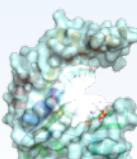




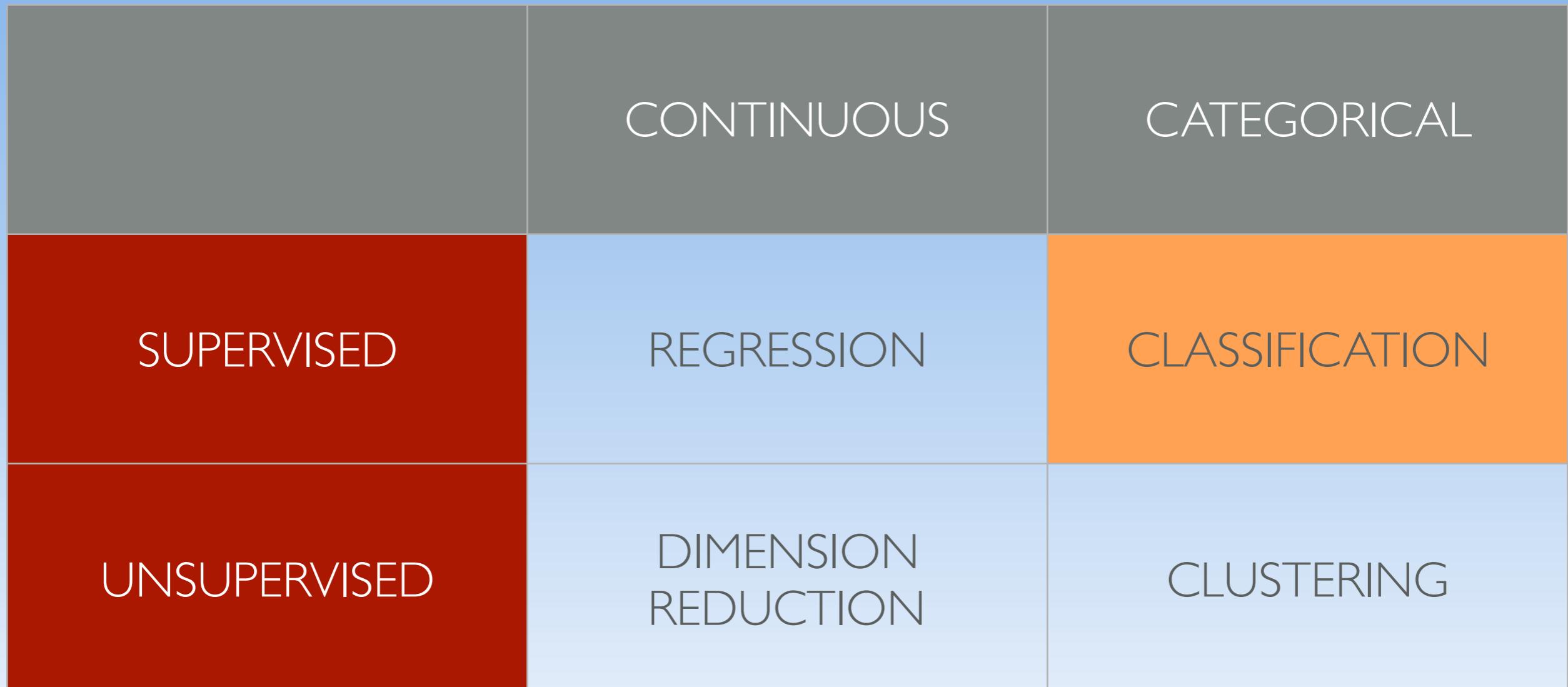
CLASSIFICATION



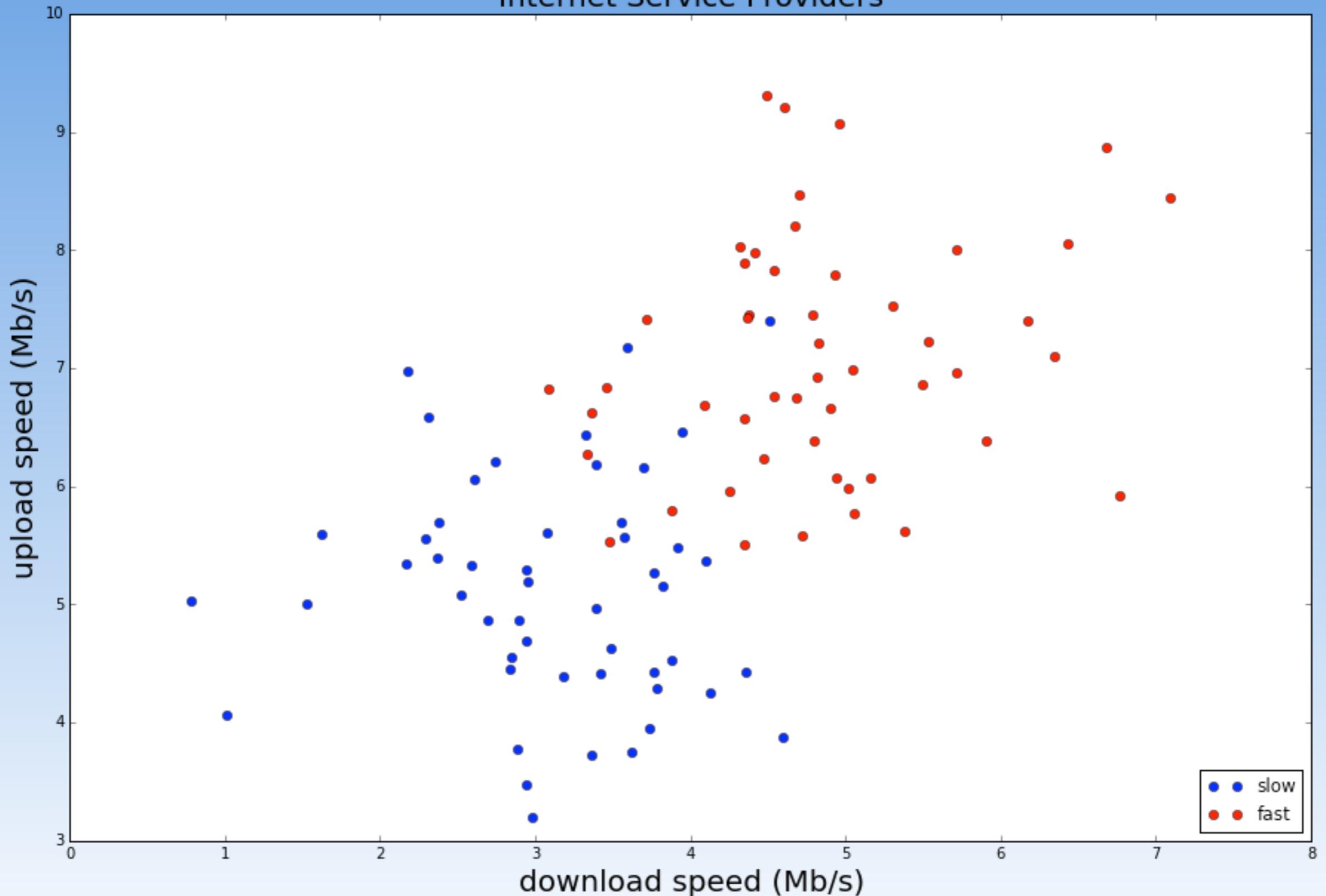
CLASSIFICATION

	CONTINUOUS	CATEGORICAL
SUPERVISED	?	?
UNSUPERVISED	?	?

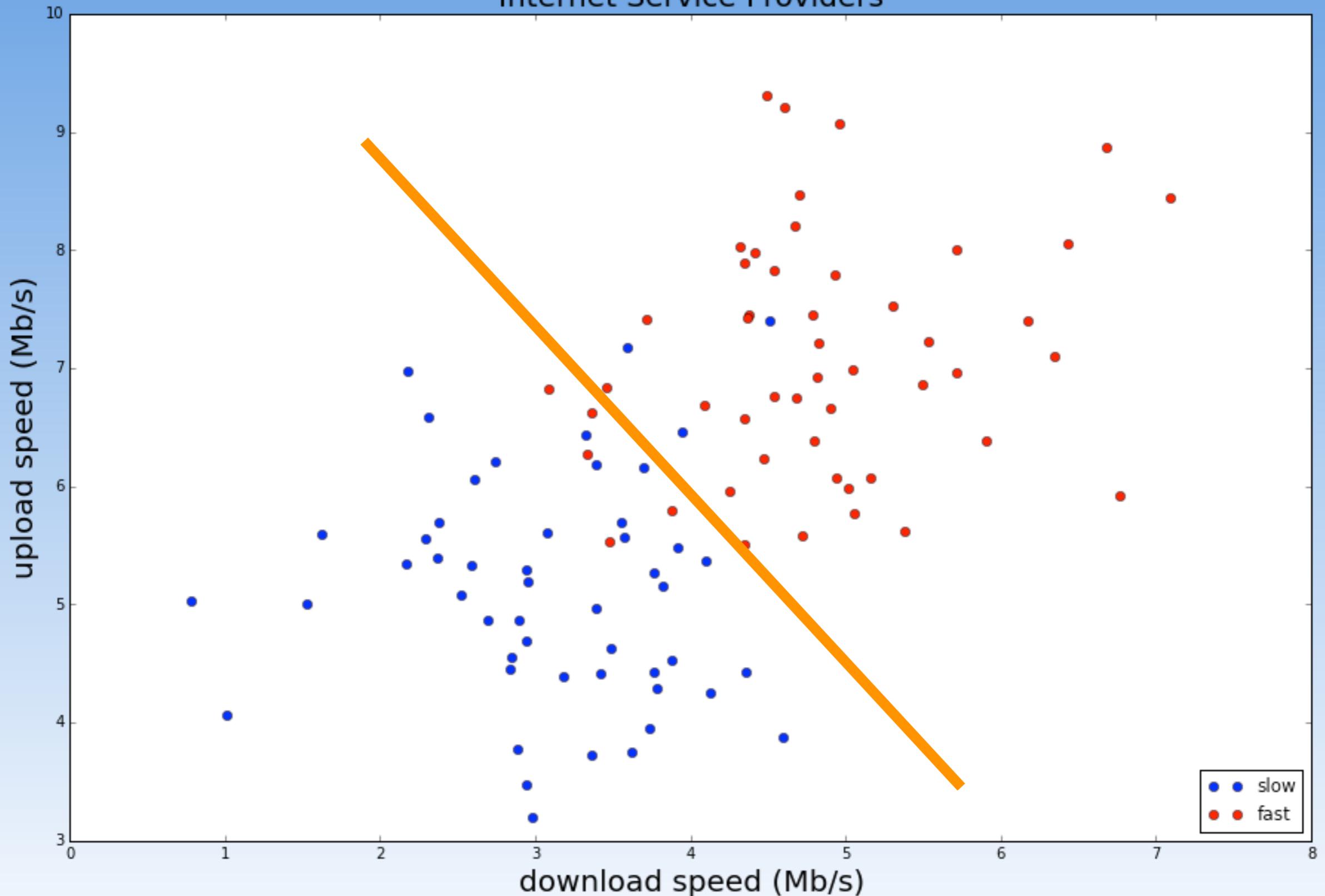
CLASSIFICATION



Internet Service Providers



Internet Service Providers



BINARY CLASSIFICATION

Features

Labels

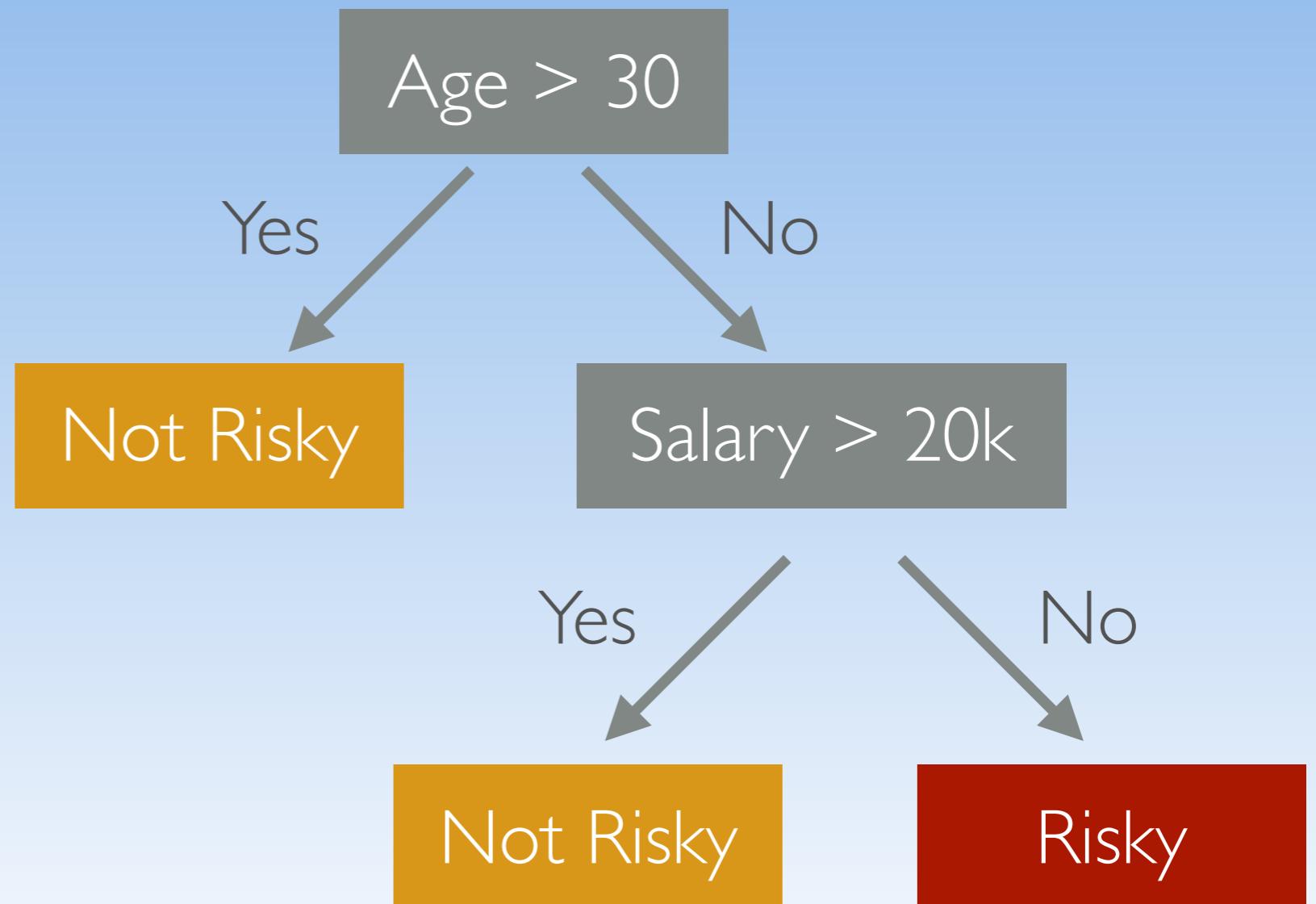
	Age	Gender	Annual Salary	Months in residence	Months in job	Current Debt	Paid off credit
Data Point	Client 1	23	M	\$30,000	36	12	\$5,000 Yes
	Client 2	30	F	\$45,000	12	12	\$1,000 Yes
	Client 3	19	M	\$15,000	3	1	\$10,000 No

DECISION TREE

- GOAL: make optimal decision at each node

- Advantages:

- easy to interpret
- fast prediction
- rules based



DECISION TREE

Binary splits

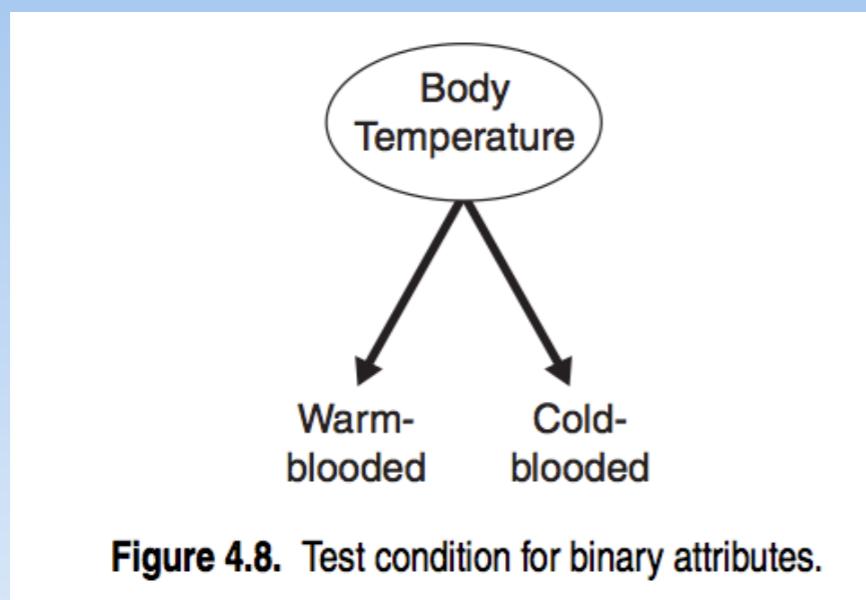
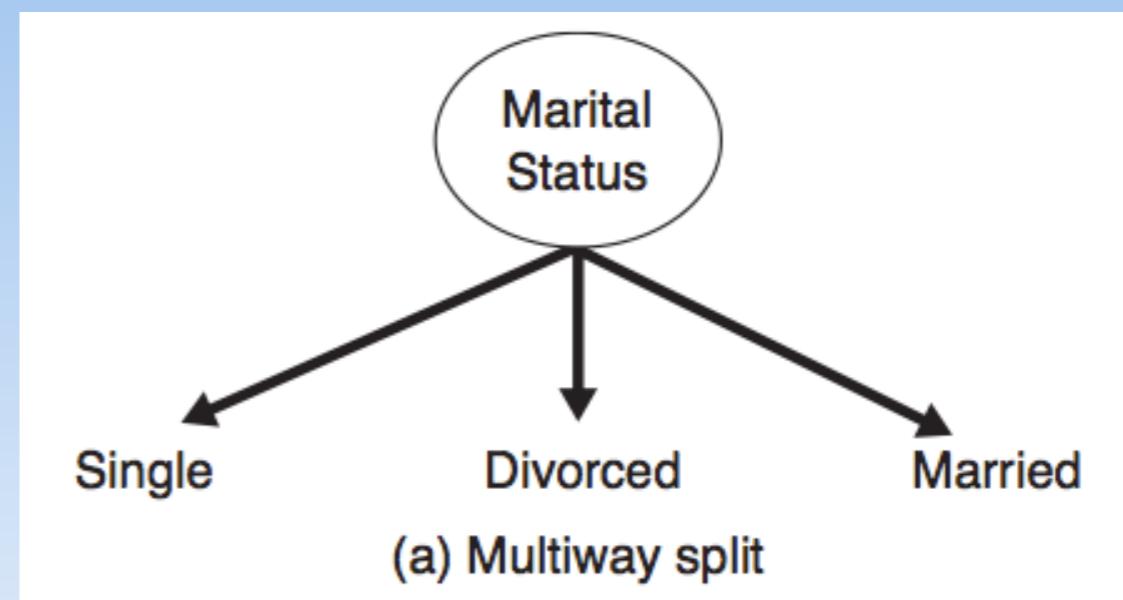
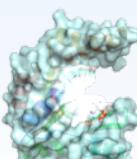


Figure 4.8. Test condition for binary attributes.

Multiway splits



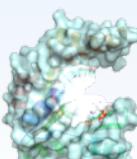
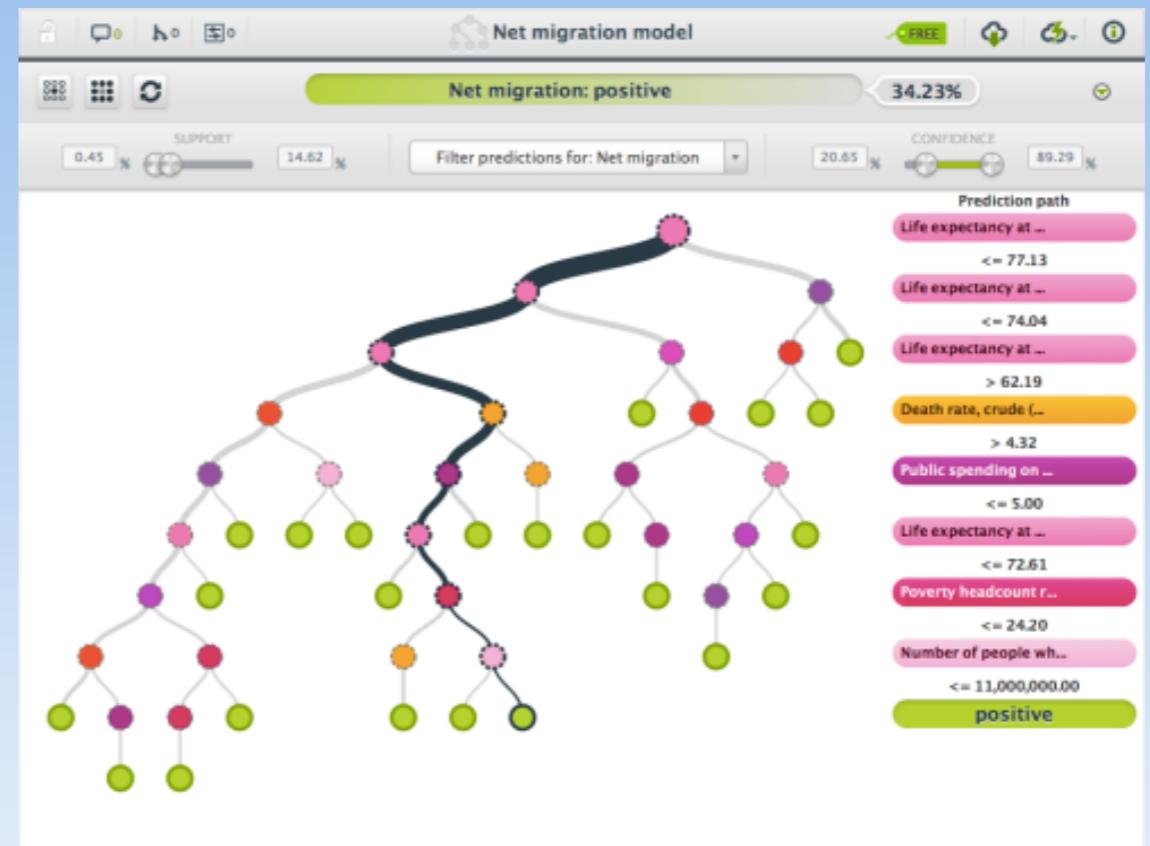
(a) Multiway split



Catalit LLC

APPLICATIONS

- Lead scoring
- Customer churn
- User segmentation

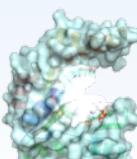
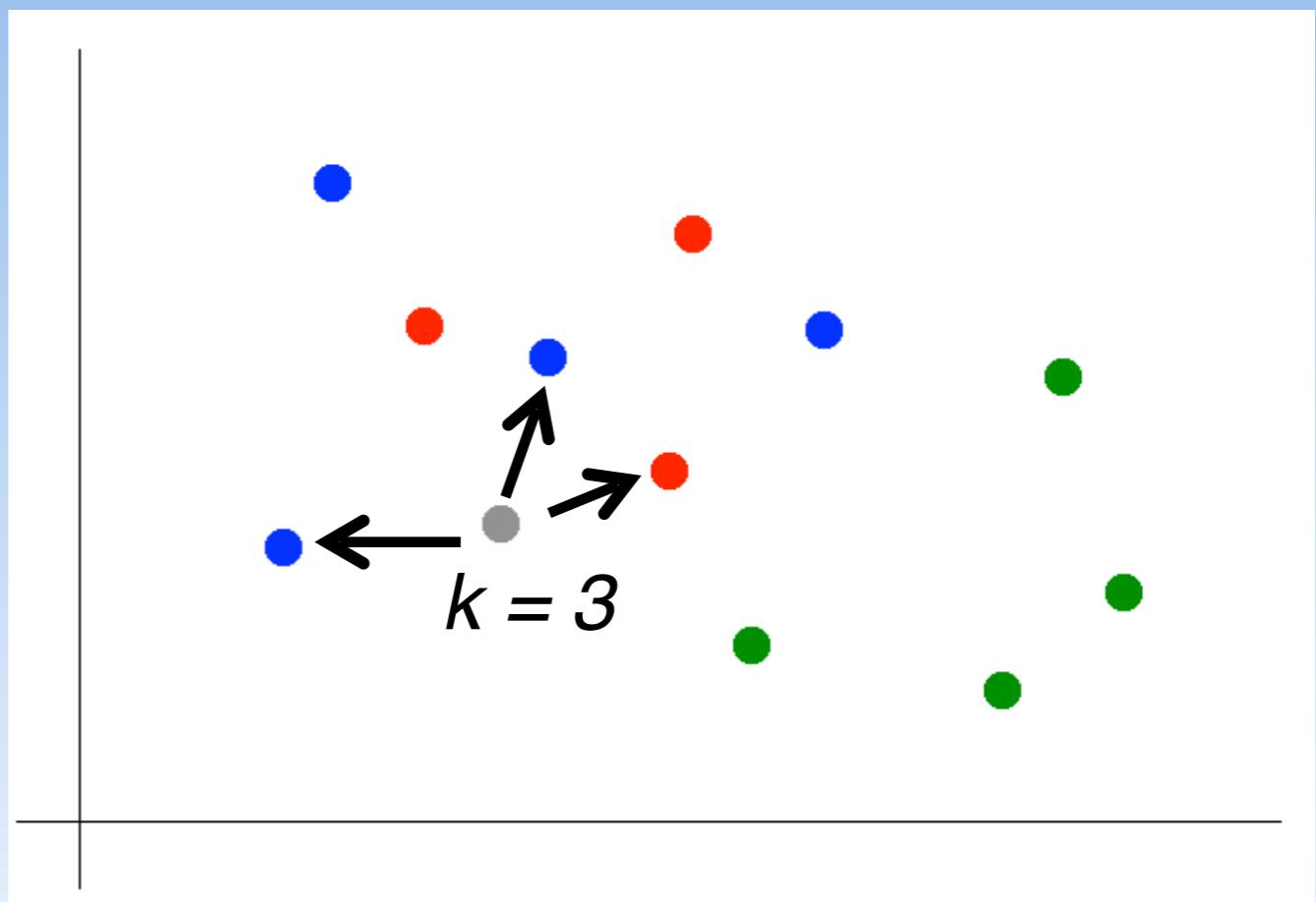


K-NEAREST NEIGHBORS

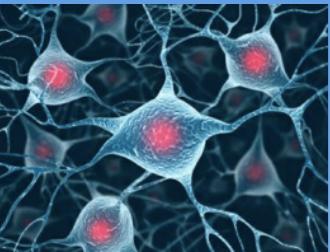
- GOAL: majority vote of nearby points

- Advantages:

- easy to understand
- any boundary



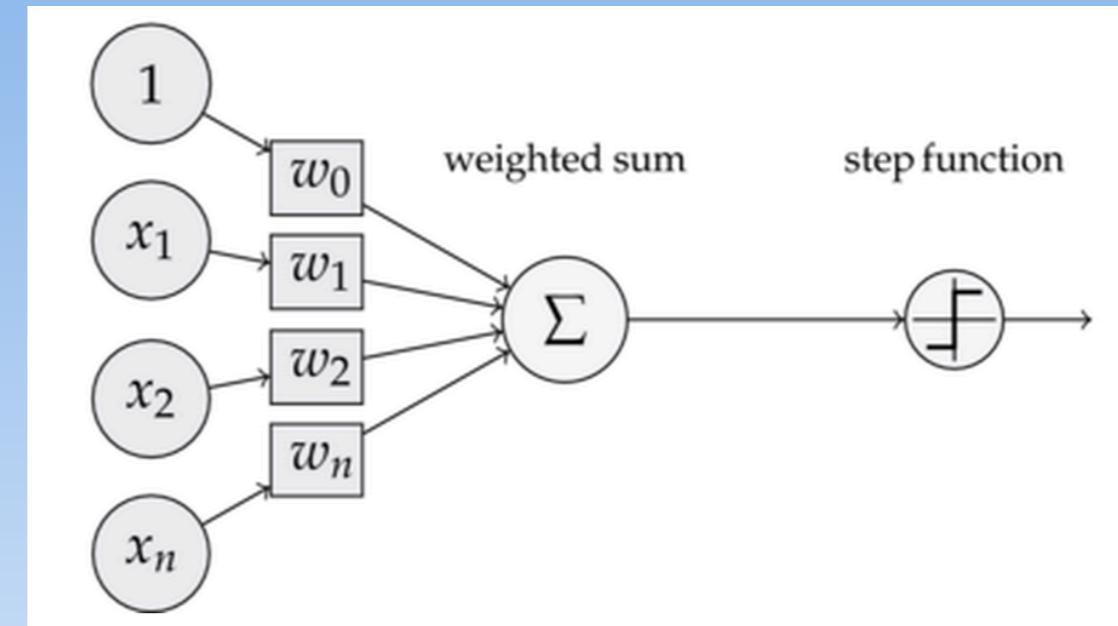
Catalit LLC



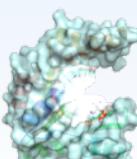
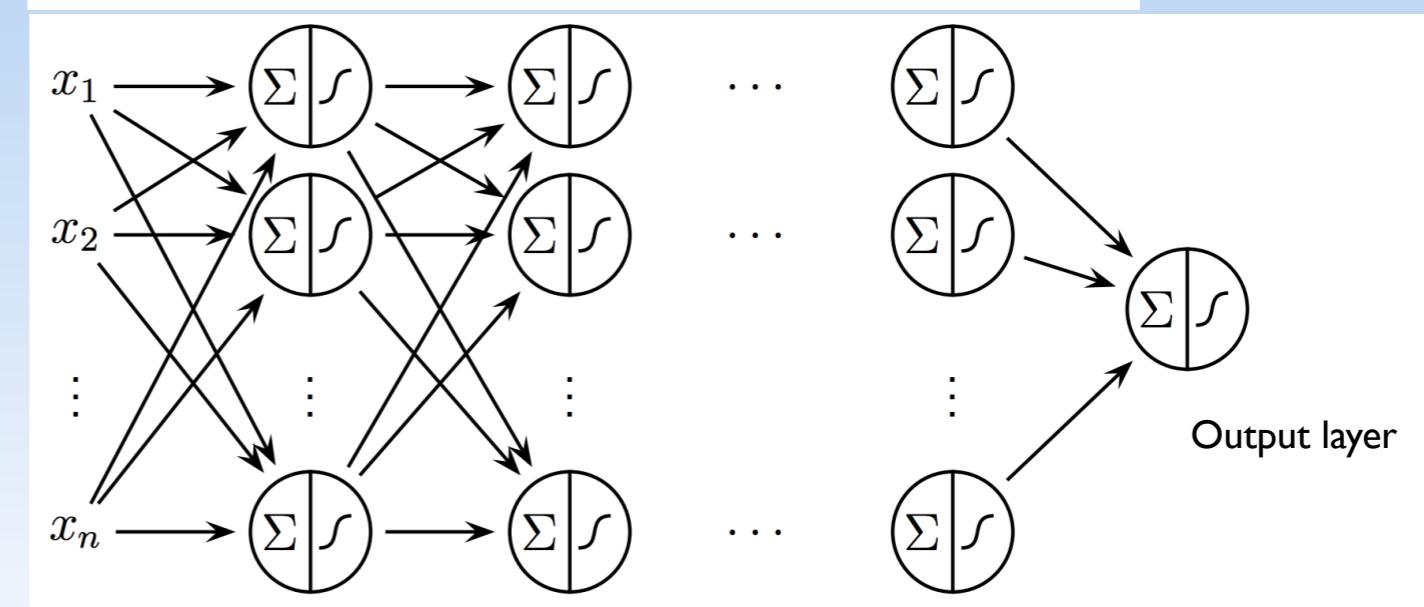
NEURAL NETWORKS

- GOAL: adjust weights to minimize output error

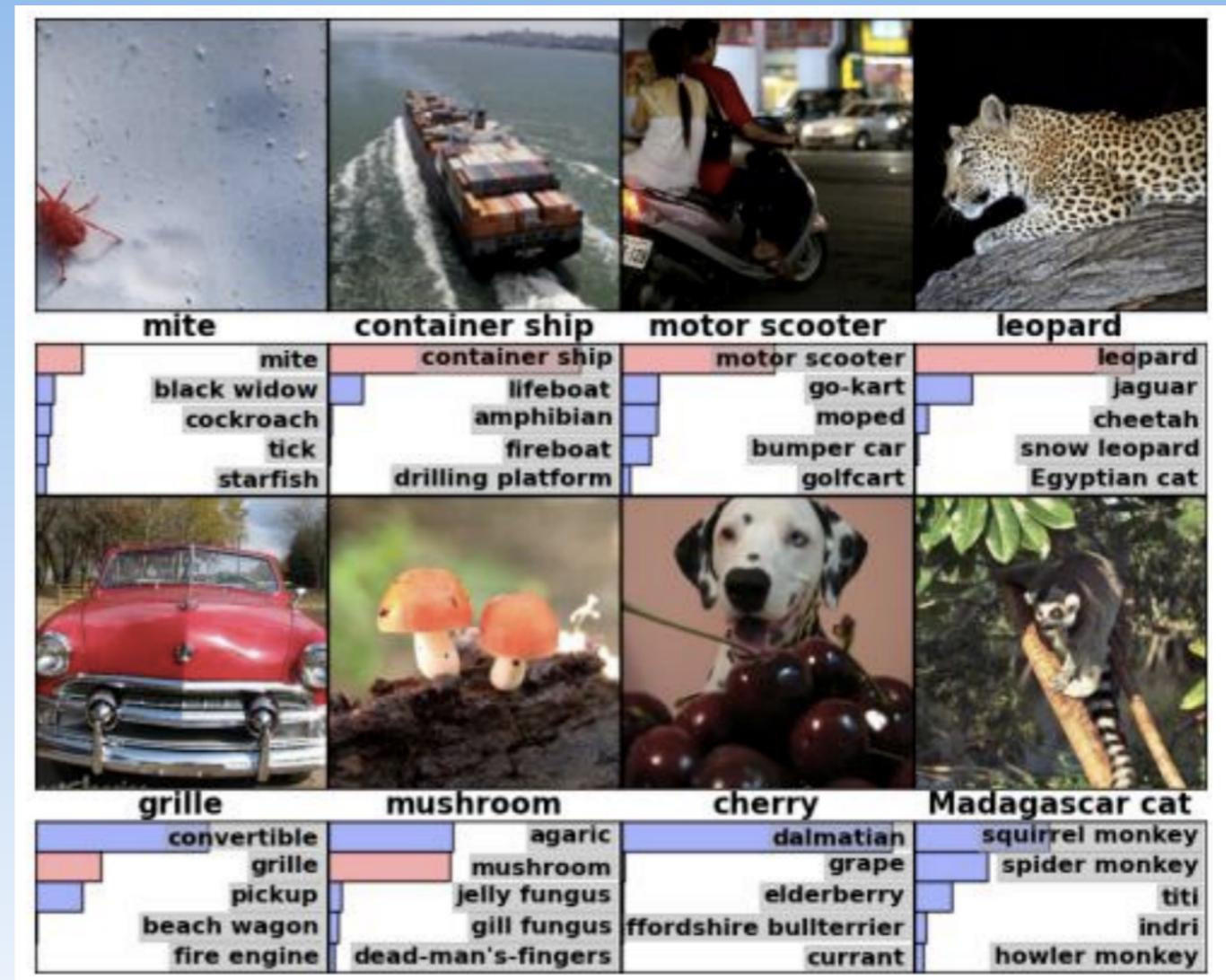
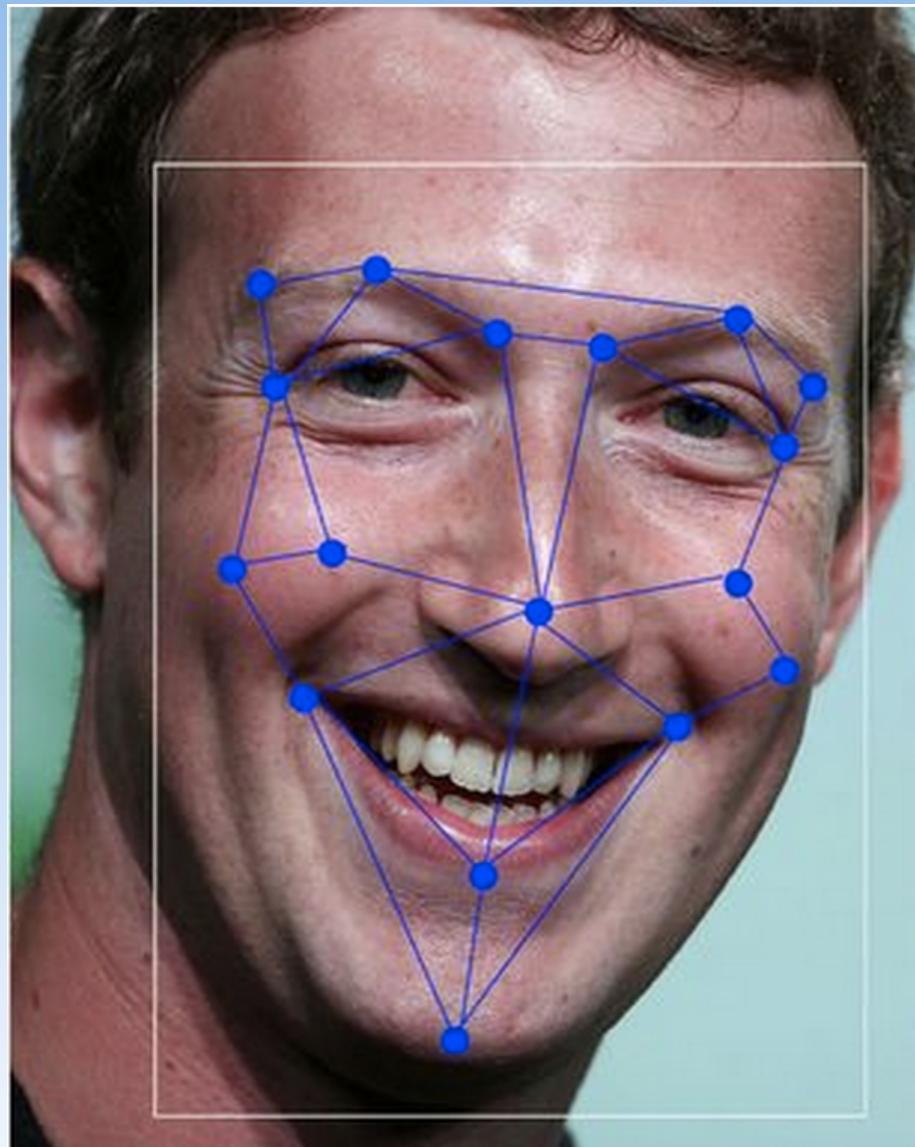
- Perceptron



- Multi-layer

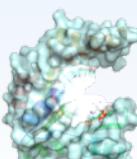
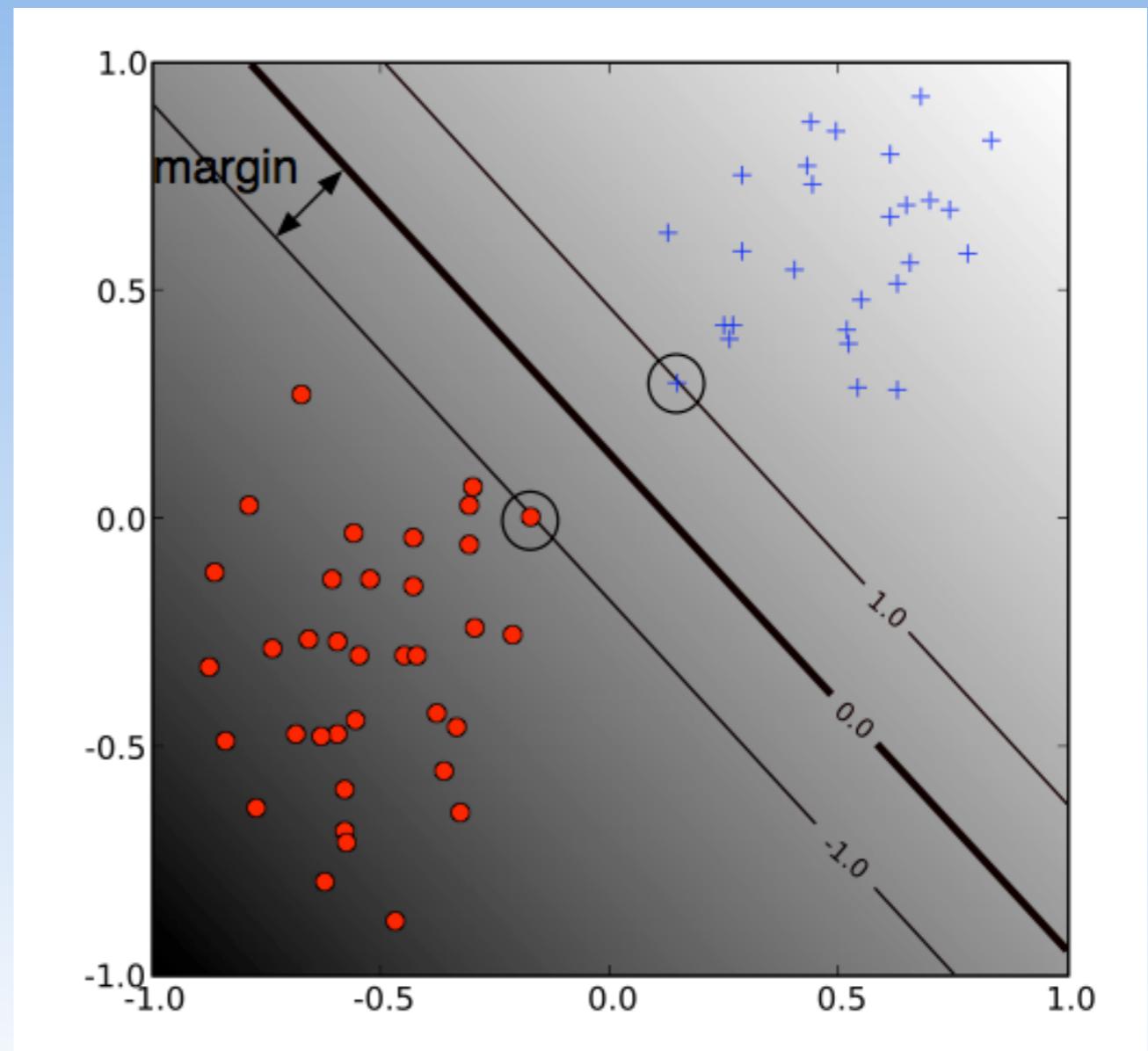


APPLICATIONS OF NEURAL NETWORKS



SUPPORT VECTOR MACHINE

- GOAL: Search for a separating hyperplane that maximizes margin
- Advantages:
 - powerful in non-linear boundaries
 - transform data to higher dimension



APPLICATIONS OF SVM

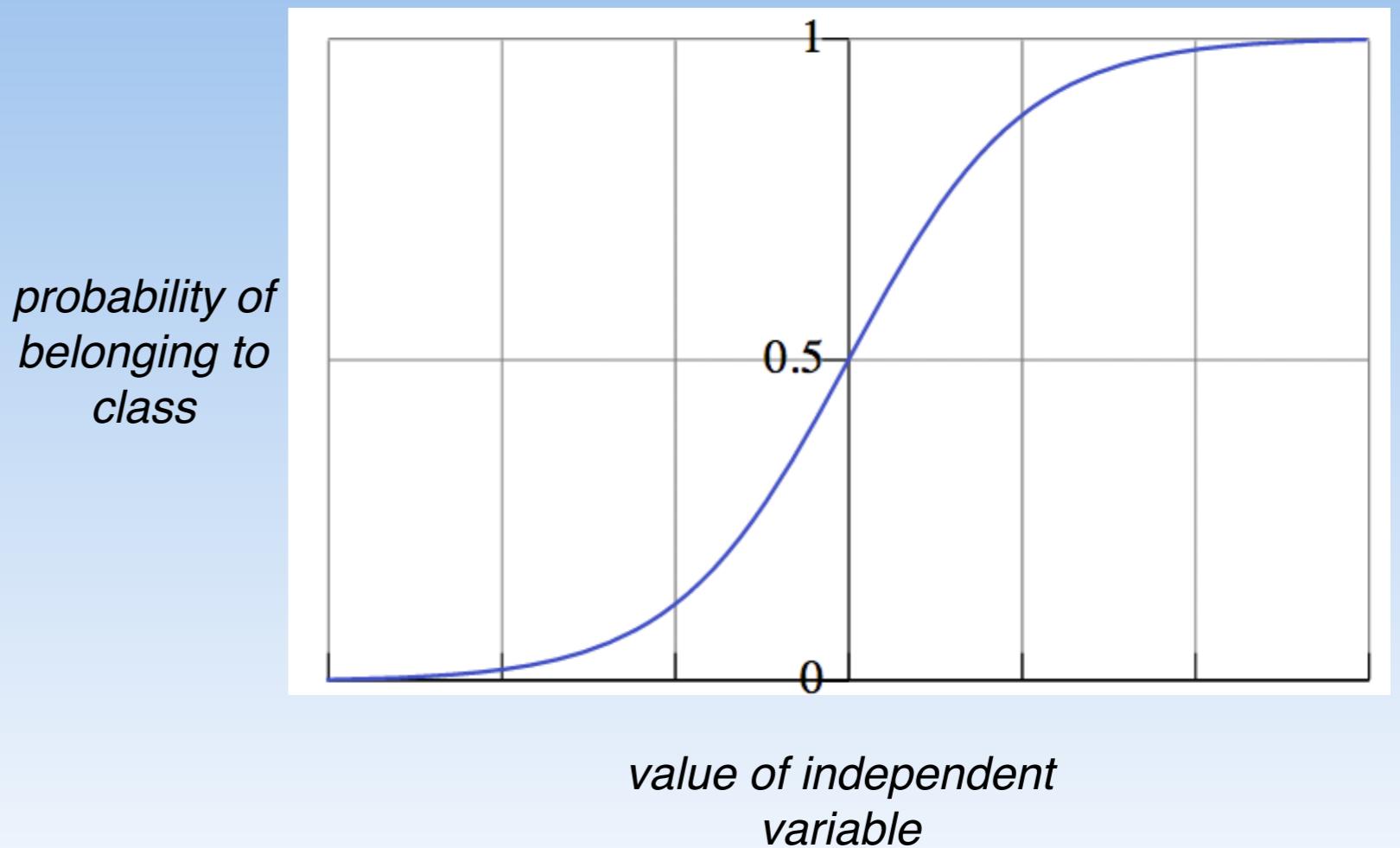
- Transmission
 - IP address — 167.12.24.55
 - Sender URL -- one-spam.com
- Email header
 - From -- “admin@one-spam.cpm”
 - To -- “undisclosed”
 - cc
- Email Body
 - # of paragraphs
 - # words
- Email structure
 - # of attachments
 - # of links



Catalit LLC

LOGISTIC REGRESSION

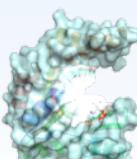
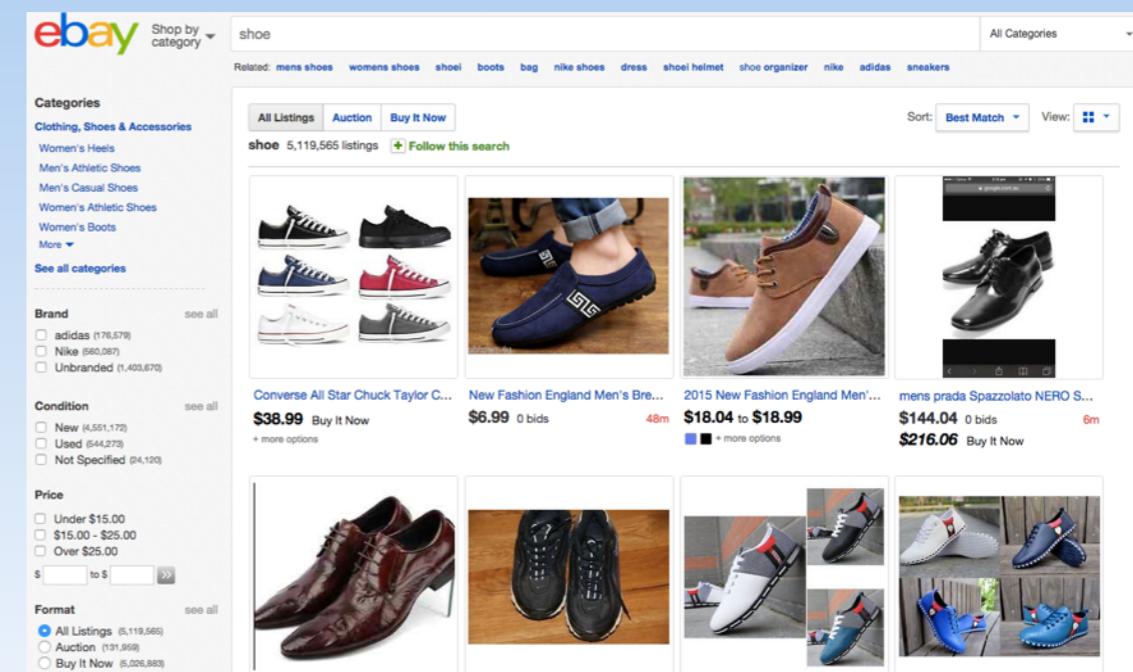
- GOAL: minimize squared error for probability
- Advantages:
 - simple function
 - can be parallelized
 - large scale



APPLICATIONS

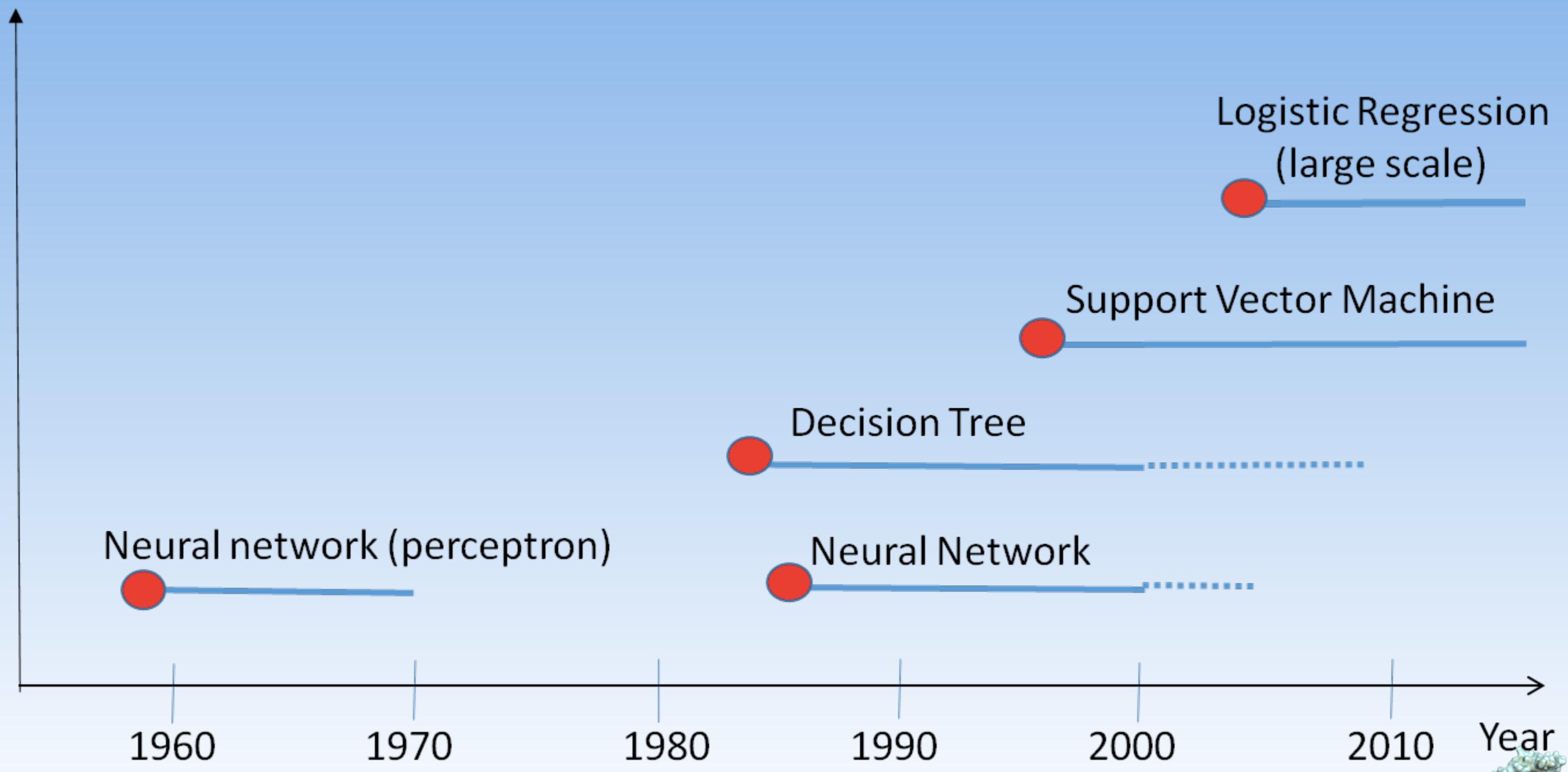
- Click prediction
 - Search ranking (web pages, products)
 - Online advertising
 - Recommendation

- The model
 - Output: Click/no click
 - Input features:
 - page content
 - search keyword
 - User information



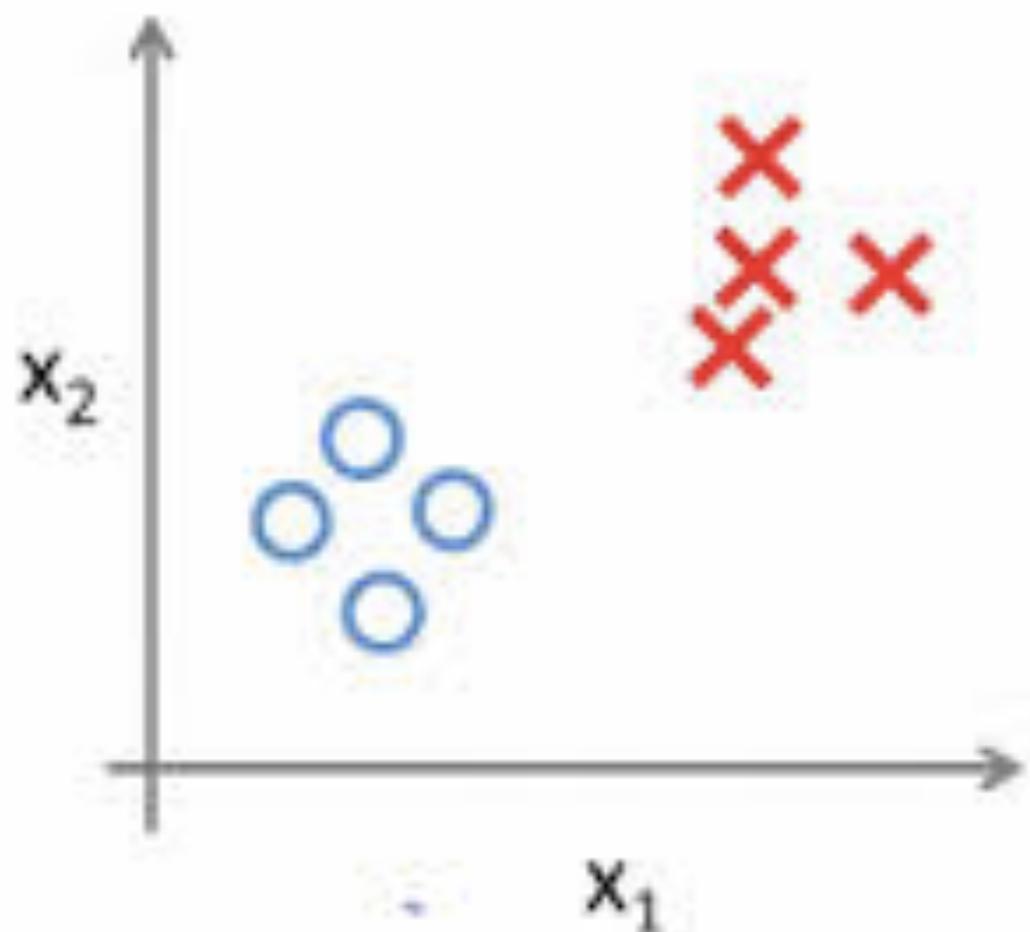
Catalit LLC

HISTORY OF SUPERVISED LEARNING

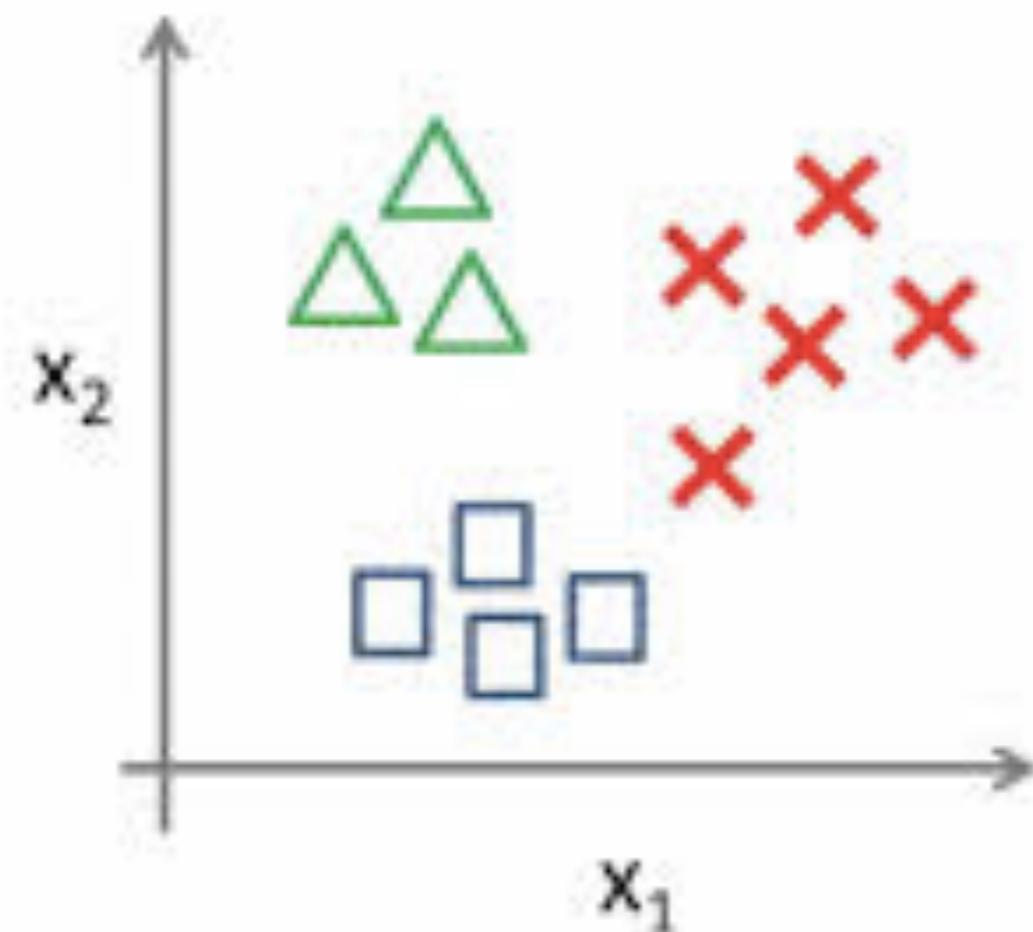


BINARY / MULTI CLASS

Binary classification:



Multi-class classification:



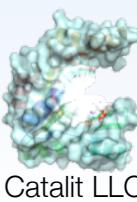
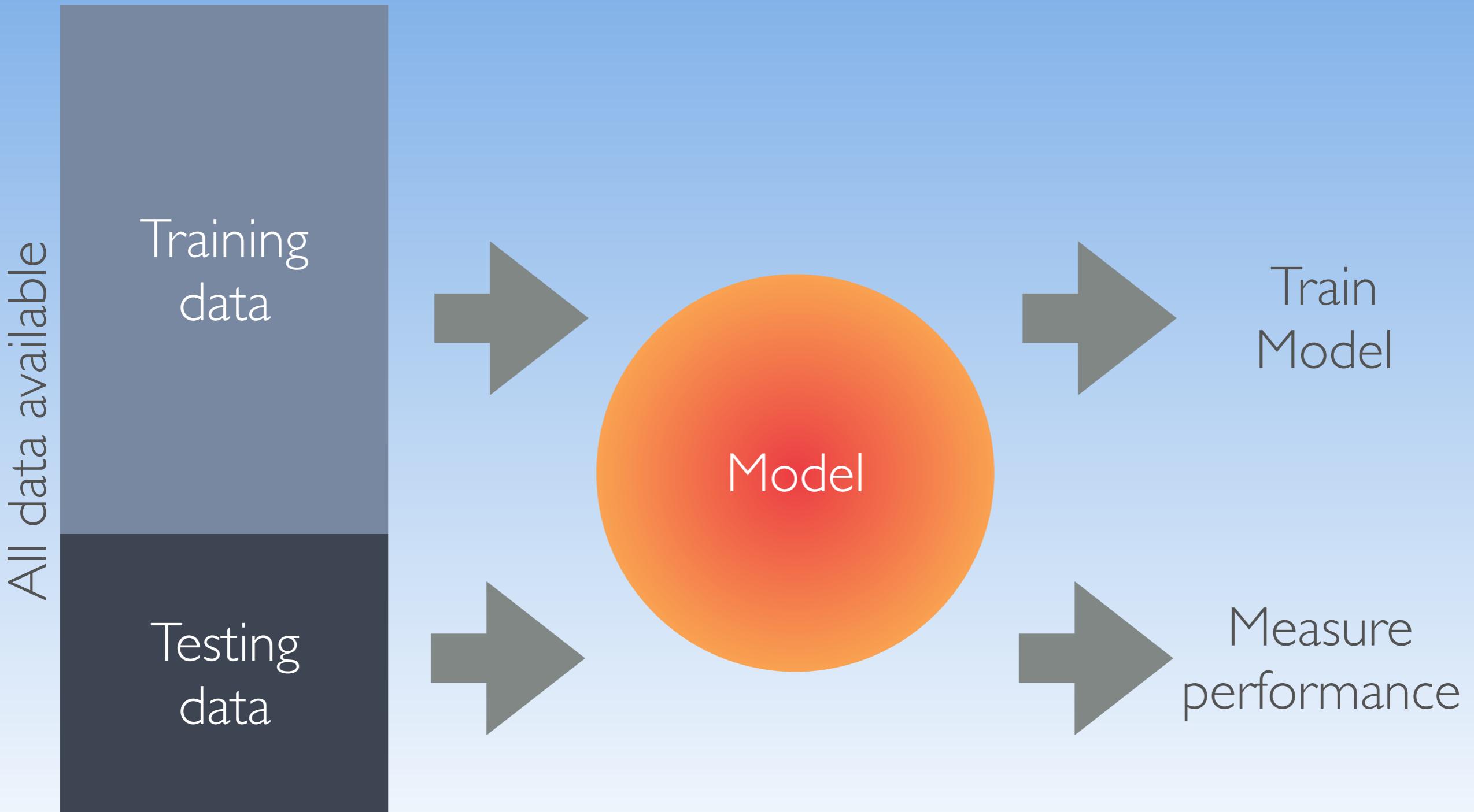
TRAIN - TEST SPLIT



http://todaysabundantliving.com/wp-content/uploads/2013/08/IMG_2221.jpg



TRAIN - TEST SPLIT



CONFUSION MATRIX

	<i>Test Negative</i>	<i>Test Positive</i>
<i>Condition Negative</i>	TRUE NEGATIVE	FALSE POSITIVE <i>(Type I error)</i>
<i>Condition Positive</i>	FALSE NEGATIVE <i>(Type II error)</i>	TRUE POSITIVE

CONFUSION MATRIX

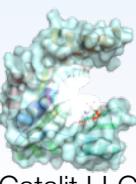
	Class A	Class B	Class C	Class D
Prediction A	84			
Prediction B		69	3	
Prediction C	2		78	
Prediction D				91

PRECISION - RECALL & ACCURACY

- **Precision:** When test is positive, how often is prediction correct?
 - TP / test yes
- **Recall:** When actual value is positive, how often is prediction correct?
 - TP / actual yes
- **Accuracy:** Overall, how often is it correct?
 - $(TP + TN) / \text{total}$

	<i>Test Negative</i>	<i>Test Positive</i>
<i>Condition Negative</i>	<i>TRUE NEGATIVE</i>	<i>FALSE POSITIVE</i>
<i>Condition Positive</i>	<i>FALSE NEGATIVE</i>	<i>TRUE POSITIVE</i>

LAB CLASSIFICATION



Catalit LLC