# Lightning Guide to Machine Learning and Deep Learning Neural Networks

(Learning about machine learning using your neural networks to discover artificial neural networks that might be useful for automatically helping others do their jobs with their own neural networks in analyzing still other people's neural networks…and other stuff.)

# Supervised vs. Unsupervised (and Semi-Supervised) Learning

- Supervised: Data classification based on examples with class labels. Use the data and class assignment to infer a mapping function for new data

- Unsupervised: Automatically learning relationships within unlabeled data

- Semi-Supervised: Maybe guess some labels for unlabeled data and/or use those to label or classify other data.

# Supervised Learning

- Data classification based on examples with class labels

| Weight (kg) | Tail Length (cm) | Color | Cat Class |
|---|---|---|---|
| 3.6 | 32 | ORANGE | DOMESTIC |
| 180 | 168 | STRIPED | TIGER |

- Use the data and class assignment to infer a mapping function for new data

- Algorithms: Aritificial Neural Networks (ANNs), Support Vector Machines (SVMs), logistic regression (LR), Naïve Bayes Classifier, K Nearest Neighbor (KNN), Decision Trees
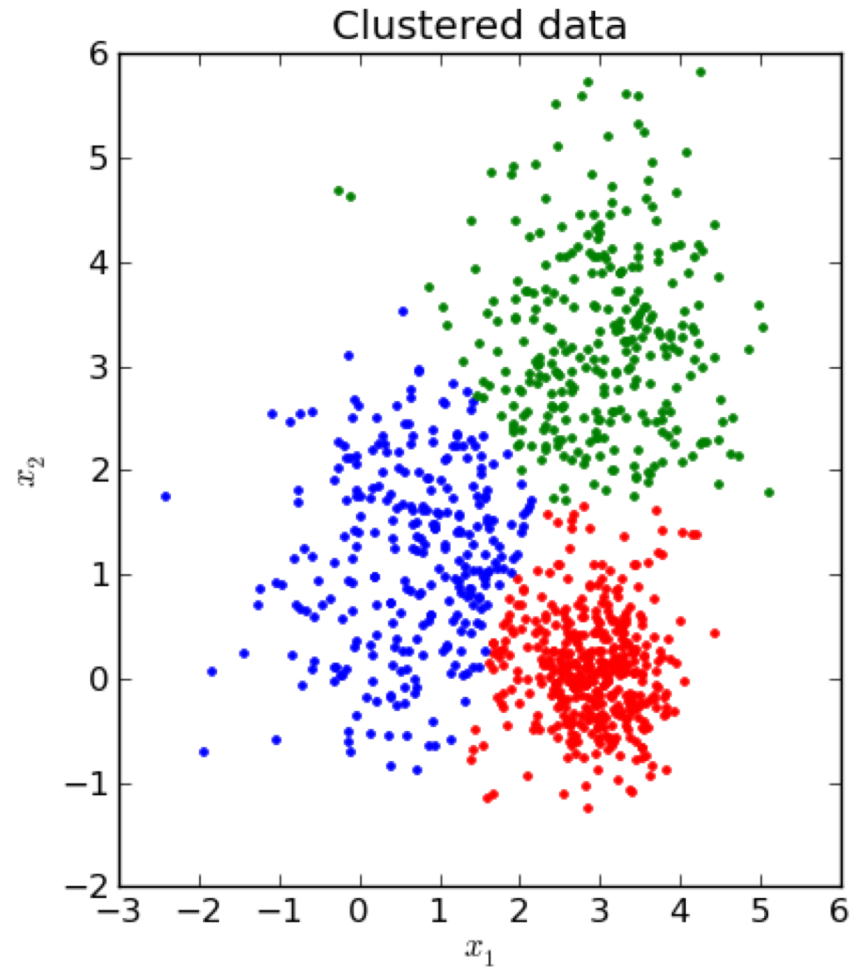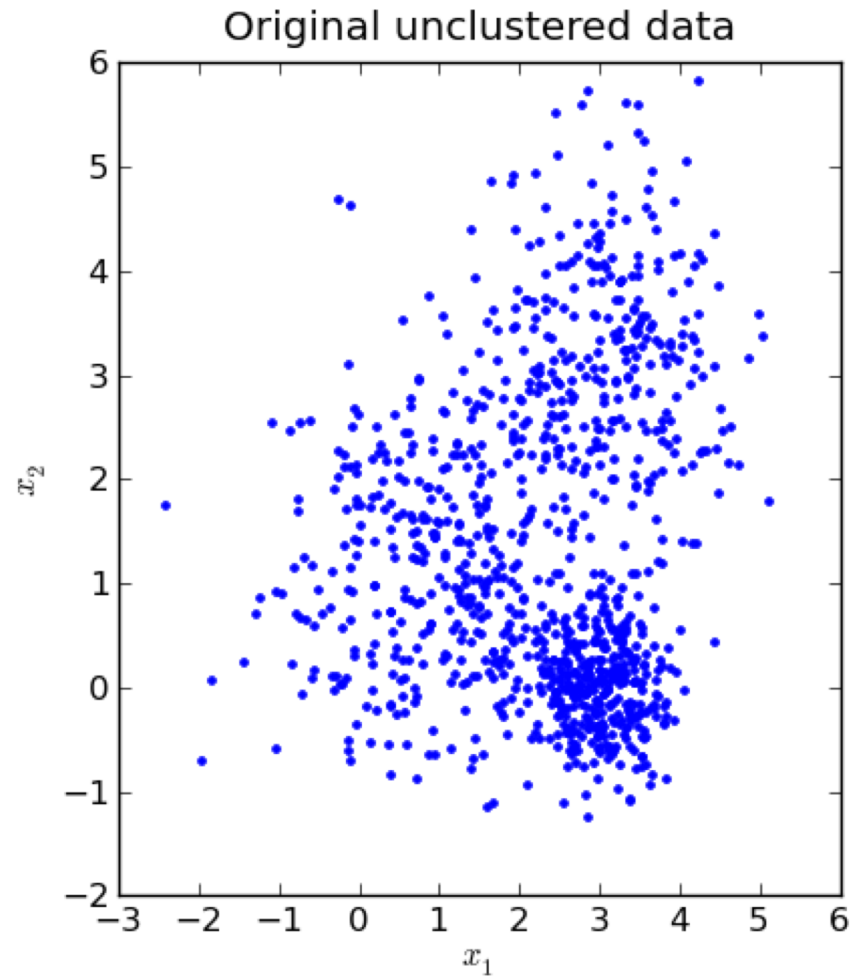
# Unsupervised Learning

- Data organization without labels

| Weight (kg) | Tail Length (cm) | Color |
|---|---|---|
| 3.6 | 32 | ORANGE |
| 180 | 168 | STRIPED |

- Organize the data based on patterns within it.
- Algorithms: Clustering (k-means, agglomerative, top-down, etc.), Self-Organizing Maps (SOMS)
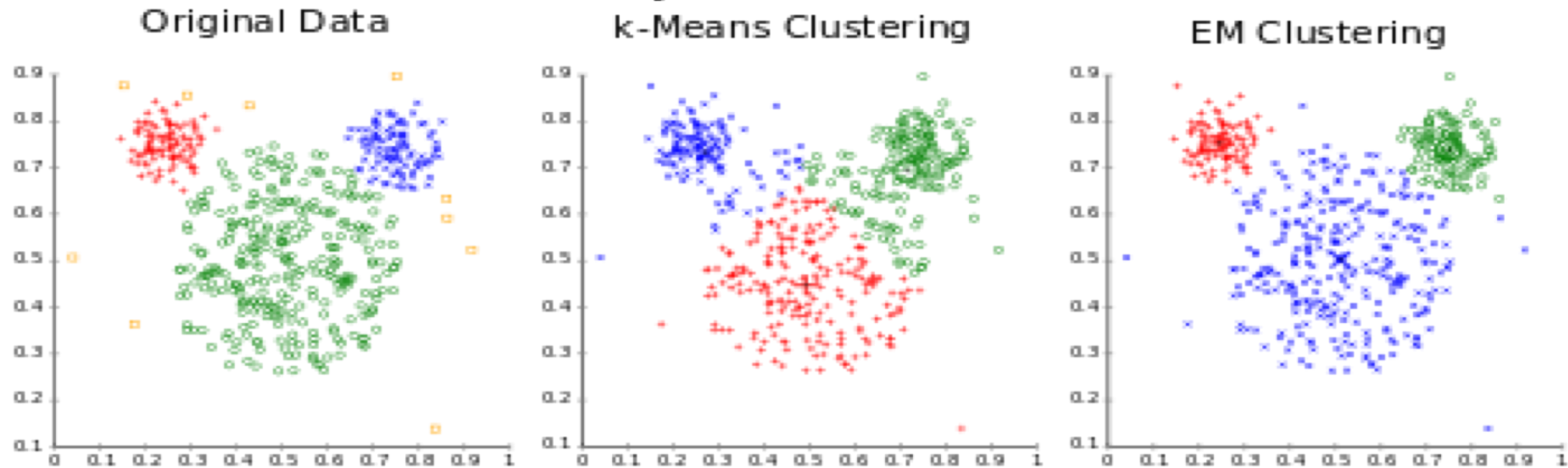
# Example Clustering

# Example Clustering

- For each data point, find the nearest data point using some distance metric (like "straight line" or Euclidean distance)

- Add those two data points to Cluster A

- Take the point midway between the data points in A and now add the nearest data point to A

- Repeat until some criteria is reached for A (like the only remaining points are too far away from the midpoint of A)

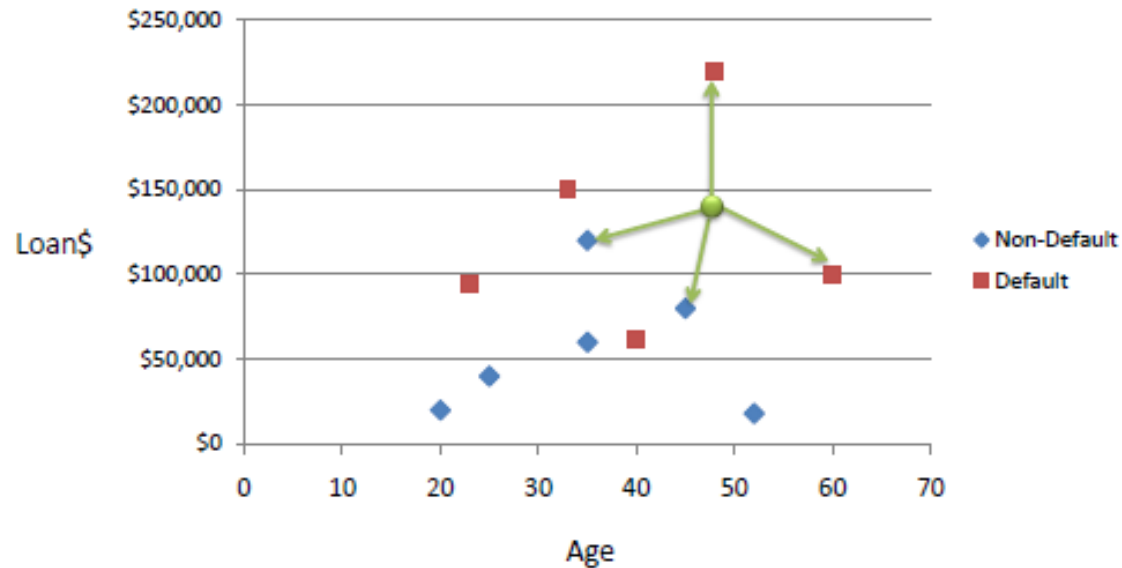- Start again on the remaining points and put them in B.

- Etc.

# More on Clustering



Different cluster analysis results on "mouse" data set:

Original Data    k-Means Clustering    EM Clustering

- Different "metrics" result in different results
- Assumes some kind of "separability in the data"
- Can't label the data very easily (the "cluster labeling problem")

# Supervised Learning
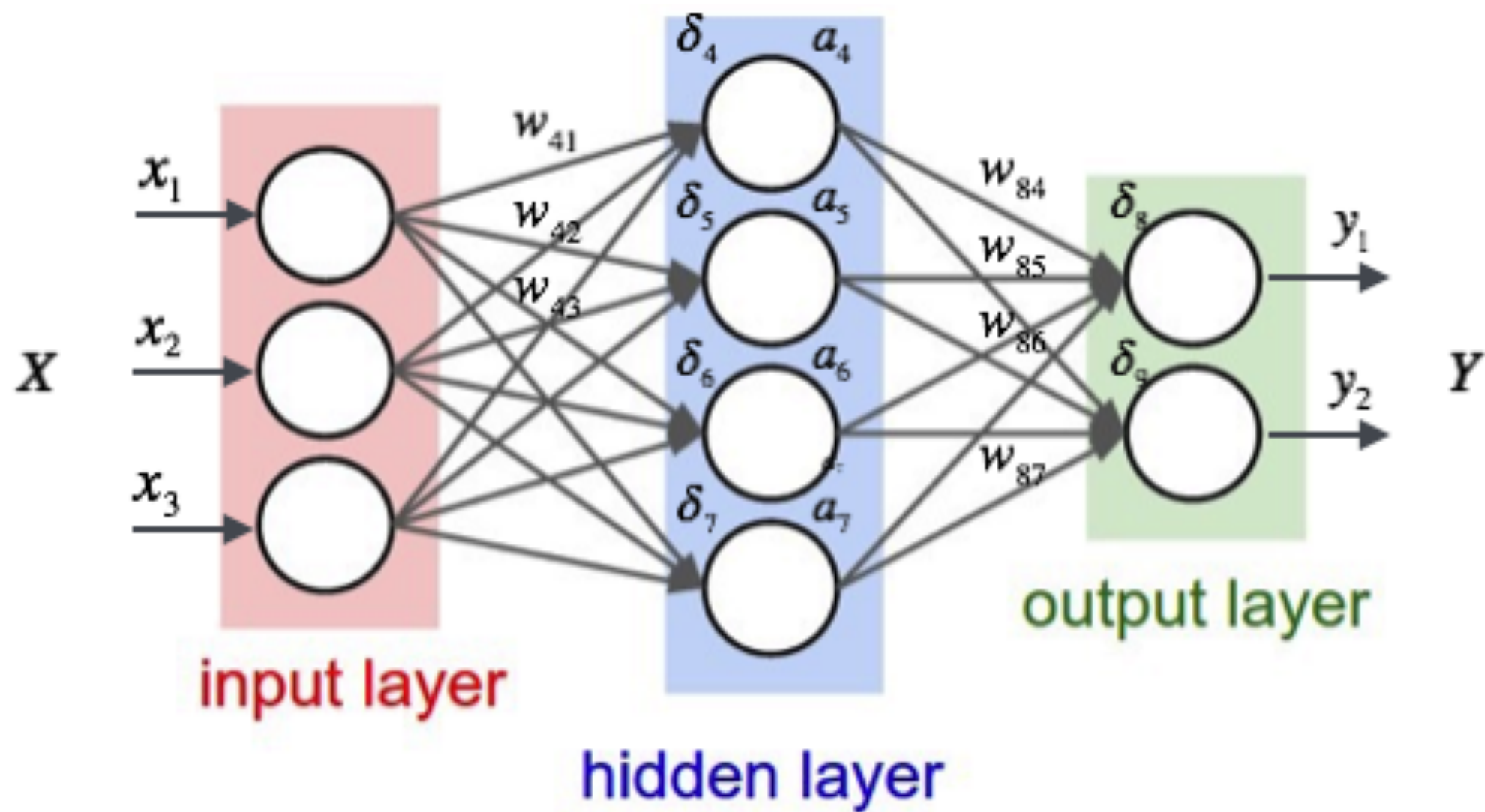
- Example: K Nearest Neighbor algorithm

## 3-KNN: Example(1)

| Customer | Age | Income | No. credit cards | Class |
|----------|-----|--------|------------------|-------|
| George | 35 | 35K | 3 | No |
| Rachel | 22 | 50K | 2 | Yes |
| Steve | 63 | 200K | 1 | No |
| Tom | 59 | 170K | 1 | No |
| Anne | 25 | 40K | 4 | Yes |
| John | 37 | 50K | 2 | **YES** |

**Distance from John**

$\sqrt{[(35-37)^2+(35-50)^2+(3-2)^2]}=15.16$

$\sqrt{[(22-37)^2+(50-50)^2+(2-2)^2]}=15$

$\sqrt{[(63-37)^2+(200-50)^2+(1-2)^2]}=152.23$

$\sqrt{[(59-37)^2+(170-50)^2+(1-2)^2]}=122$

$\sqrt{[(25-37)^2+(40-50)^2+(4-2)^2]}=15.74$
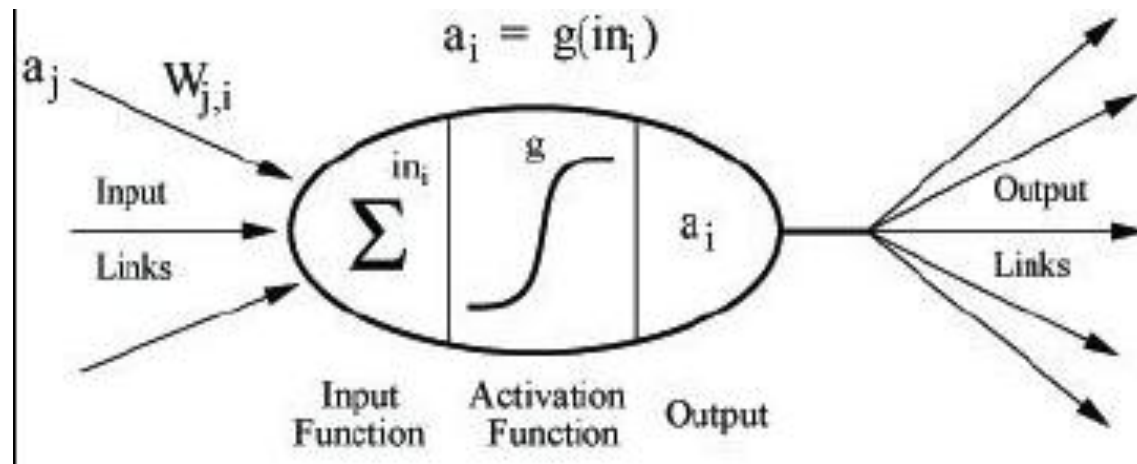


◆ Non-Default
■ Default

# Supervised Learning: Artificial Neural Networks (ANNs)

- Loosely based on early models of neurons

- The axon is connected to the dendrite

- When neuron fires, the dendrites transmit a weighted signal to the next neuron

# Training vs. Classification in ANNs

- Training phase:
  - Initialize network weights to random real values [0,1]
  - Present data vectors to input layer (suitably encode data values)
  - Multiply each value by the weights and apply some kind of non-linear (or semi-linear) firing threshold (activation function) for hidden layer nodes
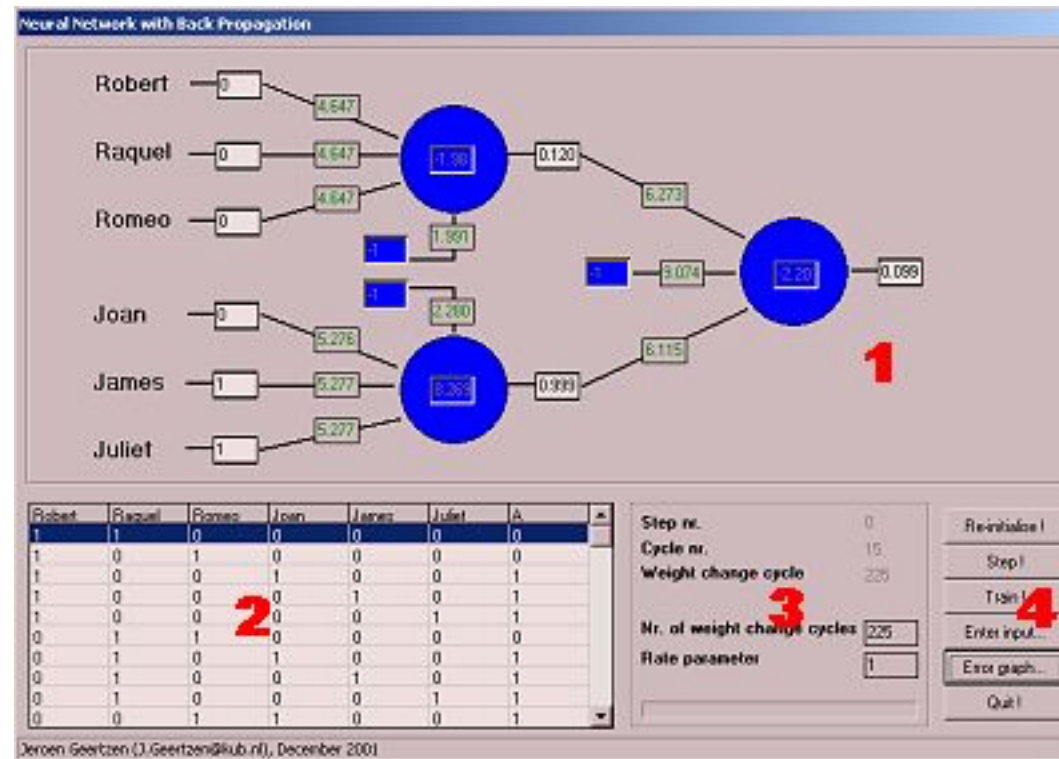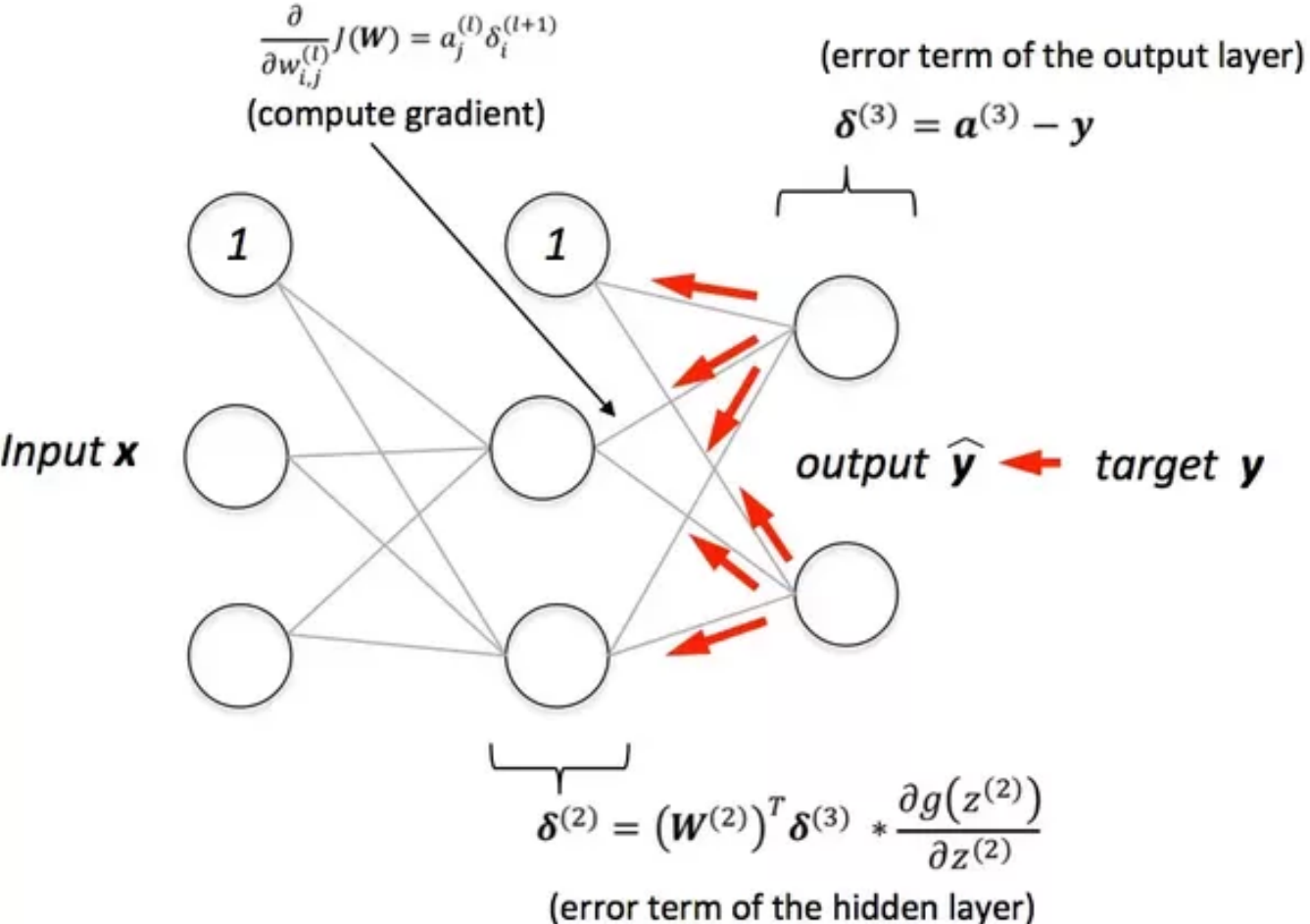
# Training vs. Classification

- Carry activation through to output layer
- Calculate the error at the output layer (the difference between the output layer activation value that should be for the input, and the actual output)
- Propagate the error back by adjusting the weight values a little bit
- Do it again through multiple random training presentations through multiple training "epochs" until the error becomes very low
- This is "backpropogation" training
- Then, you can classify new data by presenting the new data to the network and reading the output layer classes

# Backpropagation Example

• https://nlp-ml.io/jg/software/ann/

# More on Backprop

$$\frac{\partial}{\partial w_{i,j}^{(l)}} J(W) = a_j^{(l)} \delta_i^{(l+1)}$$

(compute gradient)

(error term of the output layer)

$$\delta^{(3)} = a^{(3)} - y$$

Input **x**

① ①

output $\widehat{y}$ ⬅ target **y**

$$\delta^{(2)} = \left(W^{(2)}\right)^T \delta^{(3)} * \frac{\partial g\left(z^{(2)}\right)}{\partial z^{(2)}}$$

(error term of the hidden layer)

# Problems with ANNs

- Like all machine learning, they train on whatever distinguishes the data. They are too literal. Classic tank training problem (probably apocryphal).

- Fully connected ANNs are too big: 1024x728 pixel HD frame = 745472 pixels. One input node per pixel and as many hidden layer nodes equals 745472^2 weights that need calculation (forward pass) and updating (backpropagation), just on first layer! Just too much.

- This is the "curse of dimensionality" and it applies to images or sentences (too many possible words in natural language).

- Regular ANNs don't handle sequential data very well (video scenes, sentences, stock market price fluctuations, etc.) and "recurrent" neural networks that connect output back to input are too hard to train.

# Deep Learning Neural Networks

- Bigger, better, more interesting

- Run them faster on GPUs or custom chips

- Specific architectural adjustments for handling images, sequential data, and external memory stores
  - Convolutional Neural Networks (CNN) for images
  - Treating network like an autoencoder first for images
  - Long Short-Term Memory (LSTM) for sequential data
  - Restricted Boltzman Machines (RBMs) as a variant on backpropagation
  - External memory models for things like Wikipedia-size ontologies

# Convolutional Neural Networks

- For image analysis and scene classification

- Perform different 2D discrete convolutions on images

- Pool together the results

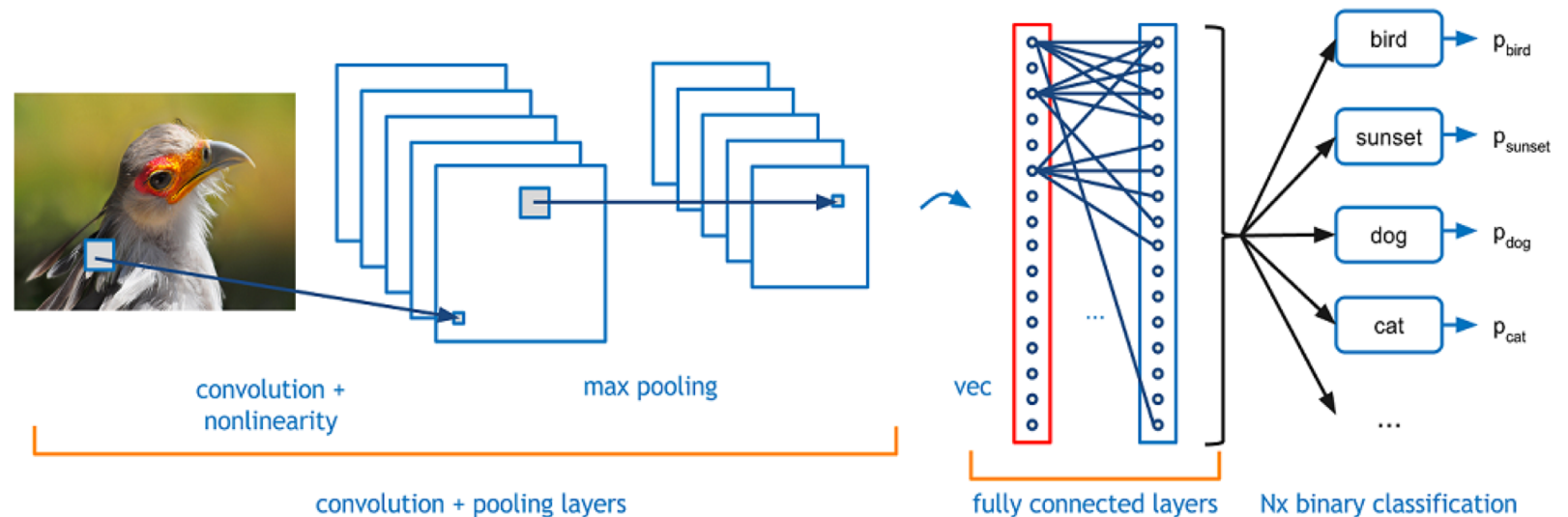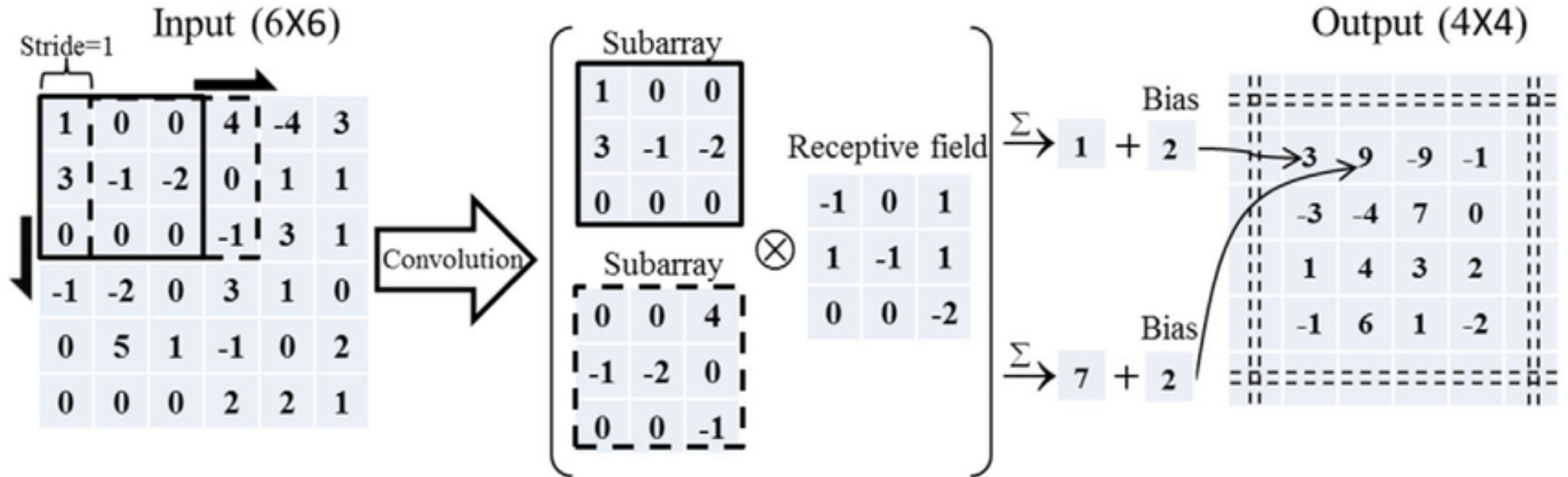- Then feed them into a backprop neural network as the input and train, train, train



convolution + nonlinearity    max pooling    vec

convolution + pooling layers    fully connected layers    Nx binary classification

bird → $p_{bird}$
sunset → $p_{sunset}$
dog → $p_{dog}$
cat → $p_{cat}$

# Image Convolution



Output size = (I − R) / S + 1,   I = Input size,   R = Receptive field size,   S = Stride size;   (6-3)/1+1=4

# Different Kernels Do Different Things



Original

Emboss

Input image

Kernel

Feature map

$$\begin{pmatrix} \dfrac{1}{16} & \dfrac{1}{8} & \dfrac{1}{16} \\[6pt] \dfrac{1}{8} & \dfrac{1}{4} & \dfrac{1}{8} \\[6pt] \dfrac{1}{16} & \dfrac{1}{8} & \dfrac{1}{16} \end{pmatrix}$$

# Some Design and Experimental Terminology

- Training/Test data splits: often we split data into ten parts (randomly), train on 9 parts and evaluate on the last part which is called the "held out data."

- If the hidden layer is too large, a single node can become assigned to each output class resulting in "overtraining" meaning that the network is not properly generalizing. By using a smaller hidden layer we "bottleneck" the system, forcing generalization.

# Next Lecture

- TensorFlow practicalities
  - Download
  - Install
  - Run
  - Understanding the code base
  - APIs and language bindings
- Programming GPUs
  - Memory models
  - A bit about CUDA