

95-865 Australia Lecture 5: Introduction to Predictive Data Analytics, Neural Nets, and Deep Learning

George Chen

Announcements

- No the quiz hasn't been graded yet
- Please make sure you have AWS set up with AWS Educate credits (bug Erick)
- Python 3.7 currently has some compatibility issues with Keras and Tensorflow — please downgrade your Python to version 3.6 if you're using 3.7!!!

conda install python=3.6
conda install keras

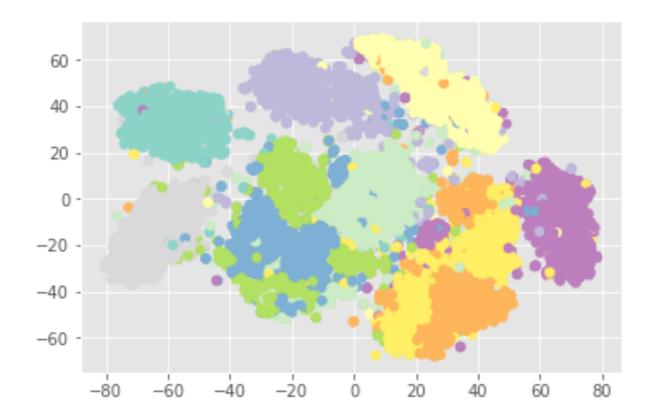
Disclaimer: unfortunately "k" means many things

Previous Lecture: Topic Modeling

- There are actually *many* topic models, not just LDA
 - Correlated topic models, Pachinko allocation, biterm topic models, anchor word topic models, ...

- Dynamic topic models: tracks how topics change over time
 - This sort of idea could be used to figure out how user tastes change over time in a recommendation system
 - Could try to see if there are existing patterns for how certain topics become really popular

What if we have labels?



Example: MNIST handwritten digits have known labels

If the labels are known...

If the labels are known...

And we assume data generated by GMM...

What are the model parameters?

k = # of colors

We can directly estimate cluster means, covariances

Flashback: Learning a GMM

Don't need this top part if we know the labels!

Step 1: Pick guesses for cluster means and covariances

Repeat until convergence:

Step 0: Pick k

Step 2: Compute probability of each point belonging to each of the k elusters

Step 3: Update **cluster means and covariances** carefully accounting for probabilities of each point belonging to each of the clusters

We don't need to repeat until convergence

If the labels are known...

And we assume data generated by GMM...

What are the model parameters?

k = # of colors

We can directly estimate cluster means, covariances

What should the label of this new point be? Whichever cluster has higher probability! Decision boundary

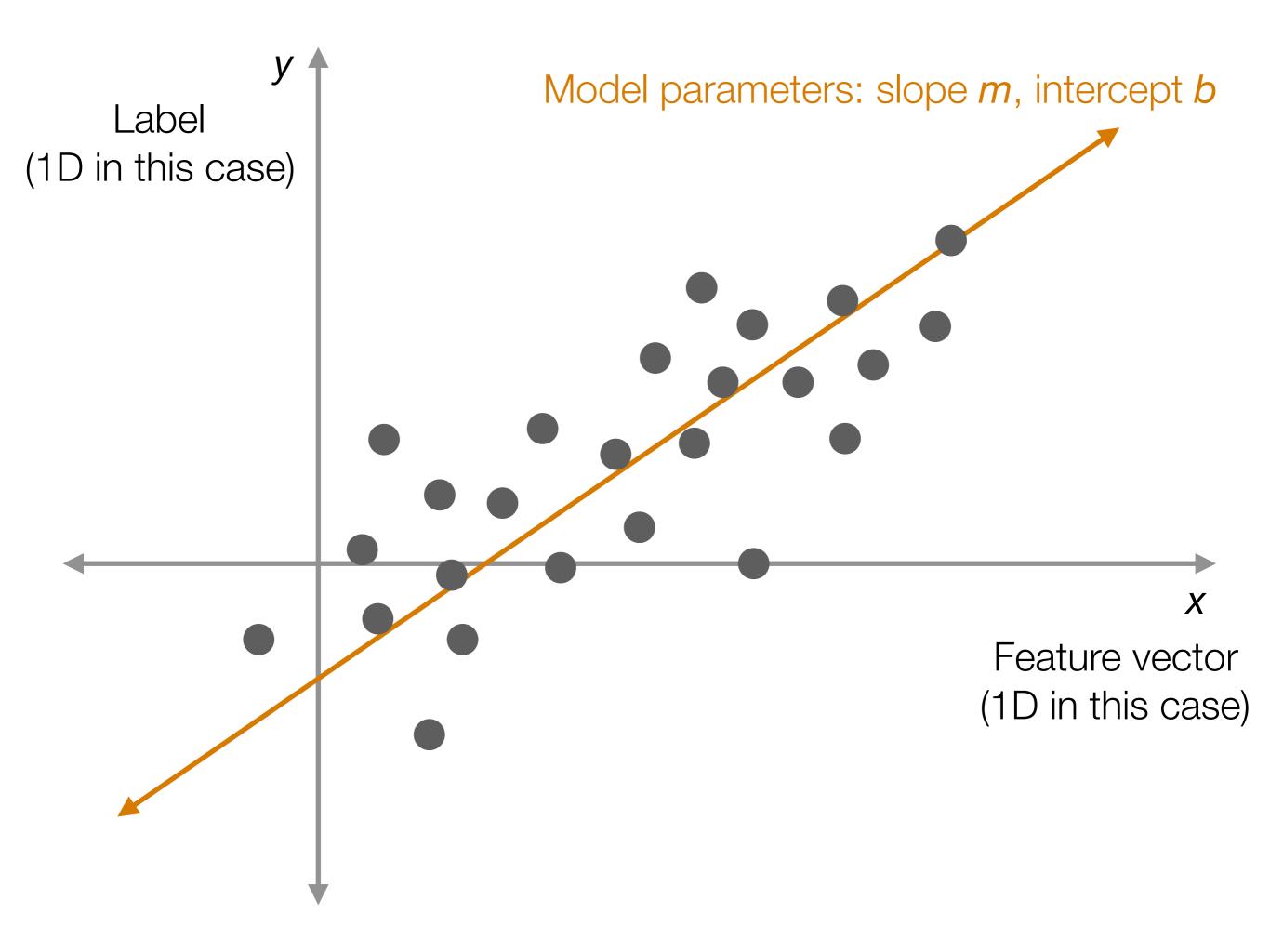
We just created a **classifier** (a procedure that given a new data point tells us what "class" it belongs to)

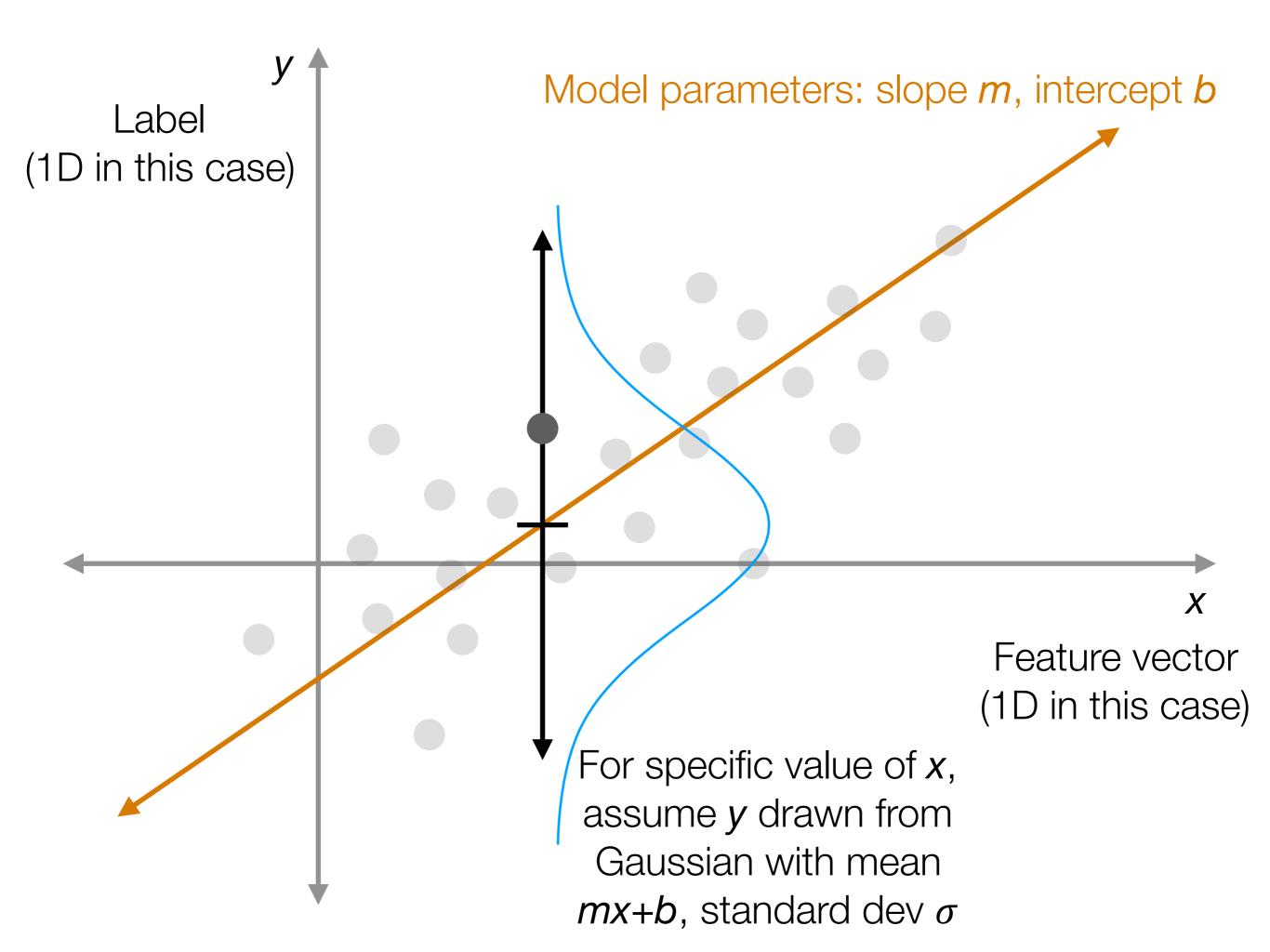
This classifier we've created assumes a *generative model*

What should the label of this new point be? Whichever cluster has higher probability!

You've seen generative models before for prediction

Linear regression!





Predictive Data Analysis

Training data

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

Goal: Given new feature vector x, predict label y

- *y* is discrete (such as colors red and blue)
 → prediction method is called a classifier
- *y* is continuous (such as a real number)
 → prediction method is called a regressor
- A giant zoo of methods

Generative Models

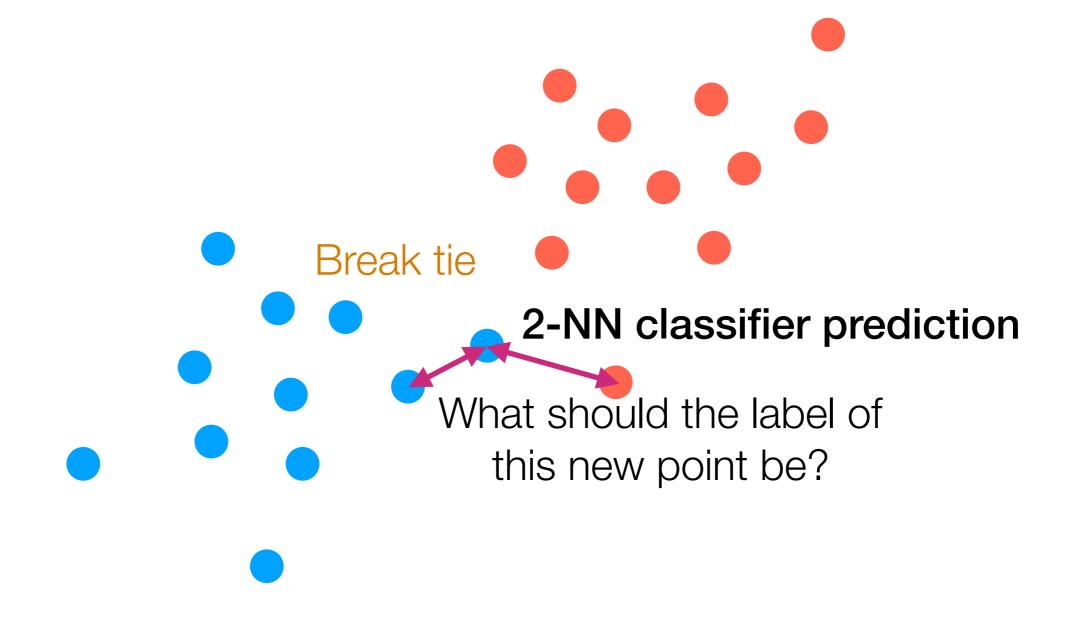
- Hypothesize a specific way in which data are generated
- After learning a generative model:
 - We can generate new synthetic data from the model
 - Usually generative models are probabilistic and we can evaluate probabilities for a new data point
- In contrast to generative models, there are *discriminative* methods that just care about learning a prediction rule

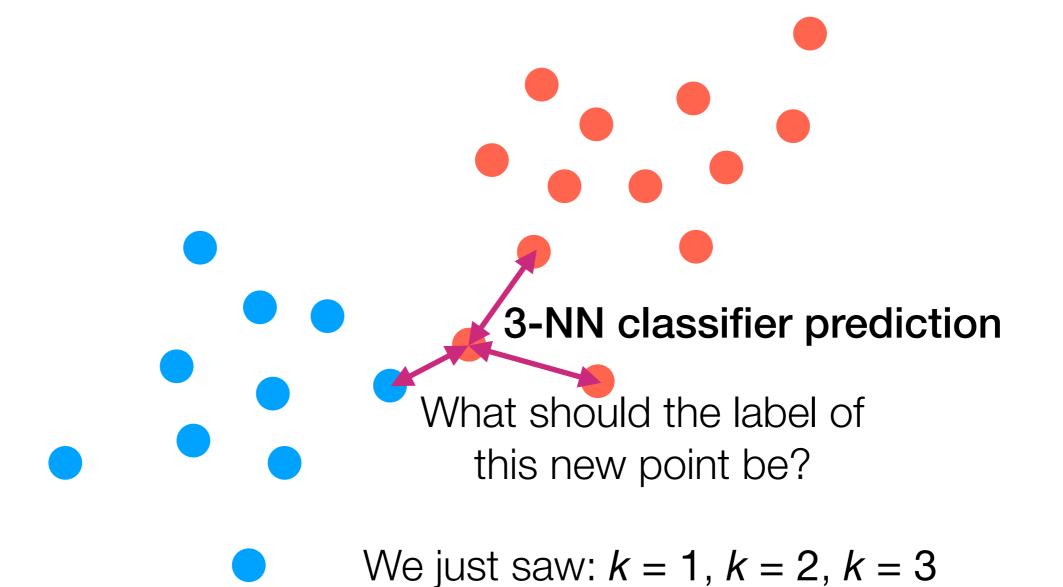
Example of a Discriminative Method: *k*-NN Classification

What should the label of this new point be?

1-NN classifier prediction

What should the label of this new point be?





What happens if k = n?

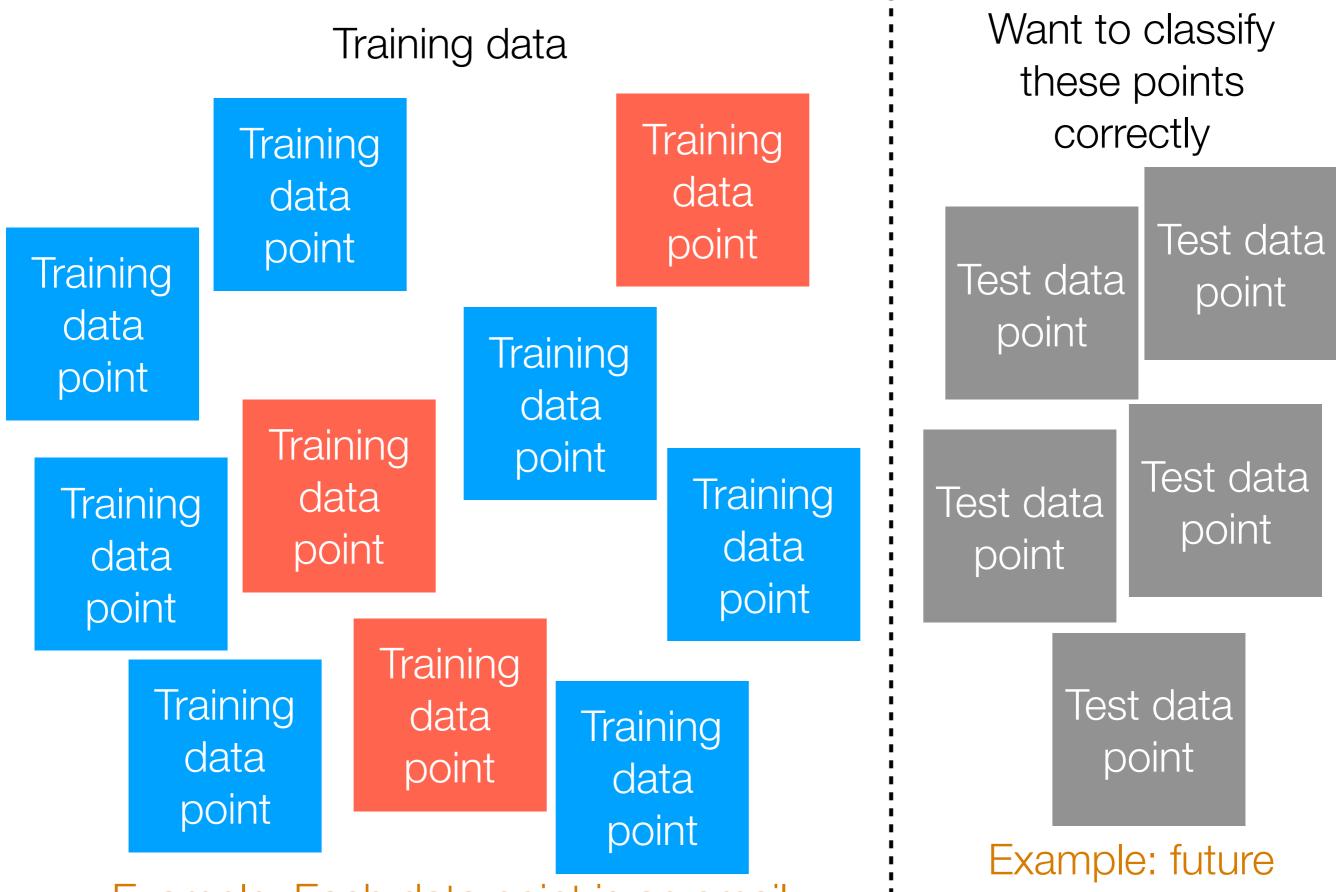
How do we choose k?

What I'll describe next can be used to select hyperparameter(s) for any prediction method

First: How do we assess how good a prediction method is?

Hyperparameters vs. Parameters

- We fit a model's parameter to training data (terminology: we "learn" the parameters)
- We pick values of hyperparameters and they do not get fit to training data
- Example: Gaussian mixture model
 - Hyperparameter: number of clusters *k*
 - Parameters: cluster probabilities, means, covariances
- Example: *k*-NN classification
 - Hyperparameter: number of nearest neighbors *k*
 - Parameters: N/A



Example: Each data point is an email and we know whether it is spam/ham

Example: future emails to classify as spam/ham

Predicted labels

Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point
Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point

Train method on data in gray

Predict on data in orange

Compute prediction error

50%

Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point
Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point

Train method on data in gray

Predict on data in orange

Compute prediction error

0% 50%

Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point
Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point

Train method on data in gray

Predict on data in orange

Compute prediction error

50% 0% 50%

Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point
Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point

Train method on data in grayPredict on data in orangeComputeCompute0%50%0%50%

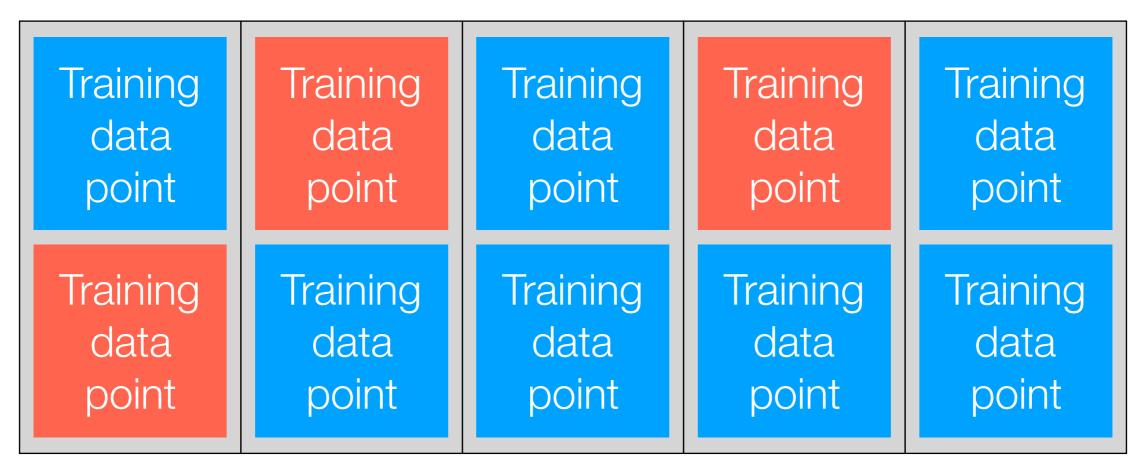
Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point
Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point

Train method on data in grayPredict on data
in orangeCompute
prediction errorCompute
prediction error0%0%50%Average error: (0+0+50+0+50)/5 = 20%

Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point
Training	Training	Training	Training	Training
data	data	data	data	data
point	point	point	point	point

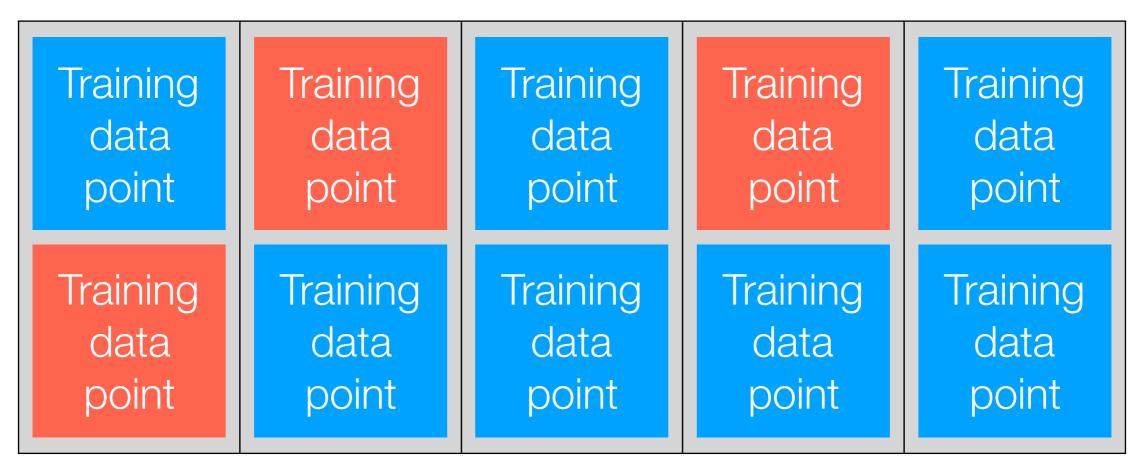
- 1. Shuffle data and put them into "folds" (5 folds in this example)
- 2. For each fold (which consists of its own train/validation sets):(a) Train on fold's training data, test on fold's validation data(b) Compute prediction error
- 3. Compute average prediction error across the folds

not the same *k* as in *k*-means or *k*-NN classification *k*-fold Cross Validation



- 1. Shuffle data and put them into "folds" (k=5 folds in this example)
- 2. For each fold (which consists of its own train/validation sets):(a) Train on fold's training data, test on fold's validation data(b) Compute prediction error
- 3. Compute average prediction error across the folds

not the same *k* as in *k*-means or *k*-NN classification *k*-fold Cross Validation



- 1. Shuffle data and put them into "folds" (k=5 folds in this example)
- 2. For each fold (which consists of its own train/validation sets):
 (a) Train on fold's training data, test on fold's validation data
 (b) Compute some sort of prediction score
- 3. Compute **average prediction score** across the folds "cross validation score"

Choosing k in k-NN Classification

Note: *k*-NN classifier has a single hyperparameter *k*

For each k = 1, 2, 3, ..., the maximum k you are willing to try:

Compute 5-fold cross validation score using *k*-NN classifier as prediction method

Use whichever k has the best cross validation score

Automatic Hyperparameter Selection

Suppose the prediction algorithm you're using has hyperparameters $\boldsymbol{\theta}$

For each hyperparameter setting θ you are willing to try:

Compute 5-fold cross validation score using your algorithm with hyperparameters θ

Use whichever θ has the best cross validation score Why 5?

People have found using 10 folds or 5 folds to work well in practice but it's just empirical — there's no deep reason

Training data

Training data

Training data

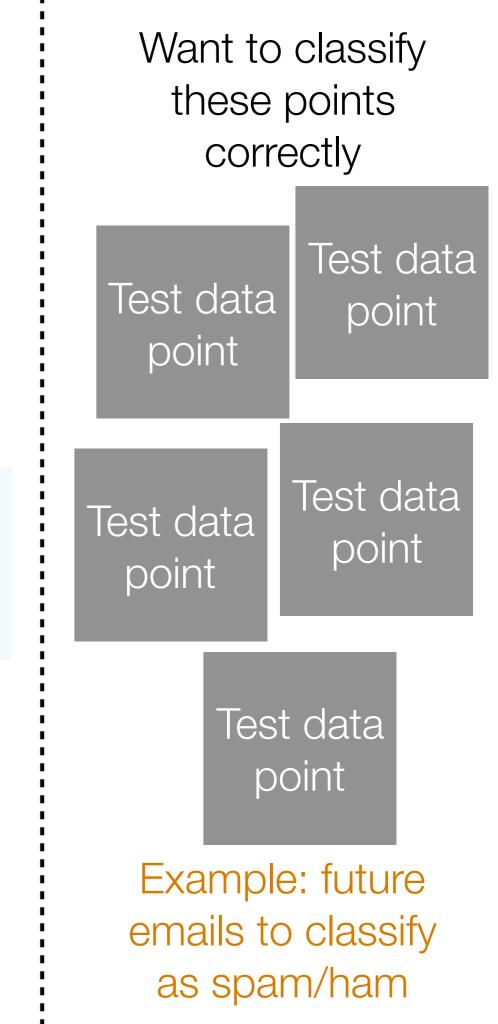
Important: the cross validation score is trying to predict what the prediction quality will be on the unseen test data

Our earlie<mark>r example</mark> had a cross validation score of 20% error

This is a guess at how well the prediction method should perform on test data

This guess is <u>not</u> always accurate

Example: Each data point is an email and we know whether it is spam/ham



Different Ways to Measure Accuracy

Simplest way:

 Raw error rate: fraction of predicted labels that are wrong (this was in our cross validation example earlier)

In "binary" classification (there are 2 labels such as spam/ham) when 1 label is considered "positive" and the other "negative":

- **Precision:** among data points predicted to be "positive", what fraction of these predictions is correct?
- **Recall:** among data points that are actually "positive", what fraction of these points is predicted correctly as "positive"? (also called **true positive rate**)
- **F1 score:** 2 × precision × recall

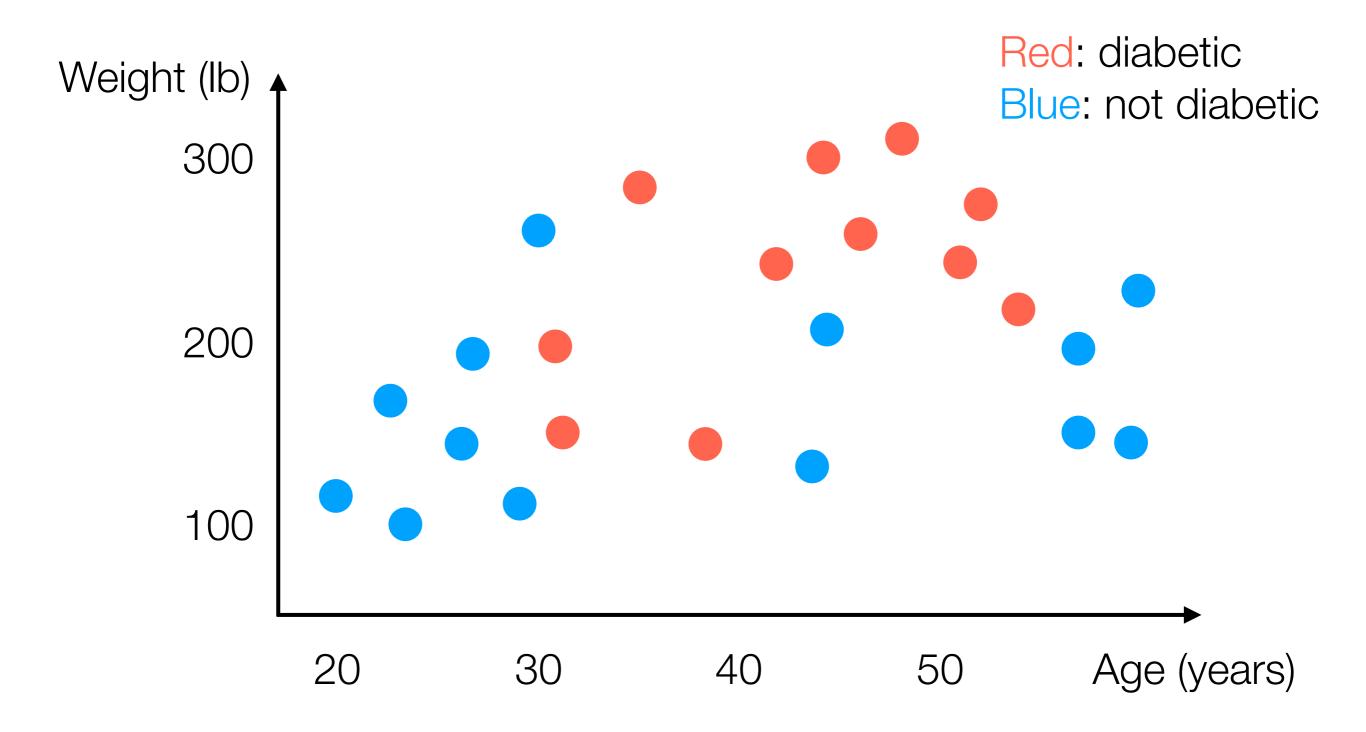
precision + recall

Prediction and Model Validation

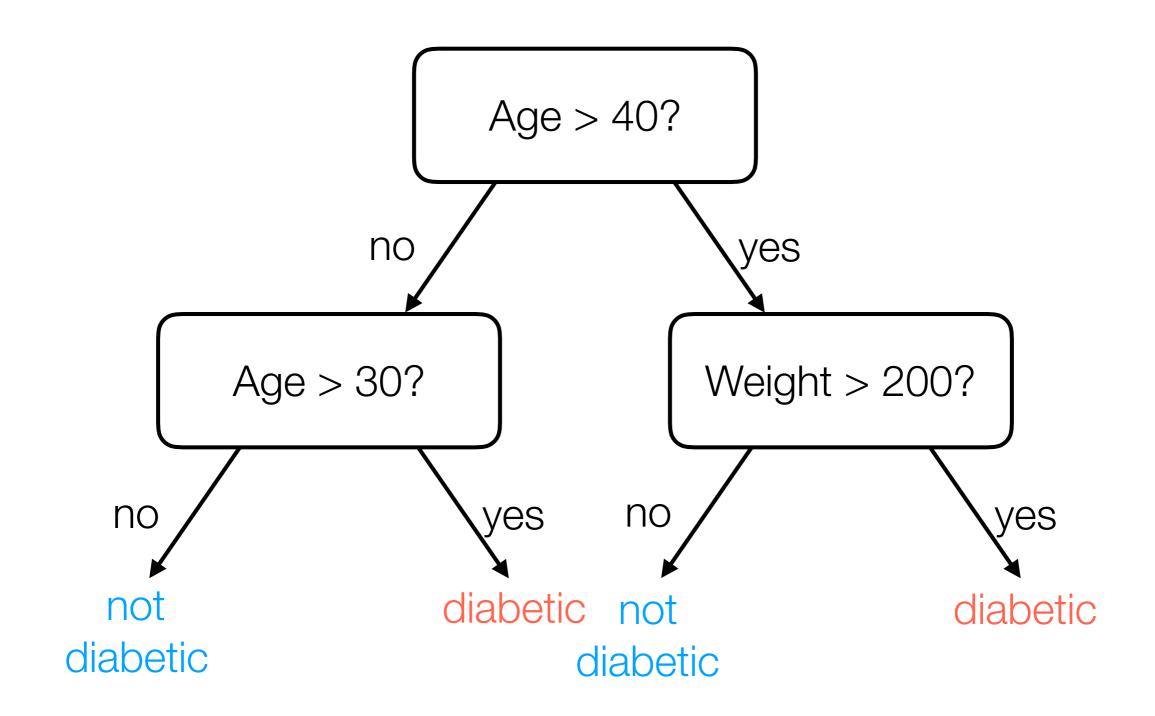
Demo

Decision Trees

Example Made-Up Data



Example Decision Tree

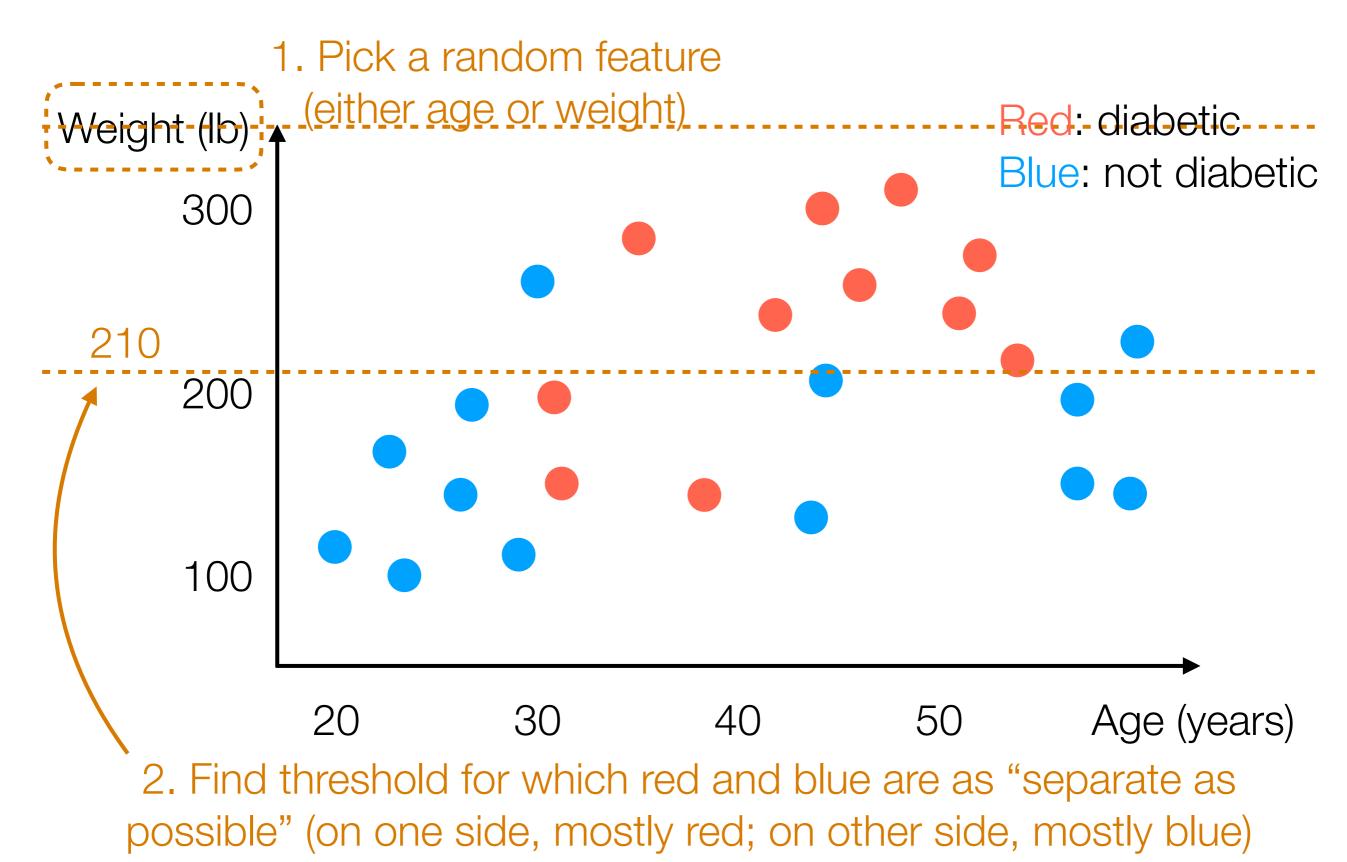


Learning a Decision Tree

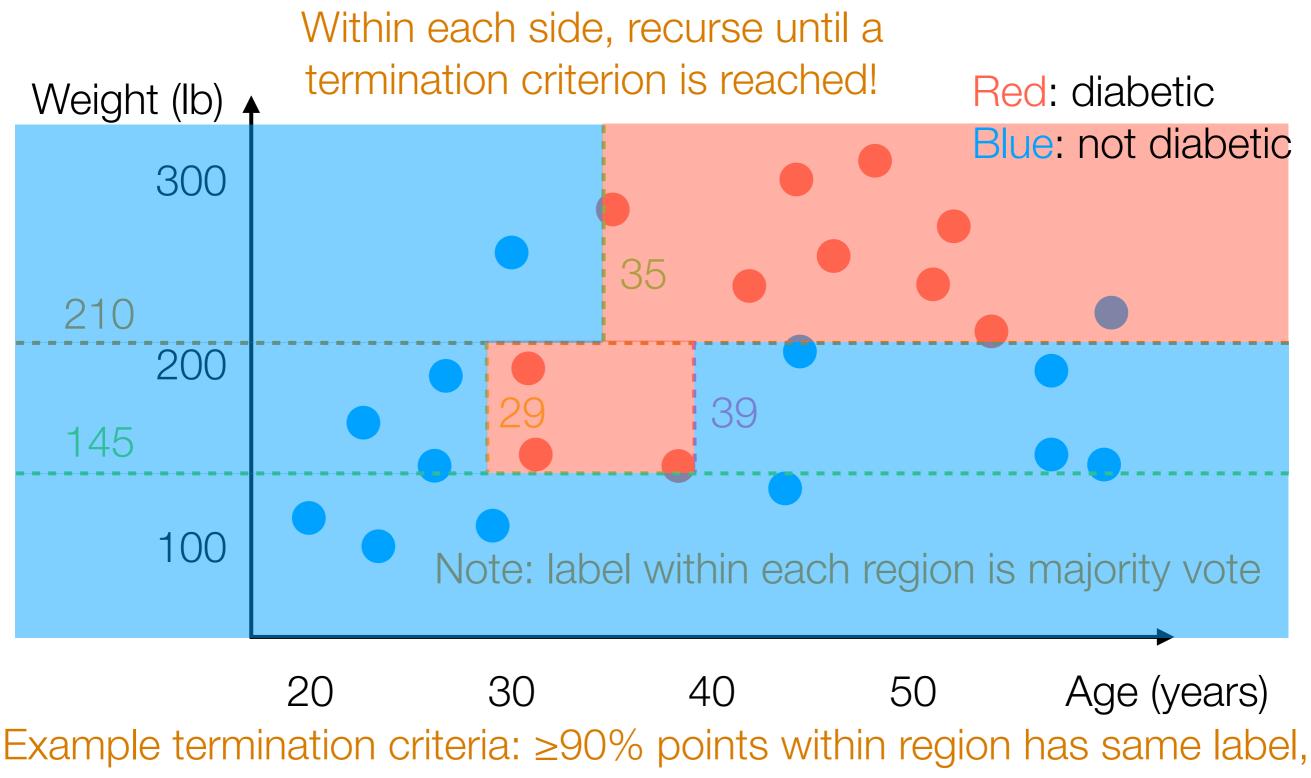
• Many ways: general approach actually looks a lot like divisive clustering *but accounts for label information*

• I'll show one way (that nobody actually uses in practice) but it's easy to explain

Learning a Decision Tree

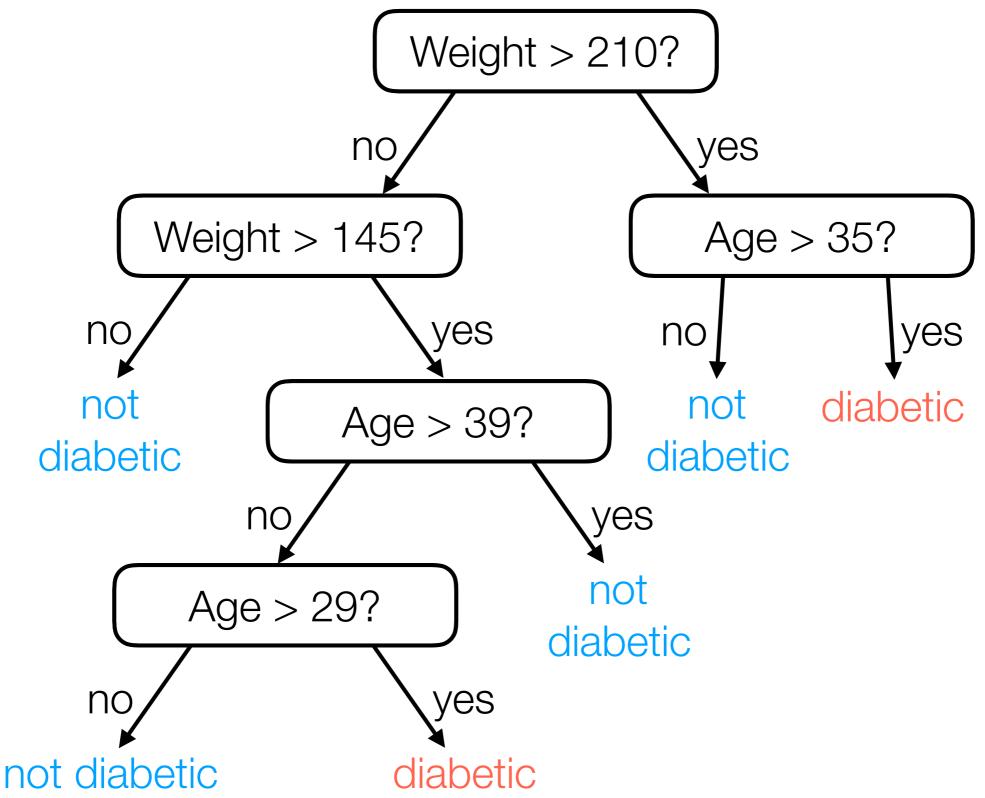


Learning a Decision Tree



number of points within region is <5

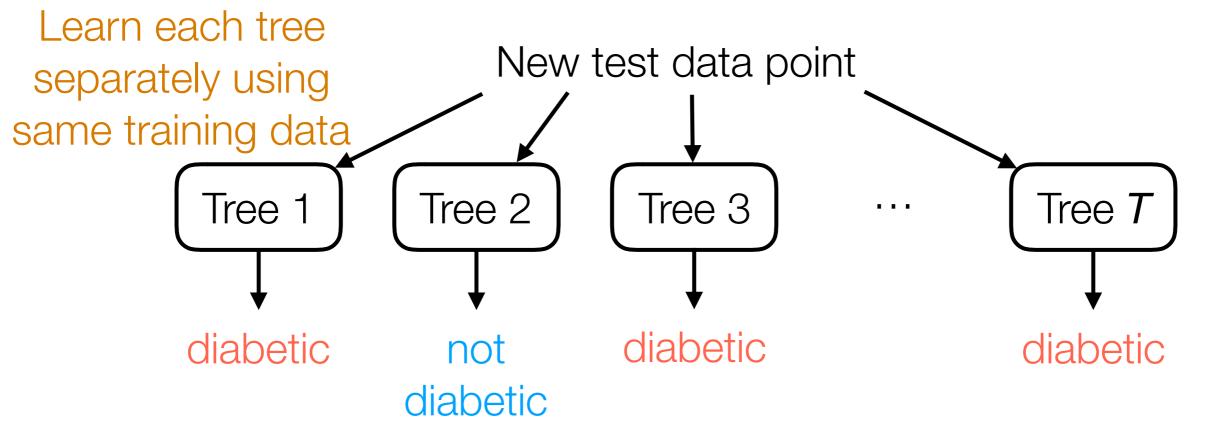
Decision Tree Learned



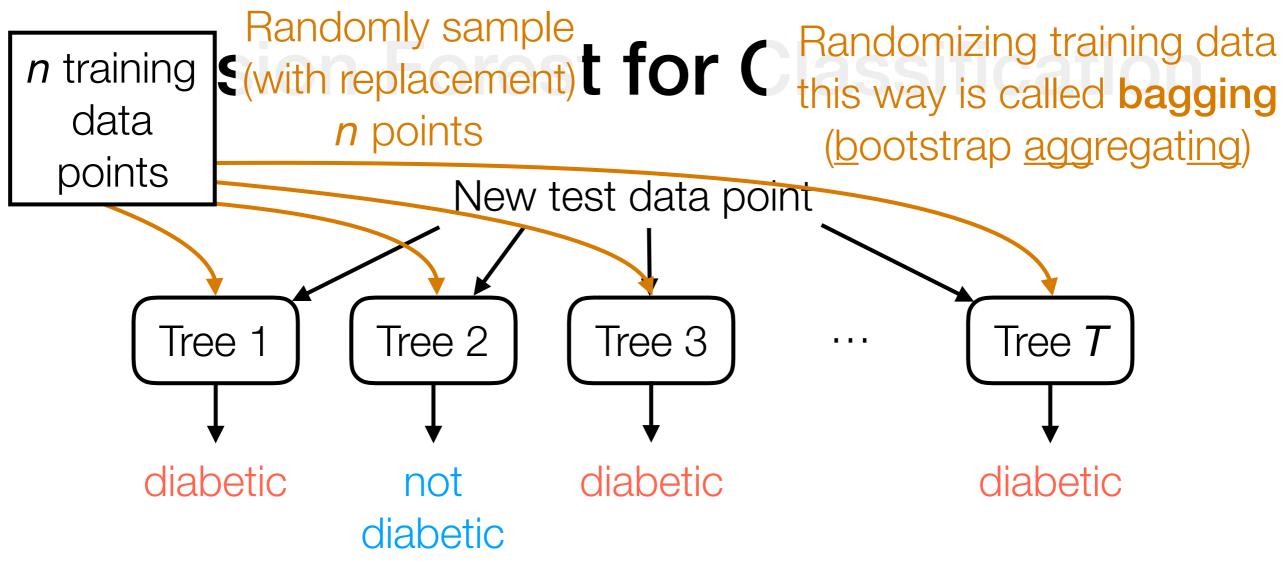
For a new person with feature vector (age, weight), easy to predict!

Decision Forest for Classification

- Typically, a decision tree is learned with randomness (e.g., we randomly chose which feature to threshold)
 - → by re-running the same learning procedure, we can get different decision trees that make different predictions!
- For a more stable prediction, use many decision trees



Final prediction: majority vote of the different trees' predictions



Question: What happens if all the trees are the same?

Adding randomness can make trees more different!

- **Random Forest:** in addition to randomly choosing features to threshold, also randomize training data used for each tree
- Extremely randomized trees: further randomize thresholds rather than trying to pick clever thresholds

Back to the demo

Neural Nets and Deep Learning



Over 10 million images, 1000 object classes



2011: Traditional computer vision achieves accuracy ~74%
2012: Initial deep neural network approach accuracy ~84%
2015 onwards: Deep learning achieves accuracy 96%+
Russakovsky et al. ImageNet Large Scale Visual Recognition Challenge. IJCV 2015.

Deep Learning Takeover

Academia:

- Top computer vision conferences (CVPR, ICCV, ECCV) are now nearly all about deep learning
- Top machine learning conferences (ICML, NeurIPS) have
- heavily been taken over by deep learning

Heavily dominated by industry now!

Google

facebook.

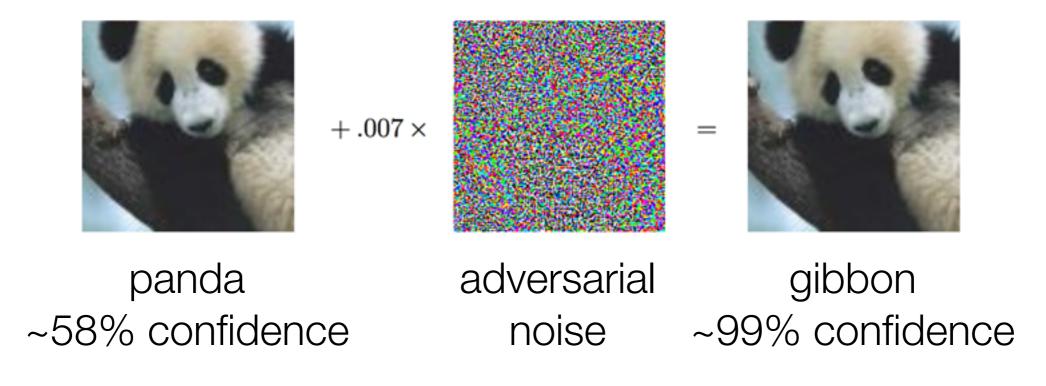
amazon

Extremely useful in practice:

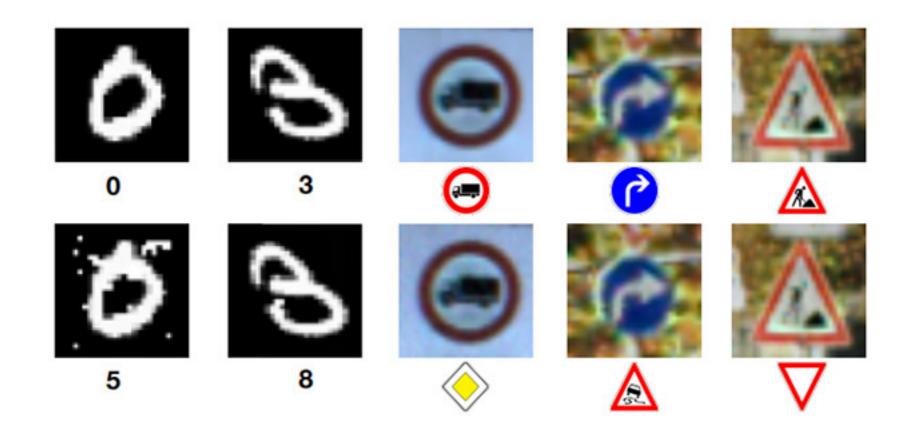
- Near human level image classification (including handwritten digit recognition)
- Near human level speech recognition
- Improvements in machine translation, text-to-speech
- Self-driving cars
- Better than humans at playing Go

Google DeepMind's AlphaGo vs Lee Sedol, 2016

Is it all hype?



Source: Goodfellow, Shlens, and Szegedy. Explaining and Harnessing Adversarial Examples. ICLR 2015.

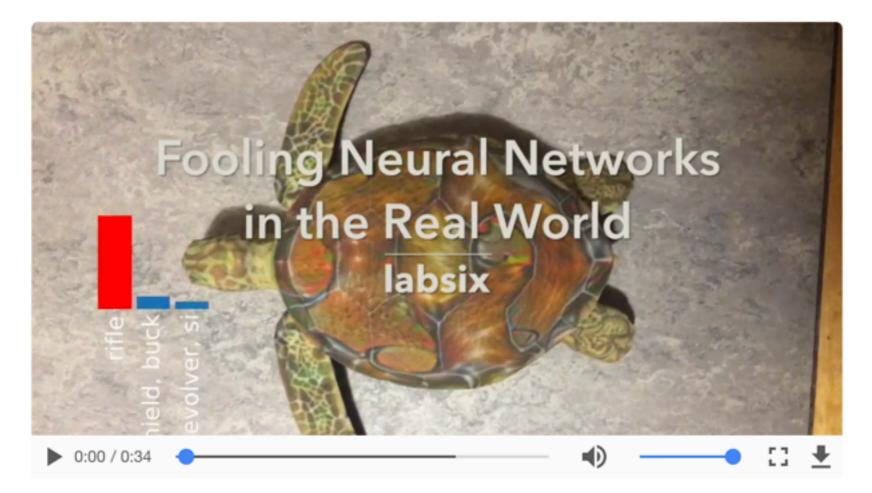


Source: Papernot et al. Practical Black-Box Attacks against Machine Learning. Asia Conference on Computer and Communications Security 2017.

Fooling Neural Networks in the Physical World with 3D Adversarial Objects

31 Oct 2017 · 3 min read — shared on Hacker News, Lobsters, Reddit, Twitter

We've developed an approach to generate *3D adversarial objects* that reliably fool neural networks in the real world, no matter how the objects are looked at.



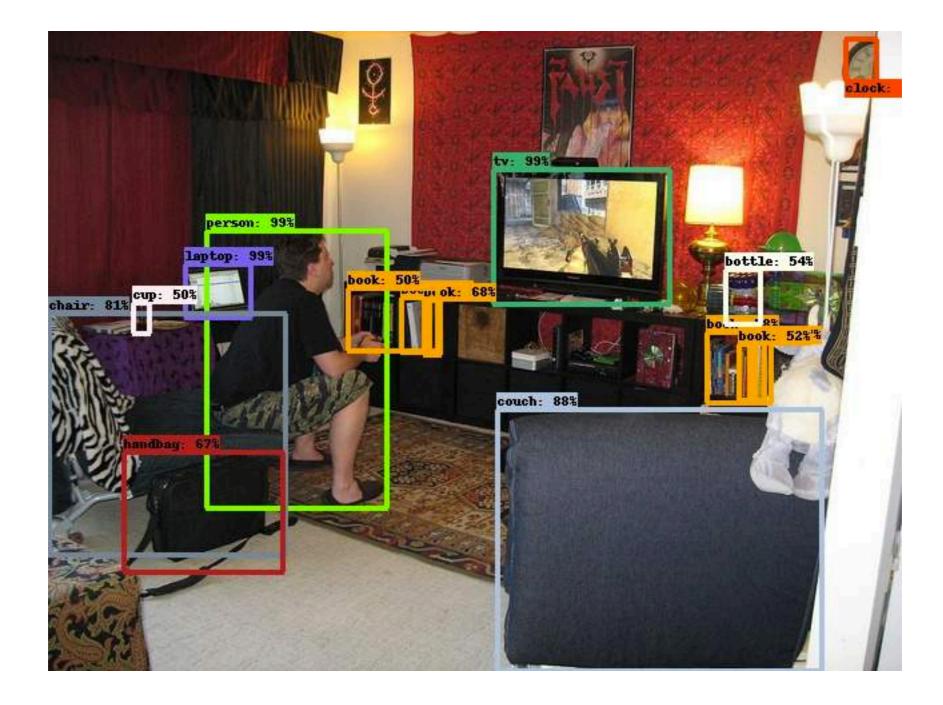
Neural network based classifiers reach near-human performance in many tasks, and they're used in high risk, real world systems. Yet, these same neural networks are particularly vulnerable to *adversarial examples*, carefully perturbed inputs that cause

Source: labsix



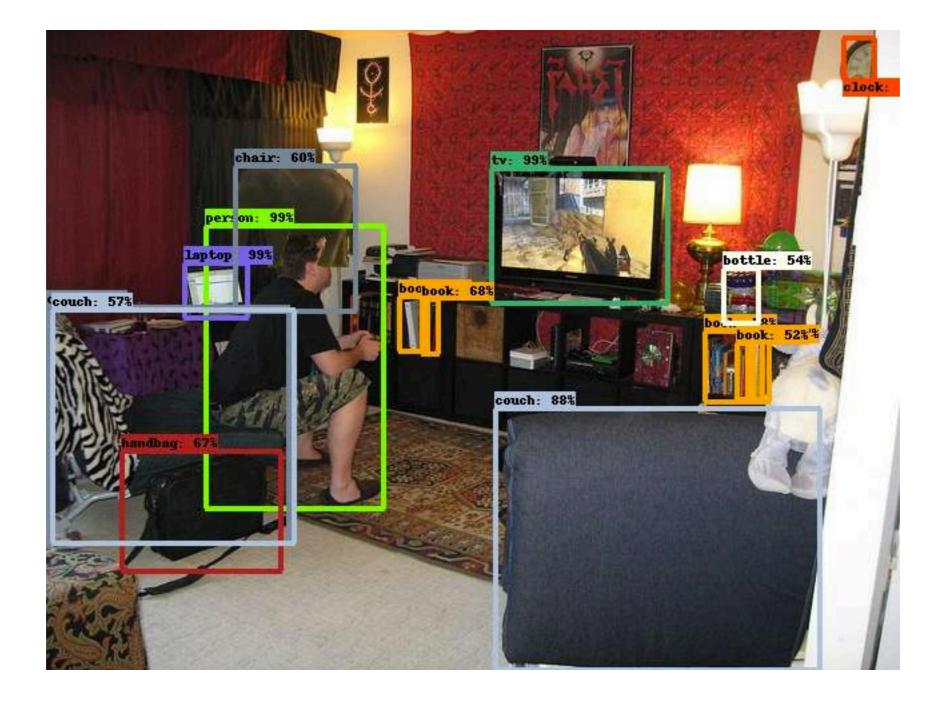
Source: Gizmodo article "This Neural Network's Hilariously Bad Image Descriptions Are Still Advanced AI". September 16, 2015. (They're using the NeuralTalk image-to-caption software.)

Slightly modifying an image results in different prediction results



Source: Quanta Magazine article "Machine Learning Confronts the Elephant in the Room". September 20, 2018.

Slightly modifying an image results in different prediction results



Source: Quanta Magazine article "Machine Learning Confronts the Elephant in the Room". September 20, 2018.

Another AI Winter?

~1970's: First AI winter over symbolic AI

~1980's: Second AI winter over "expert systems"

Every time: Lots of hype, explosion in funding, then bubble bursts

Medium





Michael Jordan Follow

Michael I. Jordan is a Professor in the Department of Electrical Engineering and Computer Sciences and the Department of Statistics at UC Berkeley. Apr 18 - 16 min read

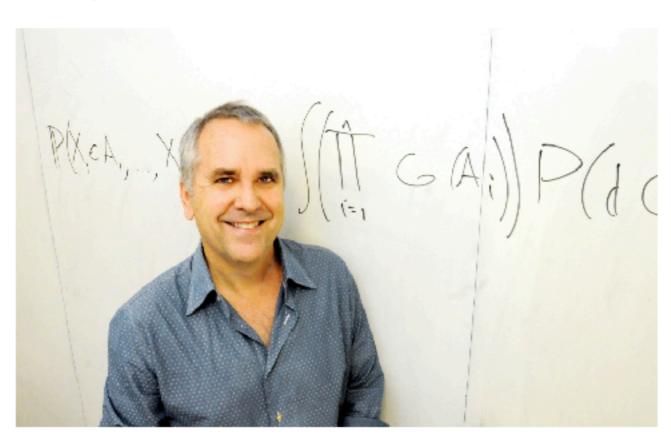


Photo credit: Peg Skorpinski

Artificial Intelligence—The Revolution Hasn't Happened Yet

Artificial Intelligence (AI) is the mantra of the current era. The phrase is intoned by technologists, academicians, journalists and venture capitalists

https://medium.com/@mijordan3/artificial-intelligence-the-revolution-hasnt-happenedyet-5e1d5812e1e7

TECHNOLOGY

How a Pioneer of Machine Learning Became One of Its Sharpest Critics

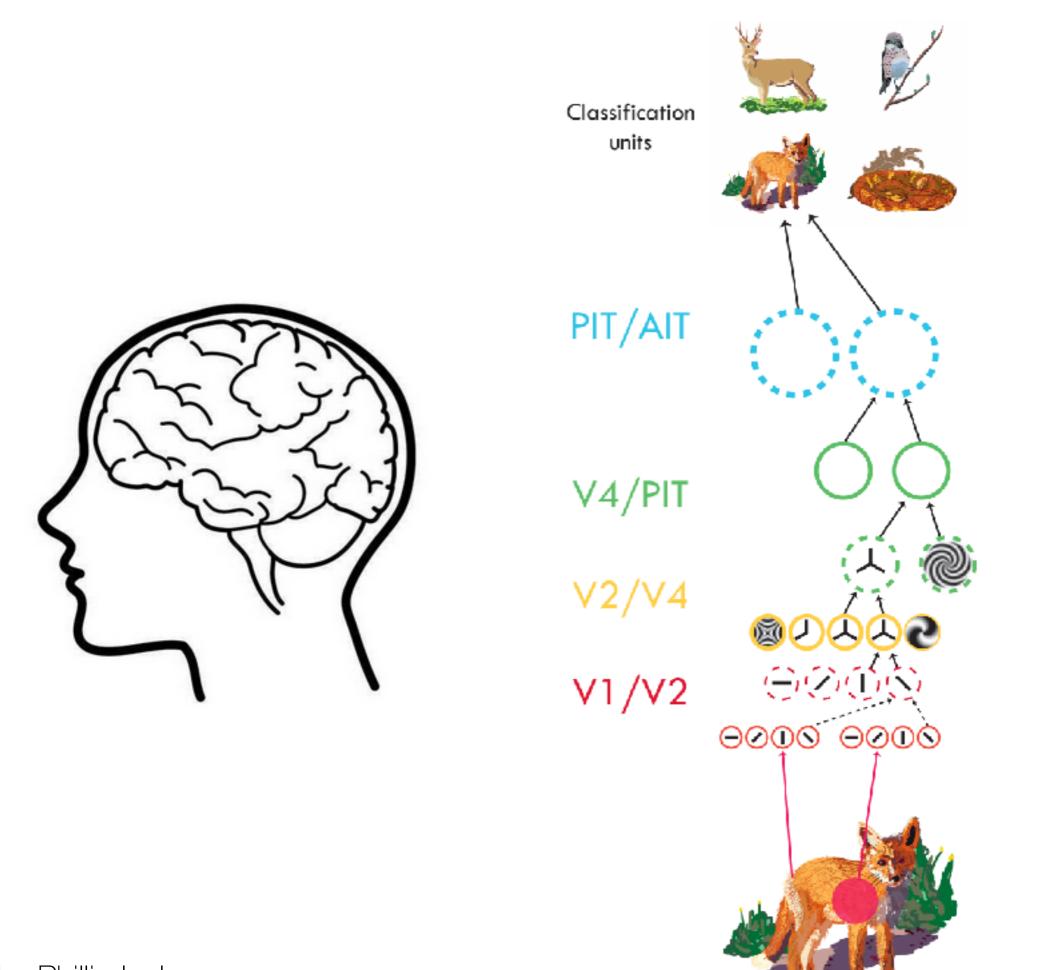
Judea Pearl helped artificial intelligence gain a strong grasp on probability, but laments that it still can't compute cause and effect.

KEVIN HARTNETT AND QUANTA MAY 19, 2018



https://www.theatlantic.com/technology/archive/2018/05/machine-learning-is-stuck-on-asking-why/ 560675/?single_page=true

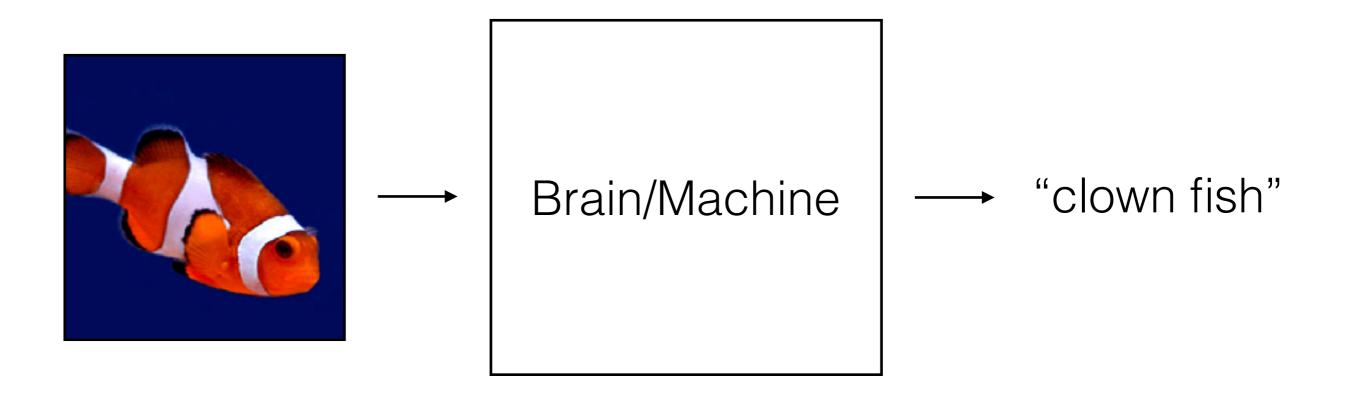
What is deep learning?



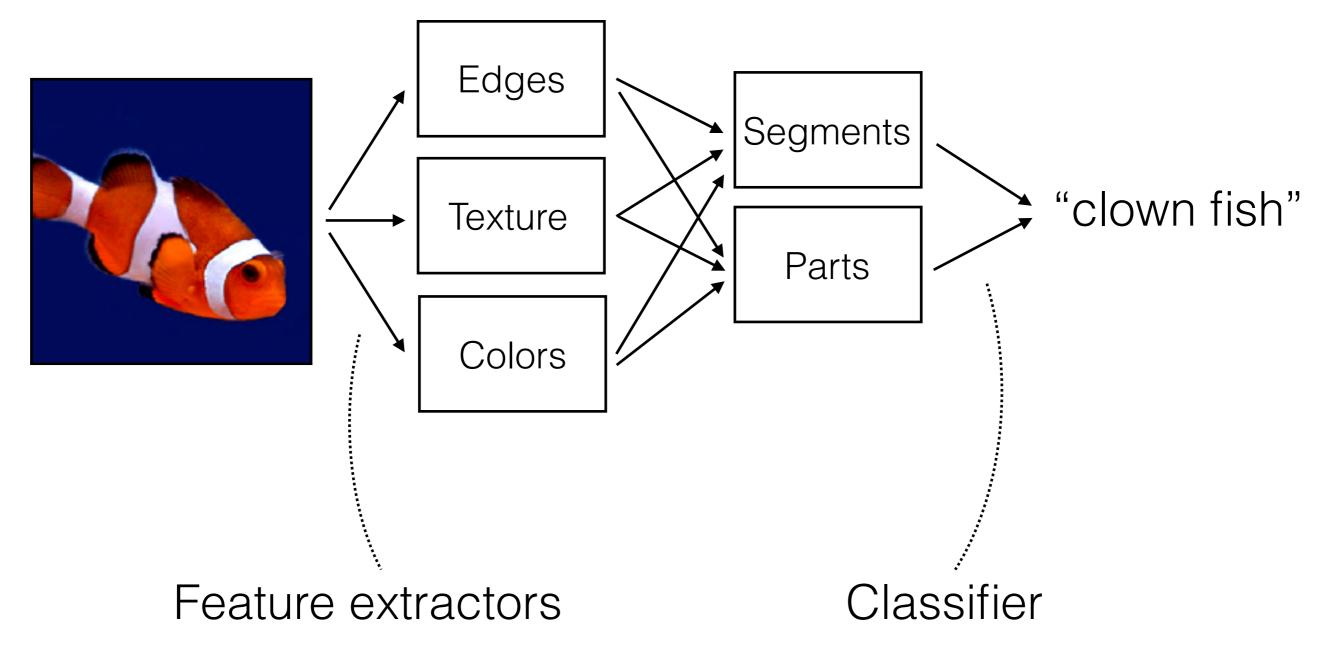
Slide by Phillip Isola

Serre, 2014

Basic Idea

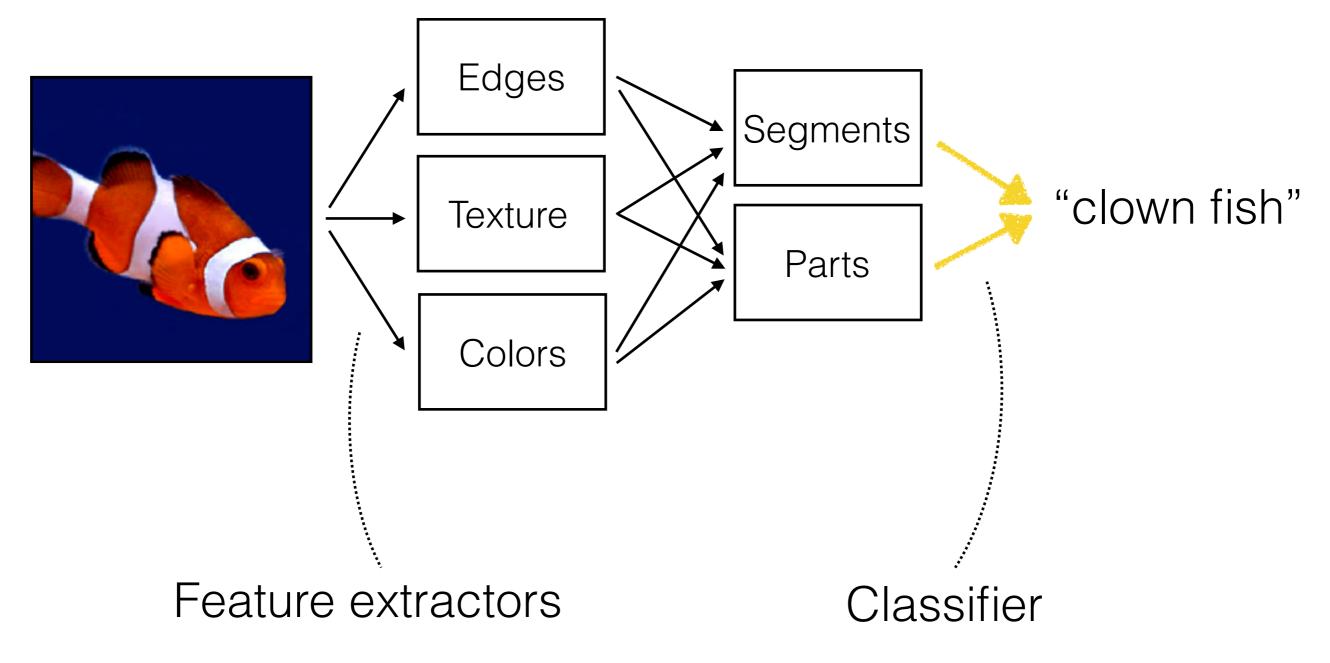


Object Recognition



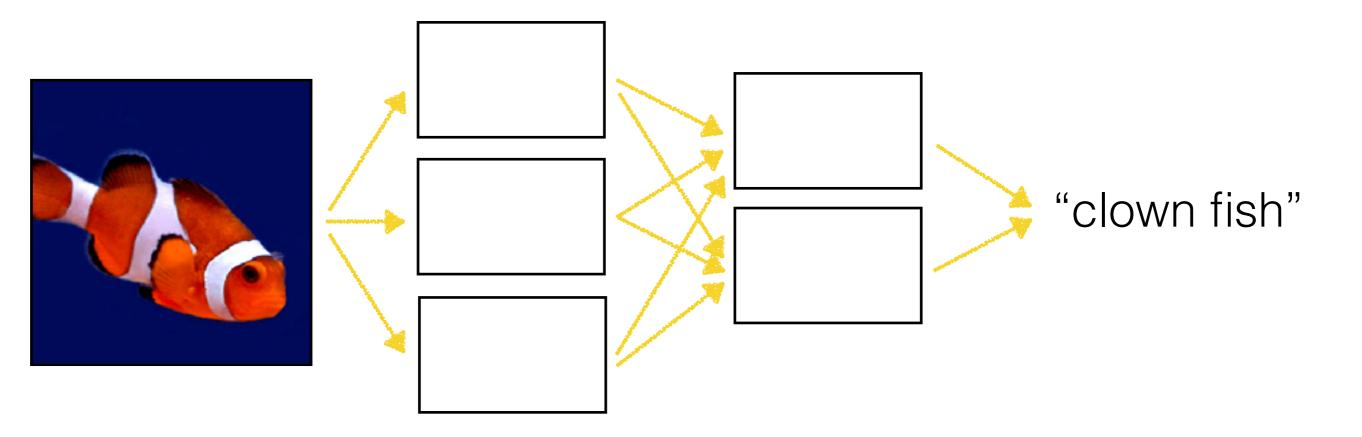
Object Recognition





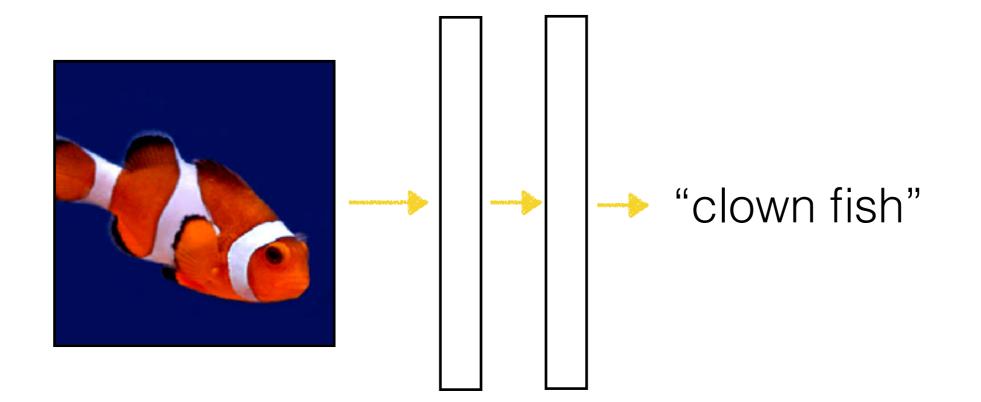
Neural Network



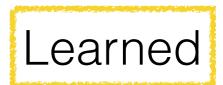


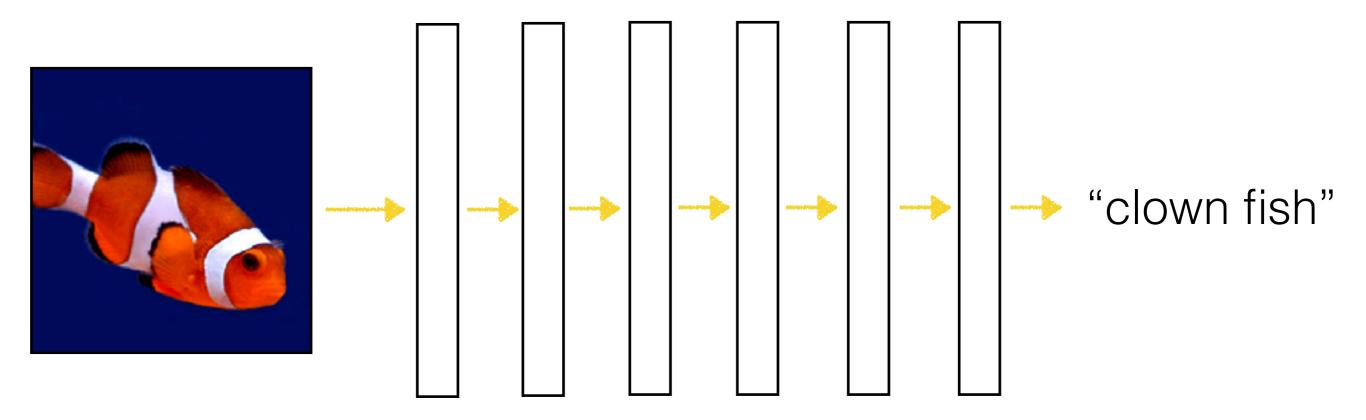
Neural Network





Deep Neural Network





Crumpled Paper Analogy

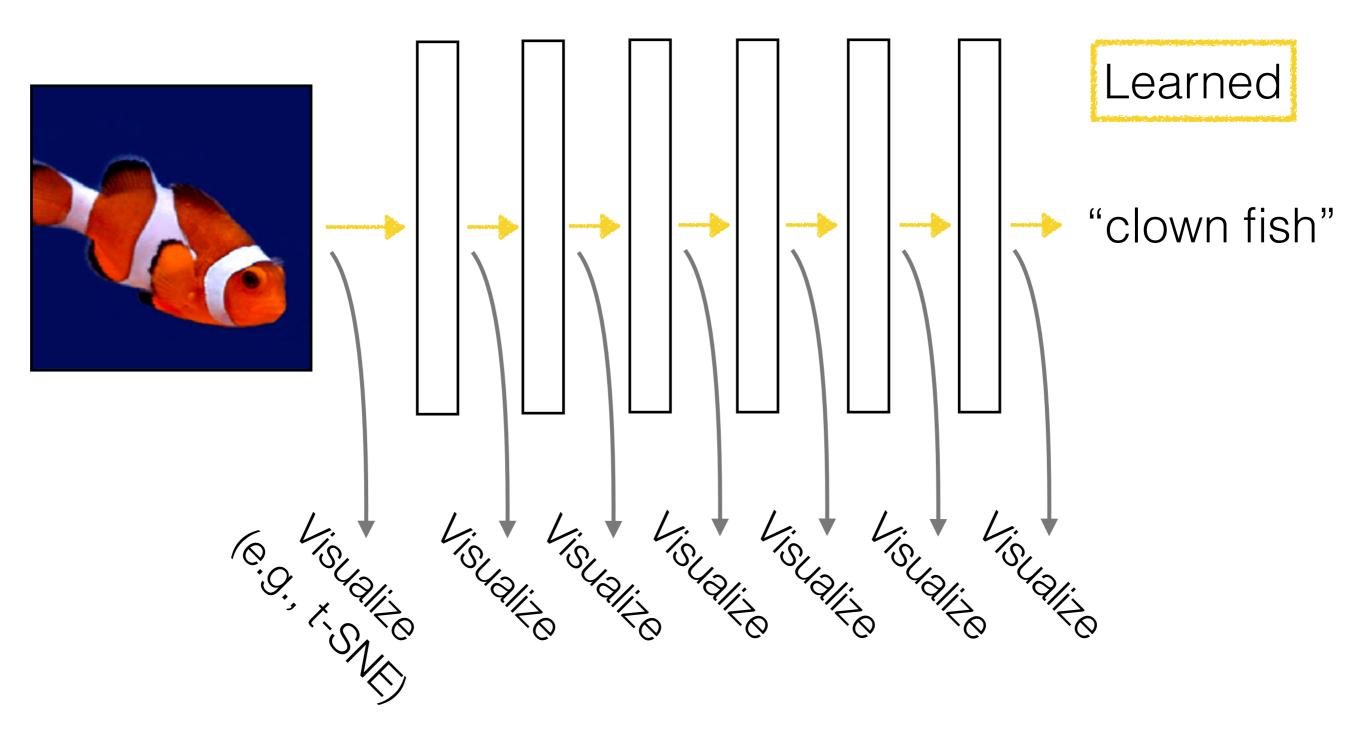
binary classification: 2 crumpled sheets of paper corresponding to the different classes

deep learning: series ("layers") of simple unfolding operations to try to disentangle the 2 sheets

Analogy: Francois Chollet, photo: George Chen

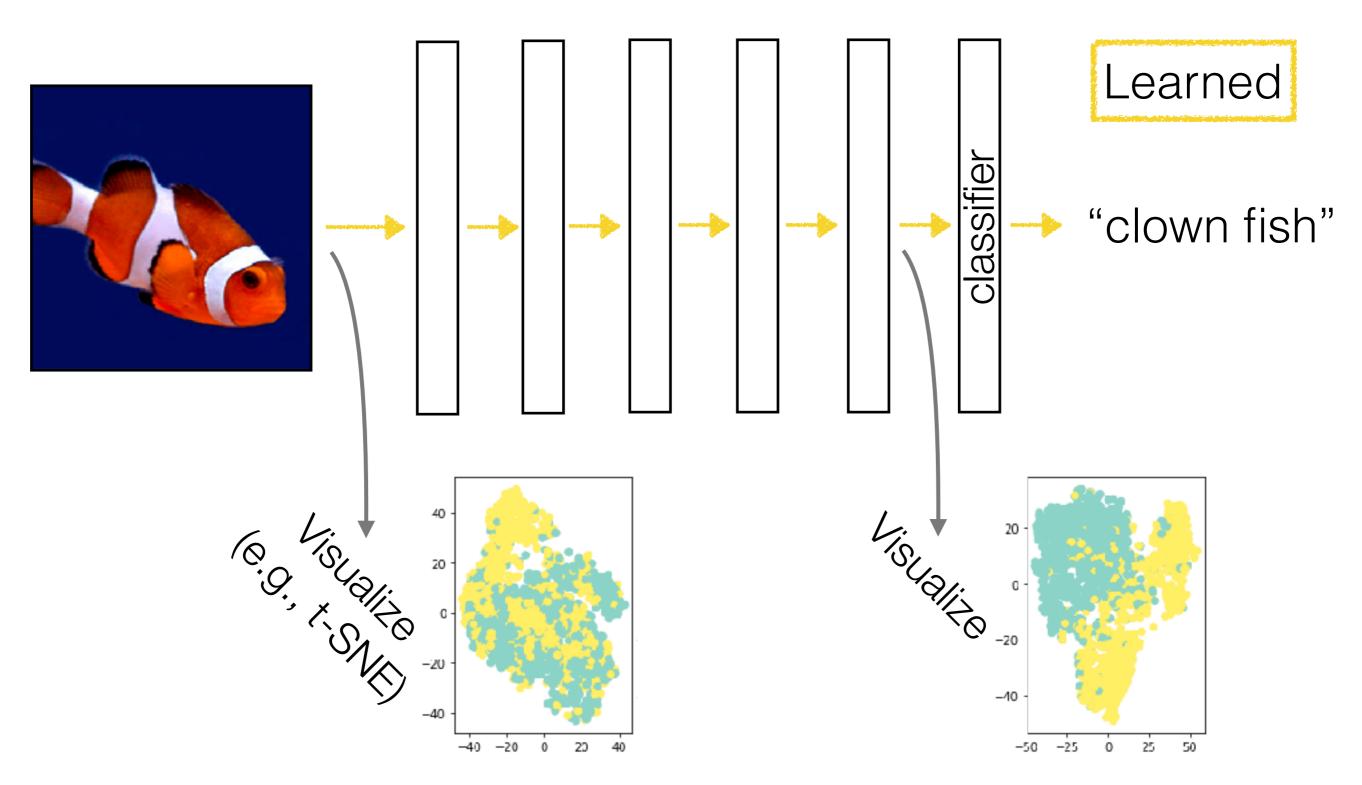
Representation Learning

Each layer's output is another way we could represent the input data



Representation Learning

Each layer's output is another way we could represent the input data



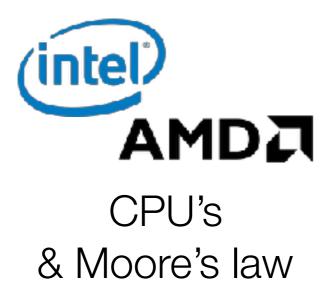
Why Does Deep Learning Work?

Actually the ideas behind deep learning are old (~1980's)

Big data



Better hardware







TPU's

• Better algorithms

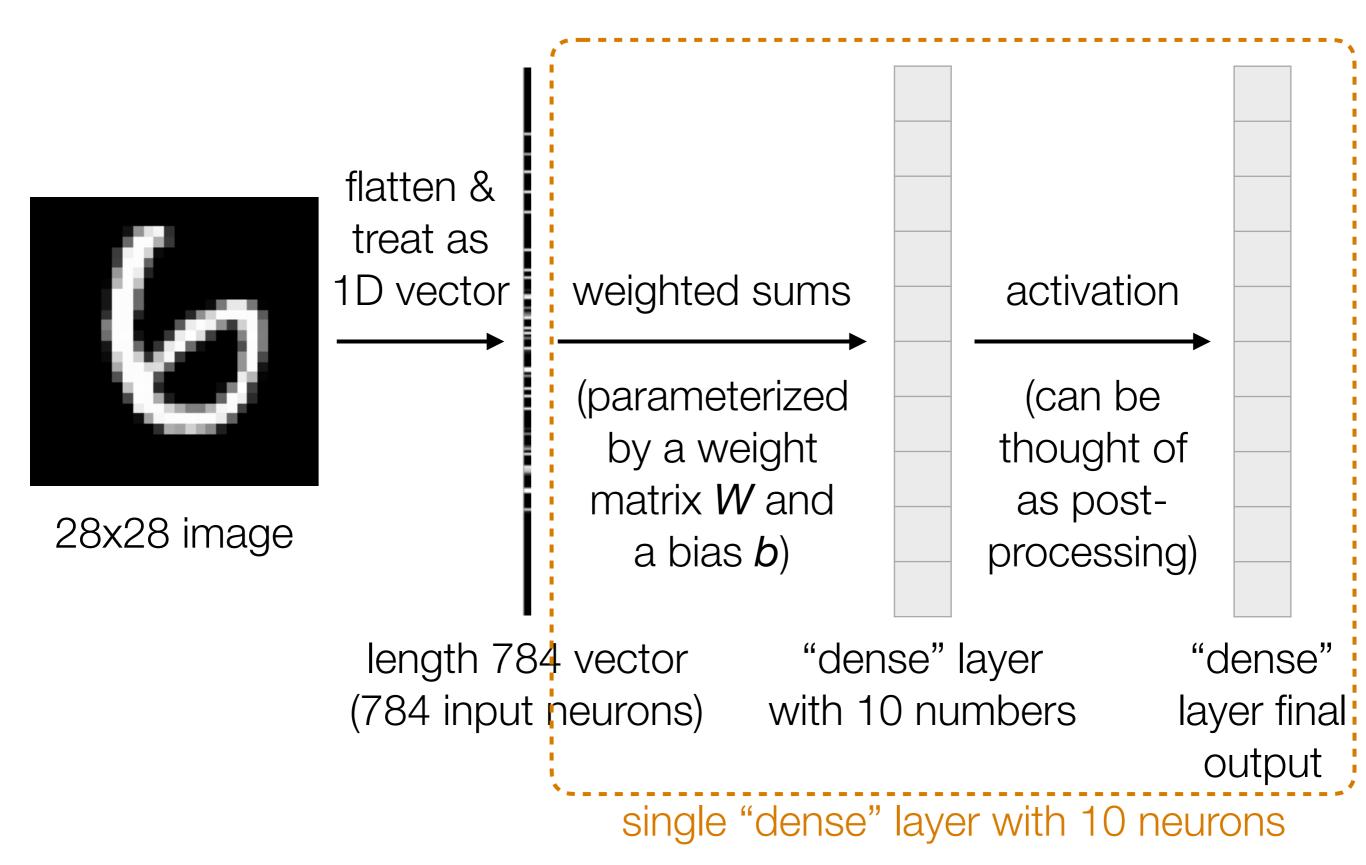
Structure Present in Data Matters

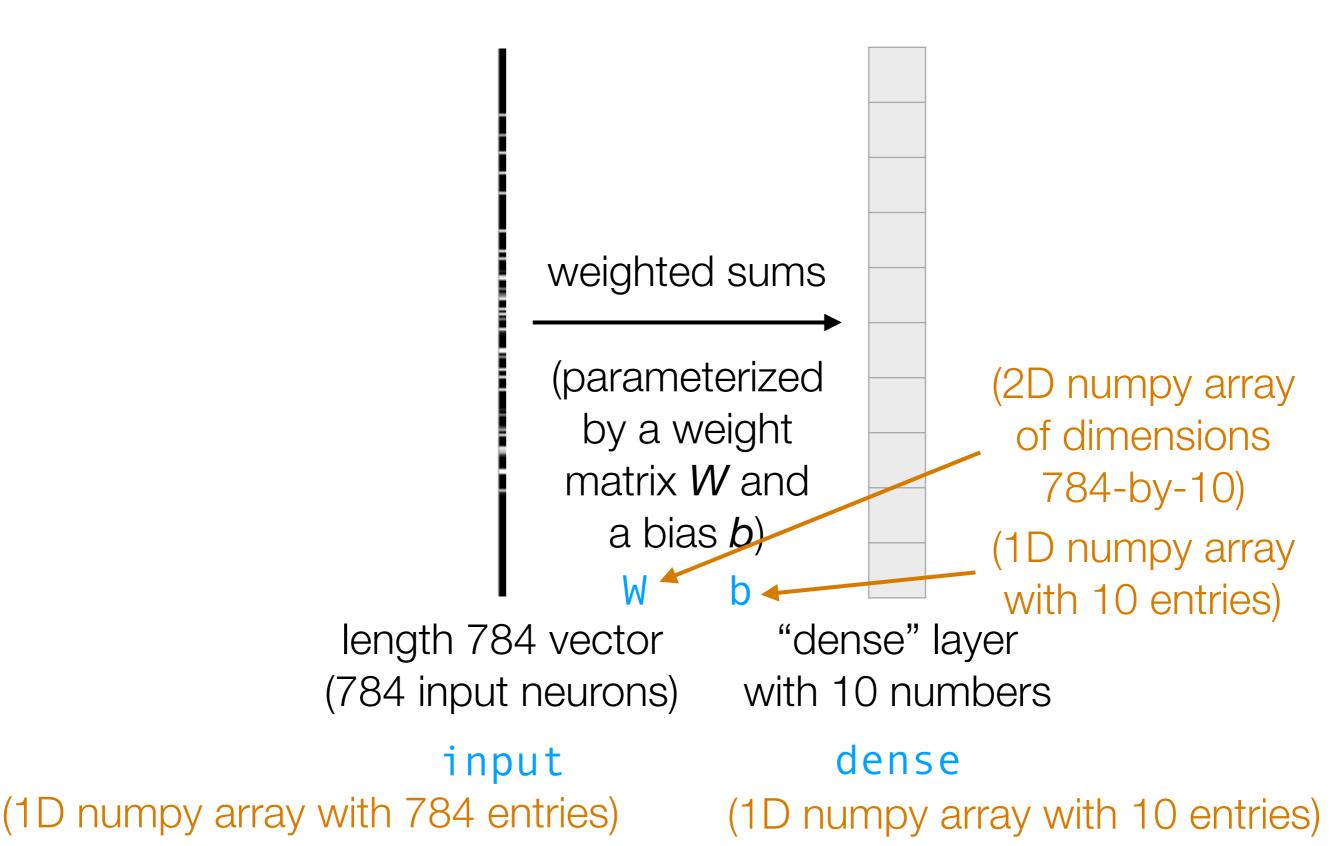
Neural nets aren't doing black magic

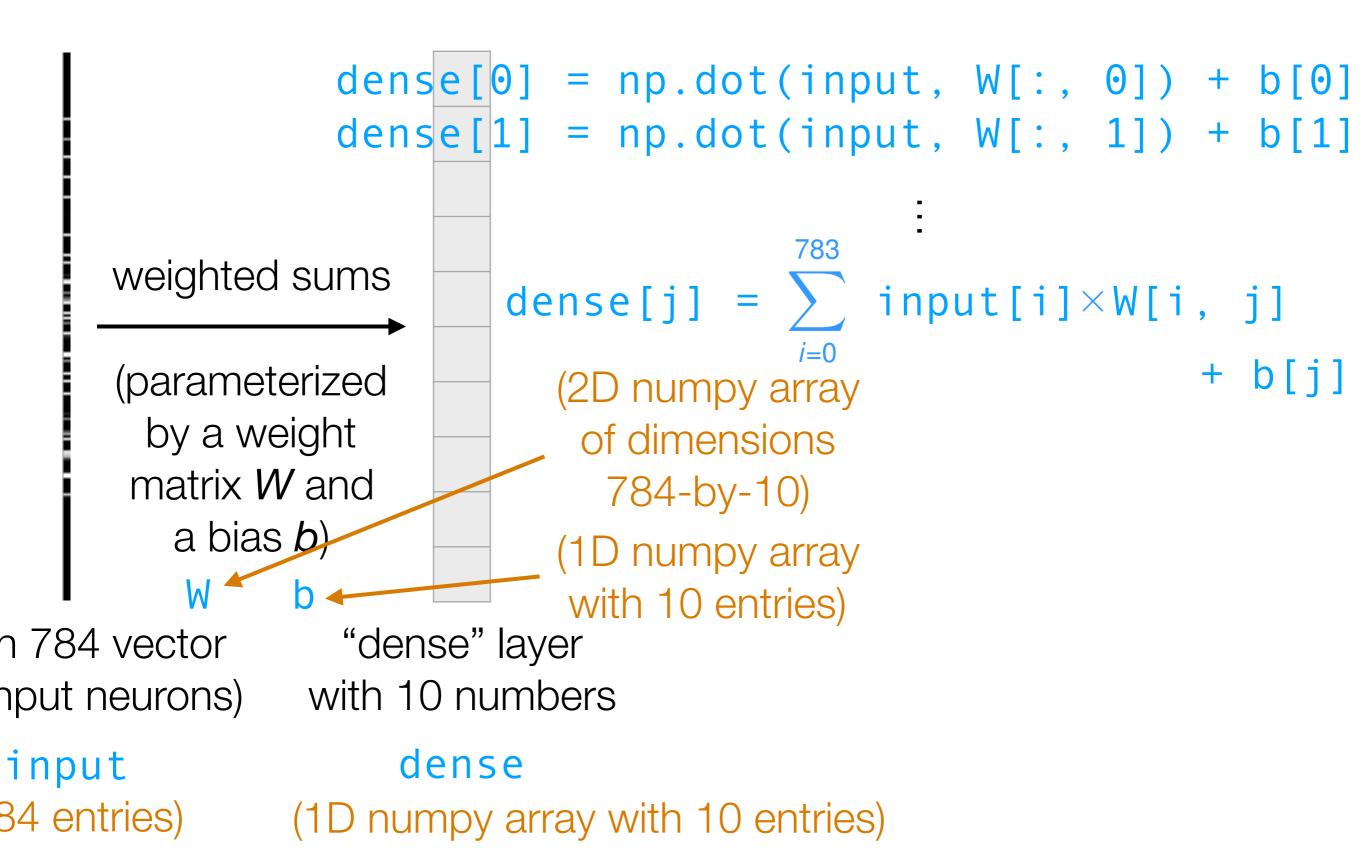
- Image analysis: convolutional neural networks (convnets) neatly incorporates basic image processing structure
- **Time series analysis:** recurrent neural networks (RNNs) incorporates ability to remember and forget things over time
 - Note: text is a time series
 - Note: video is a time series

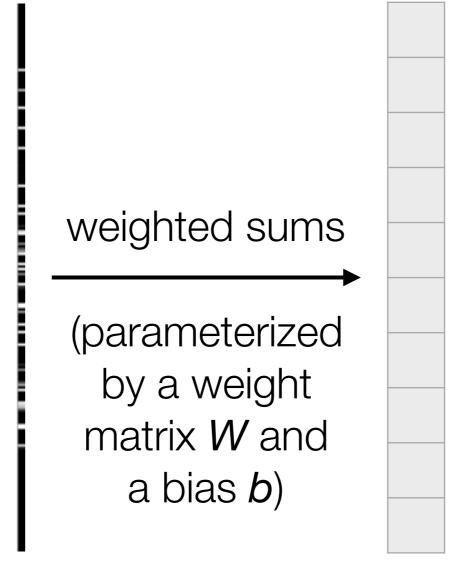
Handwritten Digit Recognition Example

Walkthrough of building a 1-layer and then a 2-layer neural net



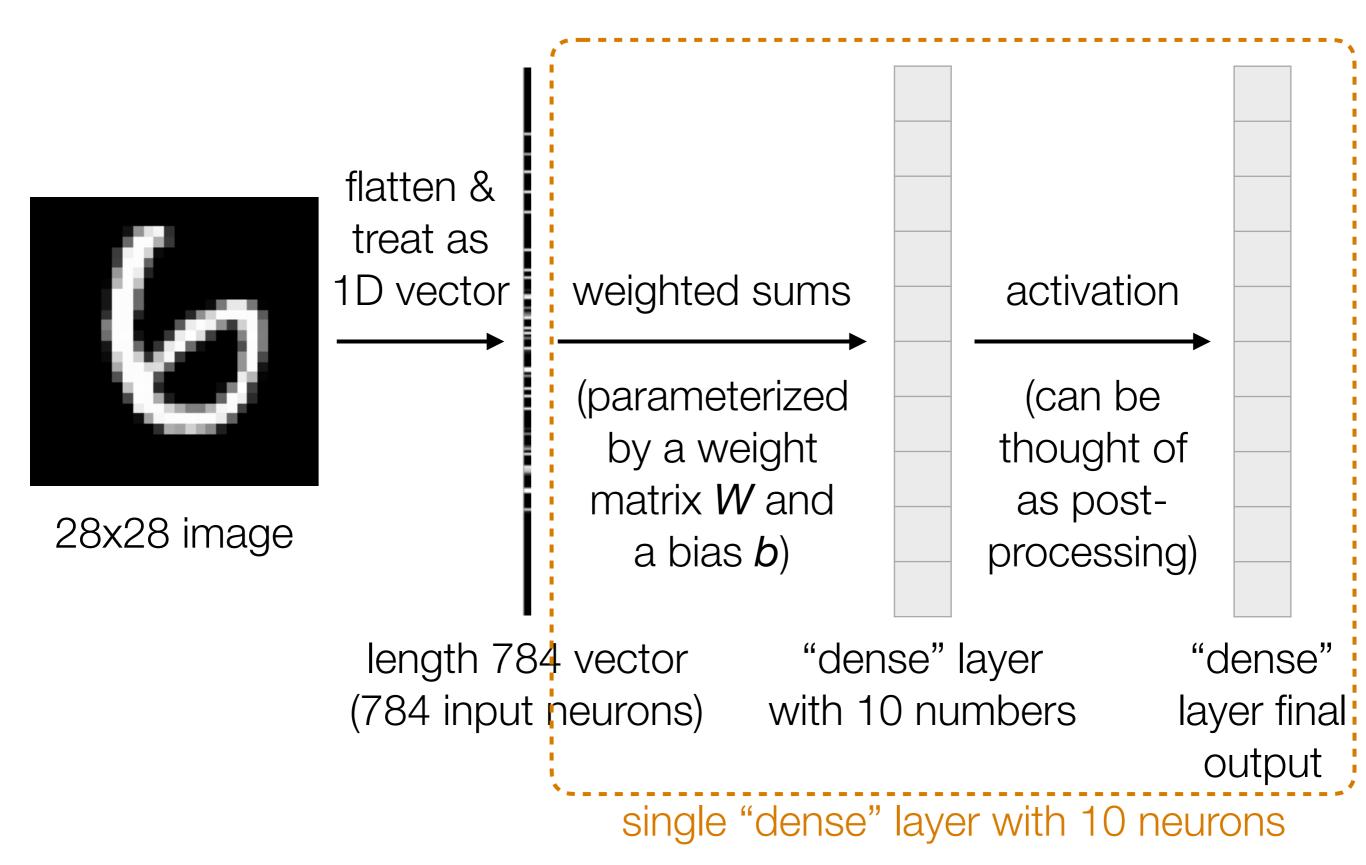


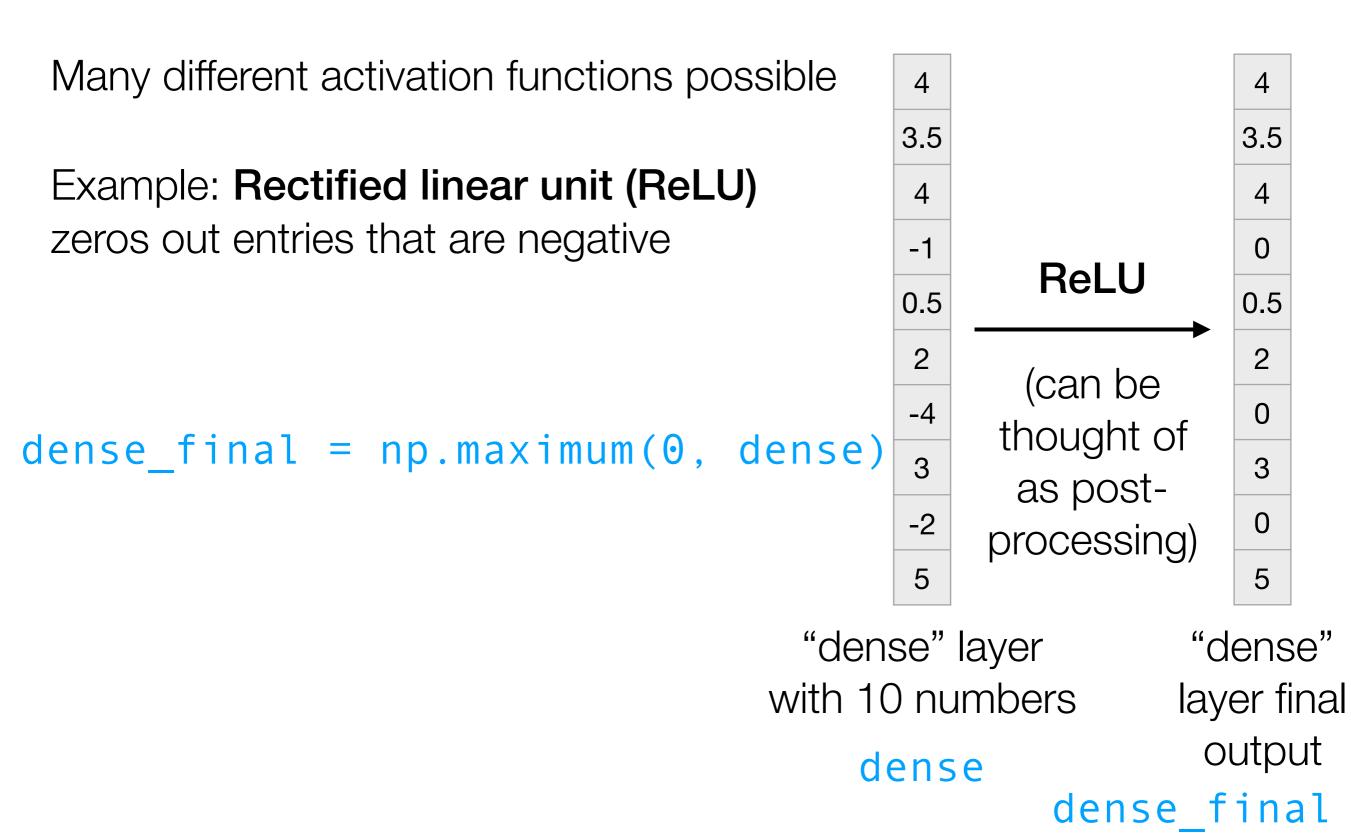




length 784 vector (784 input neurons)

"dense" layer with 10 numbers

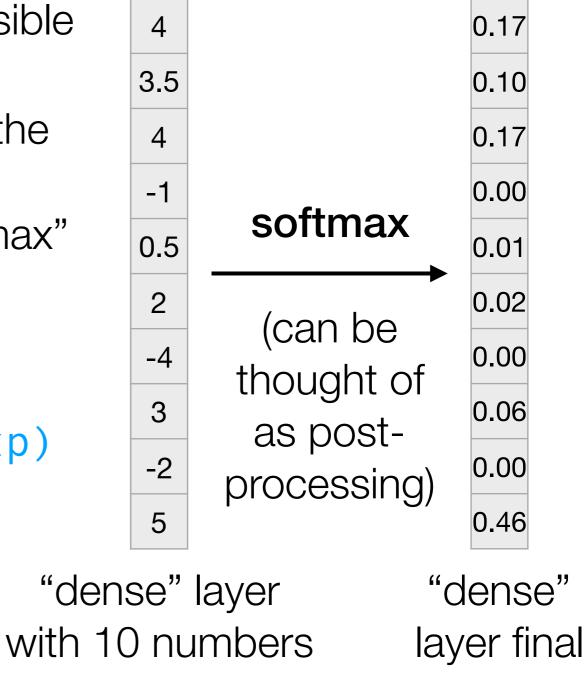




Many different activation functions possible

Example: **softmax** turns the entries in the dense layer (prior to activation) into a probability distribution (using the "softmax" transformation)

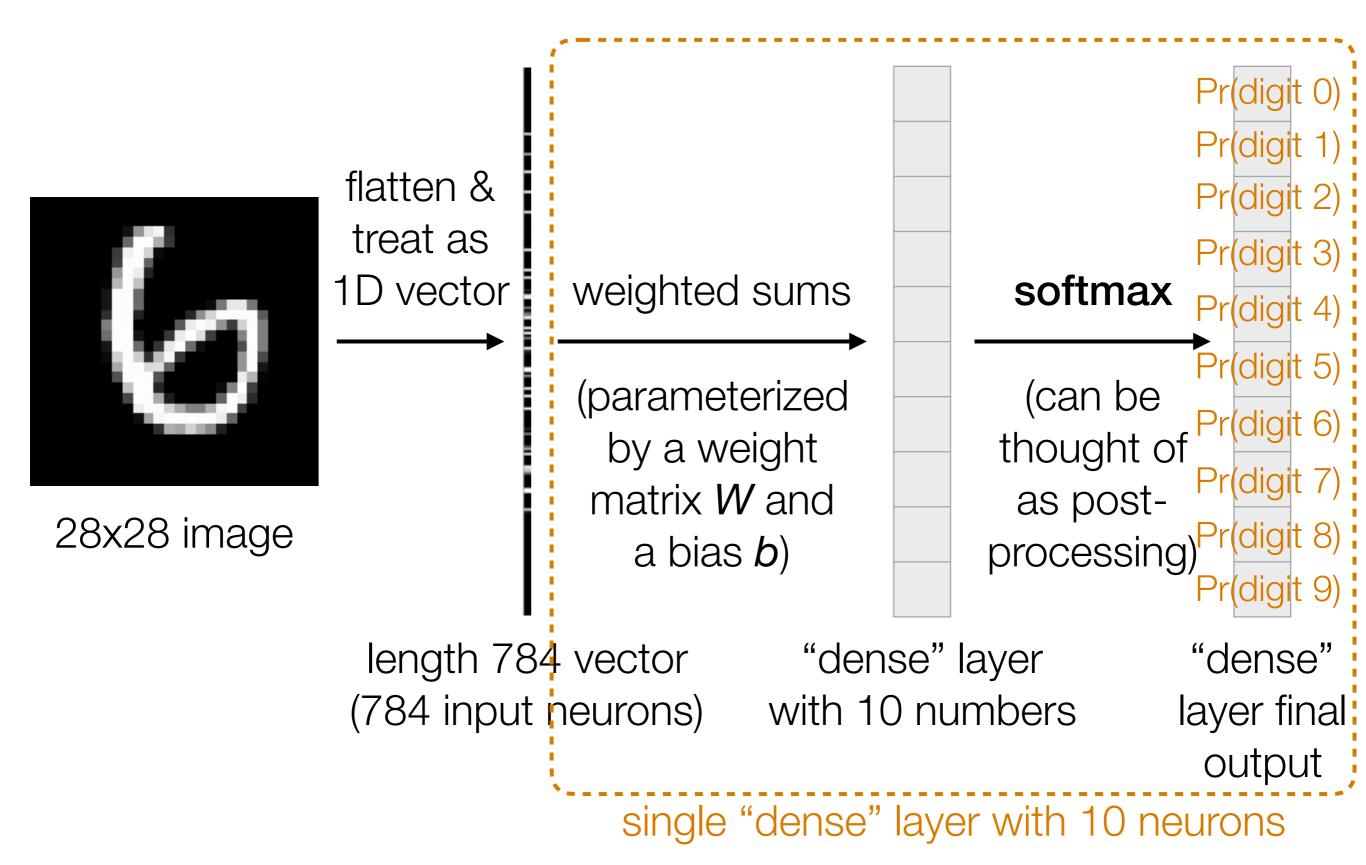
```
dense_exp = np.exp(dense)
dense_exp /= np.sum(dense_exp)
dense final = dense exp
```

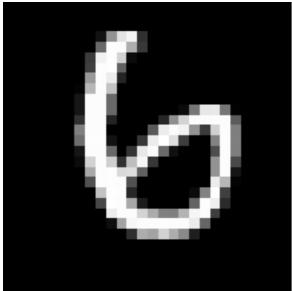


dense

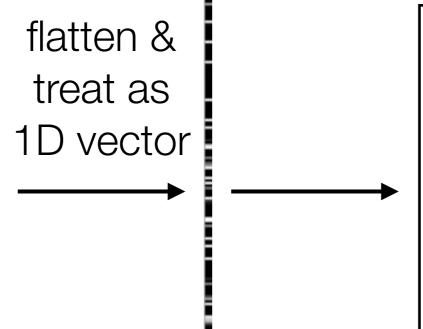
output

dense final





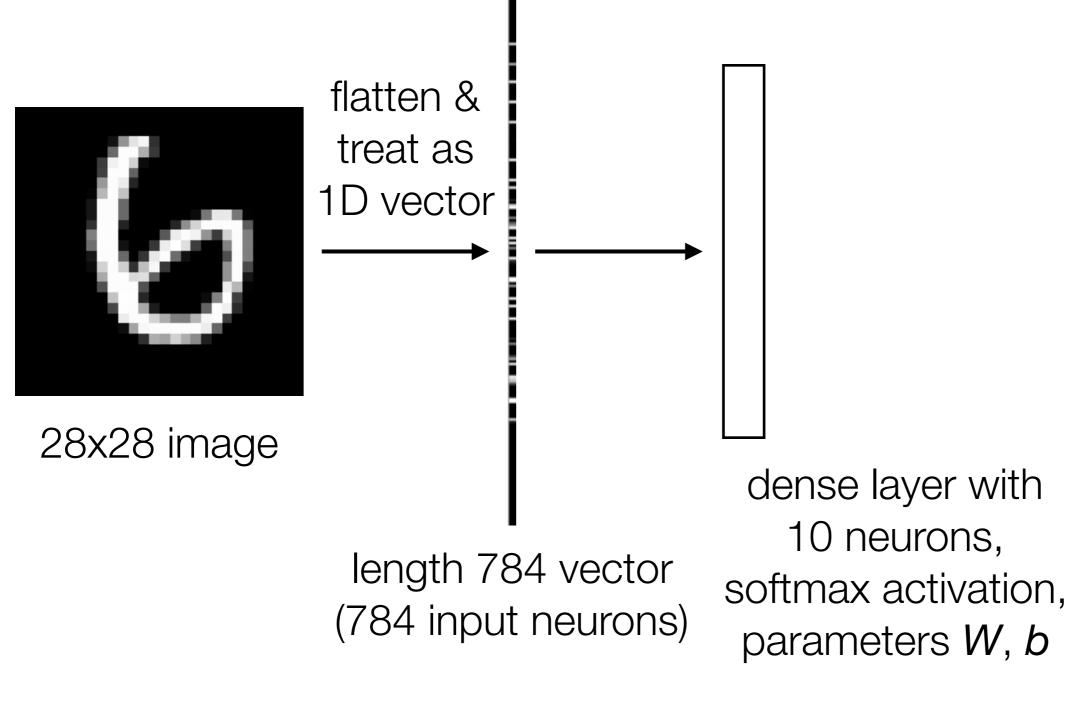
28x28 image

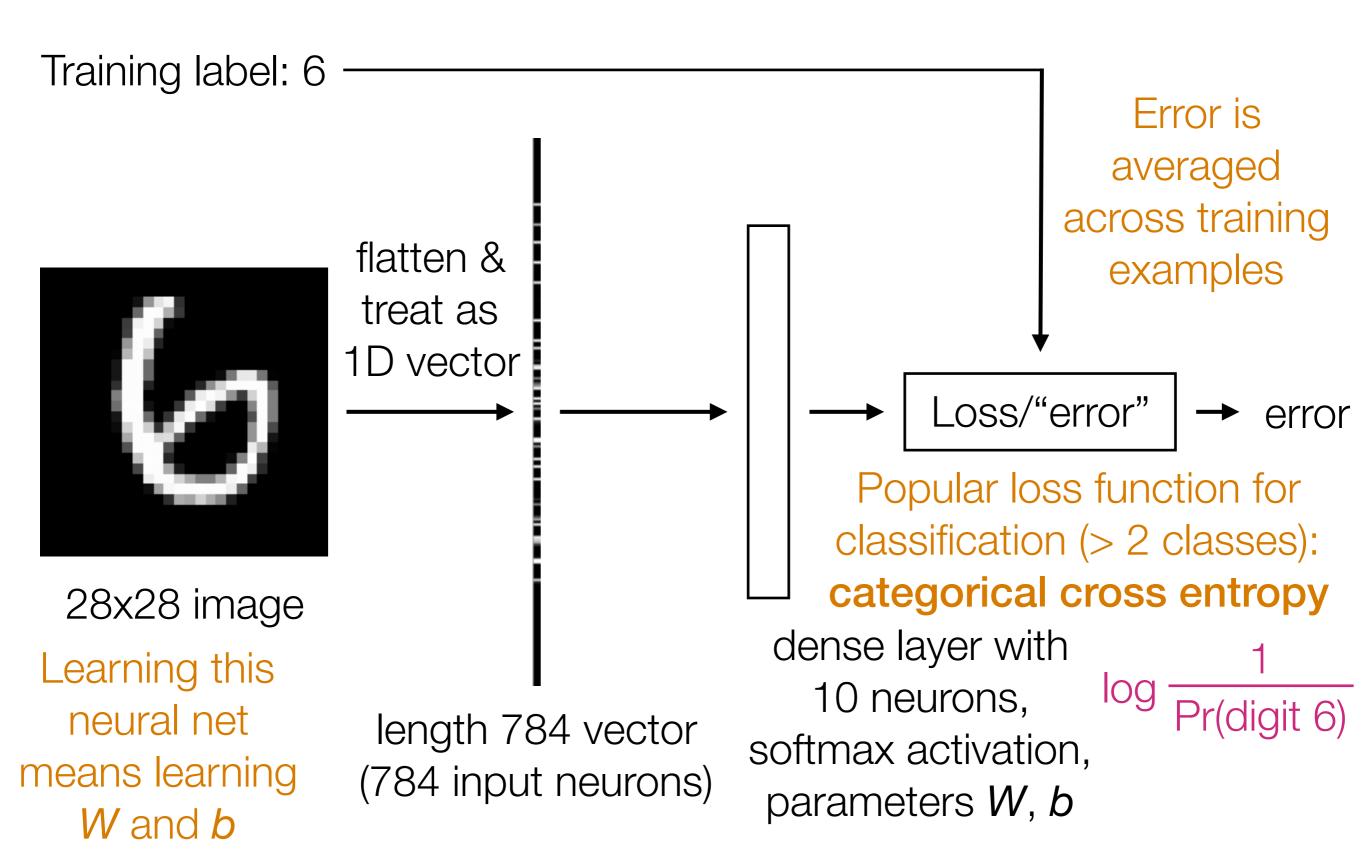


length 784 vector (784 input neurons) We want the output of the dense layer to encode probabilities for whether the input image is a 0, 1, 2, ..., 9 *but as of now we aren't providing any sort of information to enforce this*

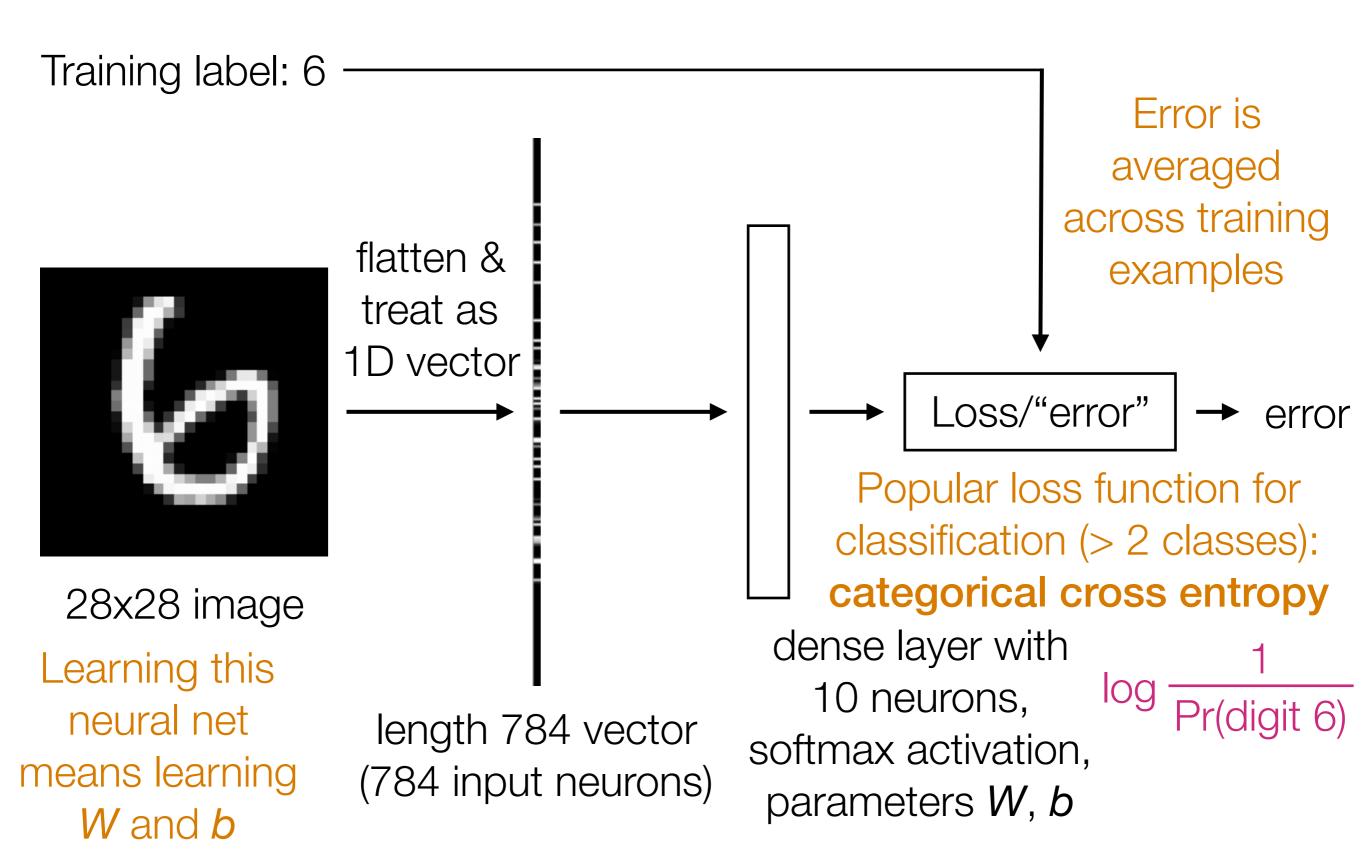
dense layer with 10 neurons, softmax activation, parameters *W*, *b*

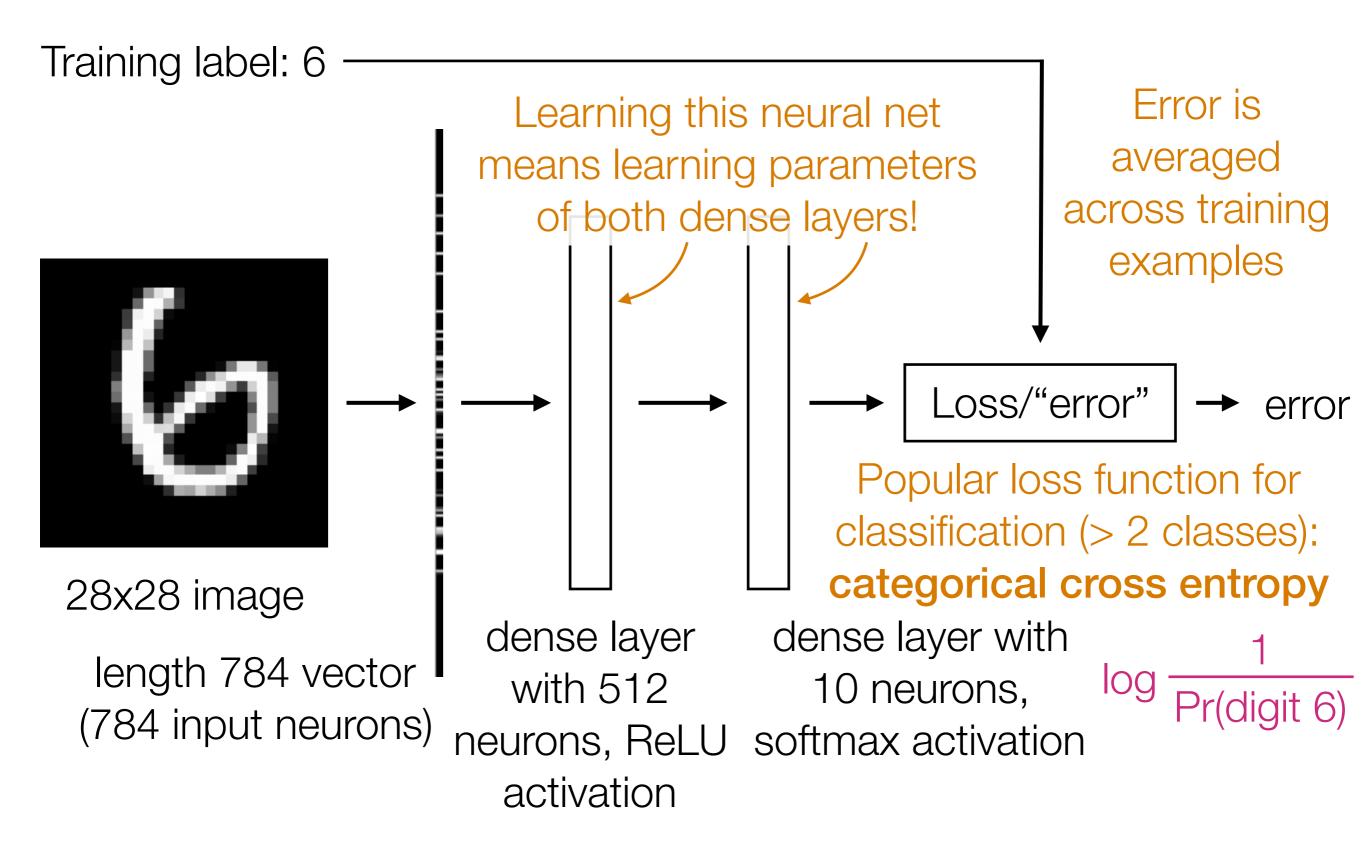
Demo part 1





Demo part 2





Demo part 3