# CPSC 340:
# Machine Learning and Data Mining

Deep Learning

Spring 2022 (2021W2)

# Admin

- Course surveys
  - Please fill them out
  - We care deeply about your education, so we take them very seriously
  - You will be able to evaluate the class overall, and then Mijung and I separately
  - Please use the text boxes to also let us know about the "lecture specialization experiment" [where we each specialized in half the lectures]
  - As always, please remember we're real people, so both praise and critical feedback are great. Please avoid personal, hurtful, or unconstructive negative comments.

- A5 deadline tonight
- A6 out soon: by Monday at the latest, due April 8 (our last class)

# End of Part 4: Key Concepts

- We discussed linear latent-factor models:

$$f(W, z) = \sum_{i=1}^{n} \sum_{j=1}^{d} (\langle w^j z_i \rangle - x_{ij})^2$$

$$= \sum_{i=1}^{n} \| W^T z_i - x_i \|^2$$

$$= \| ZW - X \|_F^2$$

- Represent 'X' as linear combination of latent factors '$w_c$'.
  - Latent features '$z_i$' give a lower-dimensional version of each '$x_i$'.
  - When k=1, finds direction that minimizes squared orthogonal distance.

- Applications:
  - Outlier detection, dimensionality reduction, data compression, features for linear models, visualization, factor discovery, filling in missing entries.

# End of Part 4: Key Concepts

- We discussed linear latent-factor models:

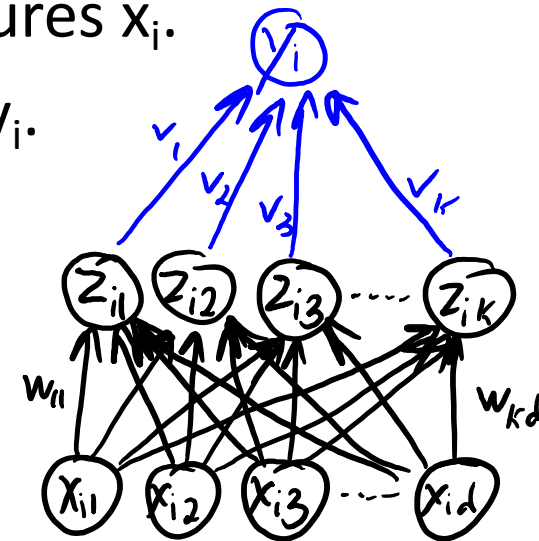$$f(W, z) = \sum_{i=1}^{n} \sum_{j=1}^{d} (\langle w_j^j z_i \rangle - x_{ij})^2$$

- Principal component analysis (PCA):
  - Often uses orthogonal factors and fits them sequentially (via SVD).
- Non-negative matrix factorization:
  - Uses non-negative factors giving sparsity.
  - Can be minimized with projected gradient.
- Many variations are possible:
  - Different regularizers (sparse coding) or loss functions (robust/binary PCA).
  - Missing values (recommender systems) or change of basis (kernel PCA).

# End of Part 4: Key Concepts

*bonus!*

- We didn't really discuss multi-dimensional scaling (MDS):
  - Non-parametric method for high-dimensional data visualization.
  - Tries to match distance/similarity in high-/low-dimensions.
    - "Gradient descent on scatterplot points".
- Main challenge in MDS methods is "crowding" effect:
  - Methods focus on large distances and lose local structure.
- Common solutions:
  - Sammon mapping: use weighted cost function.
  - ISOMAP: approximate geodesic distance using via shortest paths in graph.
  - T-SNE: give up on large distances and focus on neighbour distances.
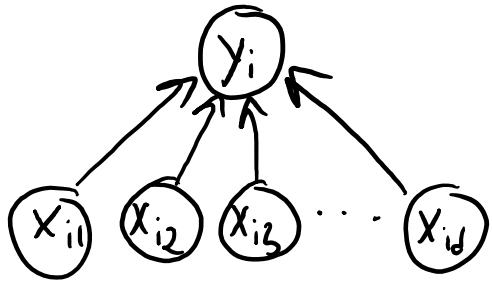- Word2vec is a recent MDS method giving better "word features".

# 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 basis $z_i$ that is good for supervised learning.**

# A Graphical Summary of CPSC 340 Parts 1-5

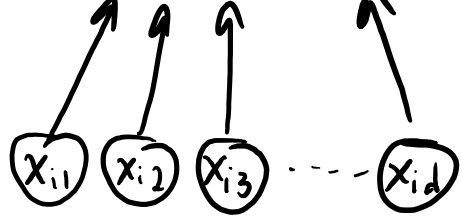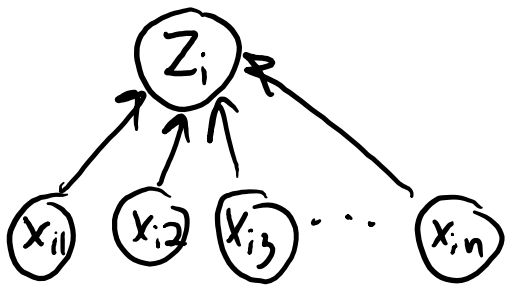# Notation for Neural Networks (MEMORIZE)

We have our usual supervised learning notation:

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

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

We have our latent features:

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

$n \times K$

We have **two** sets of parameters:

$$V = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_K \end{bmatrix} \qquad W = \begin{bmatrix} -\!\!-\; w_1 \;-\!\!- \\ -\!\!-\; w_2 \;-\!\!- \\ \vdots \\ -\!\!-\; w_K \;-\!\!- \end{bmatrix}$$

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

# Linear-Linear Model

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

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:

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

- But this is just a linear model: $\hat{y}_i = v^T z_i = v^T (W x_i) = \underbrace{(v^T W)}_{\substack{1 \times d \\ \text{some vector 'w'}}} x_i = w^T x_i$

# Introducing Non-Linearity

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

$$z_i = W_{x_i} \qquad \hat{y}_i = v^T h(z_i)$$

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

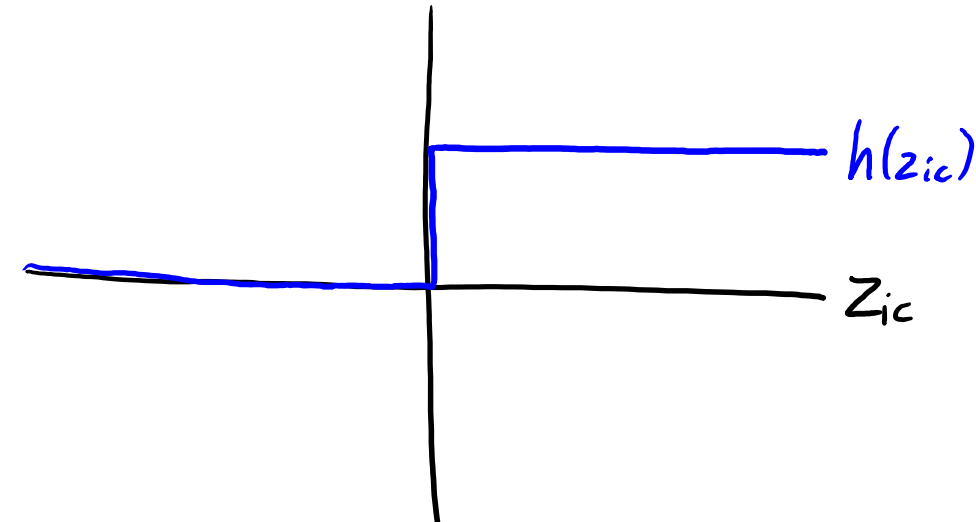- Common choice for 'h': applying sigmoid function element-wise:

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

- 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".
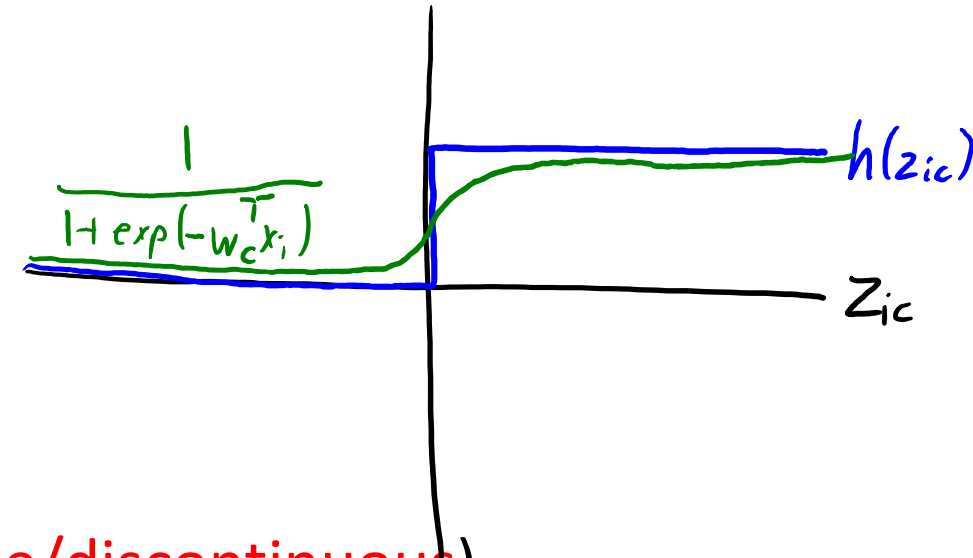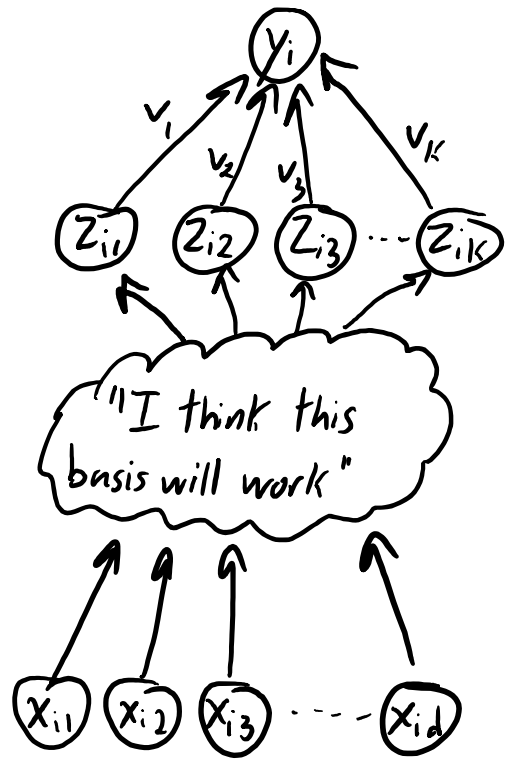


Motivation: Pixels vs. Parts

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

# Why Sigmoid?

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

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



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

  $h(z_{ic})$ , $z_{ic}$

  – Each h(zi) 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:

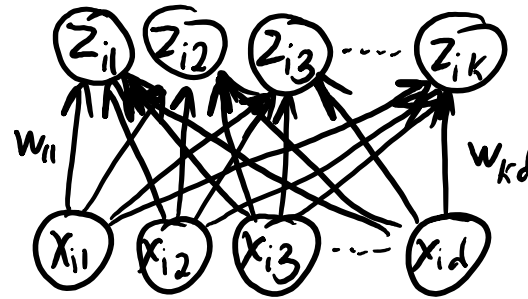    • If 'k' grows appropriately with 'n', can model any continuous function.
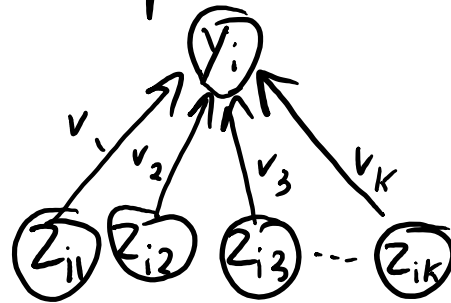
# Supervised Learning Roadmap

# 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):
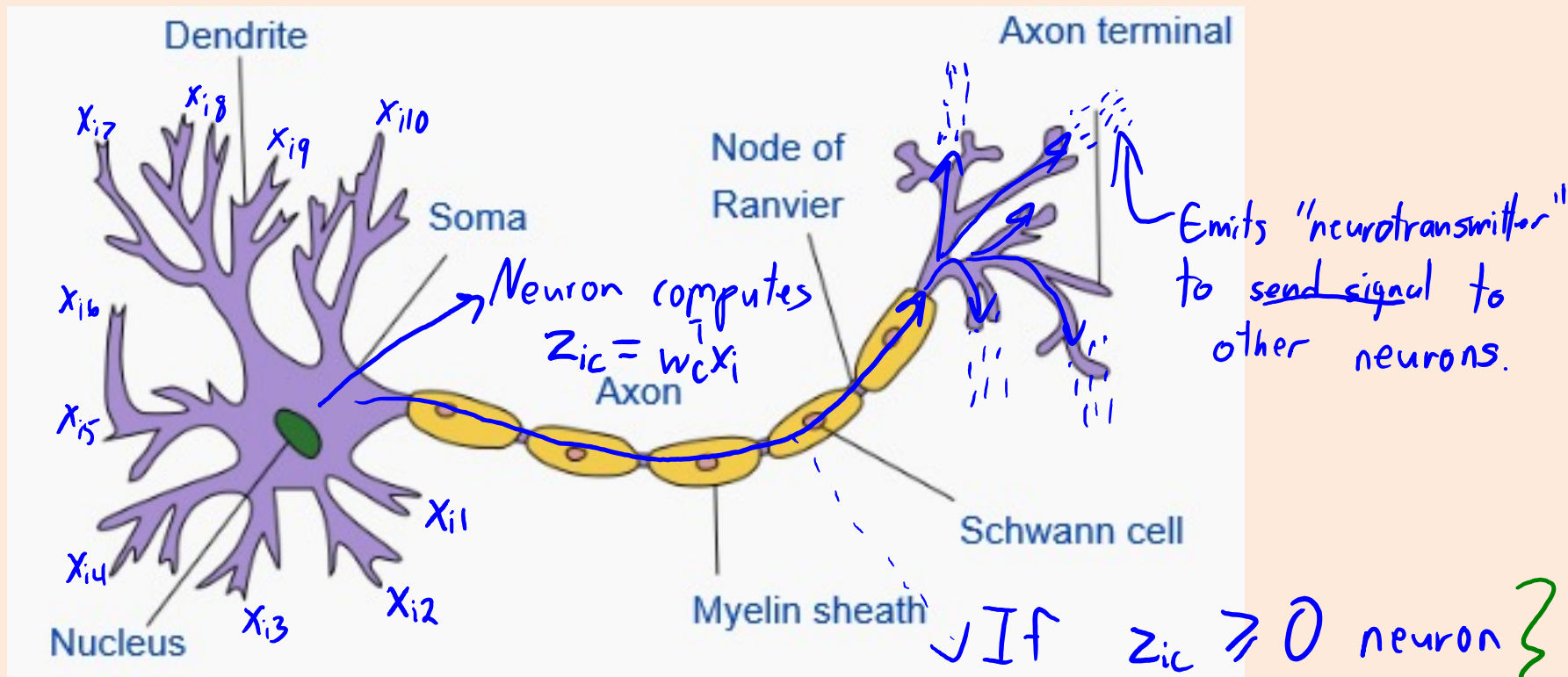
https://en.wikipedia.org/wiki/Neuron
https://en.wikipedia.org/wiki/Action_potential

# Why "Neural Network"?

Dendrite

Axon terminal

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

$x_{i6}$

Soma

Node of Ranvier

Neuron computes $z_{ic} = w_c x_i$

Emits "neurotransmitter" to send signal to other neurons.

Axon

$x_{i5}$

$x_{i1}$

$x_{i4}$

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

Nucleus

Schwann cell

Myelin sheath

If $z_{ic} \geq 0$ neuron sends signal along axon.

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

https://en.wikipedia.org/wiki/Neuron

# Why "Neural Network"?



bonus!

$h(z_{i2})$  $h(z_{i1})$

$h(z_{i3})$

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

$h(z_{i6})$

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

Axon terminal

$h(z_{i4})$

$h(z_{i4})$  $h(z_{i5})$

Emits "neurotransmitter"
to send signal to
other neurons.

Dendrite

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

Soma

$x_{i6}$

Node of Ranvier

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

Axon

$x_{i5}$

$x_{i1}$

Schwann cell

$x_{i4}$

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

Nucleus

Myelin sheath

$\sqrt{}$ If $z_{ic} \geqslant 0$ neuron
sends signal along axon.

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

# 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 recieving $x_{ij}$ values
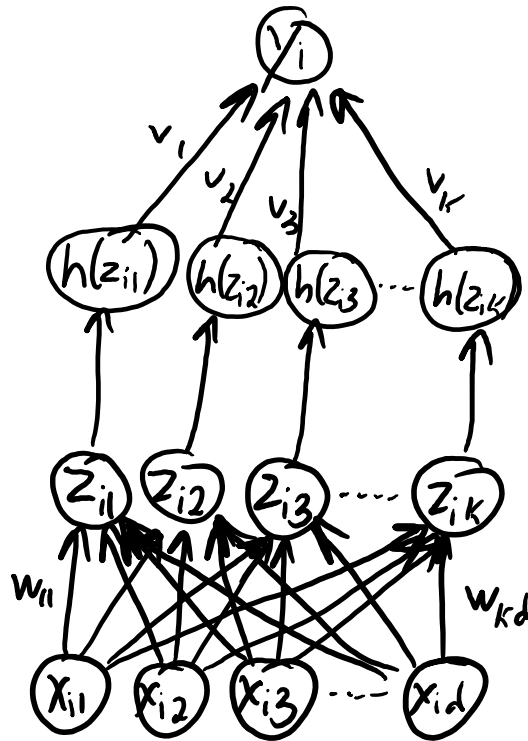
# "Artificial" Neural Nets vs. "Real" Networks Nets ~bonus!~

- 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.
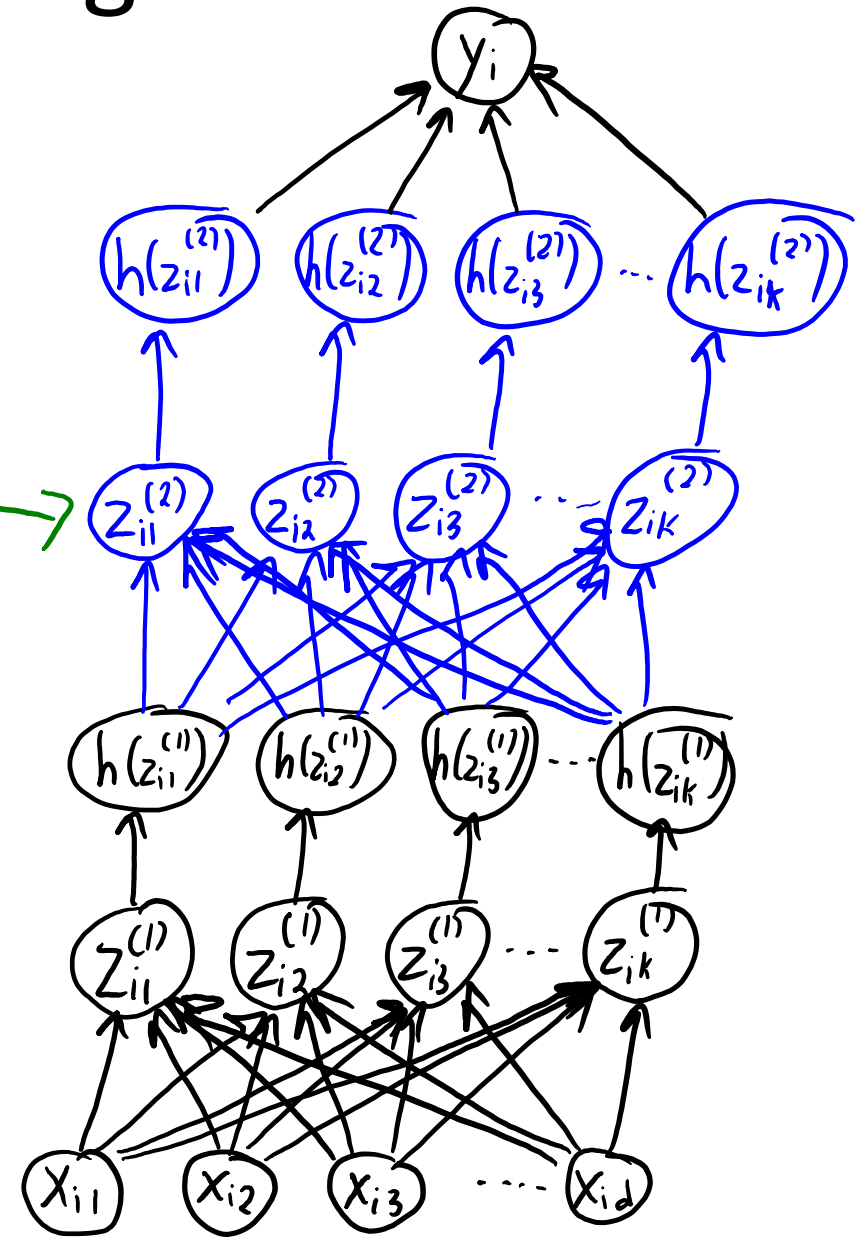  - Different types of neurotransmitters.

# Deep Learning



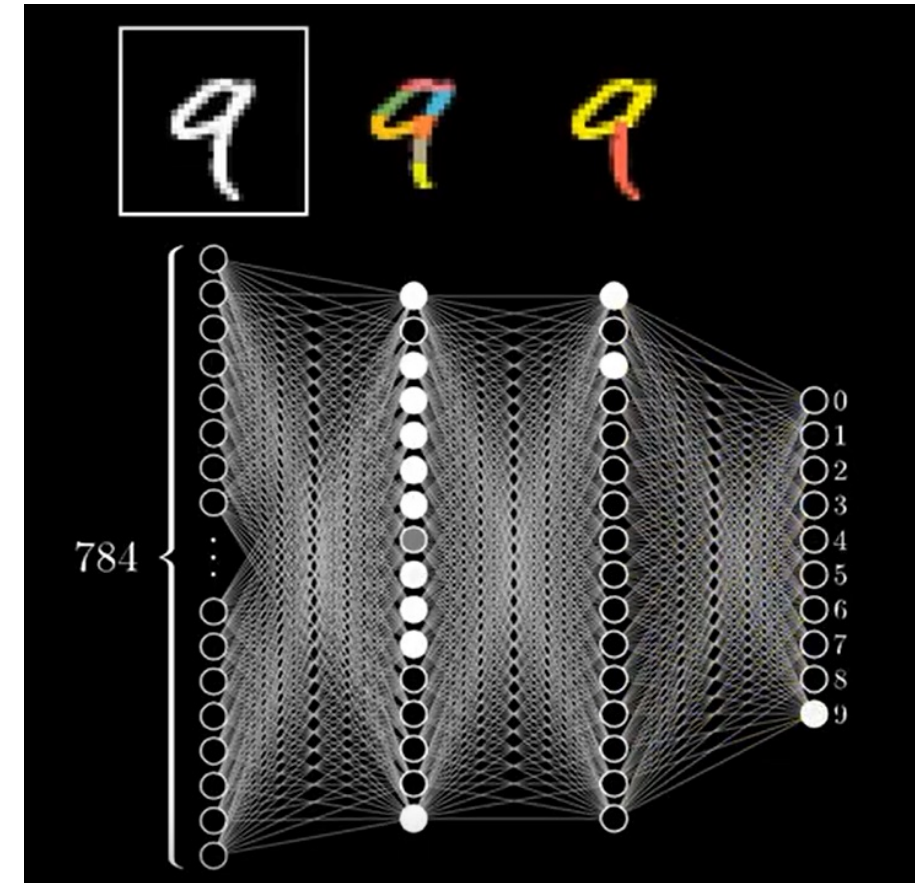Neural network:

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

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

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

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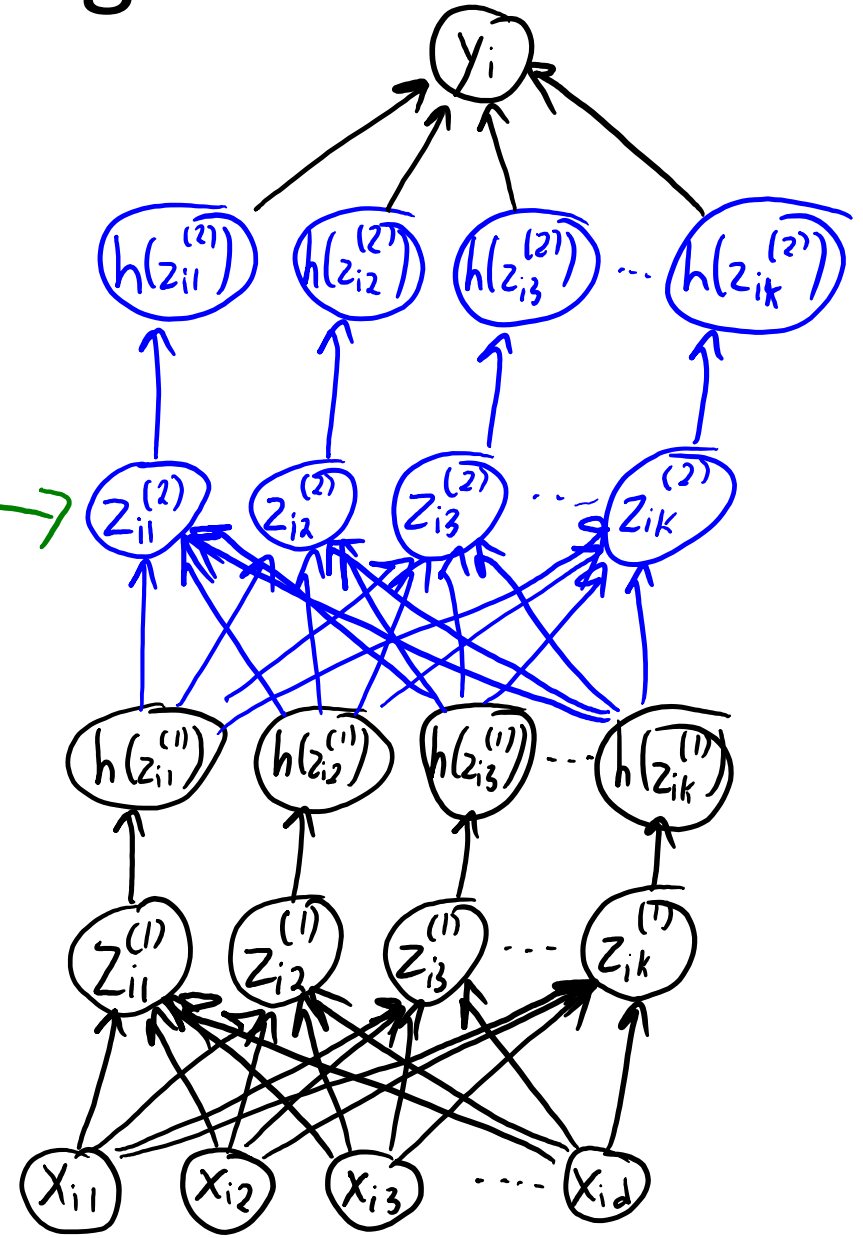
Neural network with 3 hidden layers
$$\hat{y}_i = v^T h(\underbrace{W^{(3)} h(\underbrace{W^{(2)} h(\underbrace{W^{(1)} x_i}_{z_i^{(1)}})}_{z_i^{(2)}})}_{z_i^{(3)}})$$

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:

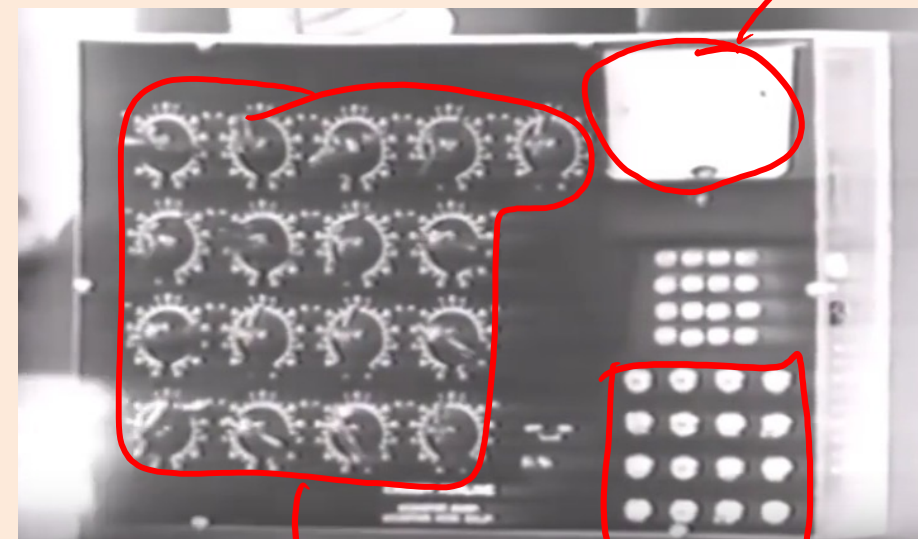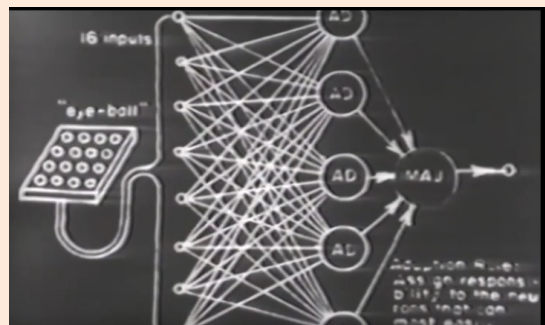$$\hat{y}_i = v^\top h\left(W^{(4)} h\left(W^{(3)} h\left(W^{(2)} h\left(W^{(1)} x_i\right)\right)\right)\right)$$

- For 'm' layers, we usually just say:

$$\hat{y}_i = v^\top h\left(W^{(m)} h\left(\cdots h\left(W^{(1)} x_i\right)\right)\right)$$
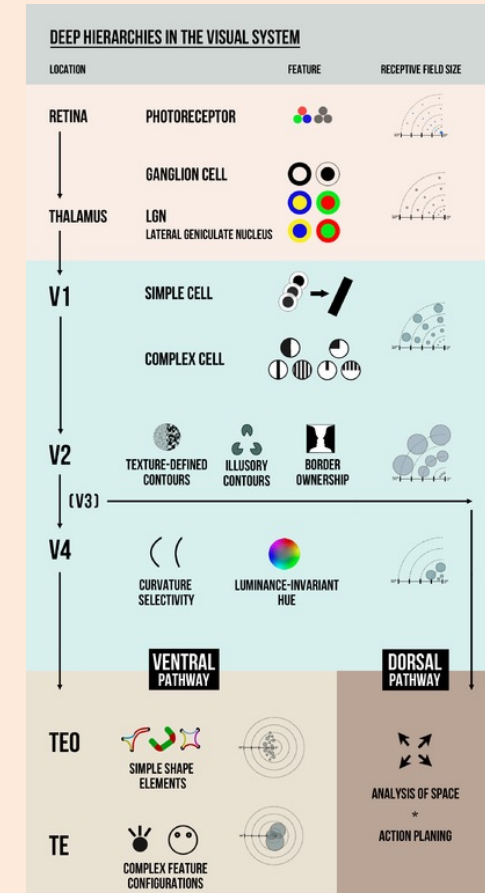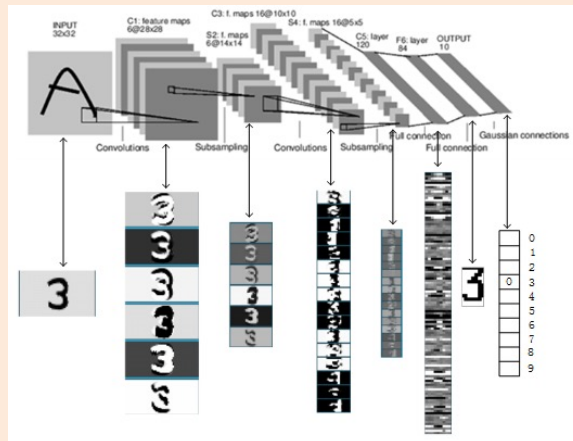
# ML and Deep Learning History

- 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

  – Object recognition assigned to students as a summer project

- Then drop in popularity:

  – Quickly realized limitations of linear models.

$w^T x_i$

$w$

$x_i$

# ML and Deep Learning History

*bonus!*

- 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.
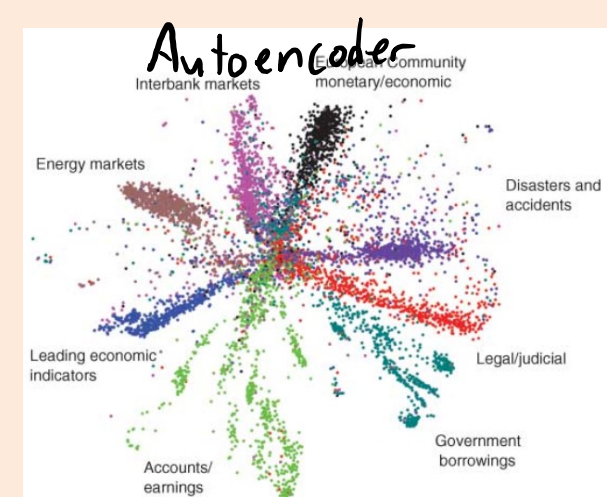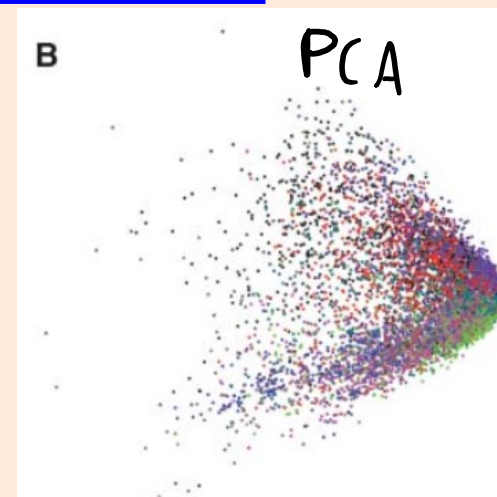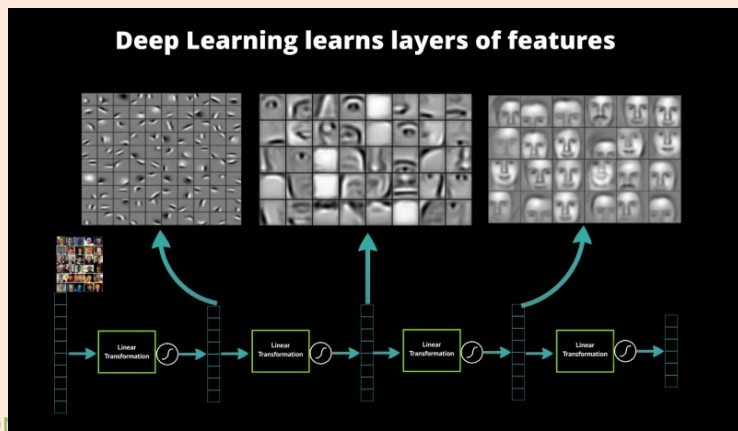
# ML and Deep Learning History

- 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.
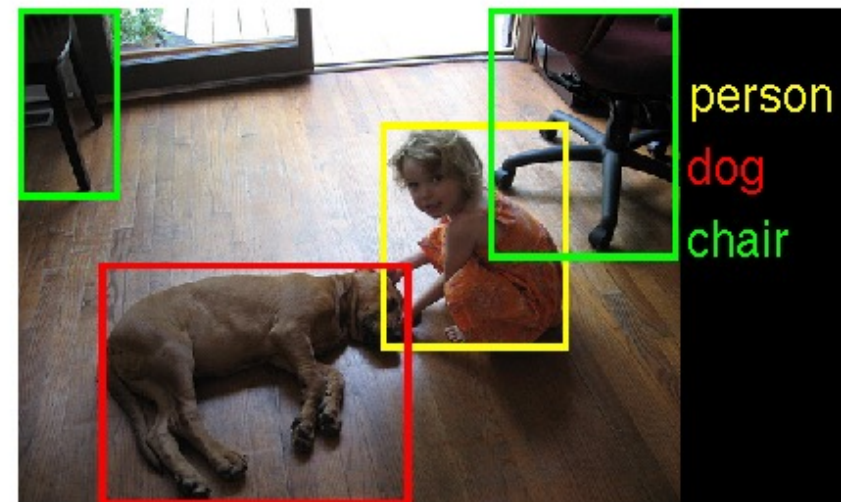
# ML and Deep Learning History

- 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
    - Unsupervised successes: "deep belief networks" and "autoencoders".
    - Could be used to initialize deep neural networks.
    - https://www.youtube.com/watch?v=KuPai0ogiHk



Deep Learning learns layers of features



PCA



Autoencoder

https://www.cs.tor

# 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.



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

# 2010s: DEEP LEARNING!!!

*bonus!*

- 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.

# 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.

# Multiple Word Prototypes

bonus!

- What about homonyms and polysemy?
  - The word vectors would need to account for all meanings.

- More recent approaches:
  - Try to cluster the different contexts where words appear.
  - Use different vectors for different contexts.

$$X_{jaguar} = \begin{bmatrix} \text{-- -- -- -- -- -- --} \\ \text{-- -- -- -- -- --} \\ \text{  } \end{bmatrix} \begin{matrix} z_{j1} \\ z_{j2} \\ z_{j3} \end{matrix}$$

# Multiple Word Prototypes

*bonus!*

# 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 *bonus!*

- Faces might be composed of different "parts":

# Cool Picture Motivation for Deep Learning *bonus!*

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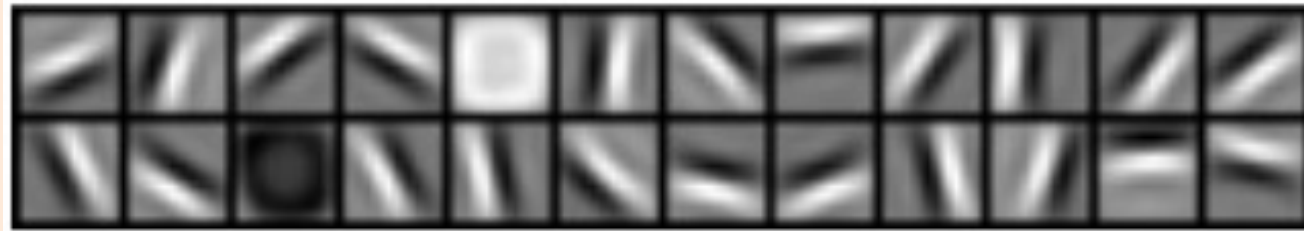


"Gabor filters"

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

- Models require many tricks to work.
  - We'll discuss these next time.
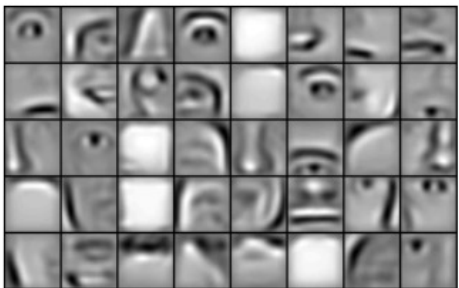
# Cool Picture Motivation for Deep Learning *bonus!*

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



{ "Gabor filters"

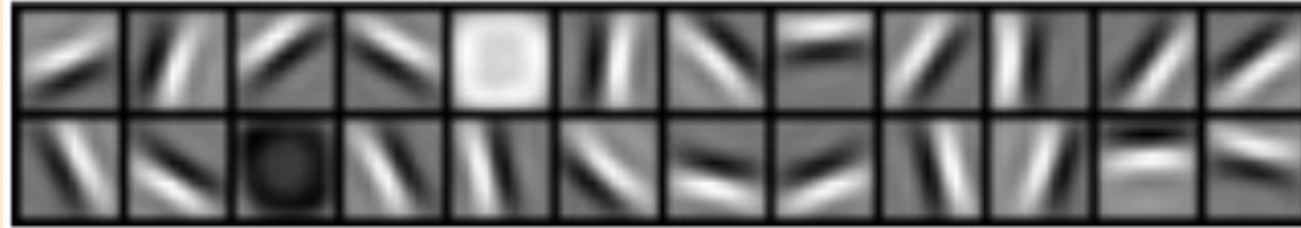- Visualization of second and third layers trained on specific objects:



faces

# Cool Picture Motivation for Deep Learning *bonus!*

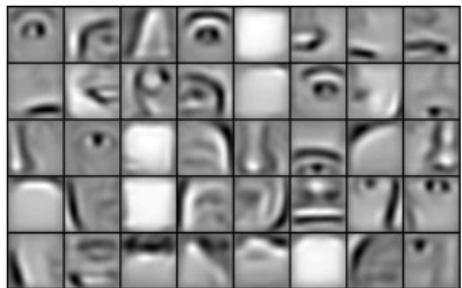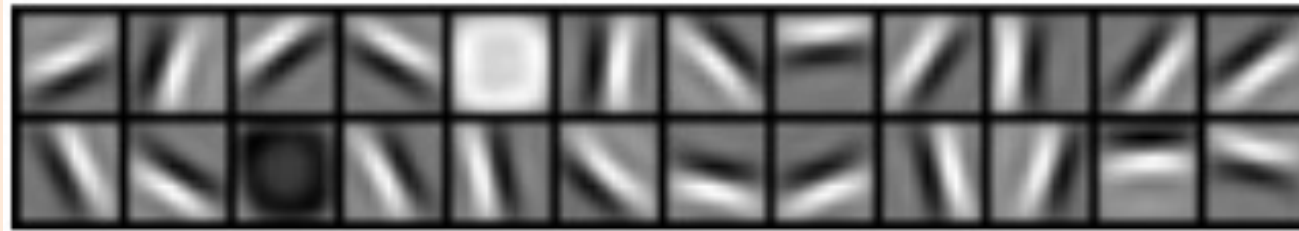- First layer of $z_i$ trained on 10 by 10 image patches:



"Gabor filters"

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

# Cool Picture Motivation for Deep Learning

_bonus!_

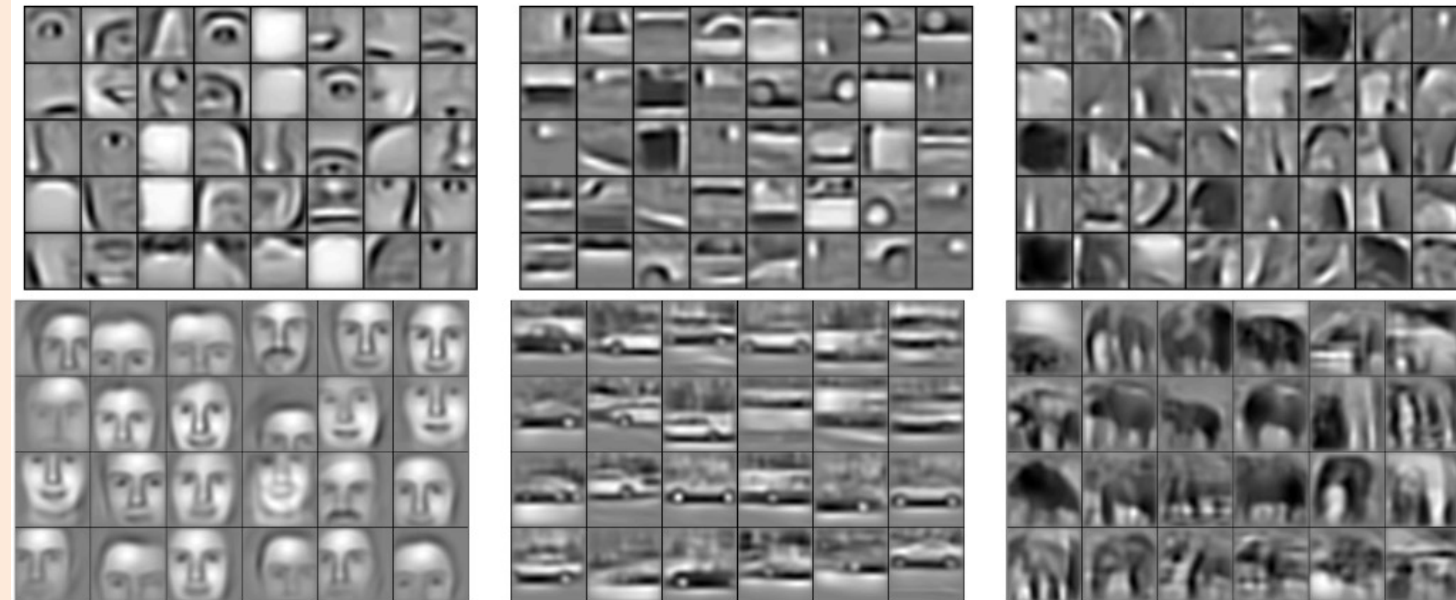- First layer of $z_i$ trained on 10 by 10 image patches:



"Gabor filters"

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

# Cool Picture Motivation for Deep Learning

*bonus!*

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



} "Gabor filters"

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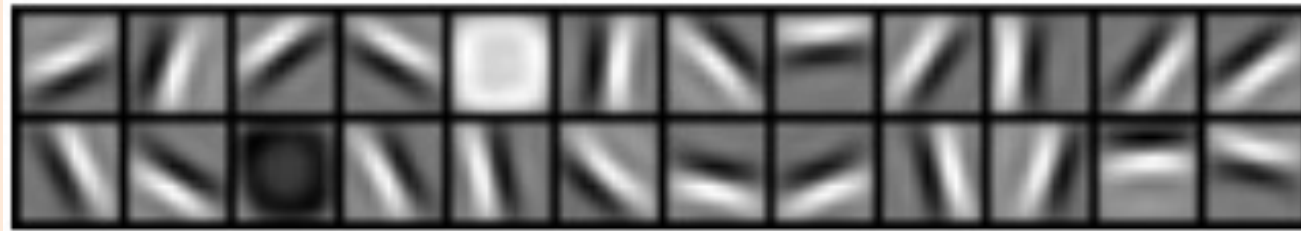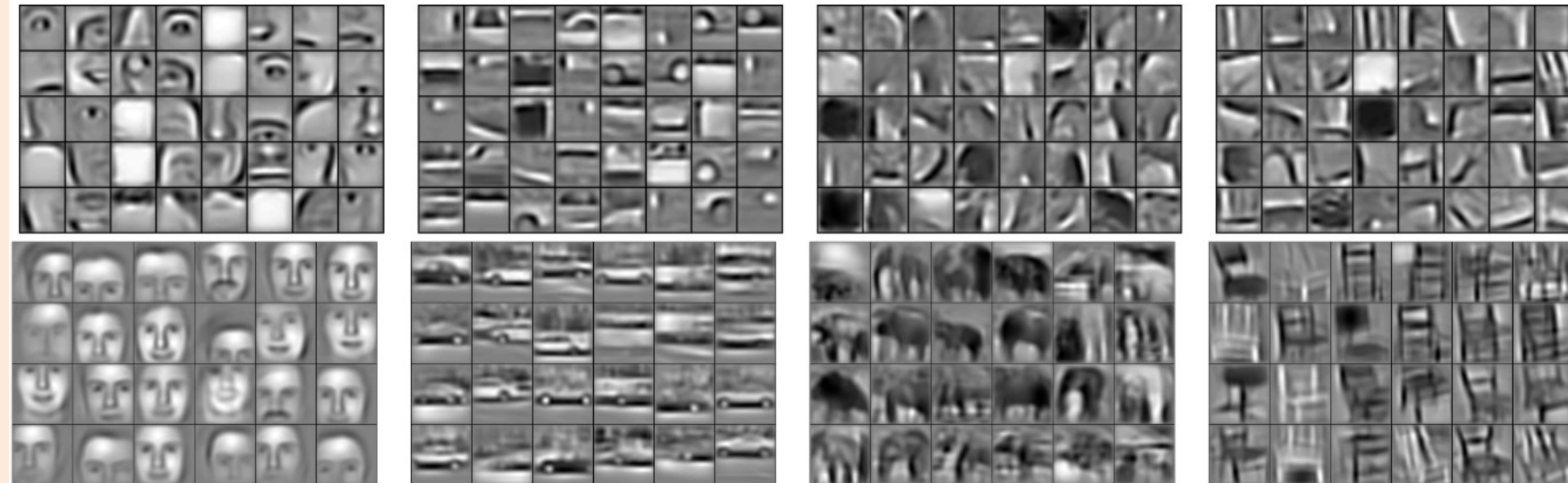


faces    cars    elephants    chairs

# Cool Picture Motivation for Deep Learning

*bonus!*

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



"Gabor filters"

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



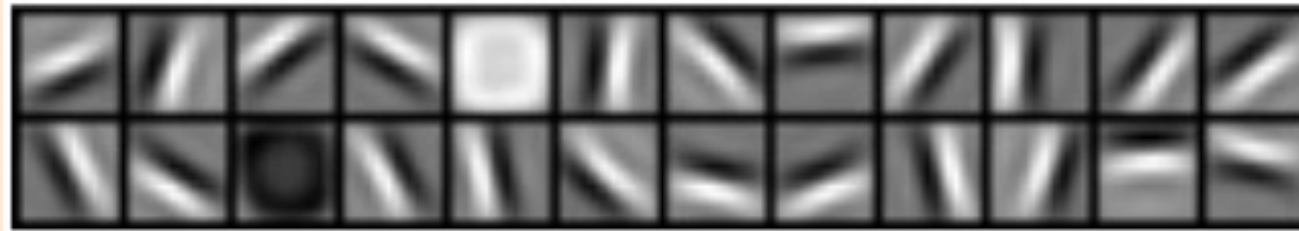faces    cars    elephants    chairs    faces, cars, airplanes, motorbikes