

# Machine Learning - Classification

CS102  
Winter 2019

# Big Data Tools and Techniques

- Basic Data Manipulation and Analysis
  - Performing well-defined computations or asking well-defined questions (“queries”)
- Data Mining
  - Looking for patterns in data
- Machine Learning
  - Using data to build models and make predictions
- Data Visualization
  - Graphical depiction of data
- Data Collection and Preparation

# Regression

Using data to build models and make predictions

- Supervised
- Training data, each example:
  - Set of predictor values - “independent variables”
  - Numerical output value - “dependent variable”
- Model is function from predictors to output
  - Use model to predict output value for new predictor values
- Example
  - Predictors: mother height, father height, current age
  - Output: height

# Classification

Using data to build models and make predictions

- Supervised
- Training data, each example:
  - Set of feature values - numeric or categorical
  - Categorical output value - “label”
- Model is method from feature values to label
  - Use model to predict label for new feature values
- Example
  - Feature values: age, gender, income, profession
  - Label: buyer, non-buyer

# Other Examples

## Medical diagnosis

- **Feature values:** age, gender, history, symptom1-severity, symptom2-severity, test-result1, test-result2
- **Label:** disease

## Email spam detection

- **Feature values:** sender-domain, length, #images, keyword<sub>1</sub>, keyword<sub>2</sub>, ..., keyword<sub>n</sub>
- **Label:** spam or not-spam

## Credit card fraud detection

- **Feature values:** user, location, item, price
- **Label:** fraud or okay

# Algorithms for Classification

Despite similarity of problem statement to regression, non-numerical nature of classification leads to completely different approaches

- K-nearest neighbors
- Decision trees
- Naïve Bayes
- ... and others

# K-Nearest Neighbors (KNN)

For any pair of data items  $i_1$  and  $i_2$ , from their feature values compute  $distance(i_1, i_2)$

Example:

Features - gender, profession, age, income, postal-code

$person_1 = (\text{male}, \text{teacher}, 47, \$25K, 94305)$

$person_2 = (\text{female}, \text{teacher}, 43, \$28K, 94309)$

$distance(person_1, person_2)$

$distance()$  can be defined as inverse of  $similarity()$

# K-Nearest Neighbors (KNN)

Features - gender, profession, age, income, postal-code

person<sub>1</sub> = (male, teacher, 47, \$25K, 94305)

person<sub>2</sub> = (female, teacher, 43, \$28K, 94309)

Remember training data has labels

# K-Nearest Neighbors (KNN)

Features - gender, profession, age, income, postal-code  
person<sub>1</sub> = (male, teacher, 47, \$25K, 94305) buyer  
person<sub>2</sub> = (female, teacher, 43, \$28K, 94309) non-buyer

Remember training data has labels

To classify a new item  $i$ : In the labeled data find the K closest items to  $i$ , assign most frequent label

person<sub>3</sub> = (female, doctor, 40, \$40K, 95123)

# KNN Example

- City temperatures - France and Germany
- Features: longitude, latitude
- Distance is Euclidean distance
$$\text{distance}([o_1, a_1], [o_2, a_2]) = \sqrt{(o_1 - o_2)^2 + (a_1 - a_2)^2}$$
= actual distance in x-y plane
- Labels: frigid, cold, cool, warm, hot

Nice (7.27, 43.72) cool

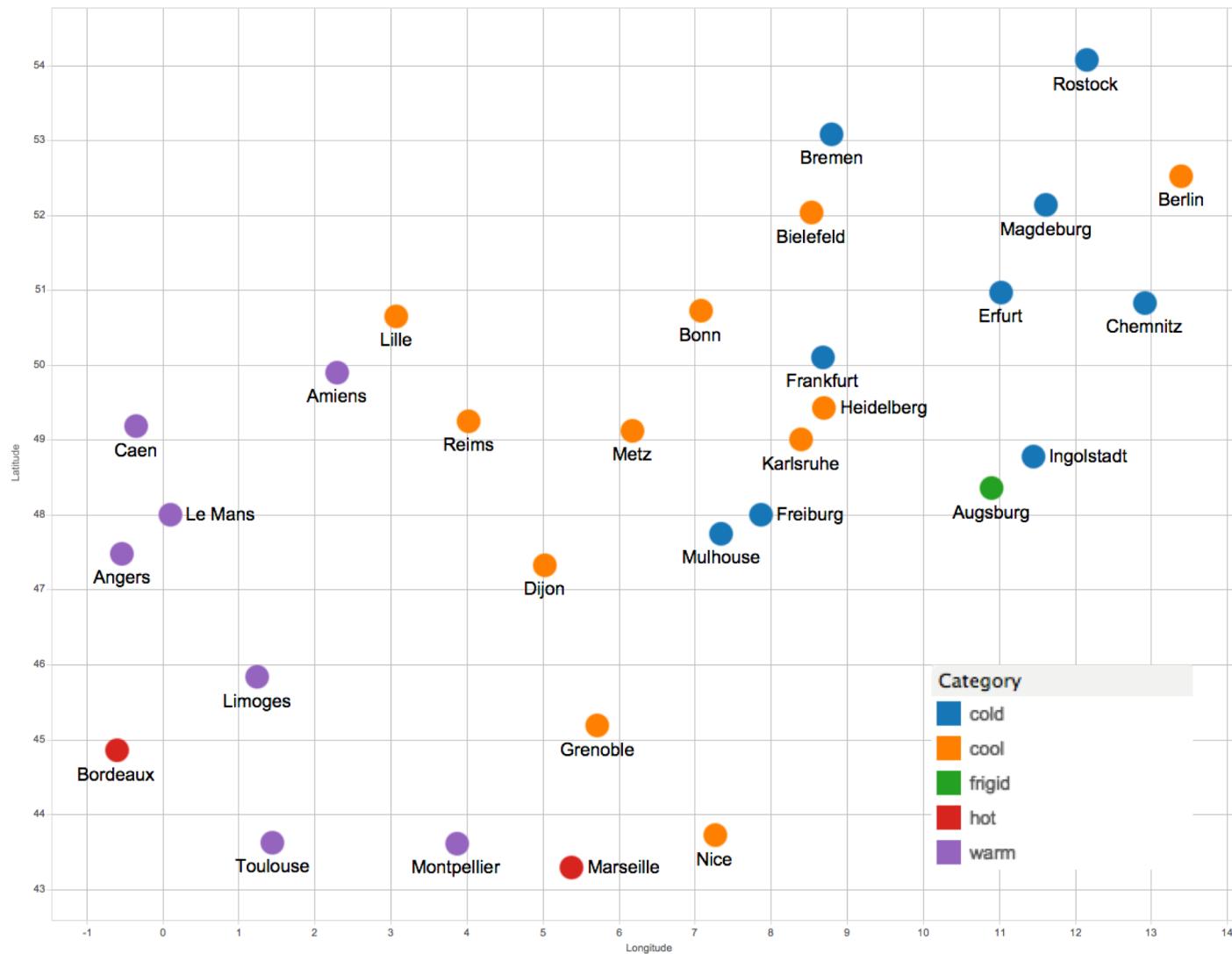
Toulouse (1.45, 43.62) warm

Frankfurt (8.68, 50.1) cold

.....

Predict temperature  
category from  
longitude and latitude

# KNN Example



# KNN Summary

To classify a new item  $i$ : find K closest items to  $i$  in the labeled data, assign most frequent label

- No hidden complicated math!
- Once distance function is defined, rest is easy
- Though not necessarily efficient

Real examples often have thousands of features

- Medical diagnosis: symptoms (yes/no), test results
- Email spam detection: words (frequency)

Database of labeled items might be enormous

# “Regression” Using KNN

Features - gender, profession, age, income, postal-code  
person<sub>1</sub> = (male, teacher, 47, \$25K, 94305) buyer  
person<sub>2</sub> = (female, teacher, 43, \$28K, 94309) non-buyer

Remember training data has labels

To classify a new item  $i$ , find K closest items to  $i$  in the labeled data, assign most frequent label

person<sub>3</sub> = (female, doctor, 40, \$40K, 95123)

# “Regression” Using KNN

Features - gender, profession, age, income, postal-code

$\text{person}_1 = (\text{male}, \text{teacher}, 47, \$25K, 94305)$  \$250

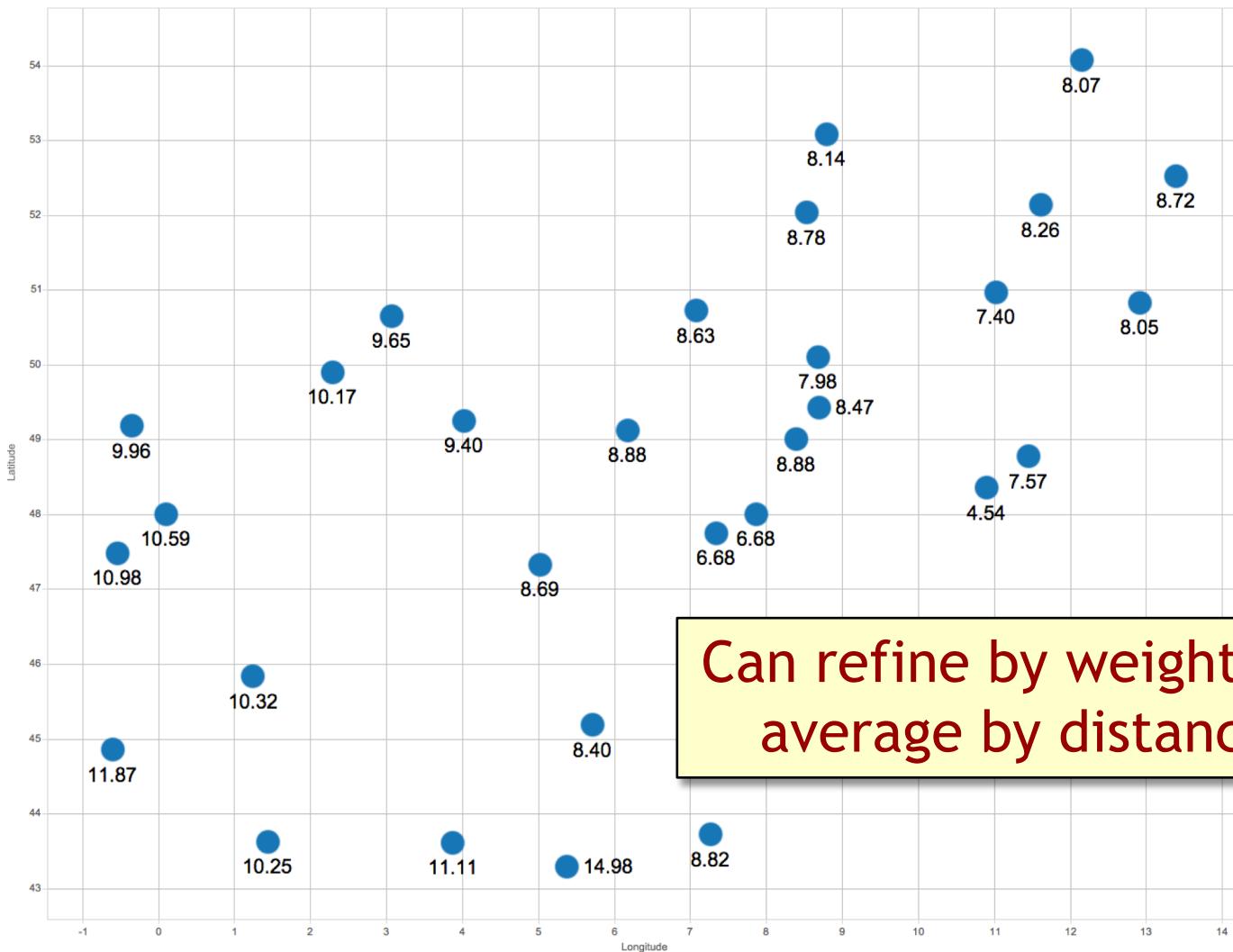
$\text{person}_2 = (\text{female}, \text{teacher}, 43, \$28K, 94309)$  \$100

Remember training data has labels

To classify a new item  $i$ , find K closest items to  $i$   
in the labeled data, assign average value of labels

$\text{person}_3 = (\text{female}, \text{doctor}, 40, \$40K, 95123)$

# Regression Using KNN - Example



# Decision Trees

- Use the training data to construct a decision tree
- Use the decision tree to classify new data

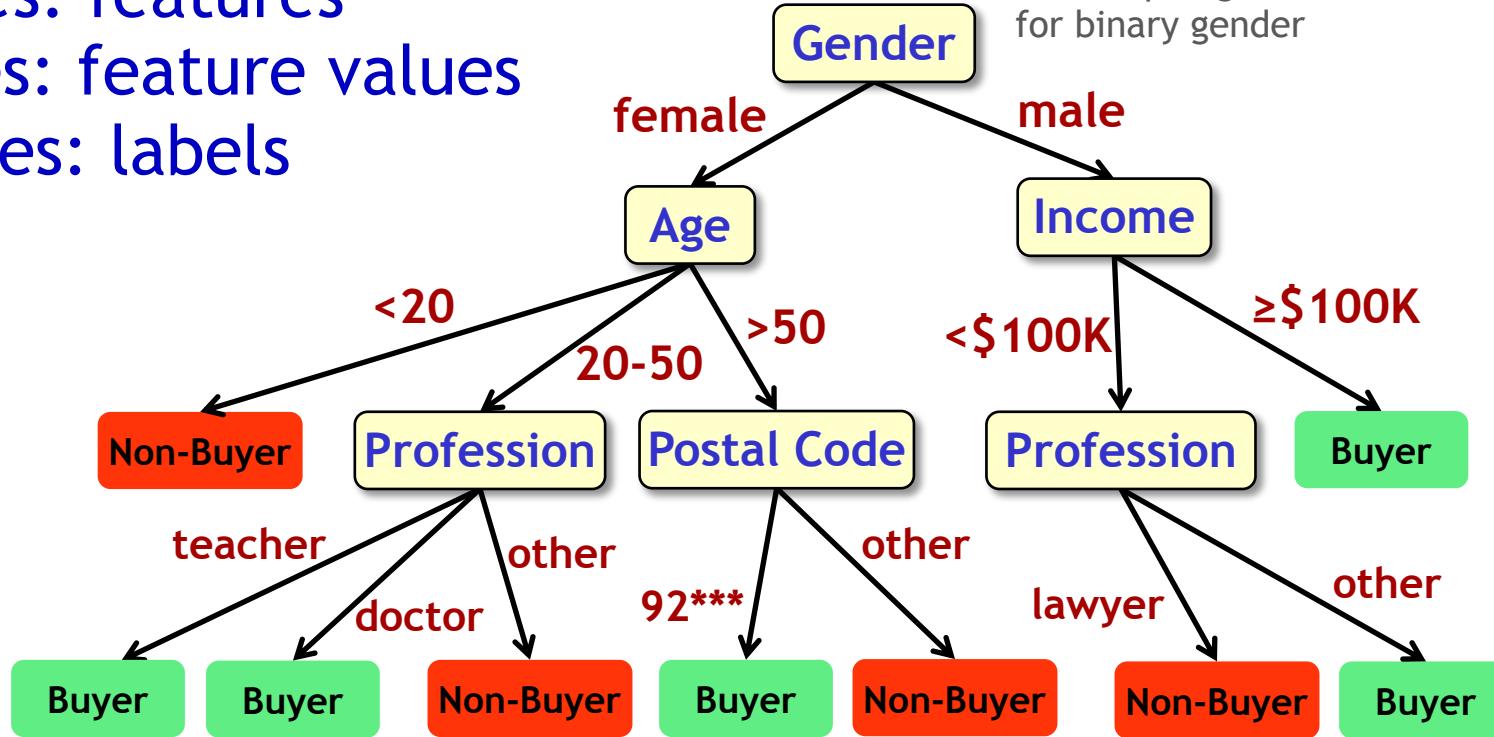
# Decision Trees

Nodes: features

Edges: feature values

Leaves: labels

with apologies  
for binary gender



New data item to classify:  
Navigate tree based on feature values

# Decision Trees

Primary challenge is building good decision trees from training data

- Which features and feature values to use at each choice point
- HUGE number of possible trees even with small number of features and values

Common approach: “forest” of many trees, combine the results

- Still impossible to consider all trees

# Naïve Bayes

Given new data item  $i$ , based on  $i$ 's feature values and the training data, compute the probability of each possible label. Pick highest one.

Efficiency relies on **conditional independence assumption**:

Given any two features  $F_1, F_2$  and a label  $L$ , the probability that  $F_1=v_1$  for an item with label  $L$  is independent of the probability that  $F_2=v_2$  for that item

Examples:

gender and age? income and postal code?

# Naïve Bayes

Given new data item  $i$ , based on  $i$ 's feature values and the training data, compute the probability of each possible label. Pick highest one.

Efficiency relies on **conditional independence assumption**:

Conditional independence assumption often doesn't hold, which is why the approach is “naive”

label L, the probability with label L is  $P(v_1=v_1 \text{ and } v_2=v_2 \text{ for that item.})$

Examples:

gender and age? income and

Nevertheless the approach works very well in practice

# Naïve Bayes Example

Predict temperature category for a country based on whether the country has coastline and whether it is in the EU

country	coastline	EU	tempAvg	category
Albania	yes	no	15.18	hot
Andorra	no	no	9.60	warm
Belarus	no	no	5.95	cool
Belgium	yes	yes	9.65	warm
Bosnia and Herzegov	no	no	9.60	warm
Bulgaria	yes	yes	10.44	warm
Croatia	yes	yes	10.87	warm
Czech Republic	no	yes	7.86	cool
Denmark	yes	yes	7.63	cool
Estonia	yes	yes	4.59	cold
Finland	yes	yes	3.49	cold
Germany	yes	yes	7.87	cool
Greece	yes	yes	16.90	hot
Hungary	no	yes	9.60	warm
Ireland	yes	yes	9.30	warm

# Naïve Bayes Preparation

Step 1: Compute fraction (probability) of items in each category

cold	.18
cool	.38
warm	.24
hot	.20

# Naïve Bayes Preparation

Step 2: For each category, compute fraction of items in that category for each feature and value

cold (.18)	coastline=yes	.83
	coastline=no	.17
	EU=yes	.67
	EU=no	.33
cool (.38)	coastline=yes	.69
	coastline=no	.31
	EU=yes	.77
	EU=no	.23

warm (.24)	coastline=yes	.5
	coastline=no	.5
	EU=yes	.5
	EU=no	.5
hot (.20)	coastline=yes	1.0
	coastline=no	.0
	EU=yes	.71
	EU=no	.29

# Naïve Bayes Prediction

New item: France, coastline=yes, EU=yes

For each category: probability of category times product of probabilities of new item's features in that category. Pick highest.

category	prob.	coastline=yes	EU=yes	product
cold	.18	.83	.67	.10
cool	.38	.69	.77	.20
warm	.24	.5	.5	.06
hot	.20	1.0	.71	.14

# Naïve Bayes Prediction

New item: France, coastline=yes, EU=yes

For each category: probability of category times product of probabilities of new item's features in that category. Pick highest.

category	prob.	coastline=yes	EU=yes	product
cold	.18	.83	.67	.10
cool	.38	.69	.77	.20
warm	.24	.5	.5	.06
hot	.20	1.0	.71	.14

# Naïve Bayes Prediction

New item: Serbia, coastline=no, EU=no

For each category: probability of category times product of probabilities of new item's features in that category. Pick highest.

category	prob.	coastline=no	EU=no	product
cold	.18	.17	.33	.01
cool	.38	.31	.23	.03
warm	.24	.5	.5	.06
hot	.20	.0	.29	.00

# Naïve Bayes Prediction

New item: Serbia, coastline=no, EU=no

For each category: probability of category times product of probabilities of new item's features in that category. Pick highest.

category	prob.	coastline=no	EU=no	product
cold	.18	.17	.33	.01
cool	.38	.31	.23	.03
warm	.24	.5	.5	.06
hot	.20	.0	.29	.00

# Naïve Bayes Prediction

New item: Austria, coastline=no, EU=yes

For each category: probability of category times product of probabilities of new item's features in that category. Pick highest.

category	prob.	coastline=no	EU=yes	product
cold	.18	.17	.67	.02
cool	.38	.31	.77	.09
warm	.24	.5	.5	.06
hot	.20	.0	.71	.0

# Naïve Bayes Prediction

New item: Austria, coastline=no, EU=yes

For each category: probability of category times product of probabilities of new item's features in that category. Pick highest.

category	prob.	coastline=no	EU=yes	product
cold	.18	.17	.67	.02
cool	.38	.31	.77	.09
warm	.24	.5	.5	.06
hot	.20	.0	.71	.0

# Naïve Bayes Prediction

New item: Austria, coastline=no, EU=yes

For example, many presentations of Naïve Bayes include an additional normalization step so the final products are probabilities that sum to 1.0. The choice of label is unchanged, so we've omitted that step for simplicity.

category	cool	warm	hot	EU
coastline	.38	.24	.20	.31
EU	.77	.5	.71	.09
temperature	.09	.06	.0	.2

# Feature Selection

Real applications often have thousands of features

- Naïve Bayes typically uses only some of the features, those most affecting the label
- Decision trees also rely on choosing features that most affect the label
- Feature selection is a key part of machine learning - an art and a science

# Training and Test

Created machine learning model from training data.  
How do you know whether it's a good model?

- Try it on known data

Training Data	Feature Values				Labels

A red bracket on the left side of the table spans all rows, labeled "Training Data". A red bracket on the right side of the table spans the bottom four rows, labeled "Test Data".

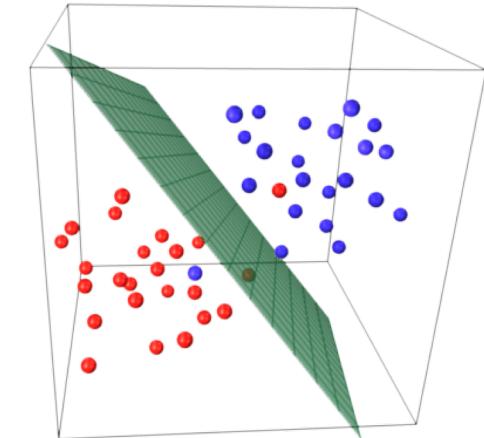
# Other Terms You Might Hear

## Logistic regression

- Recall regression model is function  $f$  from predictor values to numeric output value
- For classification: from training data obtain one regression function  $f_L$  for each label  $L$   
 $f_L(\text{feature-values})$  = probability of item having label  $L$

## Support Vector Machine

- Two labels only (“binary classifier”)
- Features = multidimensional space
- From training data SVM finds hyper-plane that best divides space according to labels



# Other Terms You Might Hear

## Deep Learning

- Complex, mysterious (the ultimate “black box” software), becoming extremely popular
- Multiple layers, each layer uses classification techniques to reduce complexity for next layer and further classification
- Important plus: identifies features from raw data

## Neural Network

- Precursor to deep learning, typically two layers
- Leap to deep learning enabled by massive amounts of data, powerful computing

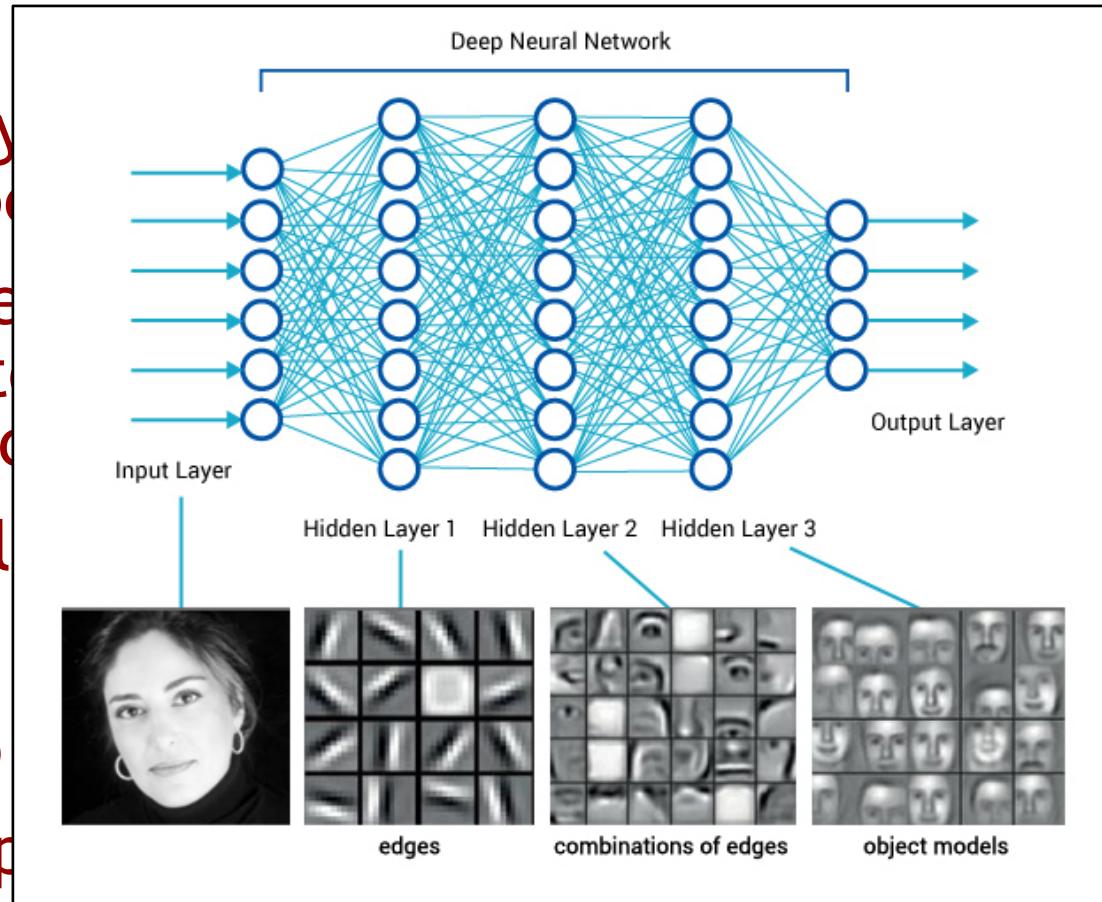
# Other Terms You Might Hear

## Deep Learning

- Complex, my software), because
- Multiple layers techniques to and further on
- Important pl

## Neural Network

- Precursor to
- Leap to deep amounts of data, powerful computing



# Classification Summary

- Supervised machine learning
- Training data, each example:
  - Set of feature values - numeric or categorical
  - Categorical output value - label
- Model is “function” from feature values to label
  - Use model to predict label for new feature values
- Approaches we covered
  - K-nearest neighbors - relies on distance (or similarity) function
  - Decision trees - relies on finding good trees/forests
  - Naïve Bayes - relies on conditional independence assumption