

Logistic Regression

Intro to ML 466/566

Fall 2033

Administrivia

- Virtual office hours
 - **(Oct. 3, 17, Nov. 14, Dec. 5 only):**
 - <https://meet.google.com/ptm-znnz-kbb>
 - Meet link is also on eclass

Homework typos 1D

- c) **(3 marks)** The perceptron algorithm will stop when it arrives at any weight vector that properly separates training data. Are all possible weight vector solutions (that separate the training data) equally likely to generalize to a larger dataset? Please explain your answer.
- d) **(7 marks)** Now consider the dataset = $\{([-1, -1], 1), ([1, 1], 1), ([1, -1], -1), ([-1, 1], -1)\}$, and the weight is initialized as $w_0 = [0, 1]^T$.
- As above, plot out the update direction of the weight using point $([-1, -1], 1)$ on the right of Figure 1. [Plot the update direction starting from the current weight point.]
 - Then, plot out the new weight vector and the separation plane in a new figure. Again, plot out the update direction of the weight using point $([1, 1], 1)$.
 - After these updates, are all training data points correctly predicted? If the algorithm uses other data points to update, will it yield a separation line that correctly predicts for all points? Why or why not?

Homework clarification 2 b

- “Use Cauchy Schwartz inequality on the result of 1 to find a lower bound of w^{k+1} .”
- “Use Cauchy Schwartz inequality on the result of 1 to find a lower bound of $\|w^{k+1}\|$.”

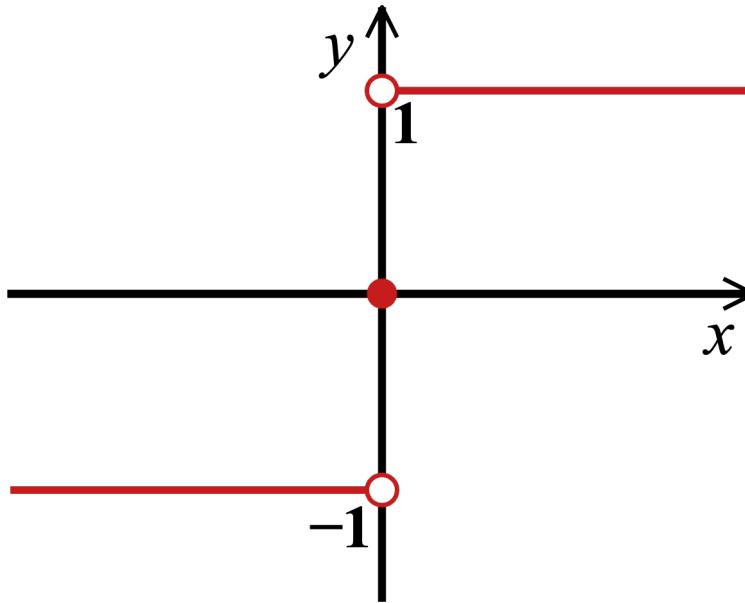
Homework clarification 2 c

- “Find an upper bound for the **inner product of the weights at $(k + 1)$ -th iteration** using the perceptron algorithm update.”
- **inner product of the weight vector at the $(k + 1)$ -th iteration with itself**

$$\| \mathbf{w}^{k+1} \|_2^2 \text{ (squared L2 norm)}$$

Perceptron

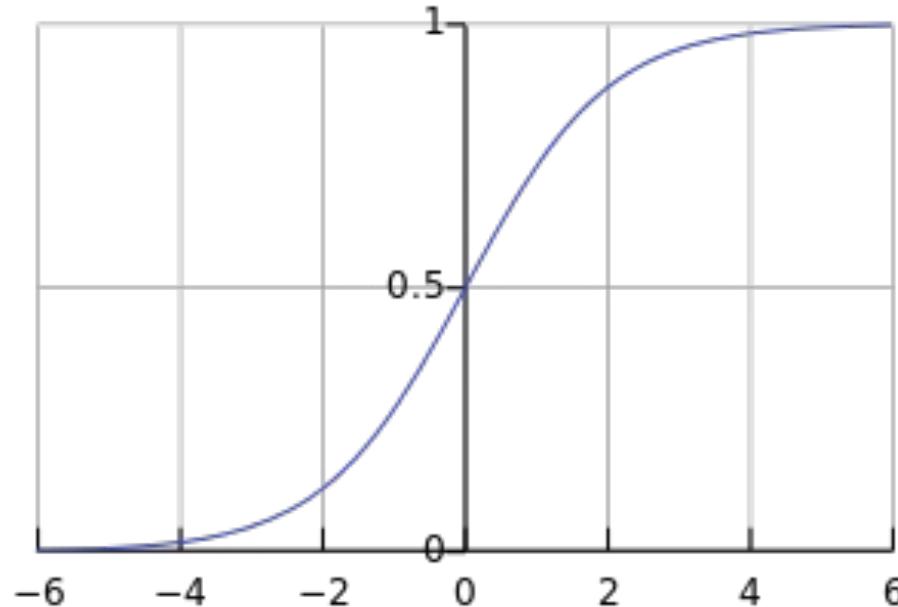
- Recall: perceptron uses the *sign* function



- Outputs either -1 or +1

Logistic Regression

- Similar to perceptron, but **sigmoid** function applied to linearity.



$$S(z) = \frac{1}{1 + e^{-z}}$$

- Error function is based on **Max Likelihood** →
as a result we get **genuine probabilities** as output of prediction (not just 1,0 discrete values).

Models so far, and the new one

$$z = \sum_{i=0}^m w_i x_i$$

Perceptron

$$h(\mathbf{x}, \mathbf{w}) = \text{sign}(z)$$

Linear regression

$$h(\mathbf{x}, \mathbf{w}) = z$$

Logistic regression

$$h(\mathbf{x}, \mathbf{w}) = S(z)$$

Output $h(\mathbf{x}, \mathbf{w})$ will be interpreted as probability.

Why? Because of the range of S , and particular cost function we'll optimize.

Probability Interpretation

- Assume tuples (\mathbf{x}, y) , where y is binary, are generated from some noisy data source according to some distribution f .

$$p(y | \mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = 0 \\ 1 - f(\mathbf{x}) & \text{for } y = 1 \end{cases}$$

For mathematical convenience we will use $[0,1]$ as our label set instead of $[-1,1]$

- We will learn $f(\mathbf{x})$ by approximating it with the sigmoid $S(z) = \frac{1}{1 + e^{-z}}$
i.e. $S(\mathbf{w} \cdot \mathbf{x})$

$$\mathbf{w} = [w_0 = b, w_1, \dots, w_m] \quad \mathbf{x} = [x_0 = 1, x_1, \dots, x_m]$$

- Makes sense as $S(z)$ is a function from 0 to 1.

Approximation

- What is the probability of $y=0$?
(according to our approximation)

$$p(y = 0 | \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}'\mathbf{x}}}$$

- What is the probability of $y=1$?
(according to our approximation)

$$\begin{aligned} p(y = 1 | \mathbf{x}) &= 1 - \frac{1}{1 + e^{\mathbf{w}'\mathbf{x}}} \\ &= \frac{1 + e^{\mathbf{w}'\mathbf{x}} - 1}{1 + e^{\mathbf{w}'\mathbf{x}}} \\ &= \frac{e^{\mathbf{w}'\mathbf{x}}}{1 + e^{\mathbf{w}'\mathbf{x}}} \end{aligned}$$

Training Logistic Regression: MCLE

- Choose parameters $W = \langle w_0, \dots, w_n \rangle$ to maximize conditional likelihood of training data

$$\text{where } p(y=0 | \mathbf{x}) = \frac{1}{1+e^{w'\mathbf{x}}}$$

$$p(y=1 | \mathbf{x}) = \frac{e^{w'\mathbf{x}}}{1+e^{w'\mathbf{x}}}$$

- Training data $D = \{ \langle X^1, Y^1 \rangle, \dots, \langle X^N, Y^N \rangle \}$

- Data conditional likelihood = $\prod_{i=1}^N P(Y^i | X^i, W)$

$$W_{MCLE} = \operatorname{argmax}_w \prod_{i=1}^N P(Y^i | X^i, W)$$

Training Logistic Regression: MCLE

- We would like to find the w that maximizes the probability of our data

$$W_{MCL\!E} = \operatorname{argmax}_w \prod_{i=1}^N P(Y^i | X^i, W)$$

Expressing Conditional Log Likelihood

$$\mathcal{L}(W) \equiv \ln \left(\prod_i P(Y^i | X^i, W) \right)$$

$$= \sum_i \ln \underbrace{P(Y^i | X^i, W)}_{\text{←}}$$

$$p(y = 0 | \mathbf{x}) = \frac{1}{1 + e^{w' \mathbf{x}}}$$

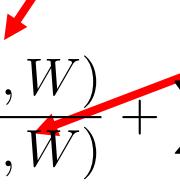
$$p(y = 1 | \mathbf{x}) = \frac{e^{w' \mathbf{x}}}{1 + e^{w' \mathbf{x}}}$$

$$\mathcal{L}(W) = \sum_i Y^i \ln \underbrace{P(Y^i = 1 | X^i, W)}_{\text{For the samples with } Y^i = 1} + \sum_i (1 - Y^i) \ln \underbrace{P(Y^i = 0 | X^i, W)}_{\text{For the samples with } Y^i = 0}$$

For the samples
with $Y^i = 1$

For the samples
with $Y^i = 0$

Expressing Conditional Log Likelihood

$$\begin{aligned}\mathcal{L}(W) &= \sum_i Y^i \ln P(Y^i = 1|X^i, W) + \sum_i (1 - Y^i) \ln P(Y^i = 0|X^i, W) \\ &= \sum_i Y^i \ln \frac{P(Y^i = 1|X^i, W)}{P(Y^i = 0|X^i, W)} + \sum_i \ln P(Y^i = 0|X^i, W)\end{aligned}$$


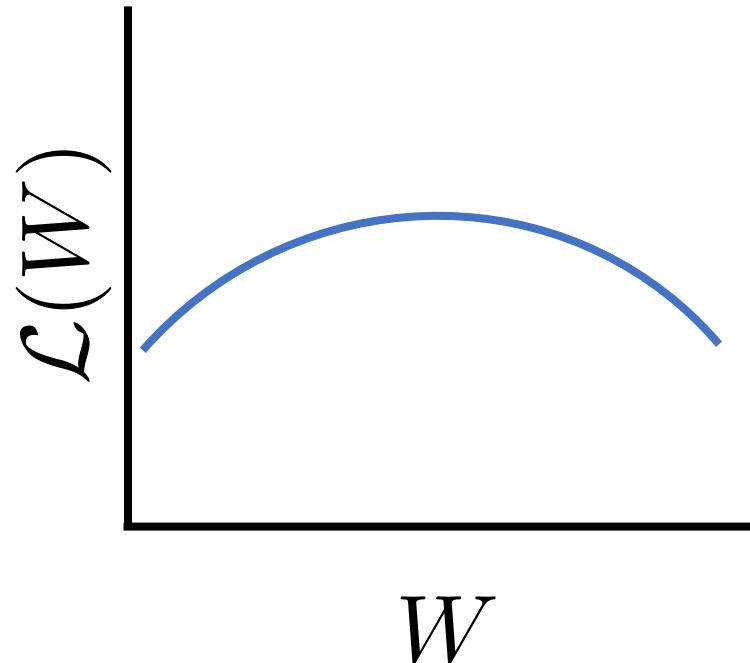
Expressing Conditional Log Likelihood

$$\begin{aligned}\mathcal{L}(W) &= \sum_i Y^i \ln P(Y^i = 1|X^i, W) + \sum_i (1 - Y^i) \ln P(Y^i = 0|X^i, W) \\ &= \sum_i Y^i \ln \frac{P(Y^i = 1|X^i, W)}{P(Y^i = 0|X^i, W)} + \sum_i \ln P(Y^i = 0|X^i, W) \\ &= \sum_i Y^i (w' x^i) - \sum_i \ln(1 + e^{w' x^i})\end{aligned}$$

$$\begin{aligned}p(y = 0 | \mathbf{x}) &= \frac{1}{1 + e^{w' \mathbf{x}}} \\ p(y = 1 | \mathbf{x}) &= \frac{e^{w' \mathbf{x}}}{1 + e^{w' \mathbf{x}}}\end{aligned}$$

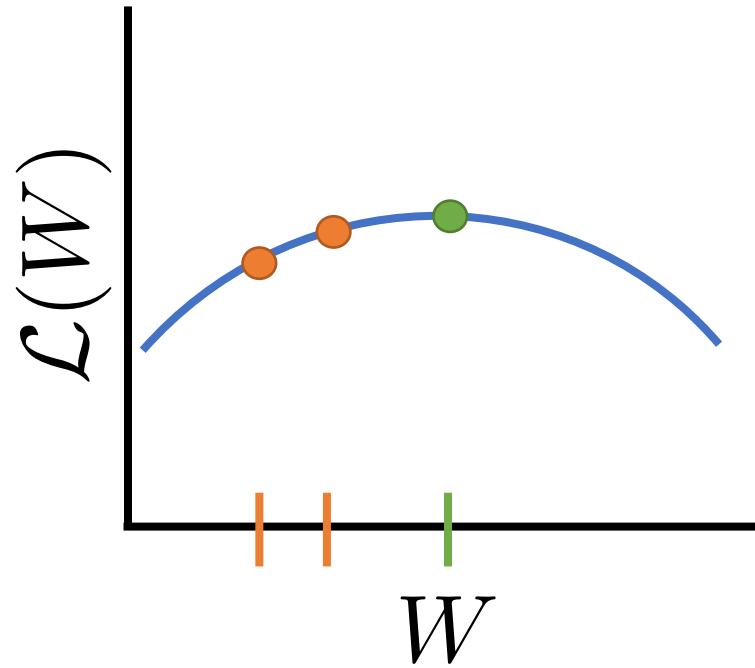
Maximizing Conditional Log Likelihood

$$\mathcal{L}(W) = \sum_i Y^i (w' x^i) - \sum_i \ln(1 + e^{w' x^i})$$



Maximizing Conditional Log Likelihood

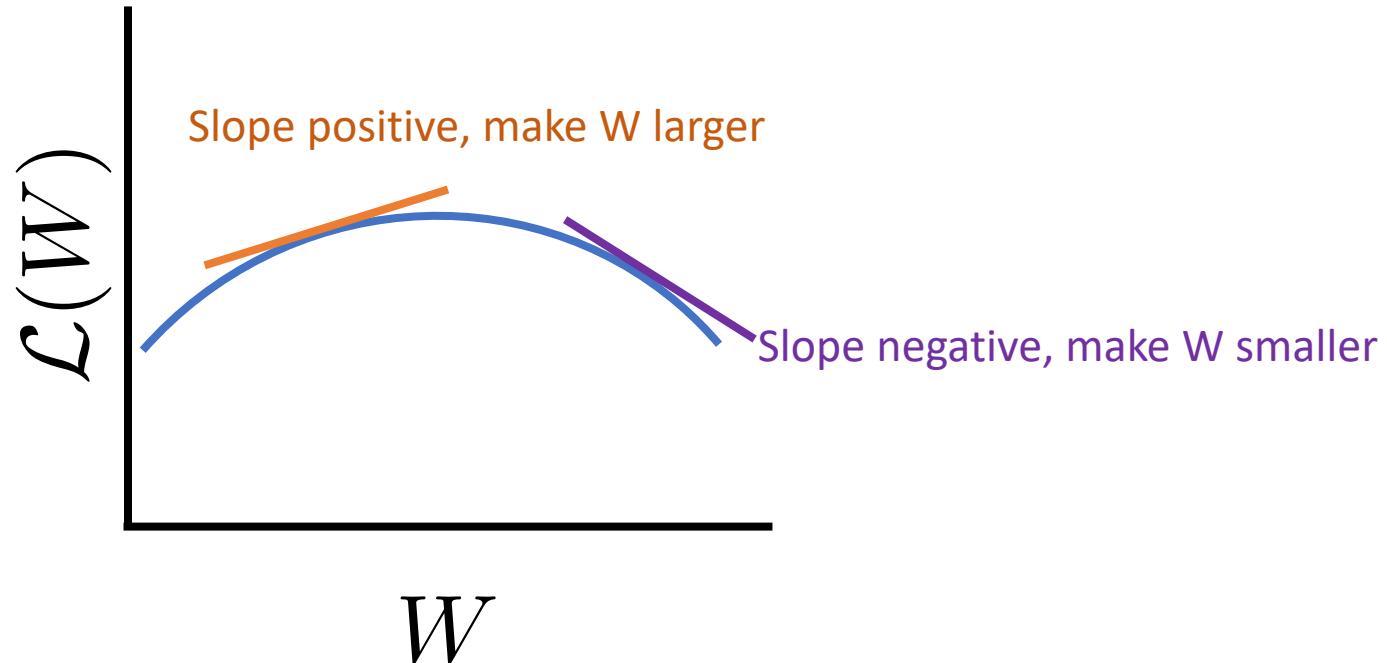
$$\mathcal{L}(W) = \sum_i Y^i (w' x^i) - \sum_i \ln(1 + e^{w' x^i})$$



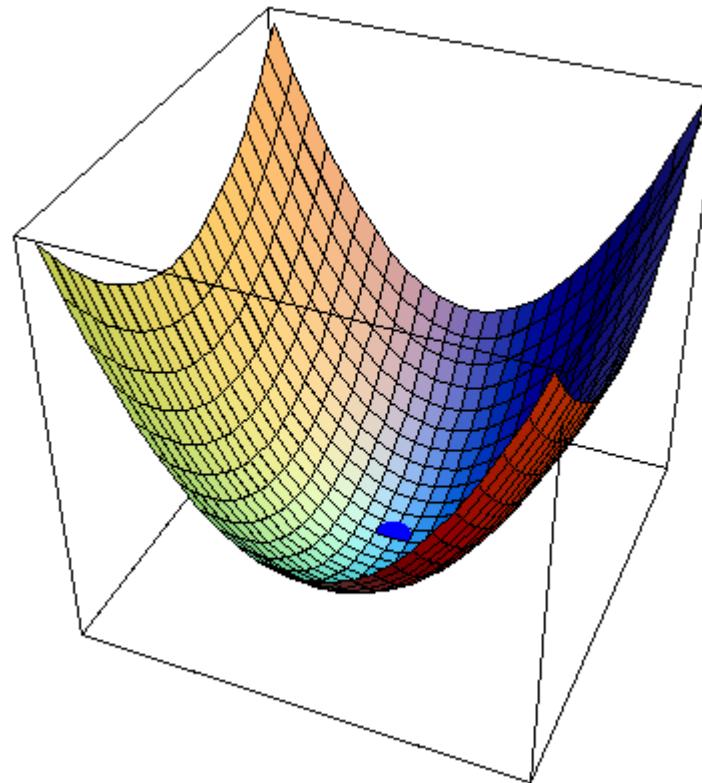
Gradient Descent/Ascent

How do we know *how* to change W ?

We want to move along the blue line in a way that makes $L(W)$ larger

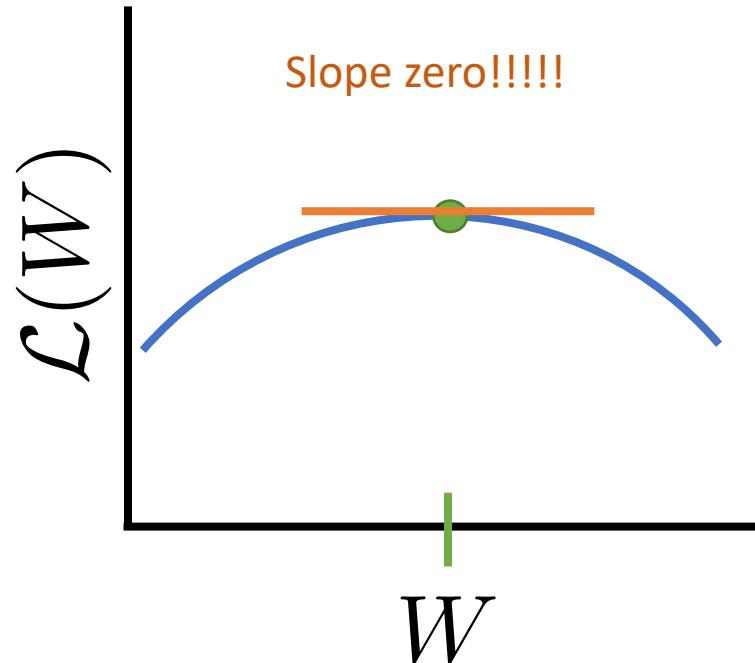


Gradient Descent/Ascent



Maximizing Conditional Log Likelihood

$$\mathcal{L}(W) = \sum_i Y^i (w' x^i) - \sum_i \ln(1 + e^{w' x^i})$$



Maximizing Conditional Log Likelihood

$$\mathcal{L}(W) = \sum_i Y^i (w' x^i) - \sum_i \ln(1 + e^{w' x^i})$$

$$\begin{aligned} \frac{d}{dw} \mathcal{L}(W) &= \sum_i Y^i x^i - \sum_i \frac{x^i e^{w' x^i}}{1 + e^{w' x^i}} \\ &= \sum_i x^i (Y^i - P(Y = 1 | W, x^i)) \end{aligned}$$

Iterative Method

- Start at $\mathbf{w}_0 = \mathbf{1}$; take a step along **steepest slope**
- Fixed step size:

$$\mathbf{w}_1 = \mathbf{w}_0 + \eta \mathbf{v}$$

- \mathbf{v} is a **vector** in the direction of the **steepest slope**.
 - What's the steepest slope?

Gradient Descent Algorithm

Kappa is the learning rate

Initialize $w_0=0$

For $t=0,1,2,\dots$ do

Compute the gradient and update the weights

$$\frac{d}{dw} \mathcal{L}(W) = \sum_i x^i (Y^i - P(Y = 1|W, x^i))$$

$$w \leftarrow w + \kappa \left(\sum_i x^i (Y^i - P(Y = 1|W, x^i)) \right)$$

Iterate with the next step until w doesn't change too much

(or for a fixed number of iterations)

Return final w .

We may not find the exact point where the slope is 0, but we will get "close enough"

Making predictions

- A new tuple comes: $(\mathbf{x}, ?)$

$$p(y = 0 | \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}' \mathbf{x}}}$$

$$p(y = 1 | \mathbf{x}) = \frac{e^{\mathbf{w}' \mathbf{x}}}{1 + e^{\mathbf{w}' \mathbf{x}}}$$

- Fix a threshold in $[0,1]$ to make predictions.

$$p(y = 1 | \mathbf{x}) > \text{threshold} \quad \text{Predict } y=1$$

$$p(y = 1 | \mathbf{x}) \leq \text{threshold} \quad \text{Predict } y=0$$

Example

GPA, GRE, and success.	Dummy	GPA	GRE	y
100, 800, 1	1	1.0	1.0	1
90, 800, 1	1	0.9	1.0	1
90, 700, 1	1	0.9	0.875	1
90, 700, 1	1	0.7	0.75	0
70, 600, 0	1	0.6	0.875	0
60, 700, 0	1	0.6	0.875	1
60, 700, 1	1	0.5	0.75	0
60, 700, 1	1	0.5	0.8125	0
50, 600, 0	1	0.5	1.0	1
50, 650, 01	1	0.5	0.875	0
50, 800, 1	1	0.5	0.875	1
50, 700, 0				
50, 700, 1				

Scaled so that max GPA/GRE is 1

It's a good idea to scale
your features!

Example

Logistic regression:

$$-7.44 + 4.4056 * \text{GPA} + 5.57 * \text{GRE}$$

$$p(y = 0 | gpa, gre) = \frac{1}{1 + e^{(-7.44 + 4.4056 * \text{GPA} + 5.57 * \text{GRE})}}$$

Fix a threshold in [0,1] to make predictions.

$$p(y = 1 | gpa, gre) > \text{threshold} \quad \text{Predict } y=1$$

$$p(y = 1 | gpa, gre) \leq \text{threshold} \quad \text{Predict } y=0$$

Odds

Definition: $odds(0 \text{ vs. } 1 \text{ given } \mathbf{x}) = \frac{p(y=0|\mathbf{x})}{p(y=1|\mathbf{x})}$

Formula: $odds(0 \text{ vs. } 1 \text{ given } \mathbf{x}) = \frac{1}{\frac{1+e^{\mathbf{w}'\mathbf{x}}}{e^{\mathbf{w}'\mathbf{x}}}}$

$$= \frac{1}{e^{\mathbf{w}'\mathbf{x}}}$$

$$= e^{-\mathbf{w}'\mathbf{x}}$$

We predict ? if this number is less than ?

Interpretation

$$odds(\text{successful vs. unsuccessful given } gpa \text{ and } gre) = e^{-7.44 + 4.4056*\text{GPA} + 5.57 * \text{GRE}}$$

If GPA increases by .1 then the odds of success will increase by 55%

$$e^{0.441} \approx 1.55$$

If GRE increases by .1 then the odds of success will increase by 75%

$$e^{0.557} \approx 1.75$$

Regularization

- Recall: What's overfitting?
- Can that happen in logistic regression?
- How can we avoid overfitting?

Regularization

- If I have two models, one with 10 parameters and one with 100, which is more complex?
- So we would like to encourage the model to “have” fewer parameters
 - In practice, this boils down to making some of the parameters close to 0
- This is what regularization can do

Regularization

- Two typical regularization forms
 - L1 (produces sparsity: many w set to 0)

$$\sum_i |w_i|$$

- L2 (weights small, but not always 0)

$$\sum_i w_i^2$$

So for L1 our likelihood becomes...

$$\mathcal{L}(W) = \sum_i Y^i (w'x^i) - \sum_i \ln(1 + e^{w'x^i}) - \lambda \sum_i |w_i|$$

- Lambda is a *hyperparameter*
- Tune lambda using nested cross validation, or the validation set.

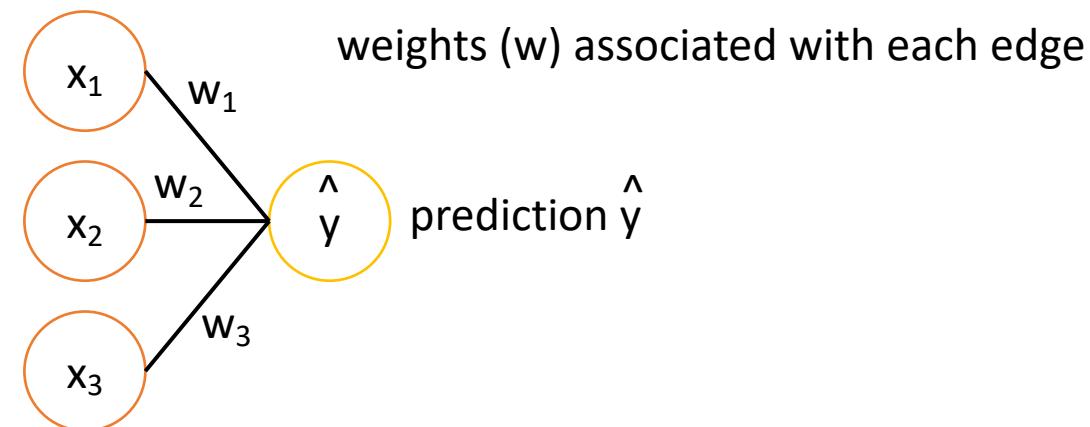
Multi-class

- What if there are more than two classes?
 - E.g. positive, negative or neutral sentiment
- Multinomial logistic regression
 - Essentially builds a logistic regression model for every class
 - Predict the highest probability class
 - In practice this can be done with the softmax function

Logistic regression is actually a simple neural net

Neural Nets

input features (x)

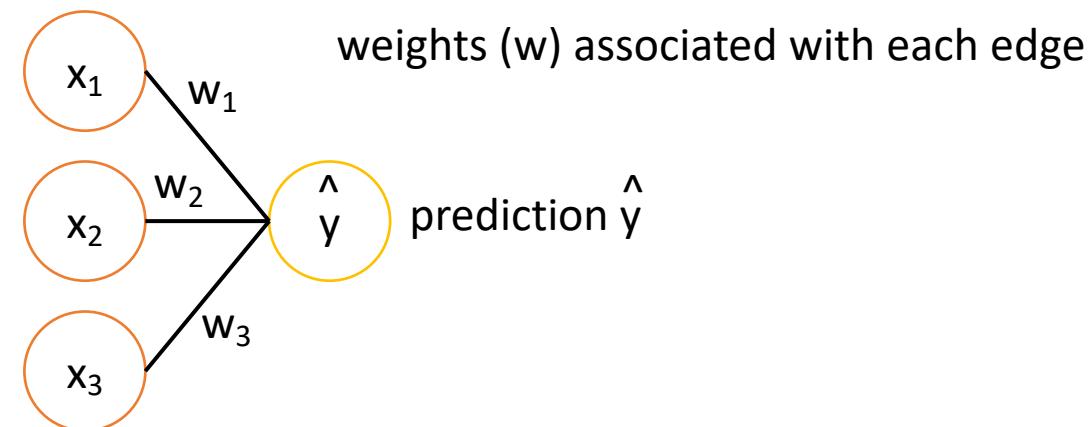


$$\hat{y} = \frac{1}{1 + e^{-\sum_i x_i w_i}}$$

output
(in the range [0...1])

Neural Nets

input features (x)



$$\hat{y} = \frac{1}{1 + e^{-\sum_i x_i w_i}}$$

output
(in the range [0...1])

Logistic regression is actually a simple neural net

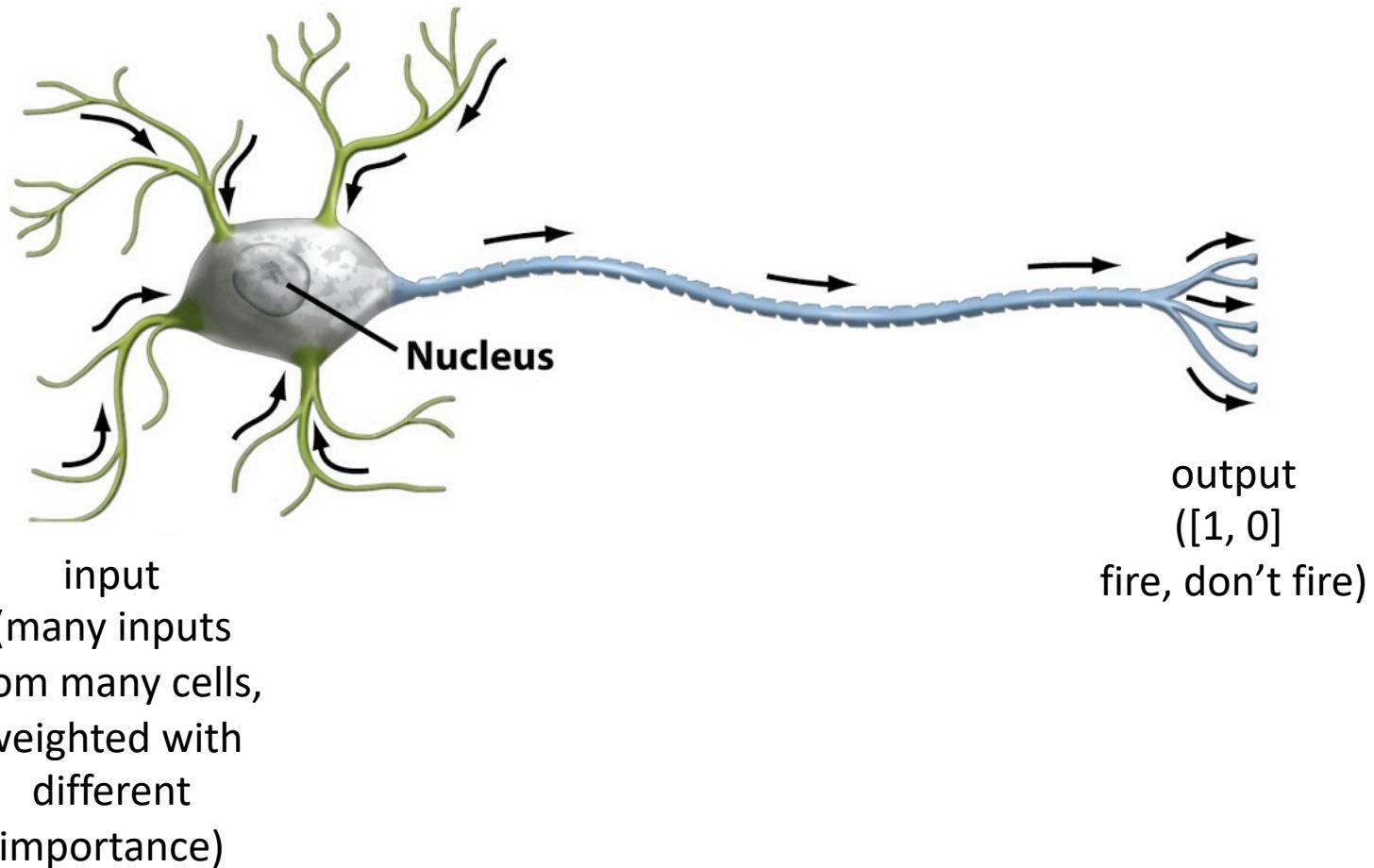
Let's play!

<https://playground.tensorflow.org/>

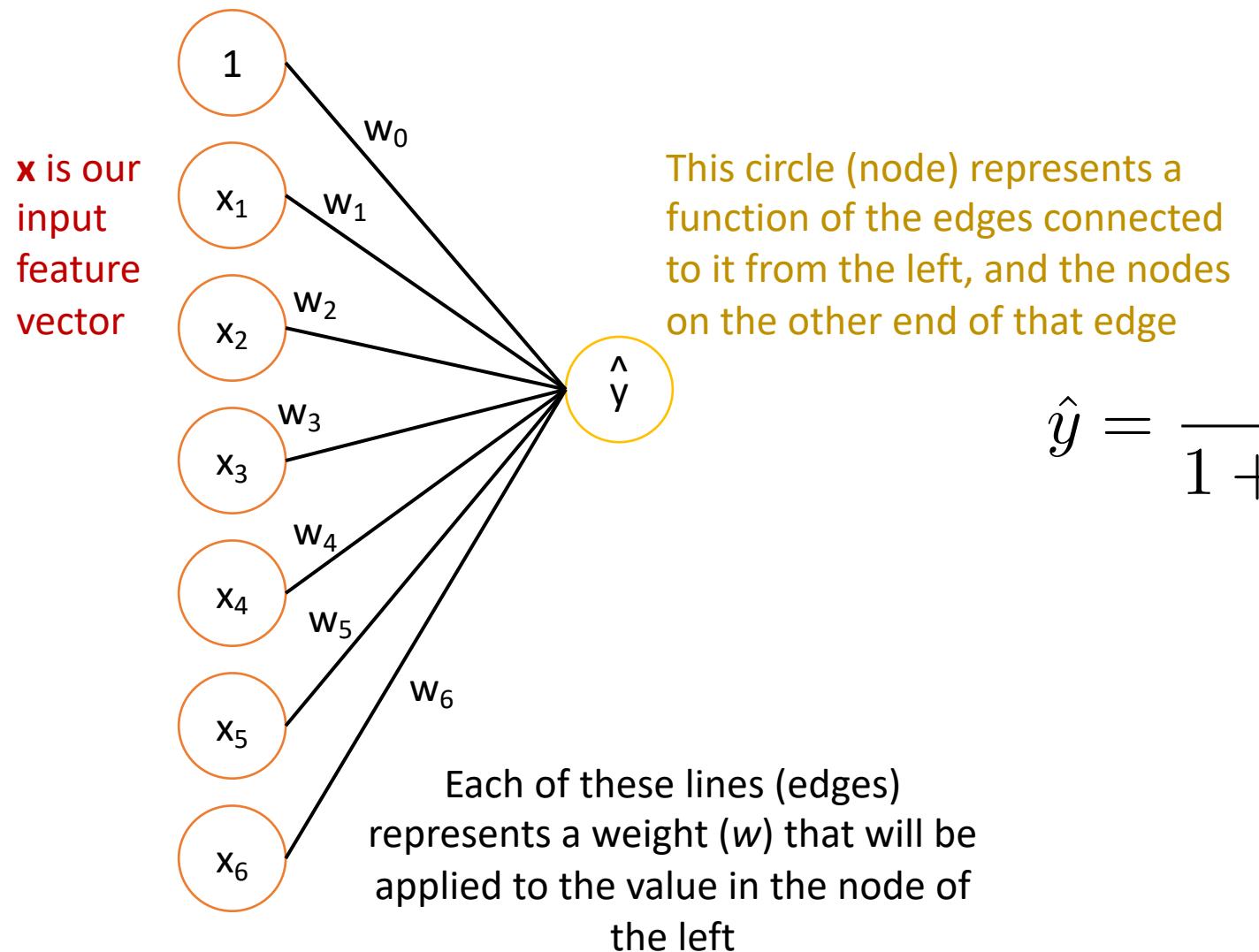
Neural Nets

- Have been around for decades
 - See <http://www.andreykurenkov.com/writing/a-brief-history-of-neural-nets-and-deep-learning/> for a brief history lesson
- Recently re-popularized due to
 - ability to process more data more quickly
 - GPUs
 - availability of large datasets (particularly in computer vision)

Inspiration

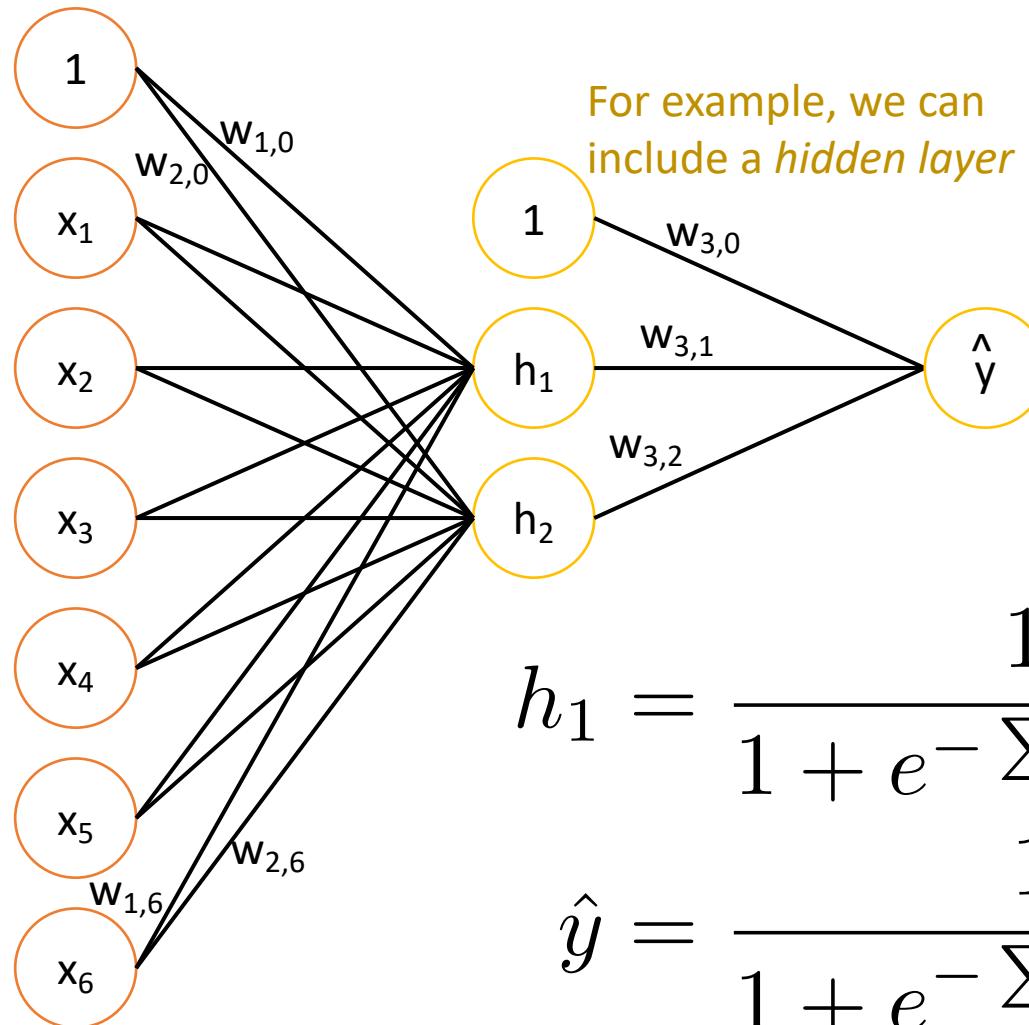


Logistic Regression Models are Neural Networks



$$\hat{y} = \frac{1}{1 + e^{-\sum_i x_i w_i}}$$

Neural nets are a generalization that includes logistic regression, but also other models

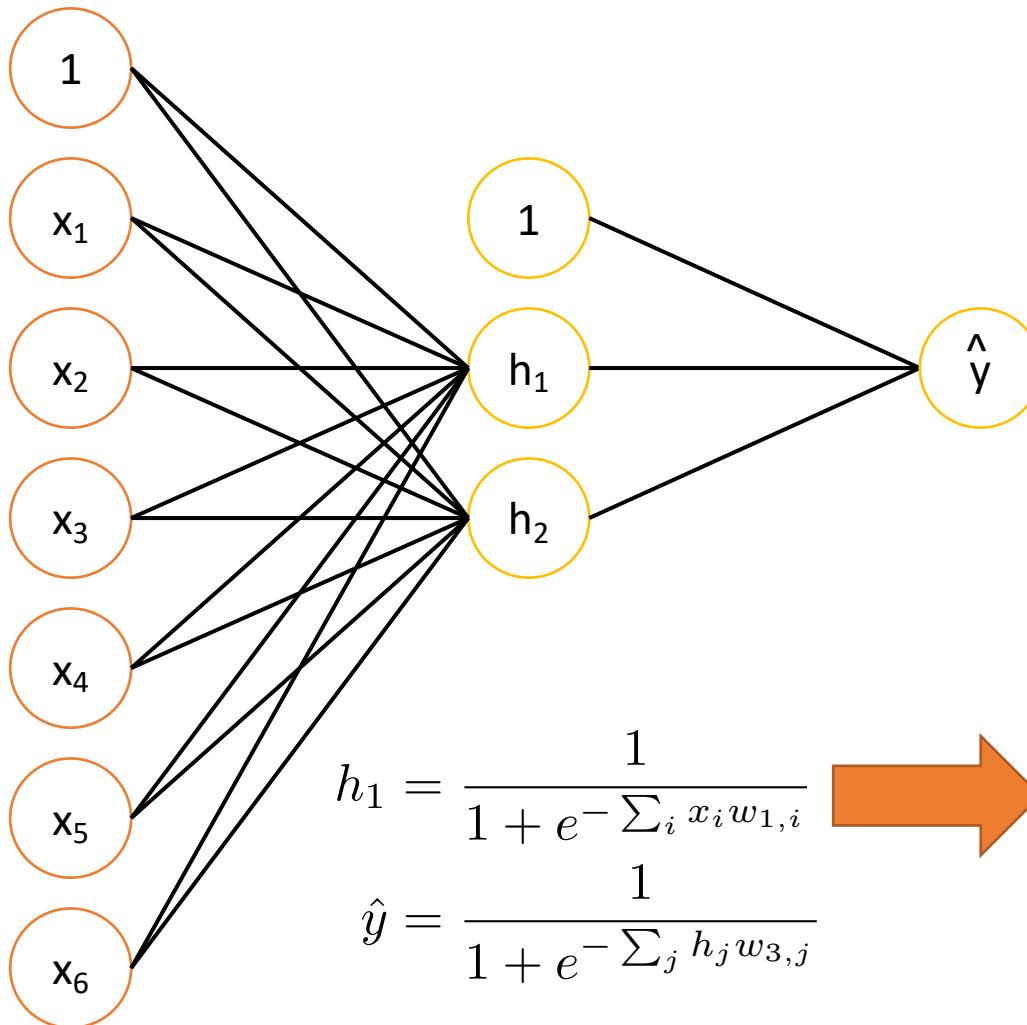


$$h_1 = \frac{1}{1 + e^{-\sum_i x_i w_{1,i}}}$$

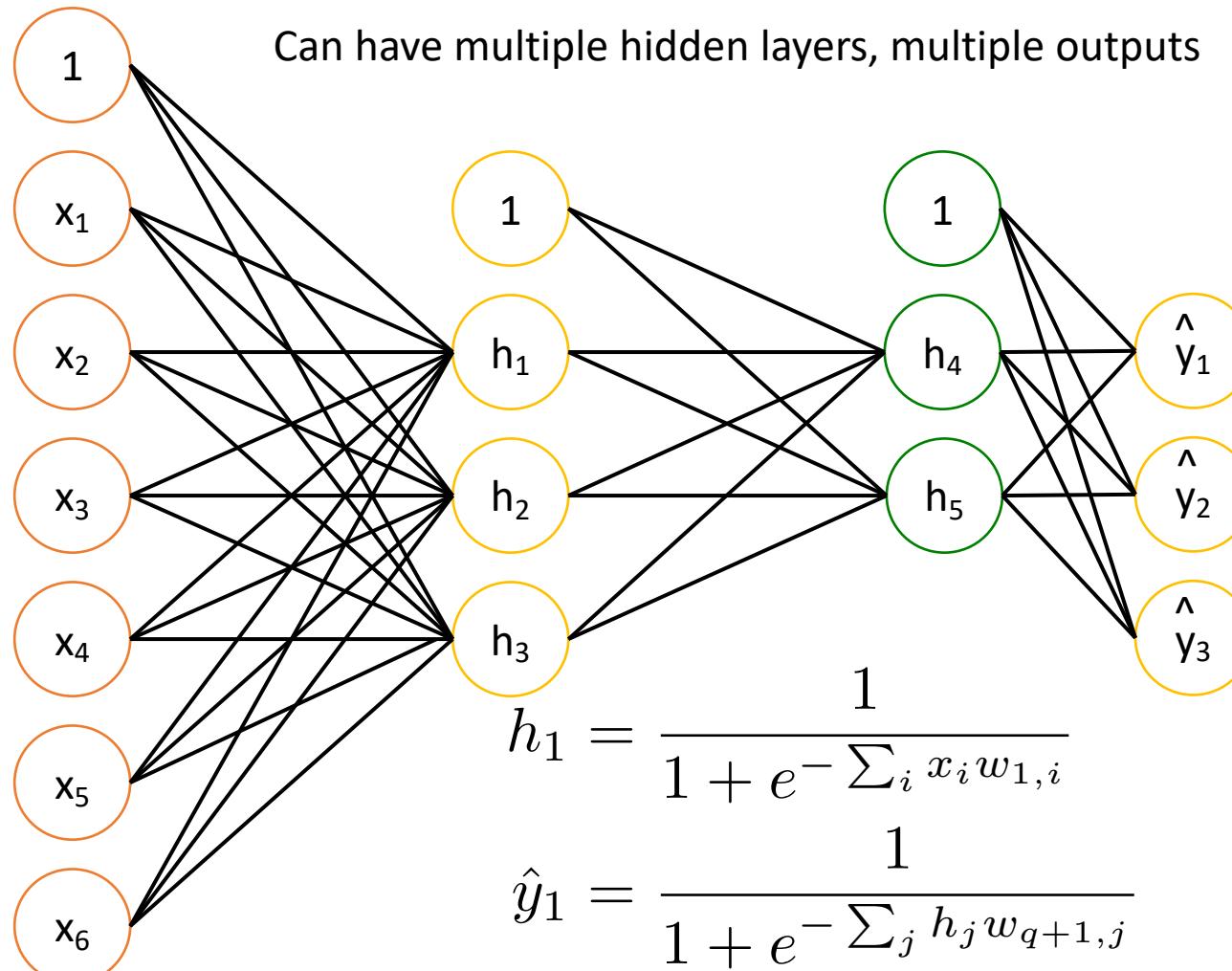
$$\hat{y} = \frac{1}{1 + e^{-\sum_j h_j w_{3,j}}}$$

(For simplicity, only some edges are labeled here)

We can also change the activation function



Neural Nets



q is total # of hidden nodes across all layers

Deep Learning

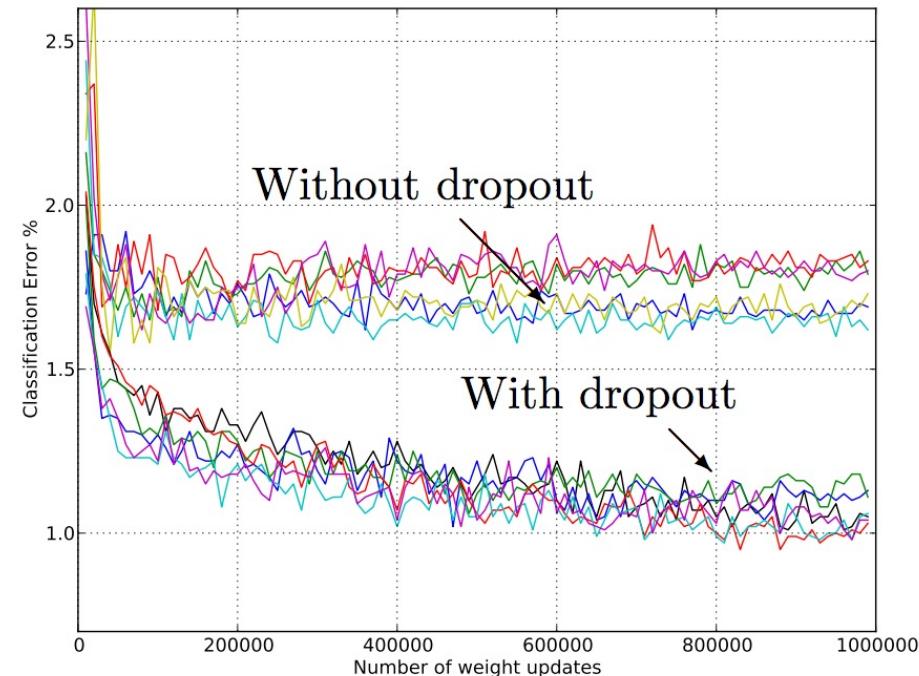
- A model with more than two hidden layers can be called “deep”
- State of the art models can have 50, 100 or more hidden layers

Regularization

- We are learning a much larger number of parameters
- We need to regularize
- Can use L1 or L2 regularizers as in logistic regression & linear regression

Regularization at the Network Level

- Dropout
 - randomly remove nodes from the network @ train
 - i.e. set to 0
 - encourages redundant connections
 - reduces overfitting



Deep Learning for Computer Vision

Convolutional neural nets (CNNs)

- Built to capture the invariances we see in images
 - objects can appear at any place in the image



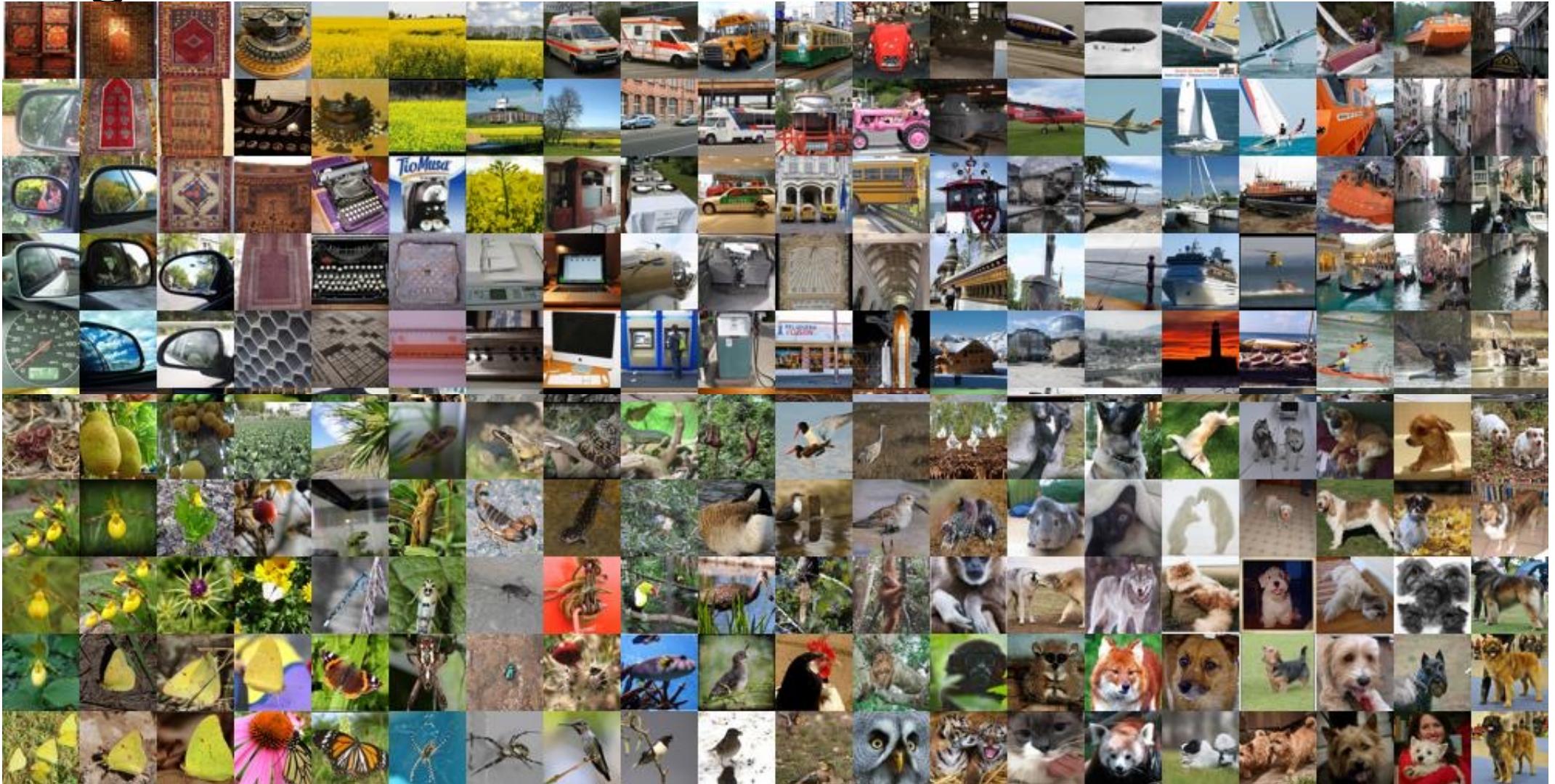
Convolutional neural nets (CNNs)

- Learn filters that respond to particular visual features
 - e.g. edges, curves, etc
- Look for those visual features anywhere in the image
 - i.e. convolve the *same* filter with many patches of the image
- Hidden layers combine these filters to create more complex shapes
 - e.g. straight edges combined to form curves, combined to form the handle on a coffee mug

CNNs

- Amazing advances in computer vision
 - ImageNet
 - First released 2009
 - 1.2 million images
 - more than 1000 concepts
 - e.g. cup, oil filter, ptarmigan
 - Deep learning
 - CNNs

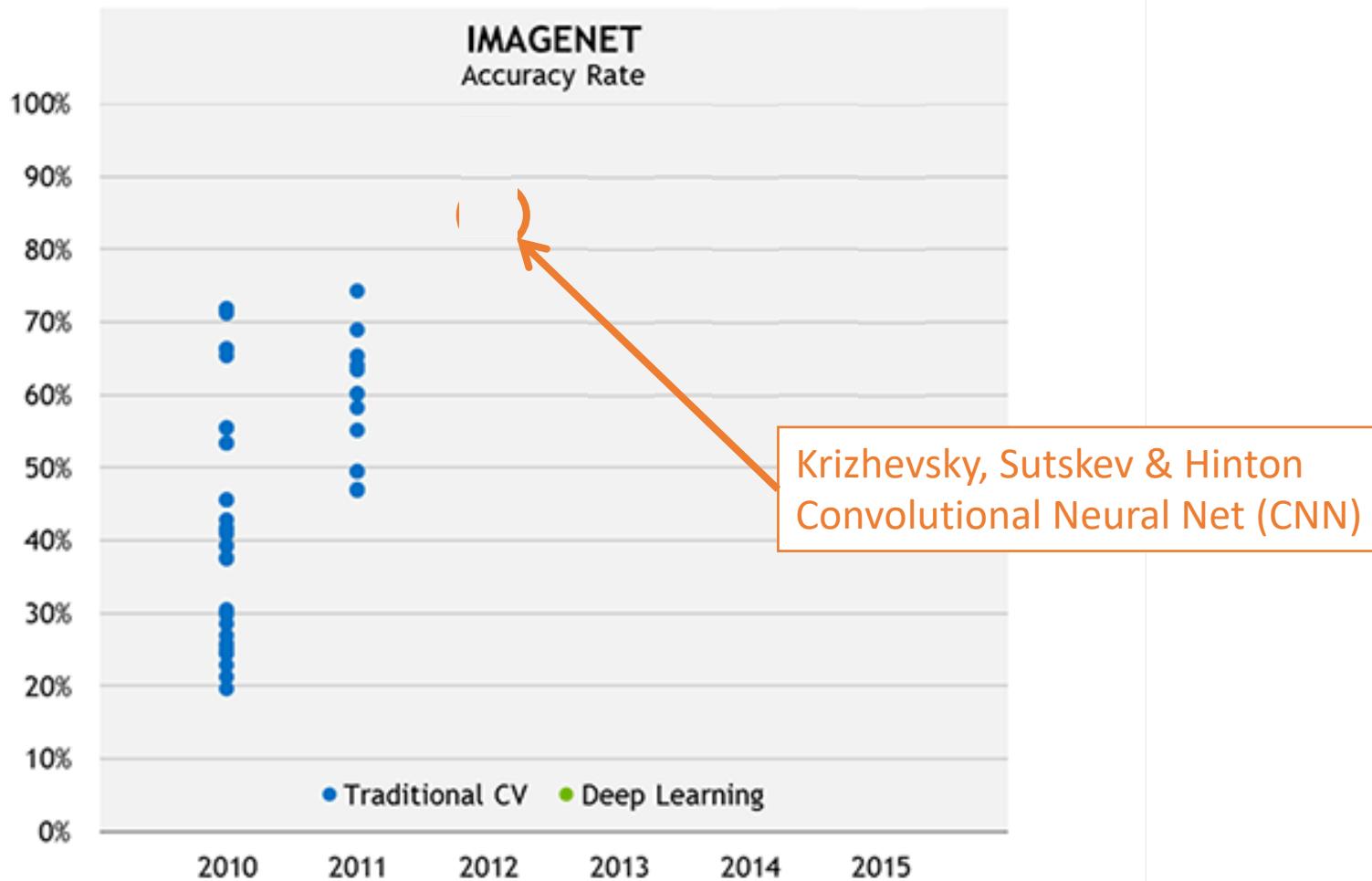
ImageNet



ImageNet

french fries mashed potato black olive face powder crab apple Granny Smith strawberry blueberry cranberry currant
blackberry raspberry persimmon mulberry orange kumquat lemon grapefruit plum fig pineapple banana jackfruit cherry
grape custard apple durian mango elderberry guava litchi pomegranate quince kidney bean soy green pea chickpea chard
lettuce cress spinach bell pepper pimento jalapeno cherry tomato parsnip turnip mustard bok choy head cabbage broccoli
cauliflower brussels sprouts zucchini spaghetti squash acorn squash butternut squash cucumber artichoke asparagus green
onion shallot leek cardoon celery mushroom pumpkin cliff lunar crater valley alp volcano promontory sandbar dune coral
reef lakeside seashore geyser bakery juniper berry gourd acorn olive hip ear pumpkin seed sunflower seed coffee bean
rapeseed corn buckeye bean peanut walnut cashew chestnut hazelnut coconut pecan pistachio lentil pea peanut okra
sunflower lesser celandine wood anemone blue columbine delphinium nigella calla lily sandwort pink baby's breath ice plant
globe amaranth four o'clock Virginia spring beauty wallflower damask violet candytuft Iceland poppy prickly poppy oriental poppy
celandine blue poppy Welsh poppy celandine poppy corydalis pearly everlasting strawflower yellow chamomile dusty miller
tansy daisy common marigold China aster cornflower chrysanthemum mistflower cosmos dahlia coneflower blue daisy
gazania African daisy male orchis butterfly orchid aerides brassavola spider orchid grass pink calypso cattleya red helleborine
coelogyne cymbid lady's slipper marsh orchid dendrobium disa helleborine fragrant orchid fringed orchis lizard orchid laelia
masdevallia odontoglossum oncidium bee orchid fly orchid spider orchid phaius moth orchid ladies' tresses stanhopea stelis
vanda cyclamen centaury gentian begonia commelina scabious achimenes African violet streptocarpus scorpionweed
calceolaria toadflax veronica bonsai star anise wattle huisache silk tree rain tree dita pandanus linden American beech
New Zealand beech live oak shingle oak pin oak cork oak yellow birch American white birch downy birch alder fringe tree
European ash fig witch elm Dutch elm cabbage tree golden shower tree honey locust Kentucky coffee tree Brazilian rosewood
logwood coral tree Japanese pagoda tree kowhai palm Arabian coffee cork tree weeping willow pussy willow goat willow China
tree pepper tree balata teak ginkgo pine ilang-ilang laurel magnolia tulip tree baobab kapok red beech cacao sorrel tree
iron tree mangrove paper mulberry Judas tree redbud mountain ash ailanthus silver maple Oregon maple sycamore box elder
Japanese maple holly dogwood truffle shiitake lichen hen-of-the-woods jelly fungus dead-man's-fingers earthstar coral
fungus stinkhorn puffball gyromitra bolete polypore gill fungus morel agaric trilobite harvestman scorpion black and gold
garden spider barn spider garden spider black widow tarantula wolf spider tick mite centipede millipede horseshoe crab





Self Driving Cars



Self Driving Cars

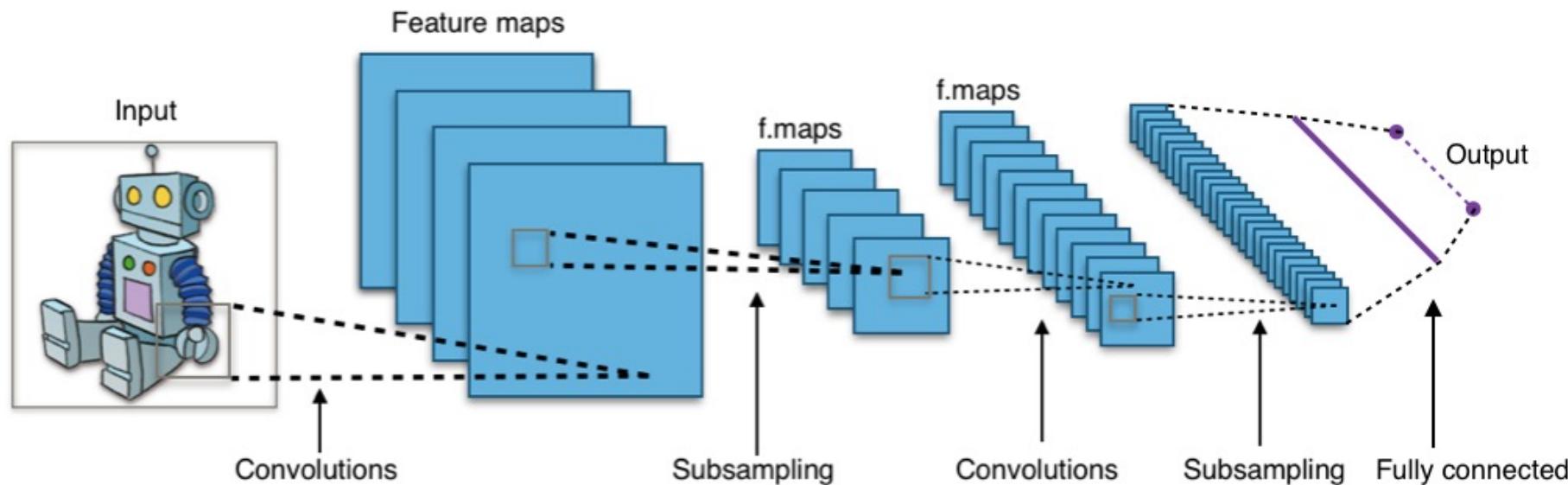


Convolutional Neural Nets

Two additional operations:

- Convolution
- Pooling/Subsampling

Convolutional Neural Nets (CNNs)



CNNs: Convolution

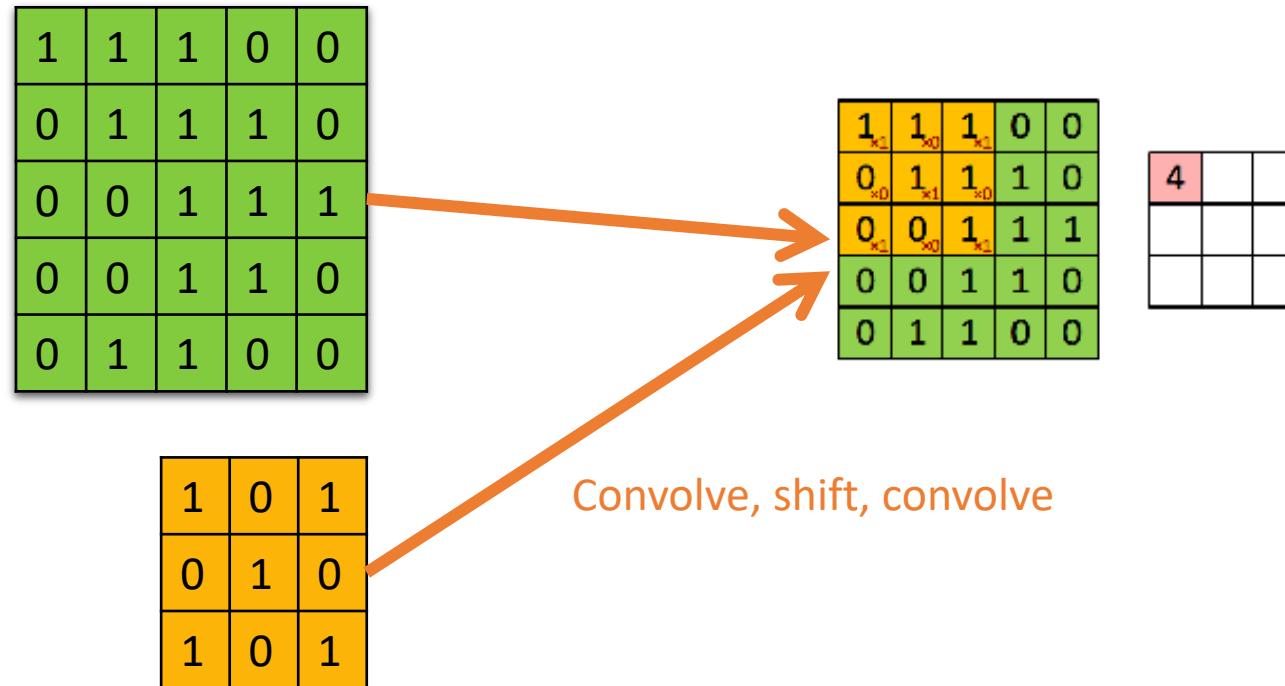
1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Input image (here binary, RGB typical)

1	0	1
0	1	0
1	0	1

Filter (here binary, but typically continuous values)

CNNs: Convolution

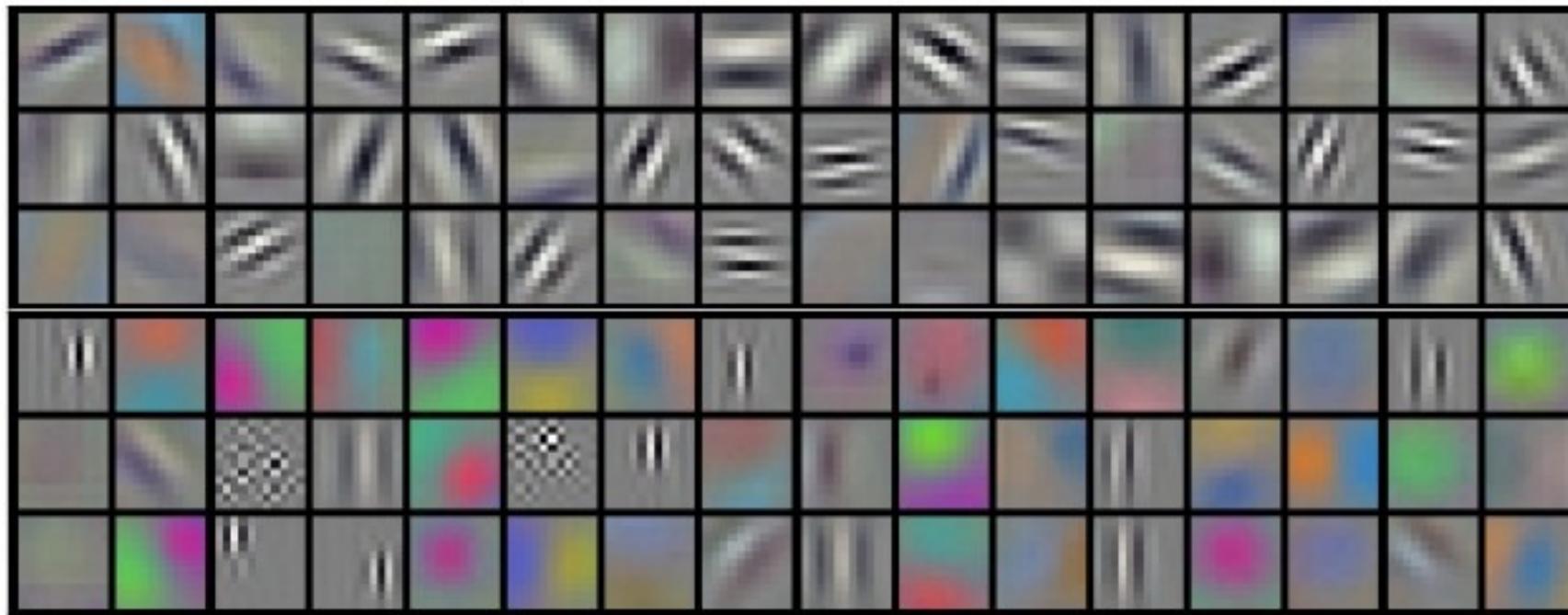


CNNs: Convolution

- Typical Filters ()

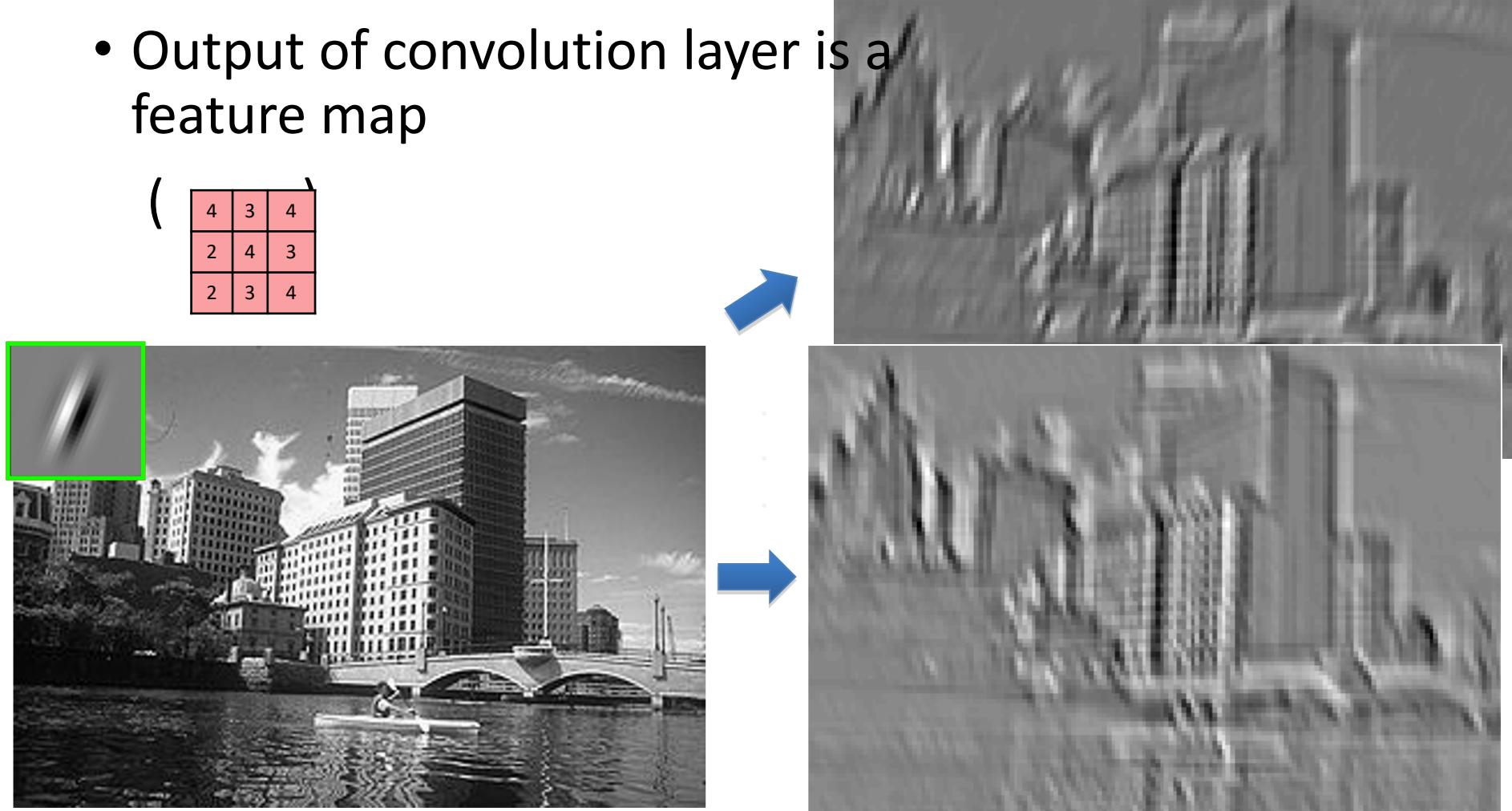
$$\begin{array}{|c|c|c|} \hline 1 & 0 & 1 \\ \hline 0 & 1 & 0 \\ \hline 1 & 0 & 1 \\ \hline \end{array}$$

Filters are learned

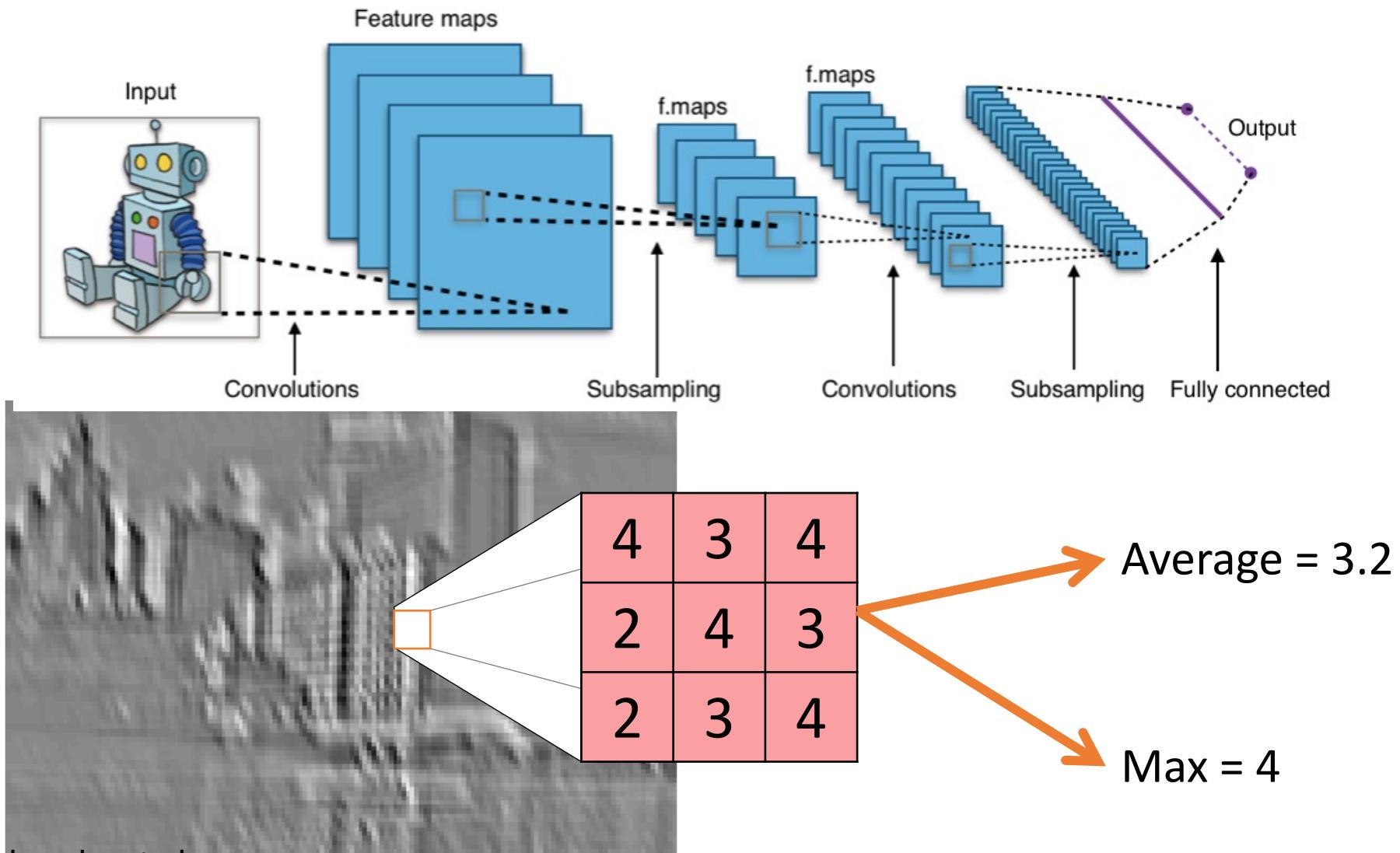


CNNs: Convolution

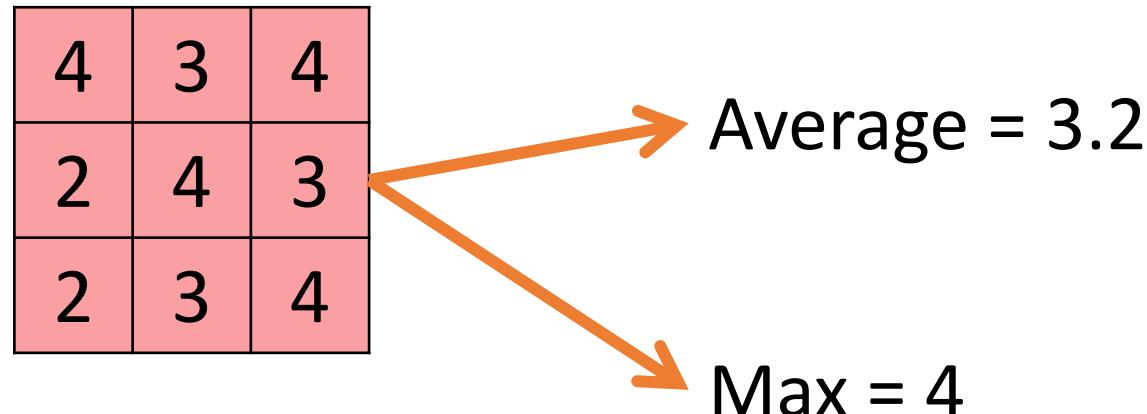
- Output of convolution layer is a feature map



CNNs: Pool/Subsample



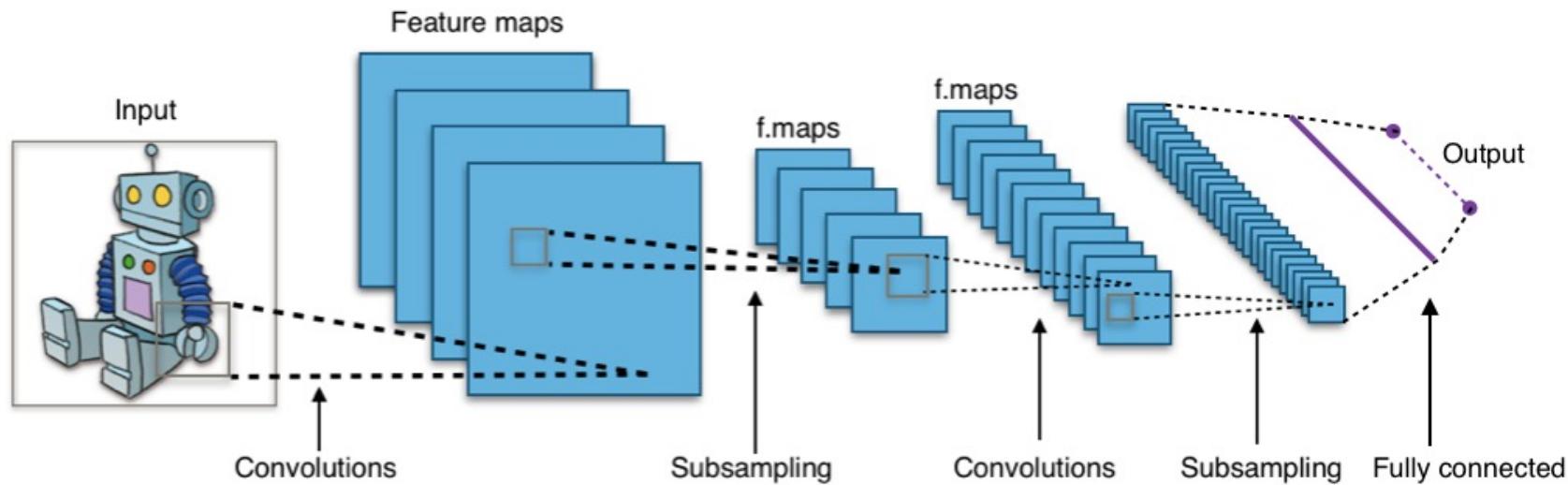
CNNs: Pool/Subsample



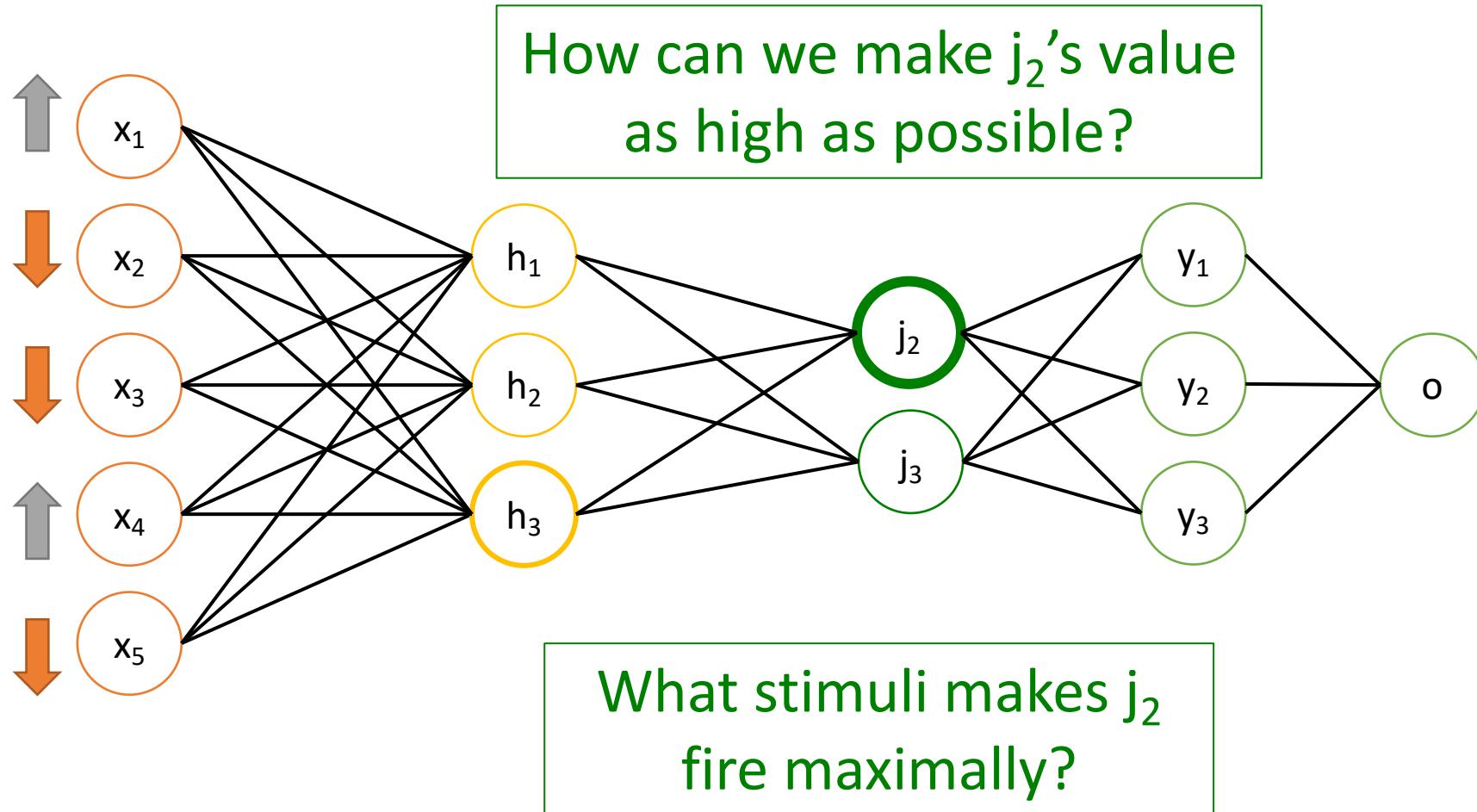
- Supplies some **invariance** to local translations
- Reduces the dimension of subsequent layers

CNNs: Hidden representations

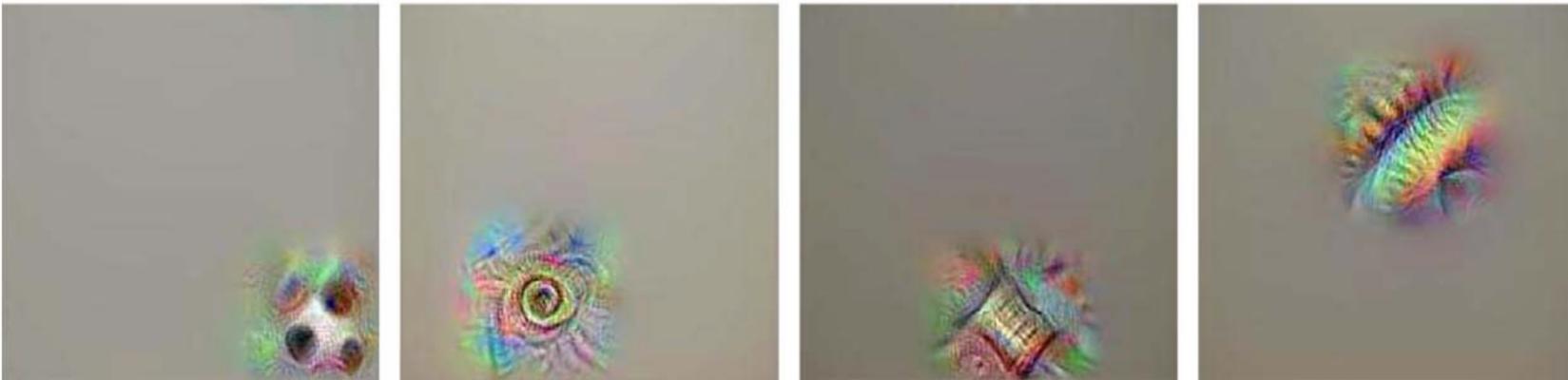
- Activations (feature maps) at each layer are a “hidden representation” of the image
- A compression of the information in image
 - not unlike PCA/SVD



What do the neurons represent?



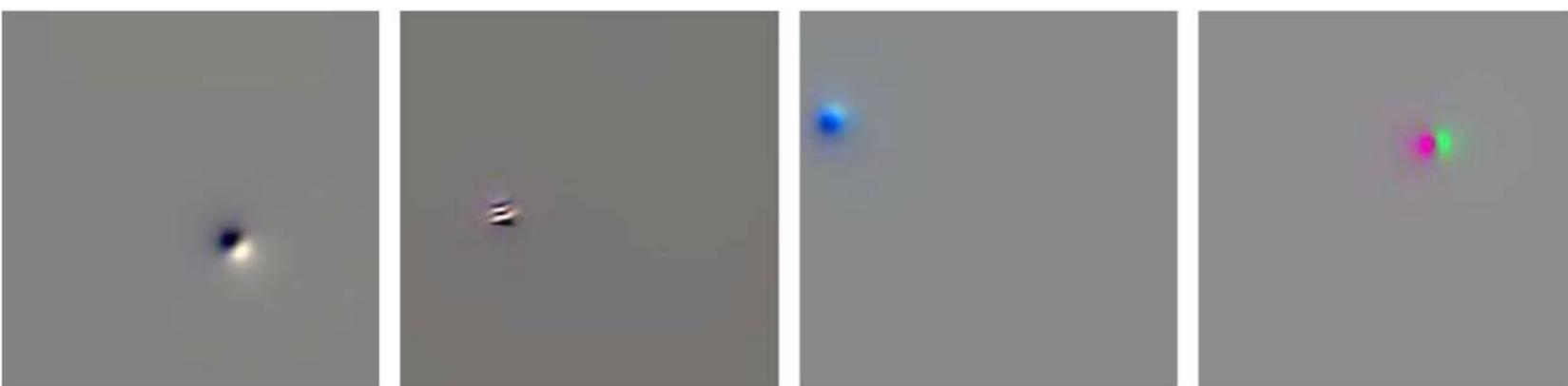
CNN3



CNN2



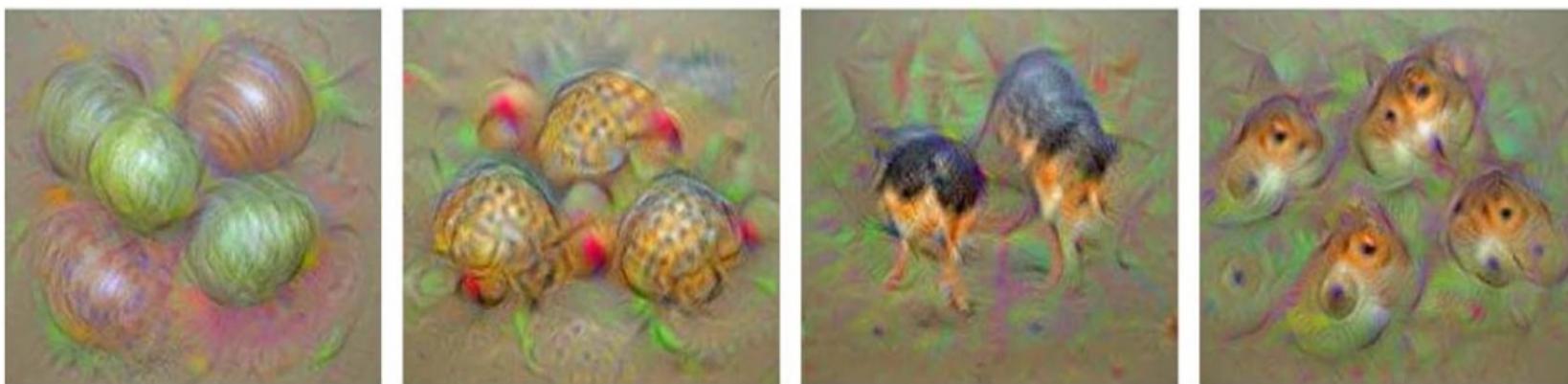
CNN1



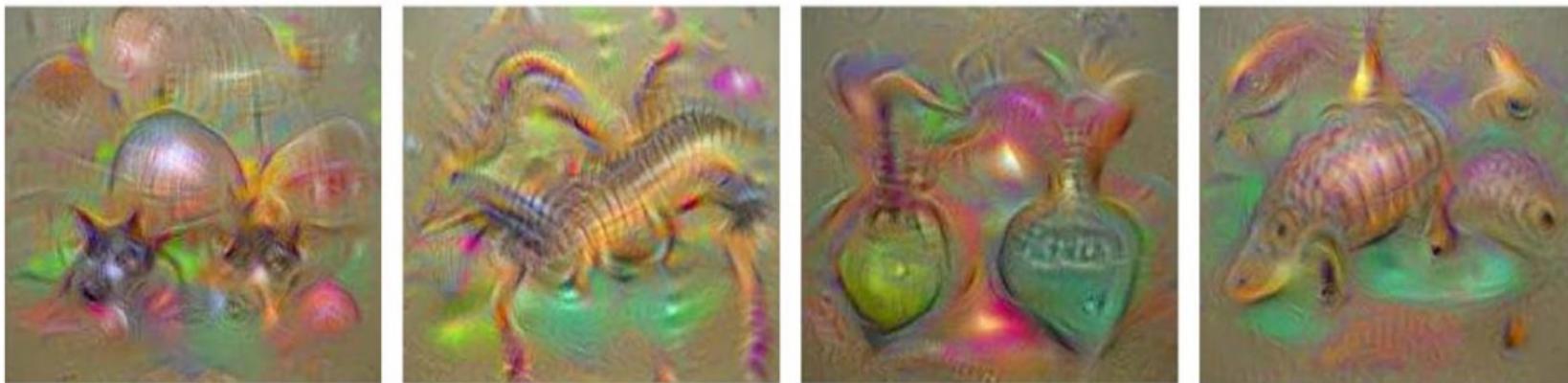
CNN8



CNN7

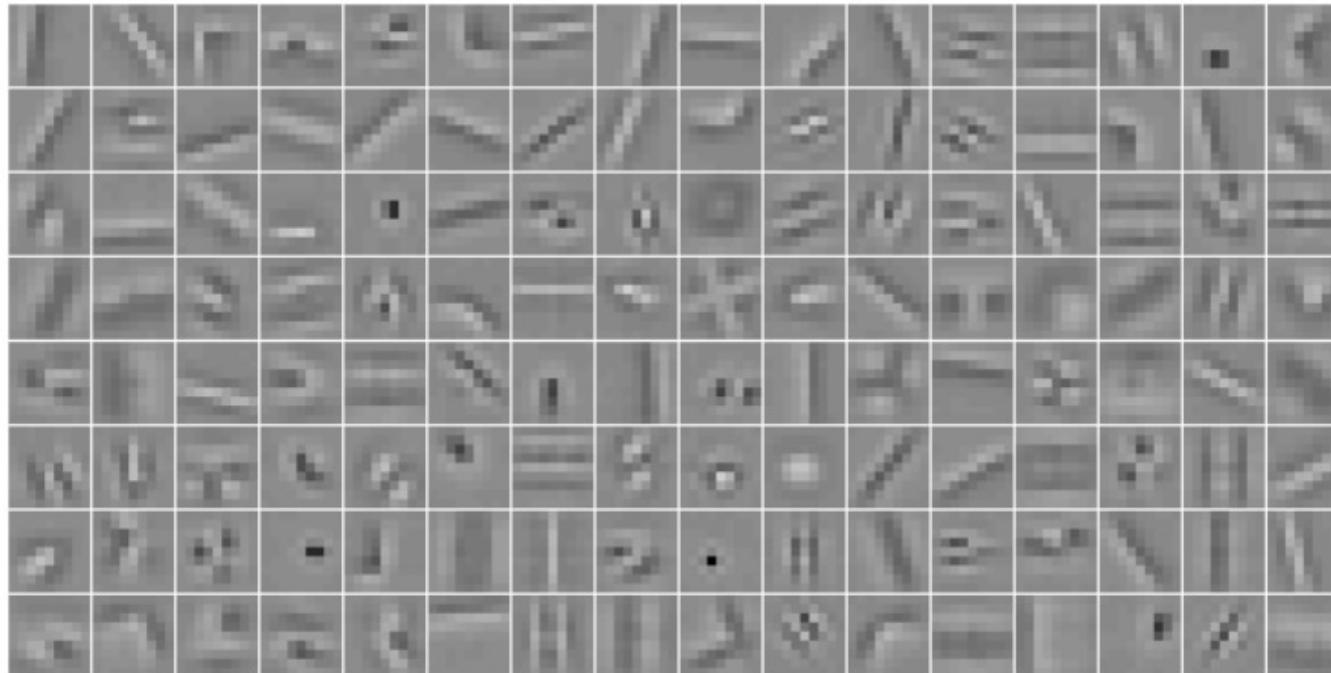


CNN6

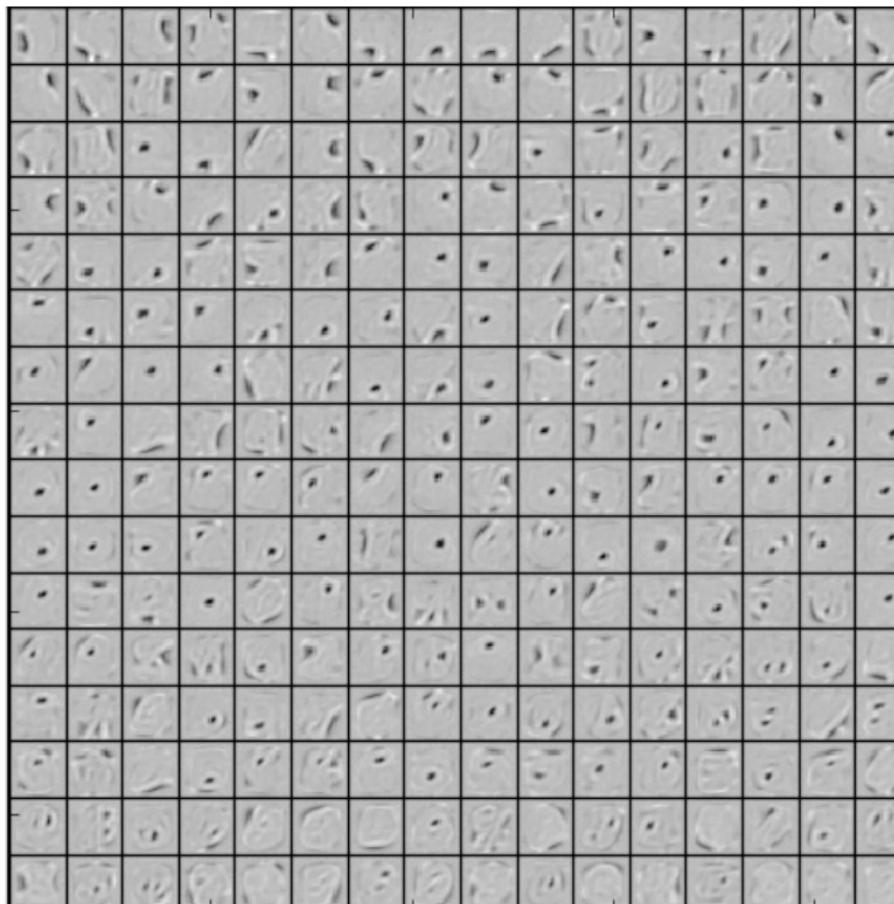


What do CNNs learn?

- When trained on images



What do CNNs Learn?



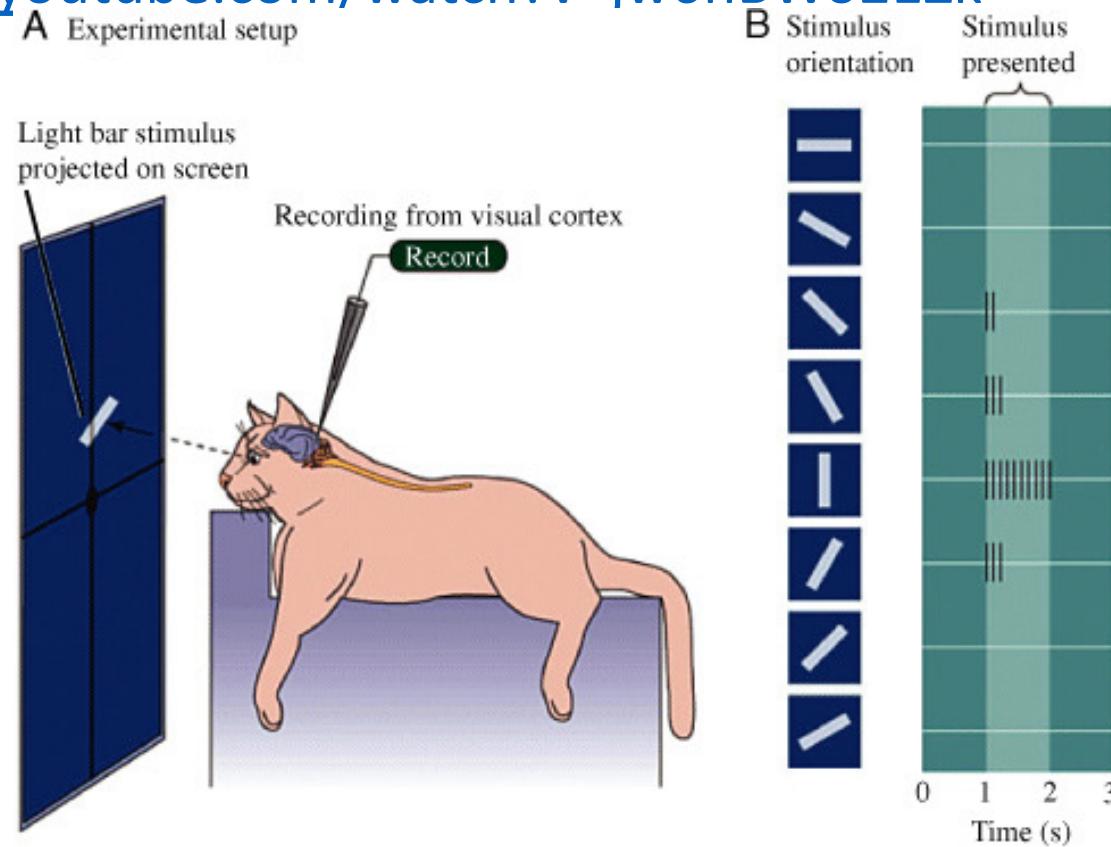
Interpretability

- Chris Olah has a great series on interpretability (the distill ones are particularly good)
 - <https://colah.github.io/>

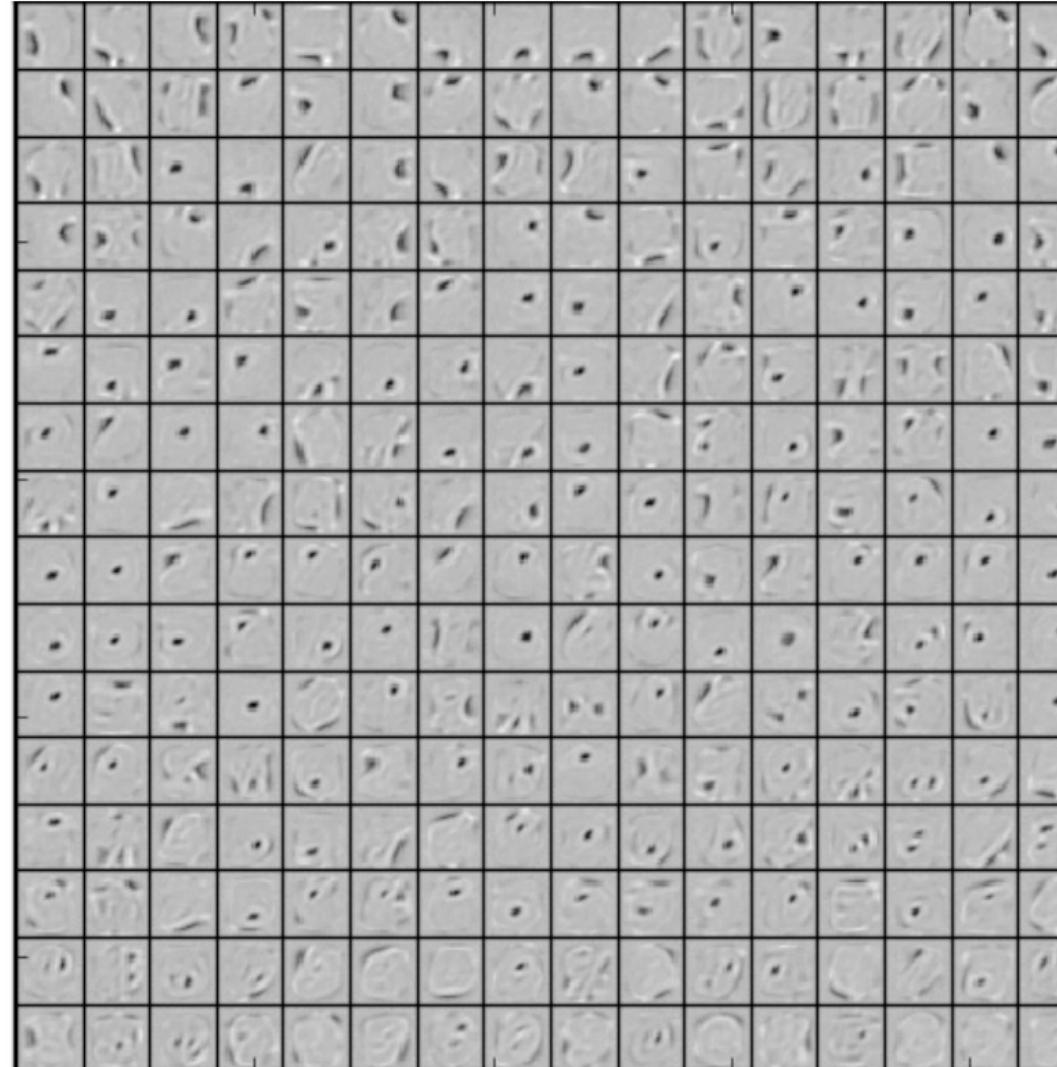
An Aside

- In this video, the static noise you hear is a representation of the neurons firing in response to the visual stimulus

- <https://www.youtube.com/watch?v=jw6nBWo21Zk>

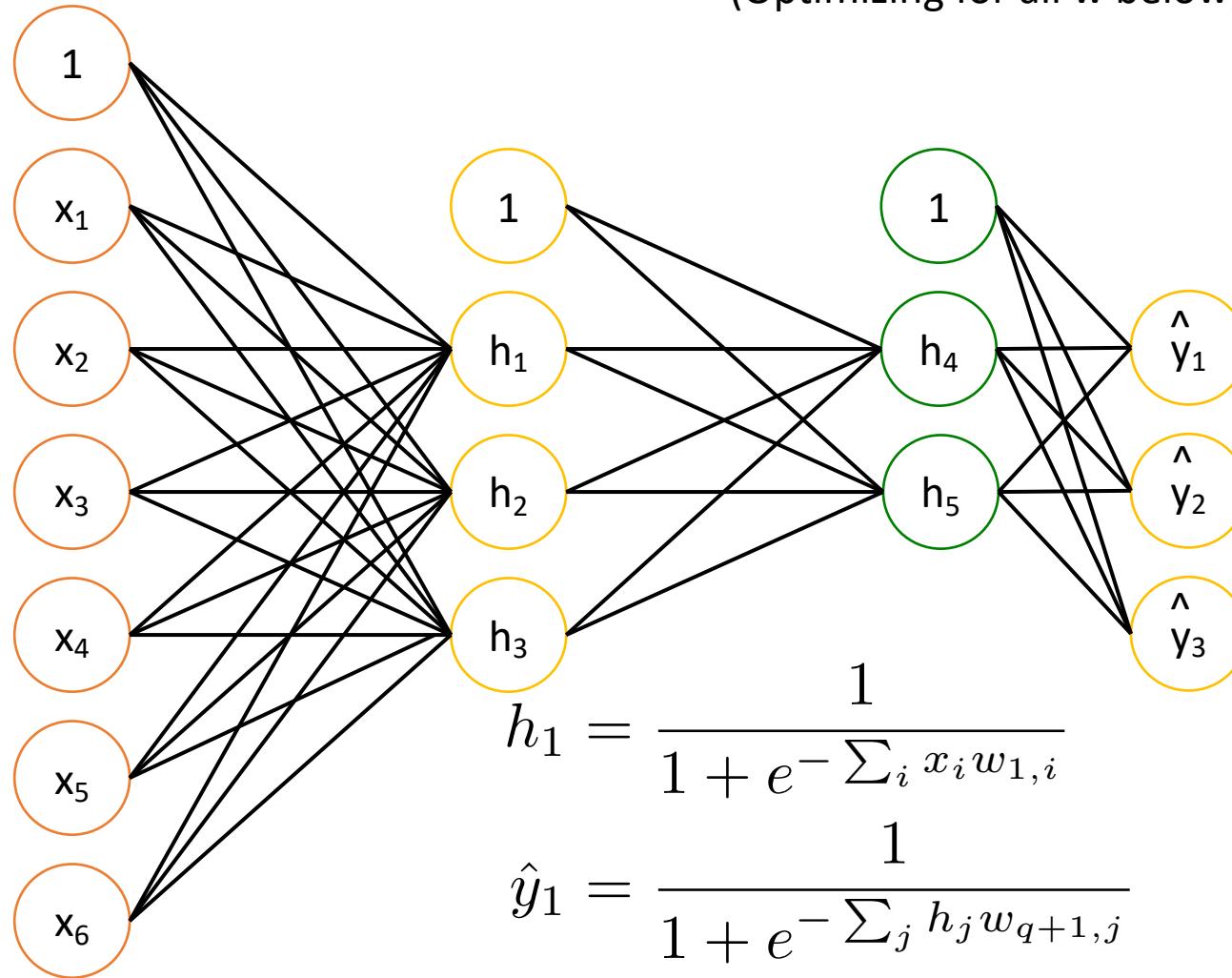


NNs learn something similar



Next time: Back Prop

(Optimizing for all w below)



q is total # of hidden nodes across all layers

Resources

- Train very simple multi-layer perceptron for classification
 - https://scikit-learn.org/stable/modules/neural_networks_supervised.html
- Programming resources for training your own NNs
 - Tensorflow <https://www.tensorflow.org/>
 - Keras <https://keras.io/>
 - Pytorch <https://pytorch.org/>
 - For intuition: <http://playground.tensorflow.org/>
- Short course on deep learning (Nando De Freitas)
 - <https://www.youtube.com/playlist?list=PLjK8ddCbDMphIMSXnw1ljyYpHU3DaUYw>
- Commentary on AlphaGo
 - <https://www.youtube.com/watch?v=UMm0XaCFTJQ>
 - <https://www.youtube.com/watch?v=g-dKXOlsf98>
- Other fun videos
 - Geoff Hinton is in this one! Neural Net stuff is towards the end
 - <https://www.youtube.com/watch?v=yxxRAHVtafl>
 - Fei Fei Li's Ted Talk (Creator of ImageNet)
 - https://www.ted.com/talks/fei_fei_li_how_we_re_teaching_computers_to_understand_pictures?language=en