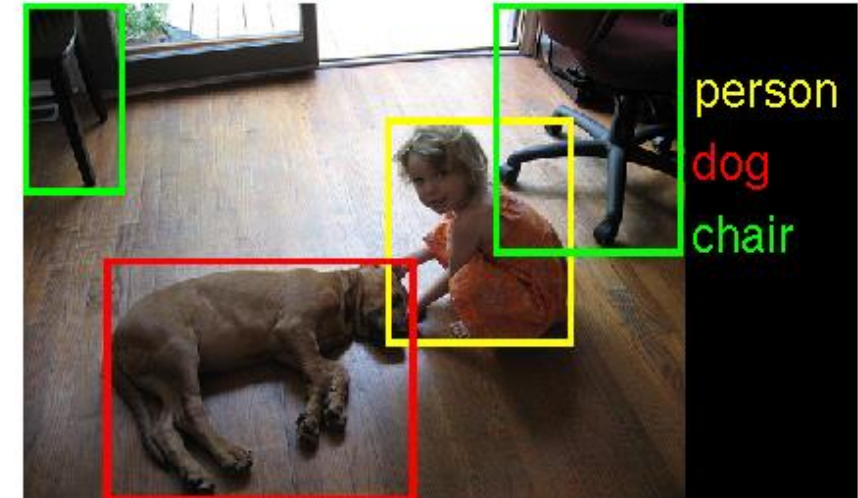
ML and Deep Learning History

- 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.

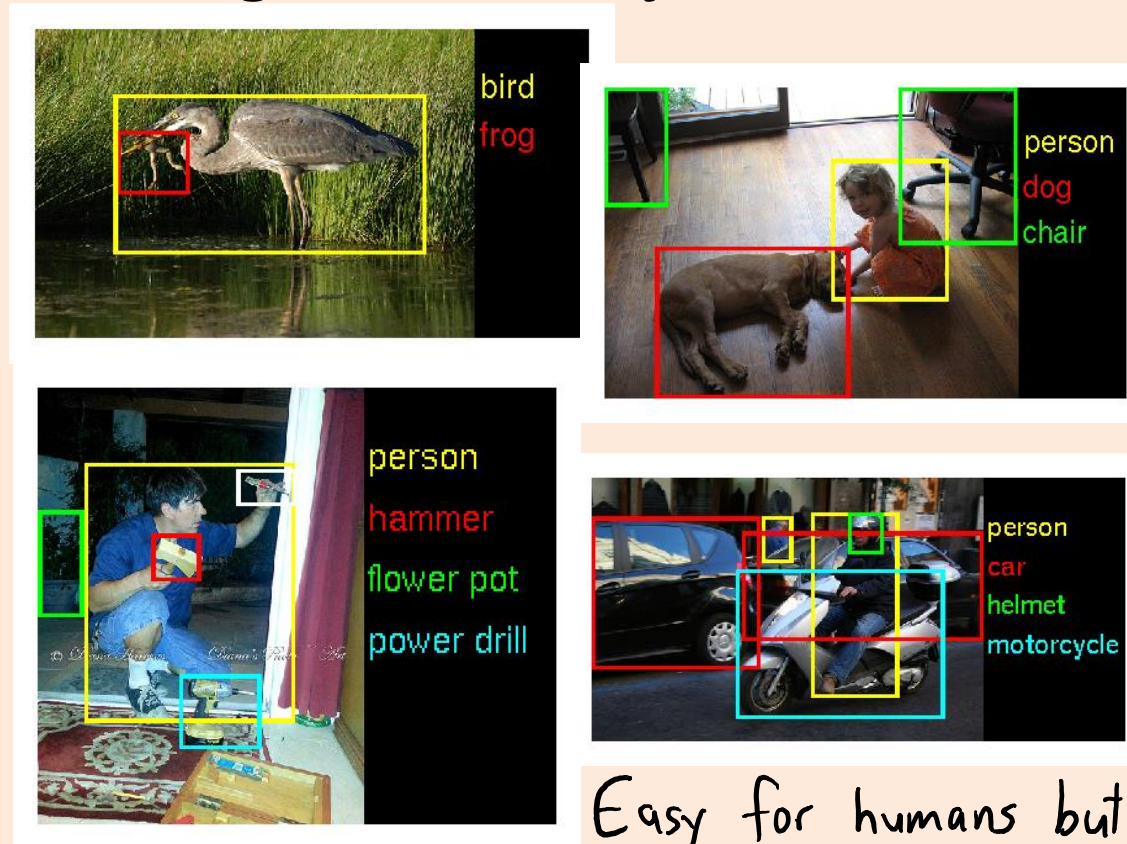


2010s: DEEP LEARNING!!!

- 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

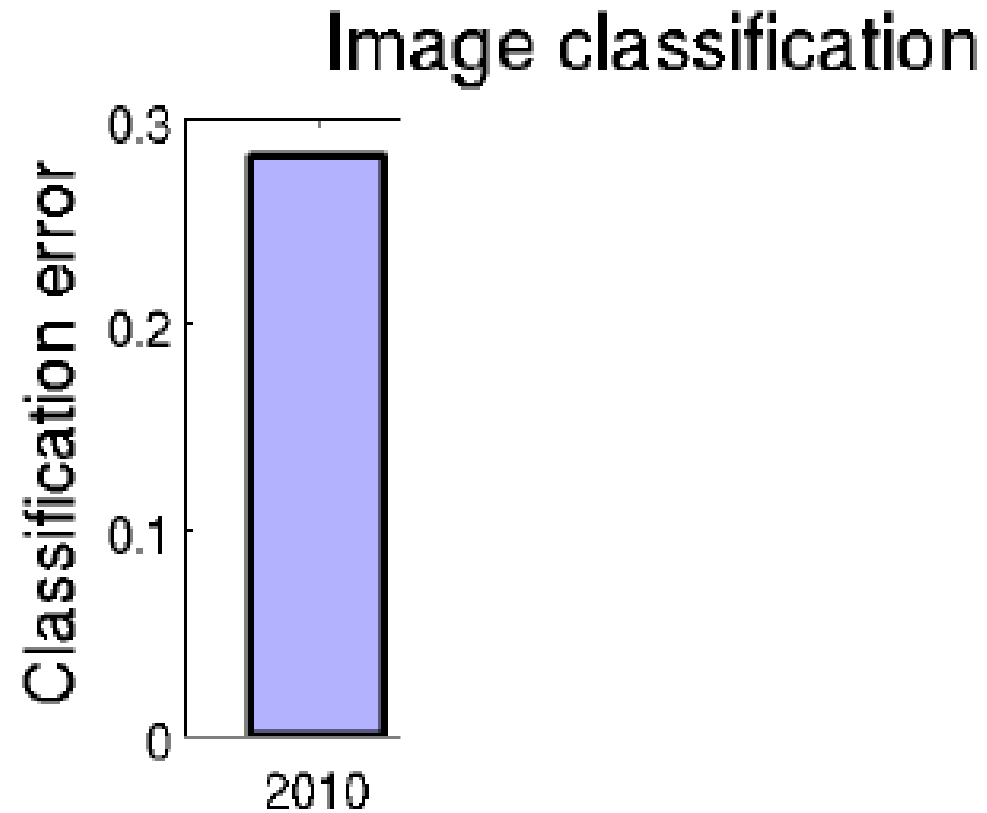
- Millions of labeled images, 1000 object classes.



Easy for humans but
hard for computers.

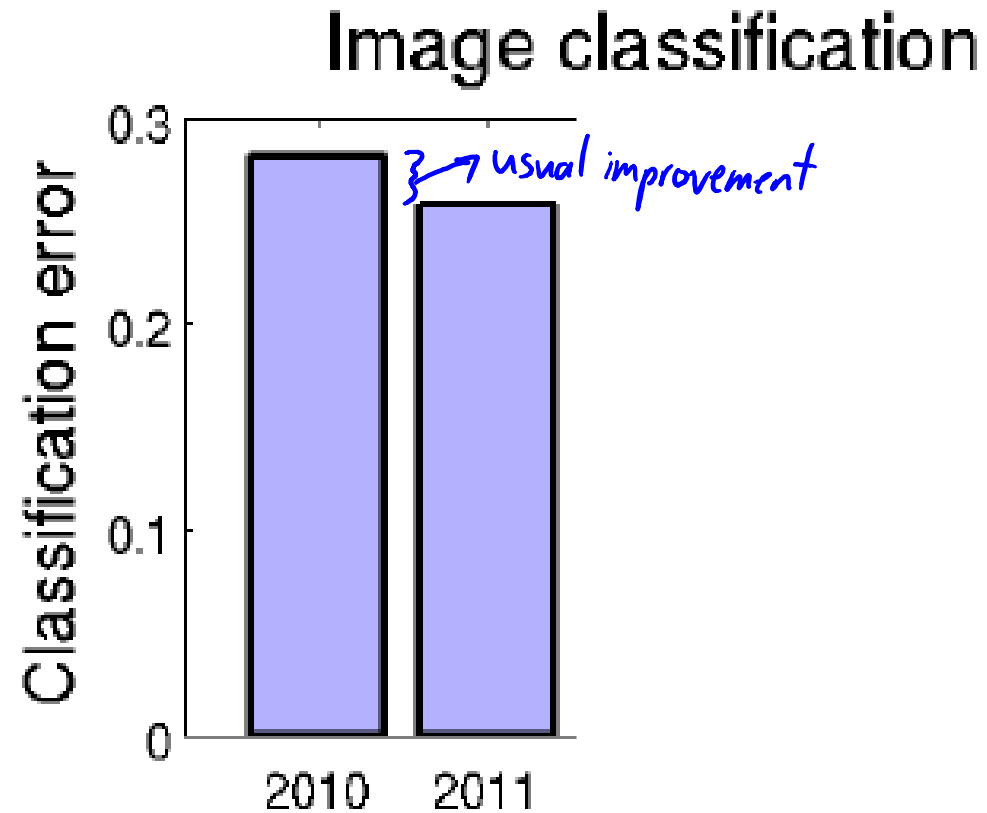
ImageNet Challenge

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



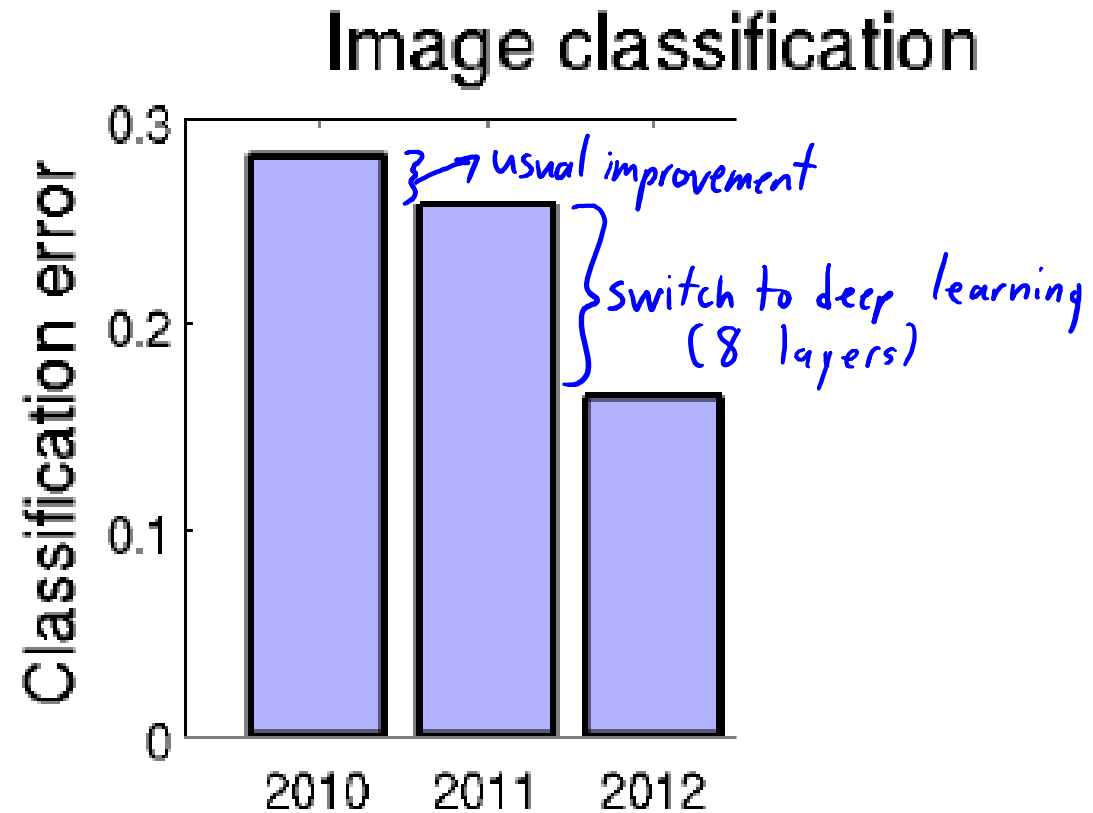
ImageNet Challenge

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



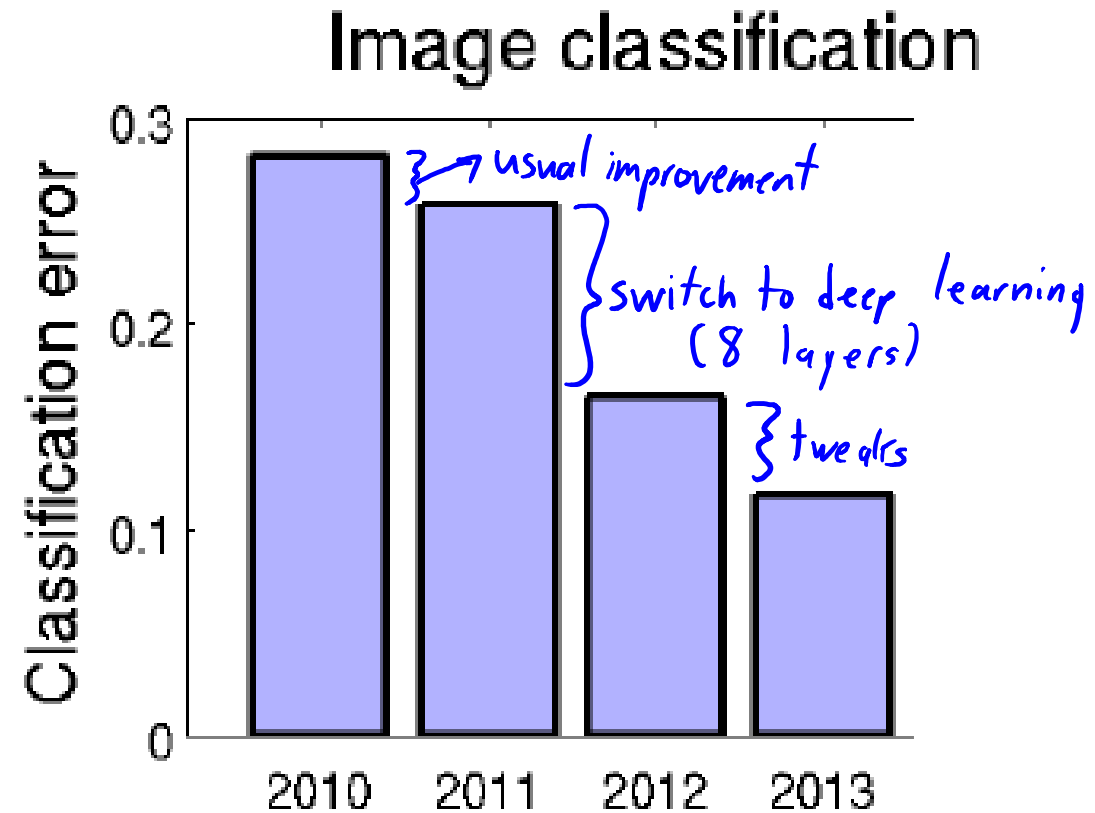
ImageNet Challenge

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



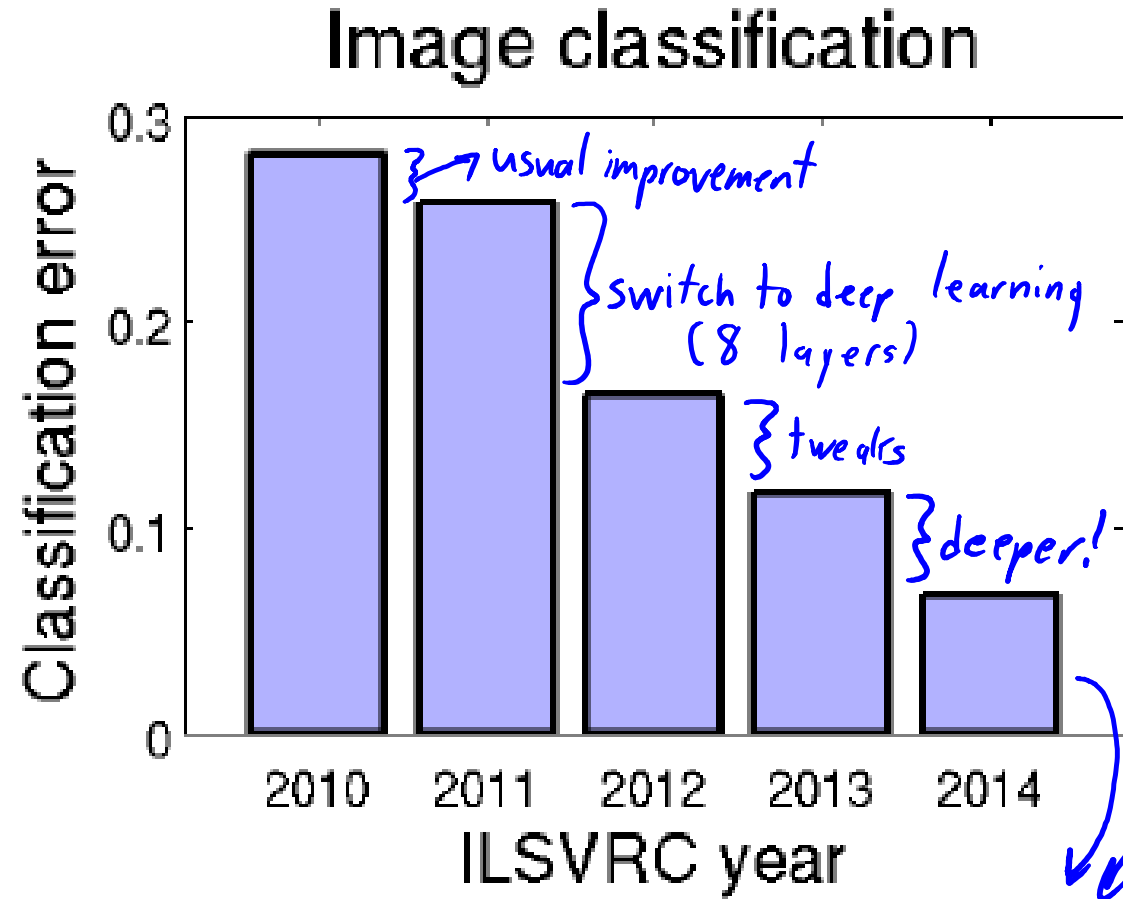
ImageNet Challenge

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



ImageNet Challenge

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



Google Net:
6.7% error
22 layers



ImageNet Challenge

- Object detection task:
 - Single label per image.
 - Humans: ~5% error.
- 2015: Won by Microsoft Research Asia
 - 3.6% error.
 - 152 layers.
- 2016: Chinese University of Hong Kong:
 - Ensembles of existing methods.
- 2017: fewer entries, organizers decided this would be last year.

Summary

- **Neural networks** learn features z_i for supervised learning.
- **Sigmoid function** avoids degeneracy by introducing non-linearity.
- **Biological motivation** for (deep) neural networks.
- **Deep learning** considers neural networks with many hidden layers.
- **Unprecedented performance** on difficult pattern recognition tasks.

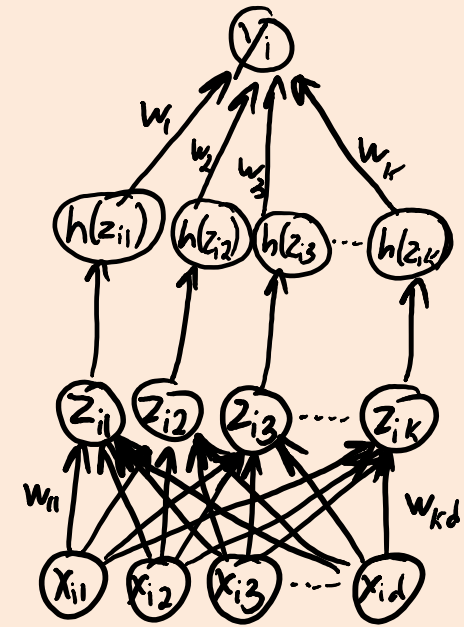
- **Next time:**
 - Training deep networks.

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.
- The value $W^T(WW^T)^{-1}$ is just “some matrix”.
 - You can think of neural networks as just directly learning this matrix.

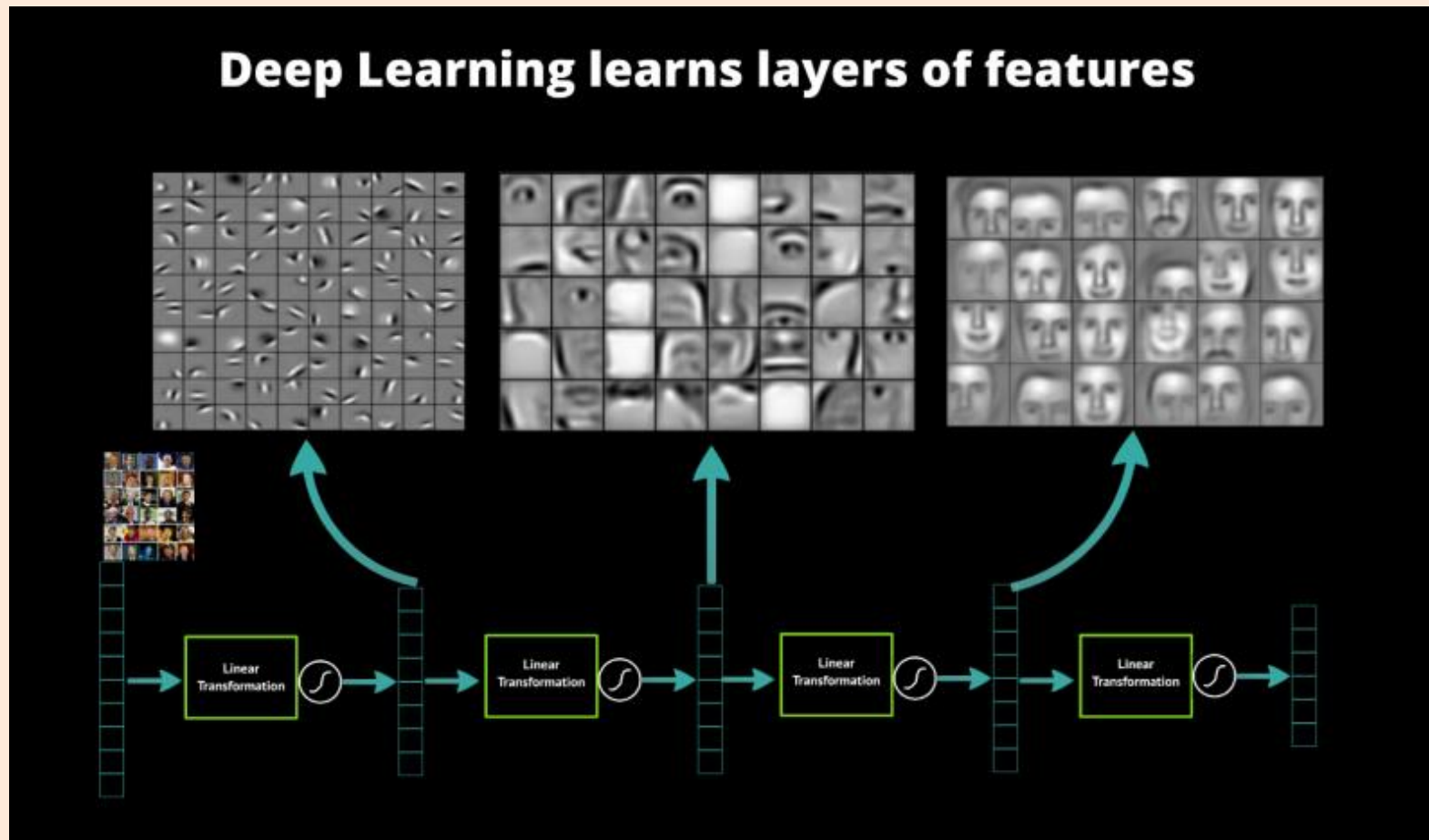
“Artificial” Neural Nets vs. “Real” Networks Nets

- 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.



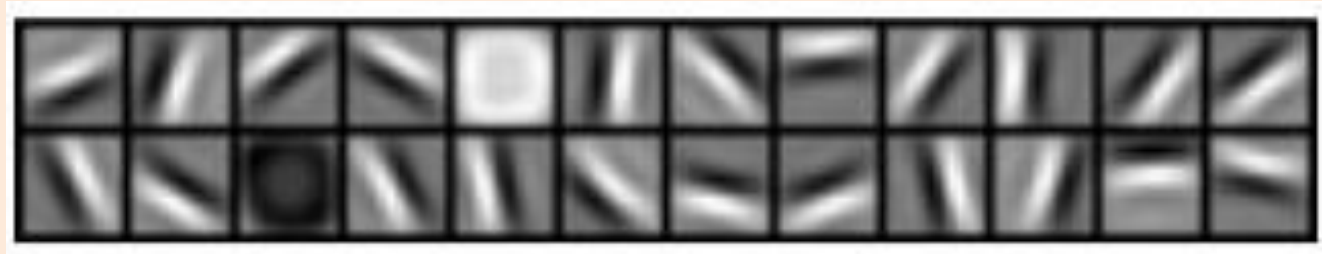
Cool Picture Motivation for Deep Learning

- Faces might be composed of different “parts”:



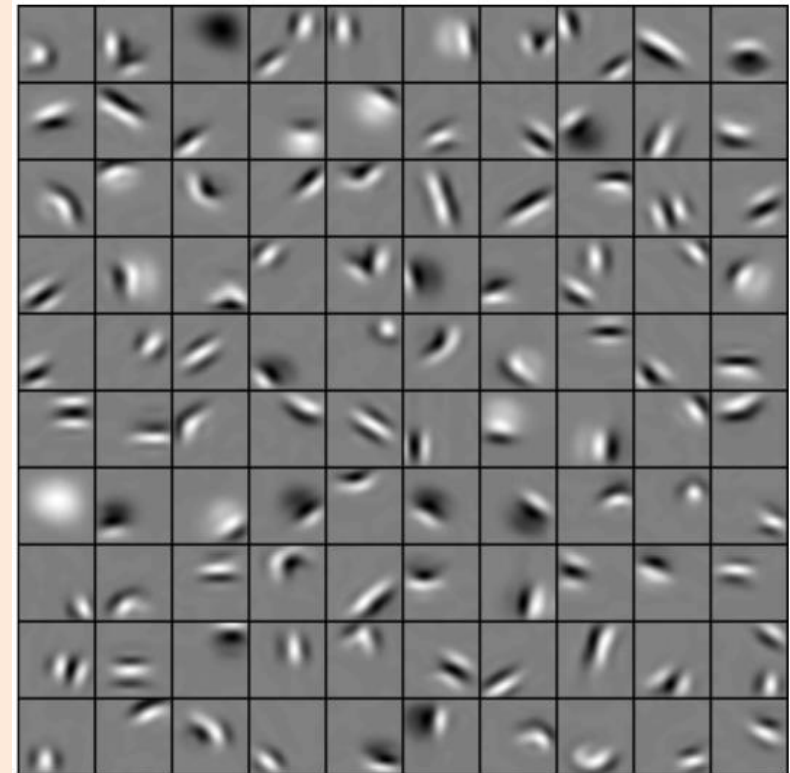
Cool Picture Motivation for Deep Learning

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



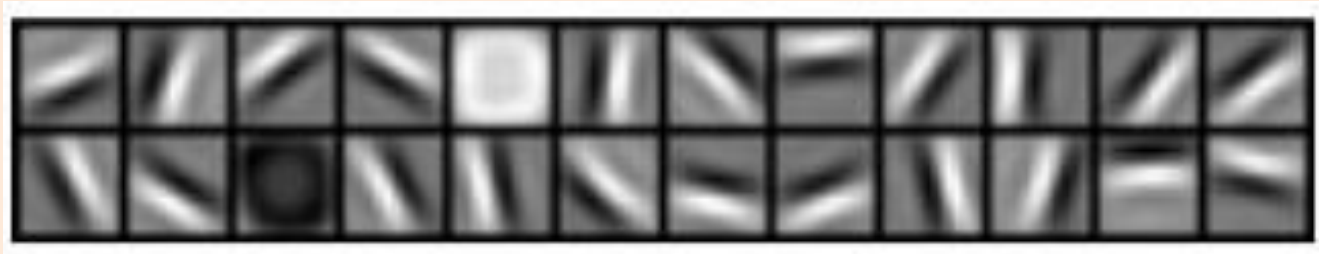
} "Gabor filters"

- 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"

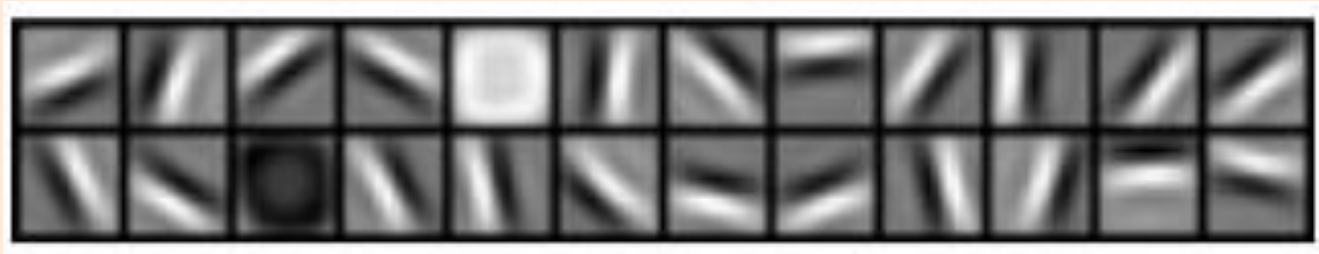
- Visualization of second and third layers trained on specific objects:

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:

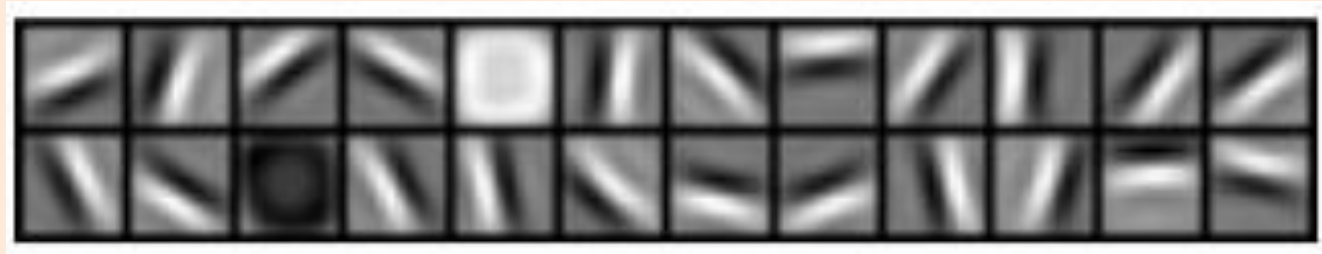
faces

cars



Cool Picture Motivation for Deep Learning

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



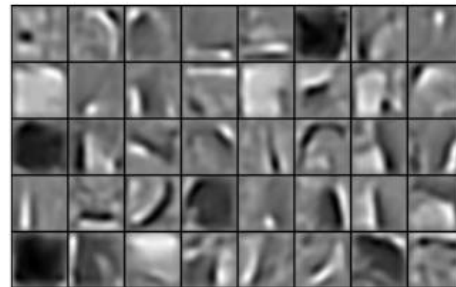
} "Gabor filters"

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

faces

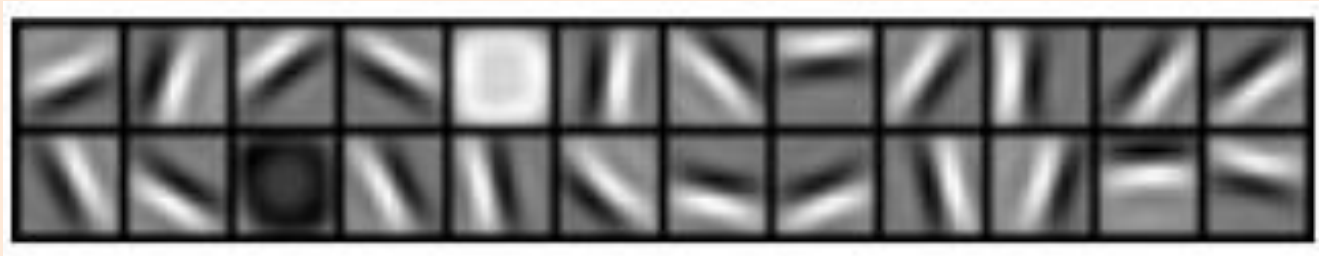
cars

elephants



Cool Picture Motivation for Deep Learning

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



} "Gabor filters"

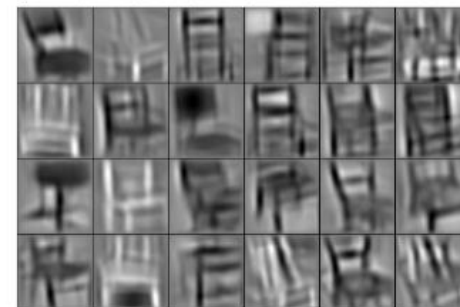
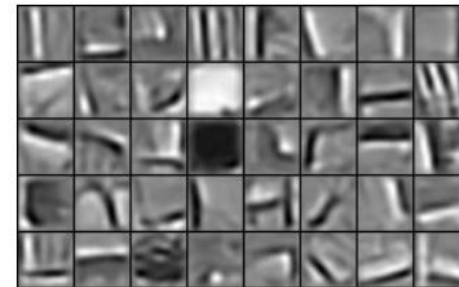
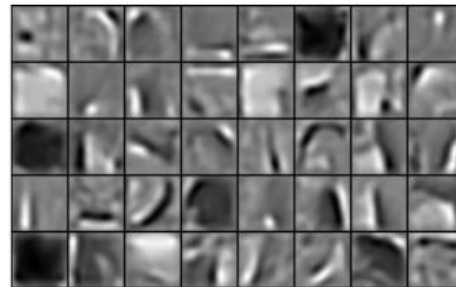
- Visualization of second and third layers trained on specific objects:

faces

cars

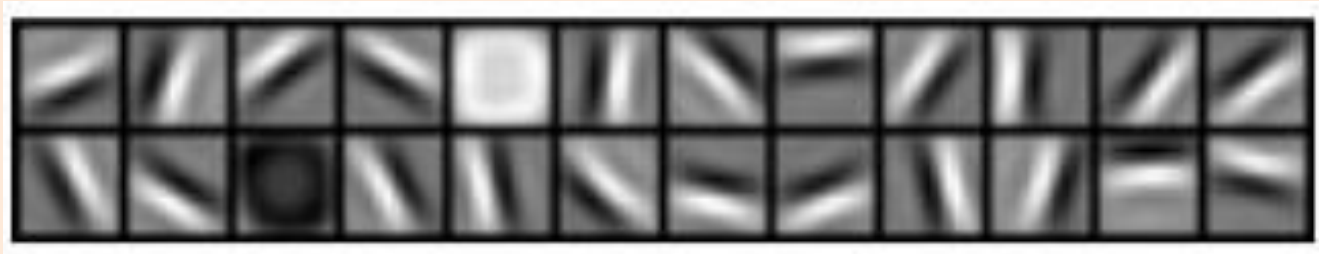
elephants

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

