

Data Mining and Analytics

Aik Choon Tan, Ph.D.

Associate Professor of Bioinformatics

Division of Medical Oncology

Department of Medicine

aikchoon.tan@ucdenver.edu

9/14/2018

<http://tanlab.ucdenver.edu/labHomePage/teaching/BSBT6111/>

Outline

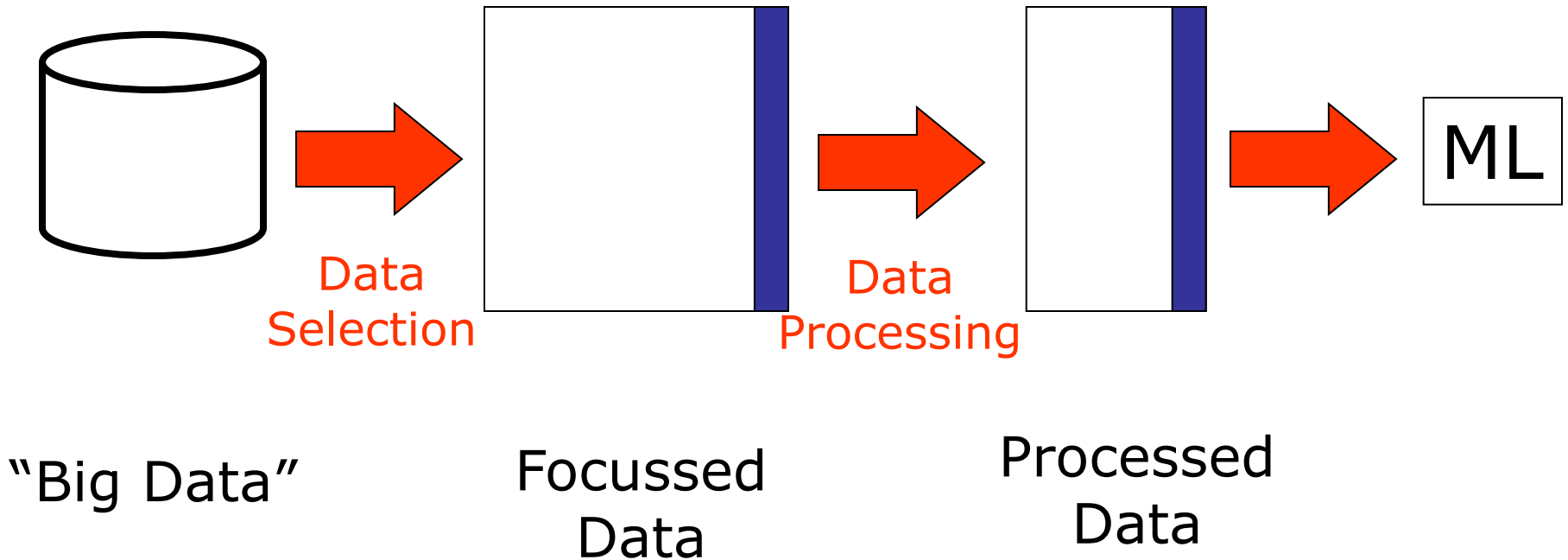


- Introduction
- Data, features, classifiers, Inductive learning
- (*Selected*) Machine Learning Approaches
 - Decision trees
 - Linear Model
 - Support Vector Machines
 - Artificial Neural Network
 - Deep Learning
 - Clustering
- Model evaluation

Steps in Class prediction problem

- Data Preparation
- Feature selection
 - Remove irrelevant features for constructing the classifier (but may have biological meaning)
 - Reduce search space in H , hence increase speed in generating classifier
 - Direct the learning algorithm to focus on “informative” features
 - Provide better understanding of the underlying process that generated the data
- Selecting a machine learning method
- Generating classifier from the training data
- Measuring the performance of the classifier
- Applying the classifier to unseen data (test)
- *Interpreting the classifier*

Data Preparation



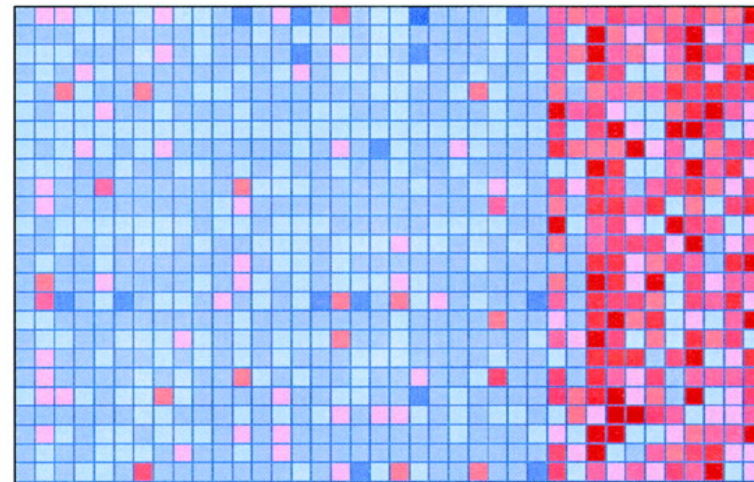
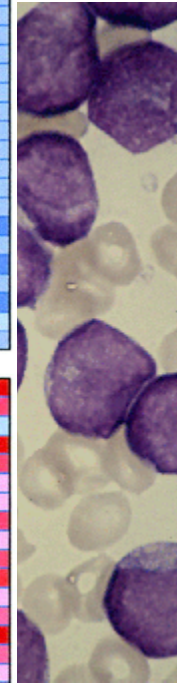
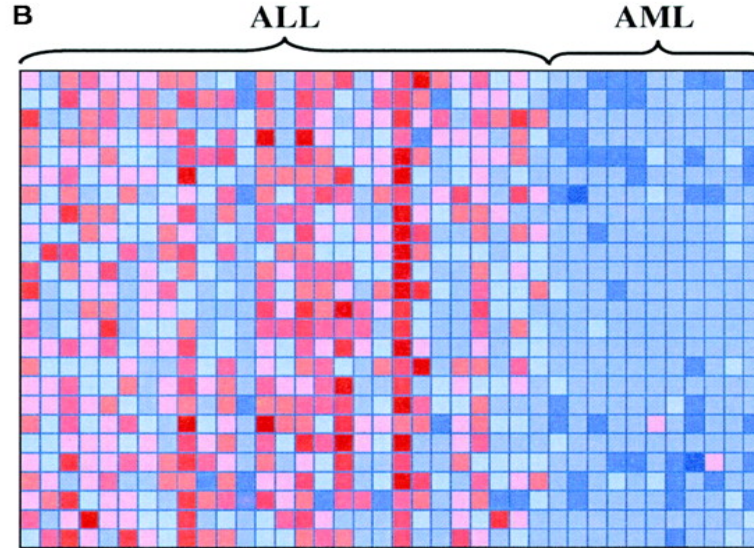
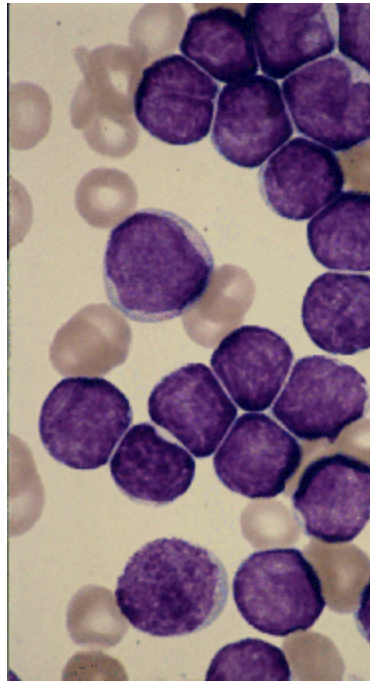
Data: Samples and Features

The diagram illustrates a data matrix with *Samples* as columns and *features* as rows. A horizontal double-headed arrow above the table is labeled *Samples*. A vertical double-headed arrow to the left of the table is labeled *features*. The table has five columns: 'Features', 'Sample 1', 'Sample 2', '...', and 'Sample *m*'. It has five rows: 'feature 1', 'feature 2', '...', 'feature *n*', and 'Feature_value' (repeated for each sample column). The cells for 'Sample 1', 'Sample 2', and 'Sample *m*' contain 'Feature_value' for each feature row. The cells for the '...' columns contain '...'.

| Features | Sample 1 | Sample 2 | ... | Sample <i>m</i> |
|------------------|---------------|---------------|-----|-----------------|
| feature 1 | Feature_value | Feature_value | ... | Feature_value |
| feature 2 | Feature_value | Feature_value | ... | Feature_value |
| ... | ... | ... | ... | ... |
| feature <i>n</i> | Feature_value | Feature_value | ... | Feature_value |

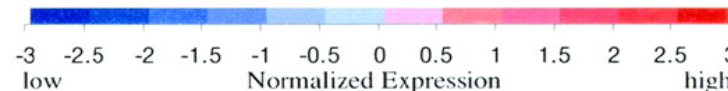
Cancer Classification Problem

(Golub et al 1999)




ALL
acute lymphoblastic
(lymphoid precursor)

AML
acute myeloid leukemia
(myeloid precursor)



Gene Expression Profile



m samples

n genes

| Geneid | Condition 1 | Condition 2 | ... | Condition m |
|----------|-------------|-------------|-----|---------------|
| Gene1 | 103.02 | 58.79 | ... | 101.54 |
| Gene2 | 40.55 | 1246.87 | ... | 1432.12 |
| ... | ... | ... | ... | ... |
| Gene n | 78.13 | 66.25 | ... | 823.09 |

A (very) Brief Introduction to Machine Learning



To Learn

“ ... to acquire *knowledge* of (a subject) or skill in (an art, etc.) as a result of *study*, *experience*, or *teaching*... ”
(OED)

What is Machine Learning?

“ ... a computer program that can learn from *experience* with respect to some class of *tasks* and *performance measure* ... ”
(Mitchell, 1997)

Key Steps of Learning

- Learning task
 - what is the learning task?
- Data and assumptions
 - what data is available for the learning task?
 - what can we assume about the problem?
- Representation
 - how should we represent the examples to be classified
- Method and estimation
 - what are the possible hypotheses?
 - how do we adjust our predictions based on the feedback?
- Evaluation
 - how well are we doing?
- Model selection
 - can we rethink the approach to do even better?

Learning Tasks

- Classification – Given positive and negative examples, find hypotheses that distinguish these examples. It can extend to multi-class classification.
- Clustering – Given a set of unlabelled examples, find clusters for these examples (unsupervised learning)

Learning Approaches

- Supervised approach – given *predefined* class of a set of positive and negative examples, construct the classifiers that distinguish between the classes $\langle \mathbf{x}, y \rangle$
- Unsupervised approach – given the *unassigned* examples, group together the examples with similar properties $\langle \mathbf{x} \rangle$

Concept Learning

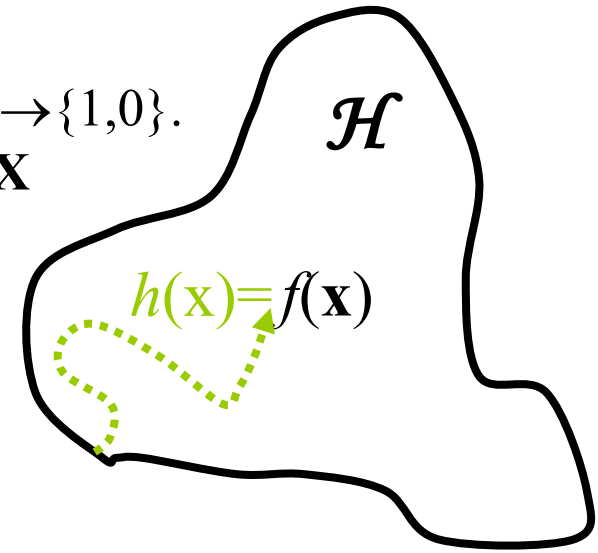
Given a set of training examples $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ where \mathbf{x} is the instances usually in the form of tuple $\langle x_1, \dots, x_n \rangle$ and y is the class label, the function $y = f(\mathbf{x})$ is unknown and finding the $f(\mathbf{x})$ represent the essence of concept learning.

For a binary problem $y \in \{1, 0\}$, the unknown function $f: \mathbf{X} \rightarrow \{1, 0\}$.
The learning task is to find a hypothesis $h(\mathbf{x}) = f(\mathbf{x})$ for $\mathbf{x} \in \mathbf{X}$

Training examples $\langle \mathbf{x}, f(\mathbf{x}) \rangle$ where:

$f(\mathbf{x}) = 1$ are Positive examples,

$f(\mathbf{x}) = 0$ are Negative examples.



\mathcal{H} is the set of all possible hypotheses, where $h: \mathbf{X} \rightarrow \{1, 0\}$

A machine learning task:

Find hypothesis, $h(\mathbf{x}) = c(\mathbf{x})$; $\mathbf{x} \in \mathbf{X}$.

(in reality, usually ML task is to approximate $h(\mathbf{x}) \cong c(\mathbf{x})$)

Inductive Learning

- Given a set of observed examples
- Discover concepts from these examples
 - class formation/partition
 - formation of relations between objects
 - patterns

Learning paradigms

- Discriminative (model $\Pr(y|\mathbf{x})$)
 - only model decisions given the input examples; no model is constructed over the input examples
- Generative (model $\Pr(\mathbf{x}|y)$)
 - directly build class-conditional densities over the multidimensional input examples
 - classify new examples based on the densities

Decision Trees

- Widely used - simple and practical
- Quinlan - ID3 (1986), C4.5 (1993) & See5/C5 (latest)
- Classification and Regression Tree (CART by Breiman et.al., 1984)
- Given a set of instances (with a set of properties/attributes), the learning system constructs a tree with internal *nodes* as an *attribute* and the *leaves* as the *classes*
- Supervised learning
- Symbolic learning, give interpretable results

Information Theory - Entropy

Entropy – a measurement commonly used in information theory to characterise the (im)purity of an arbitrary collection of examples

$$Entropy(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

where S is a collection of training examples with c classes and p_i is the proportion of examples S belonging to class i .

Example:

If S is a set of examples containing positive (+) and negative (-) examples ($c \in \{+, -\}$), the entropy of S relative of this boolean classification is:

$$Entropy(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

$$Entropy(S) = \begin{cases} 0 & \text{if all members of } S \text{ belong to the same class} \\ 1 & \text{if } S \text{ contains an equal number of positive (+) and negative (-) examples} \end{cases}$$

Note*: Entropy \downarrow Purity \uparrow

ID3 (Induction of Decision Tree)

- Average entropy of attribute A

$$\hat{E}_A = \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

v = all the values of attribute A, S = training examples, S_v = training examples of attribute A with value v

$$E_A = \begin{cases} 0 & \text{if all members of } S \text{ belong to the same value } v \\ 1 & \text{if } S \text{ contains an equal number of value } v \text{ examples} \end{cases}$$

Note: Entropy ↓ Purity ↑*

Splitting rule of ID3 (Quinlan, 1986)

- Information Gain

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Note: Gain ↑ Purity ↑*

Decision Tree Algorithm

Function *Decision_Tree_Learning* (*examples*, *attributes*, *target*)

Inputs: *examples* = set of training examples

attributes = set of attributes

target = class label

1. **if** *examples* is empty **then return** *target*
2. **else if** all *examples* have the same *target* **then return** *target*
3. **else if** *attributes* is empty then return most common value of *target* in *examples*
4. **else**
5. *Best* \leftarrow the attribute from *attributes* that best classifies *examples*
6. *Tree* \leftarrow a new decision tree with root attribute *Best*
7. **for** each value v_i of *Best* **do**
8. $examples_i \leftarrow \{\text{elements with } Best = v_i\}$
9. $subtree \leftarrow \text{Decision_Tree_Learning}(examples_i, attributes - best, target)$
10. add a branch to *Tree* with label v_i and subtree *subtree*
11. **end**
12. **return** *Tree*

Training Data

Decision attributes
(dependent)

Independent condition attributes

| Day | outlook | temperature | humidity | windy | play |
|-----|----------|-------------|----------|-------|------|
| 1 | sunny | hot | high | FALSE | no |
| 2 | sunny | hot | high | TRUE | no |
| 3 | overcast | hot | high | FALSE | yes |
| 4 | rainy | mild | high | FALSE | yes |
| 5 | rainy | cool | normal | FALSE | yes |
| 6 | rainy | cool | normal | TRUE | no |
| 7 | overcast | cool | normal | TRUE | yes |
| 8 | sunny | mild | high | FALSE | no |
| 9 | sunny | cool | normal | FALSE | yes |
| 10 | rainy | mild | normal | FALSE | yes |
| 11 | sunny | mild | normal | TRUE | yes |
| 12 | overcast | mild | high | TRUE | yes |
| 13 | overcast | hot | normal | FALSE | yes |
| 14 | rainy | mild | high | TRUE | no |

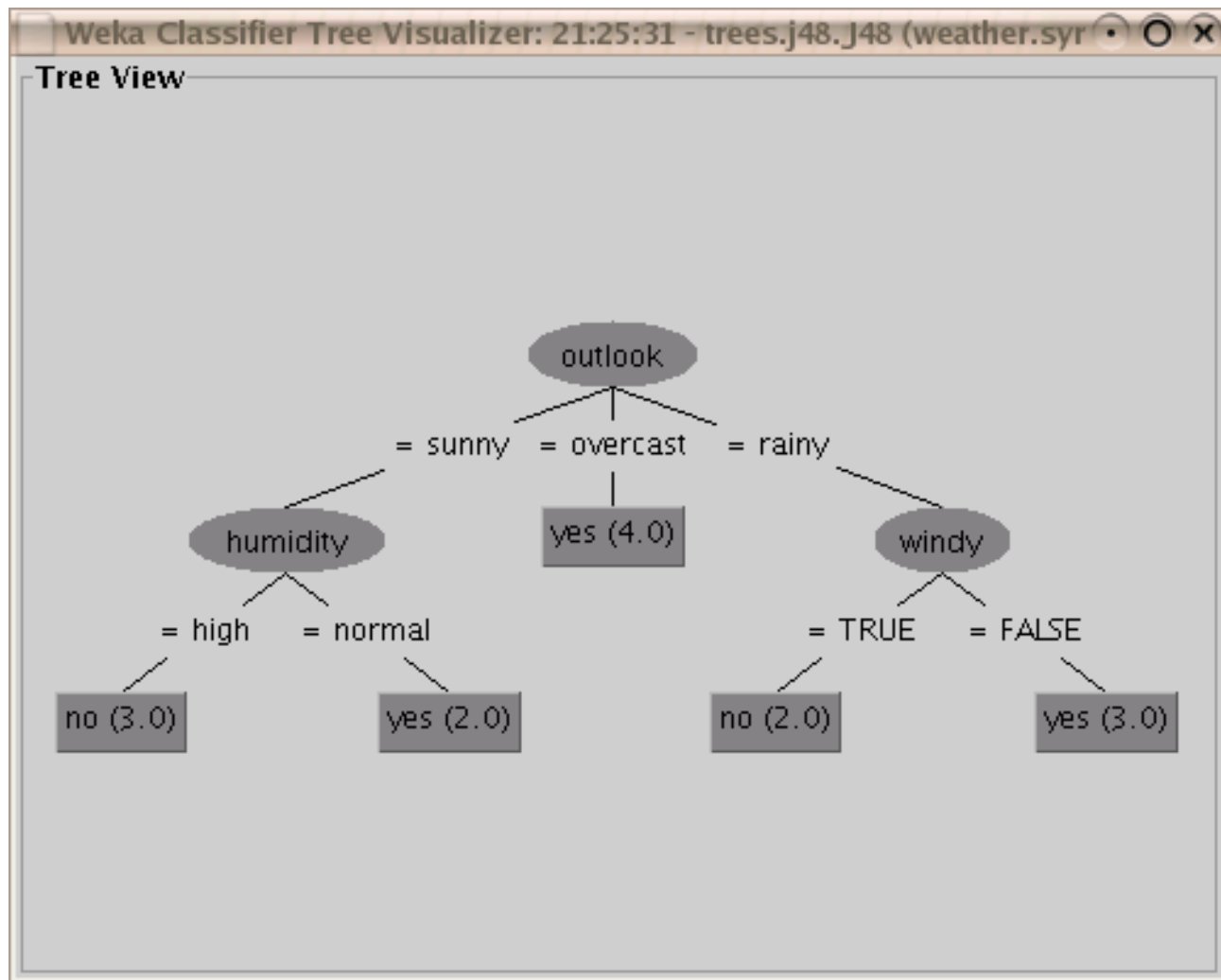
| | | | | | |
|-------|-------|------|------|------|---|
| Today | sunny | cool | high | TRUE | ? |
|-------|-------|------|------|------|---|



Entropy(S)



Decision Tree



Decision Trees (Quinlan, 1993)

J48 pruned tree

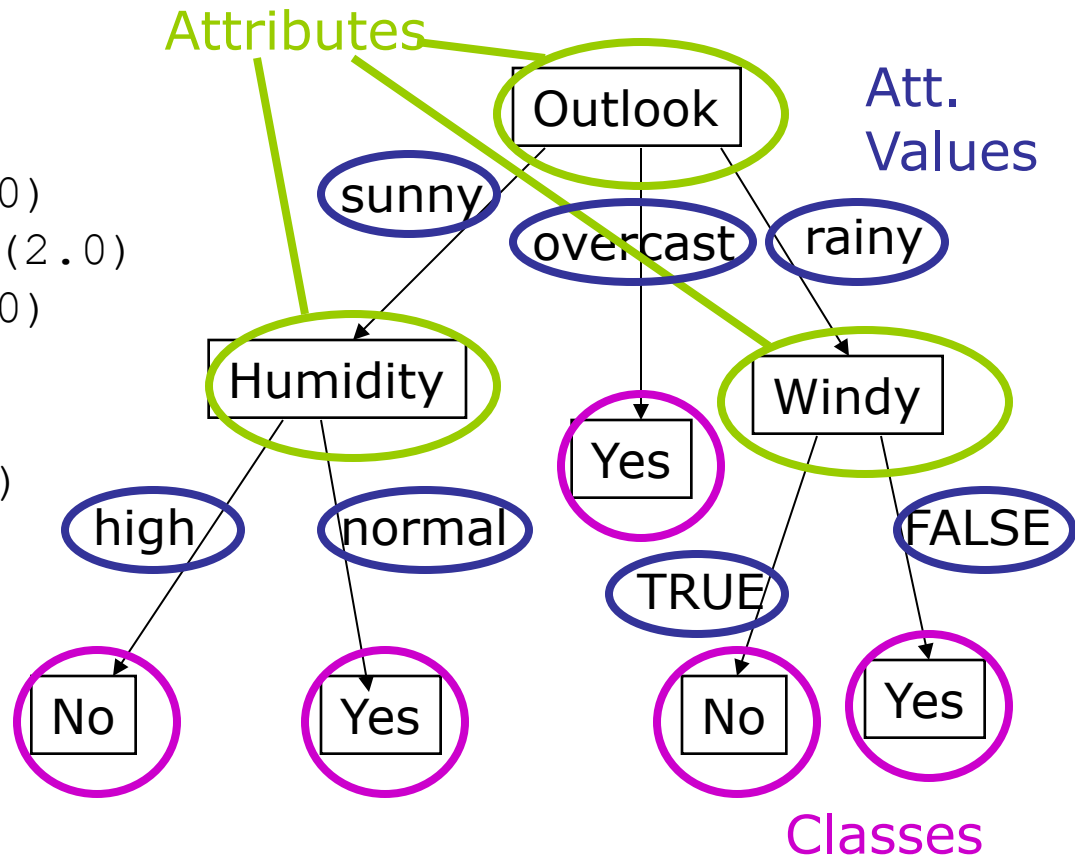
outlook = sunny
| humidity = high: no (3.0)
| humidity = normal: yes (2.0)
outlook = overcast: yes (4.0)
outlook = rainy
| windy = TRUE: no (2.0)
| windy = FALSE: yes (3.0)

Number of Leaves : 5

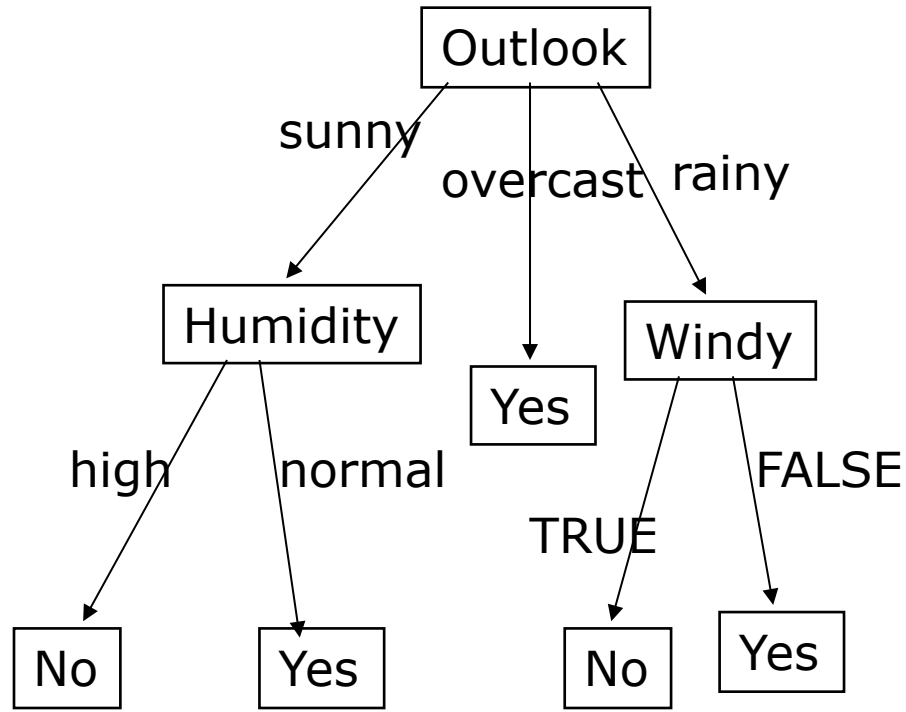
Size of the tree : 8

Time taken to build model: 0.05 seconds

Time taken to test model on training data: 0 seconds



Converting Trees to Rules



R1: IF Outlook = sunny \wedge Humidity = high THEN play = No

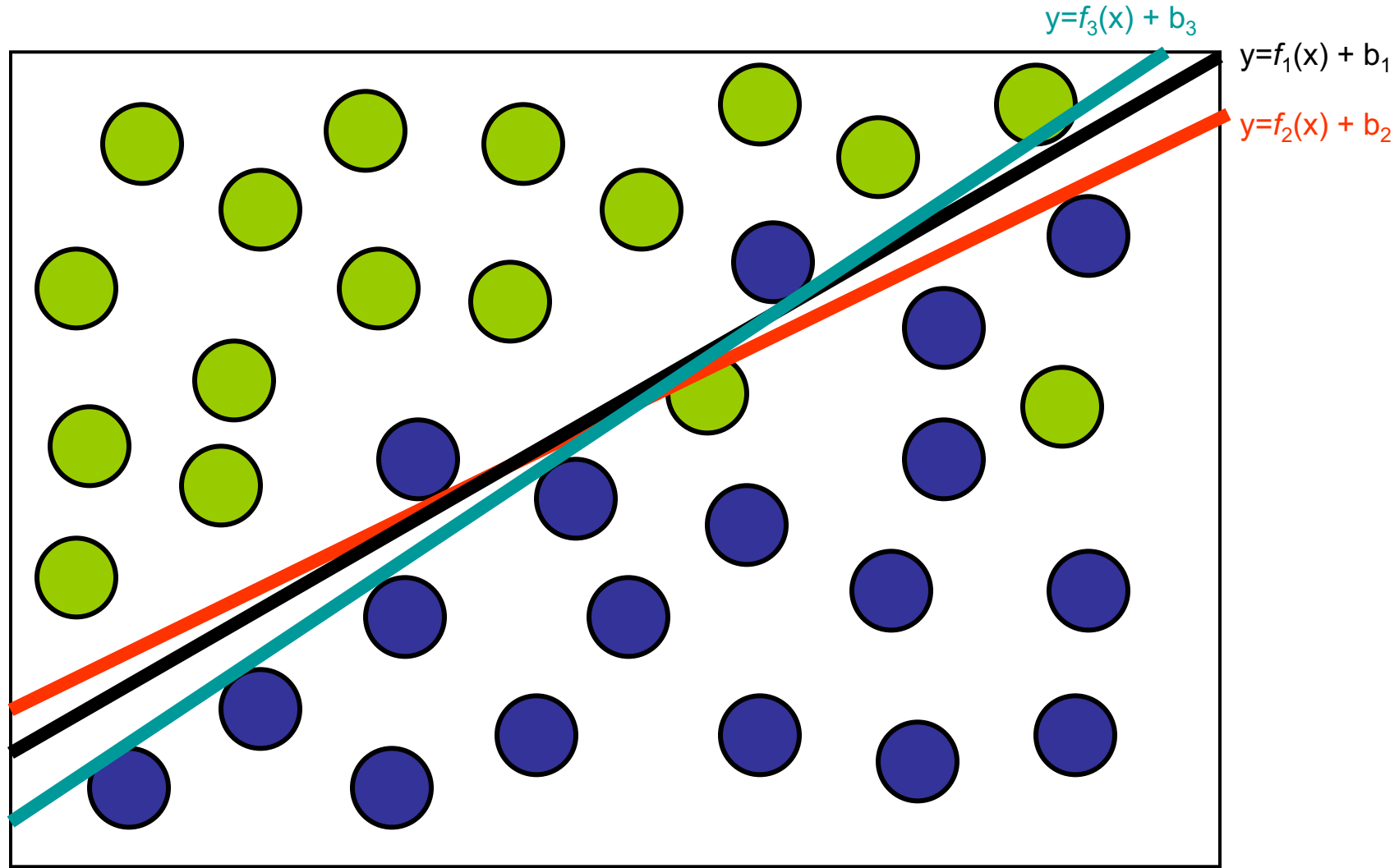
R2: IF Outlook = sunny \wedge Humidity = normal THEN play = Yes

R3: IF Outlook = overcast THEN play = Yes

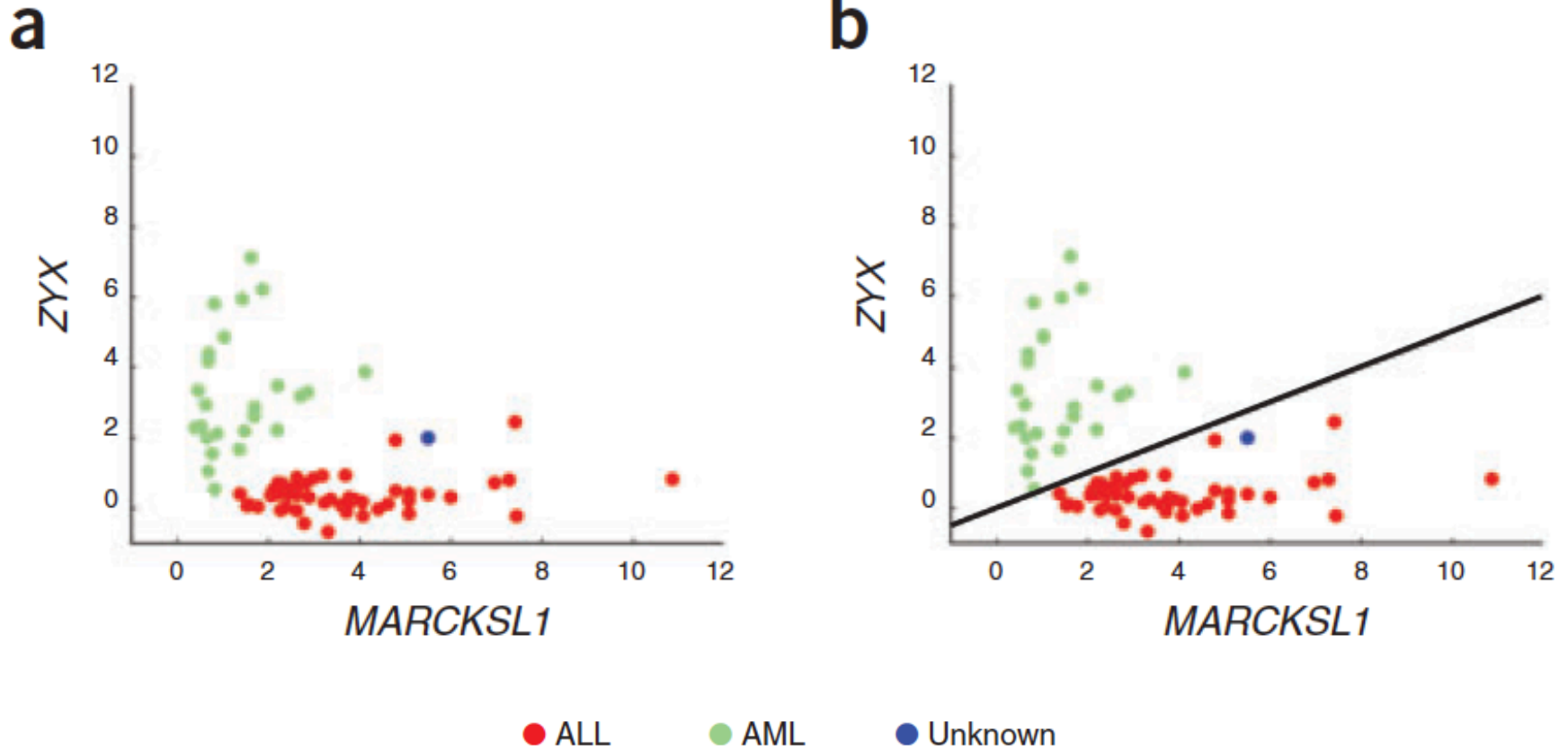
R4: IF Outlook = rainy \wedge Windy = TRUE THEN play = No

R5: IF Outlook = rainy \wedge Windy = FALSE THEN play = Yes

Linear Model



Straight Line as Classifier in 2D Space



Support Vector Machines (SVM)

Key concepts:

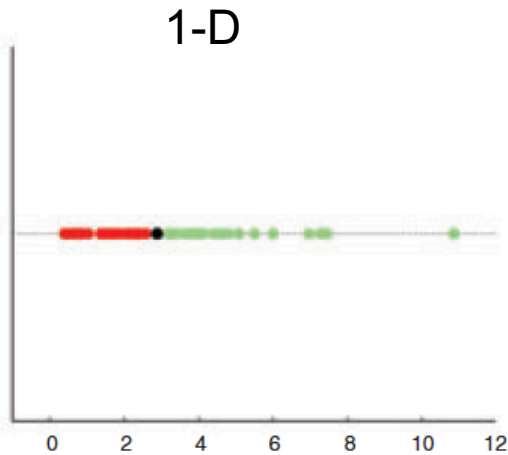
Separating hyperplane – straight line in high-dimensional space

Maximum-margin hyperplane - the distance from the separating hyperplane to the nearest expression vector as the margin of the hyperplane
Selecting this particular hyperplane maximizes the SVM's ability to predict the correct classification of previously unseen examples.

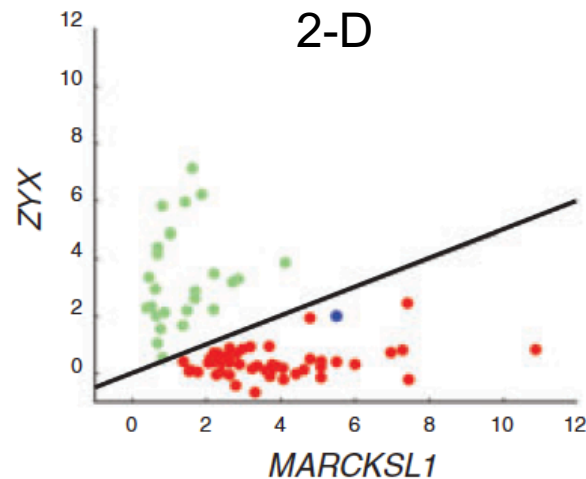
Soft margin - allows some data points (“soften”) to push their way through the margin of the separating hyperplane without affecting the final result.
User-specified parameter.

Kernel function - mathematical trick that projects data from a low-dimensional space to a space of higher dimension. The goal is to choose a good kernel function to separate data in high-dimensional space.

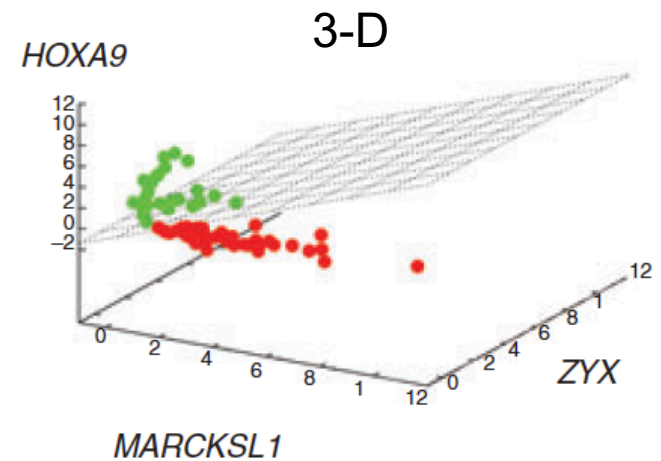
Separating Hyperplane



Separating
hyperplane =
dot



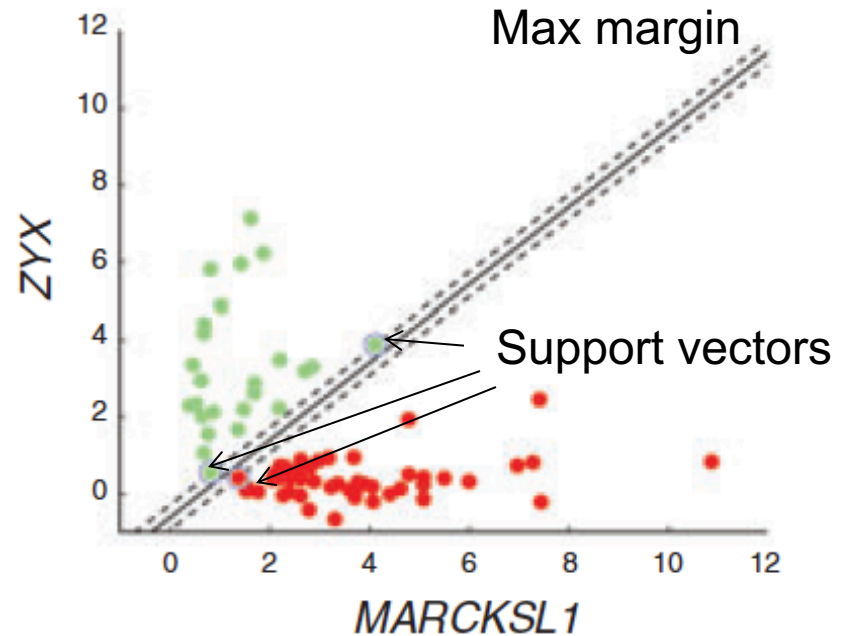
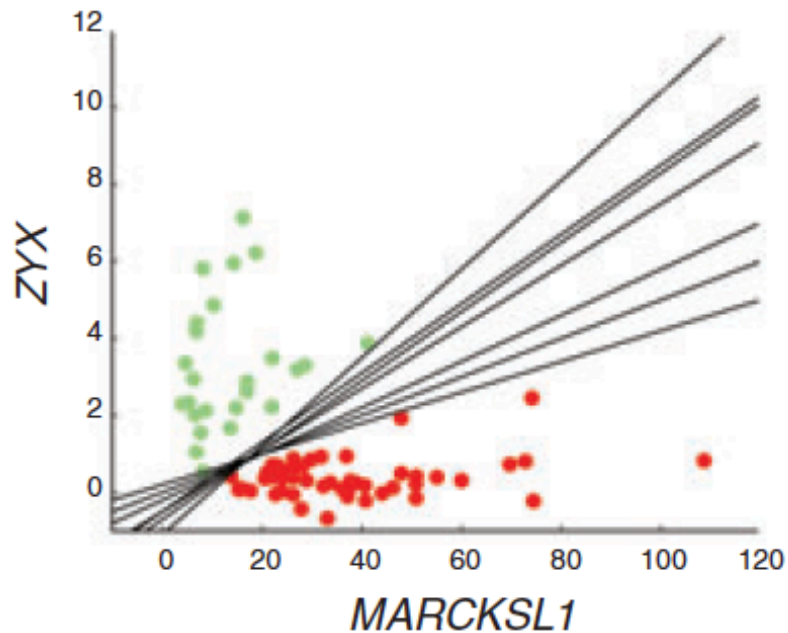
Separating
hyperplane =
line



Separating
hyperplane =
hyperplane

● ALL ● AML ● Unknown

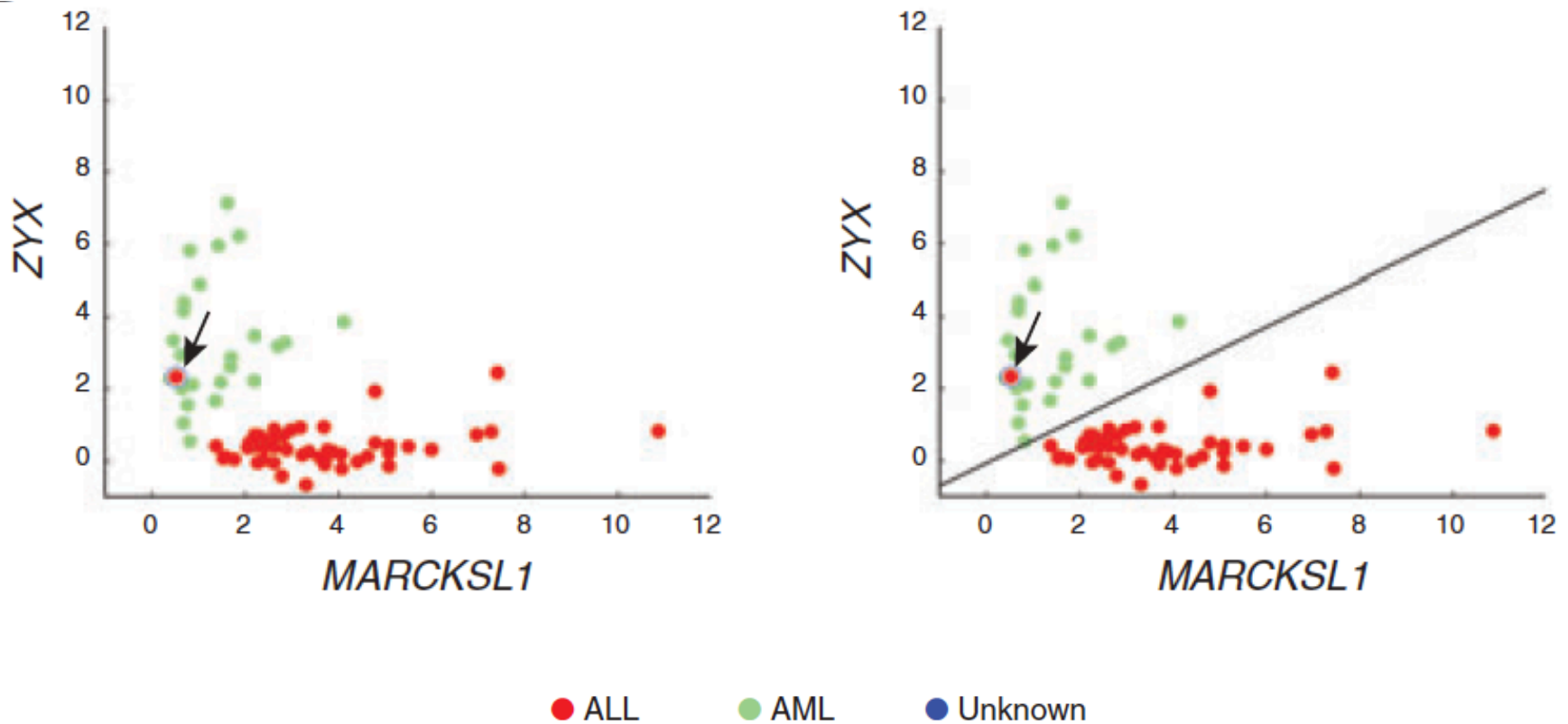
Maximum-margin Hyperplane



● ALL ● AML ● Unknown

(Adapted from Noble 2006)

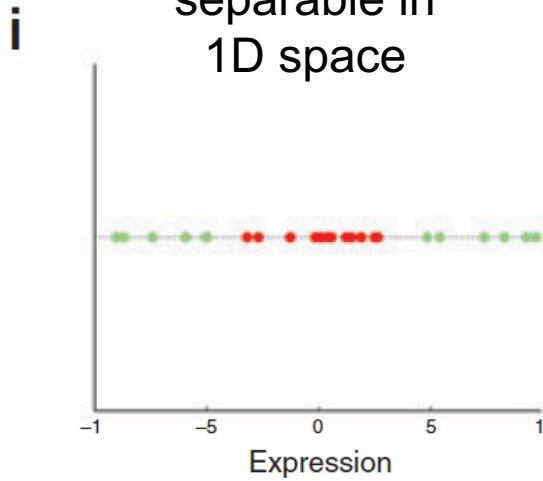
Soft Margin



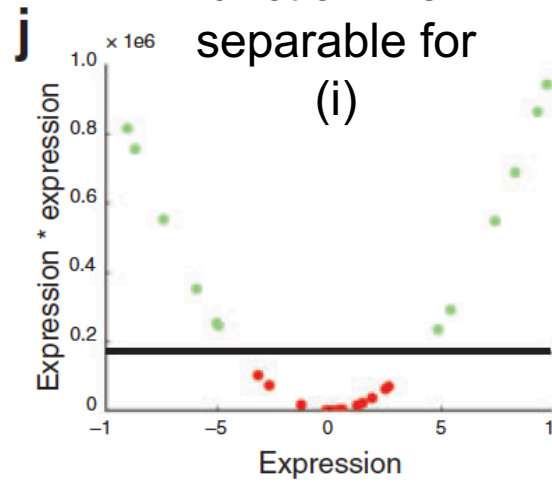
(Adapted from Noble 2006)

Kernel Function

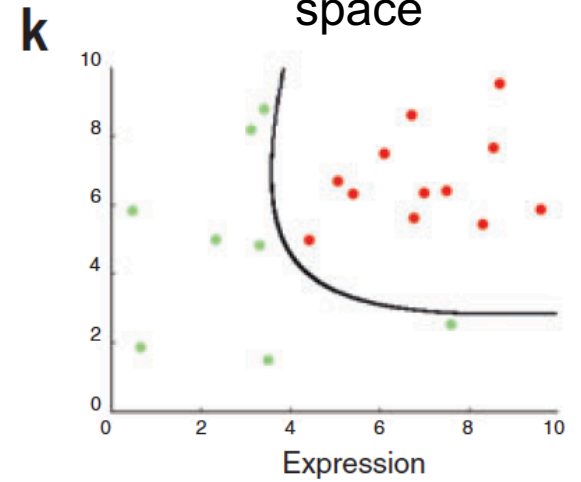
Not
separable in
1D space



Kernel
function: now
separable for
(i)



Kernel
function: 4-D
space



● ALL ● AML ● Unknown

Kernal: Linear SVM = Linear Regression

Predictive
Accuracy
= 50%

Support
Vectors = 0

```
SVM
Classifier for classes: yes, no
BinarySVM
Machine linear: showing attribute weights, not support vectors.

  0.8440785904115866 * outlook=sunny
+ -0.9533207559861846 * outlook=overcast
+ 0.10924216557459787 * outlook=rainy
+ 0.5276359628579281 * temperature
+ 0.7712122046533554 * humidity
+ -0.8907578344254022 * windy
- 0.8688305080362968

Number of kernel evaluations: 66

=== Stratified cross-validation ===

Correctly Classified Instances      7           50      %
Incorrectly Classified Instances    7           50      %
Kappa statistic                    -0.2564
Mean absolute error                 0.5
Root mean squared error             0.7071
Relative absolute error             105      %
Root relative squared error         143.3236 %
Total Number of Instances          14

=== Confusion Matrix ===

 a b  <-- classified as
 7 2 | a = yes
 5 0 | b = no
```


Kernal: Polynomial (Quadratic Function)

Predictive
Accuracy
= 78.6%

Support
Vectors = 10

SMD

Classifier for classes: yes, no

BinarySMD

```
-1 * 0.8100635668551557 * K[X(1) * X]
+ -1 * 0.019817568367163058 * K[X(3) * X]
+ -1 * 0.8887836783080866 * K[X(4) * X]
+ -1 * 1.0 * K[X(6) * X]
+ -1 * 0.3534785049716326 * K[X(7) * X]
+ 1 * 0.3727234704263585 * K[X(9) * X]
+ 1 * 0.1796126505860263 * K[X(10) * X]
+ 1 * 1.0 * K[X(11) * X]
+ 1 * 0.7245265999700636 * K[X(12) * X]
+ 1 * 0.7952805975195893 * K[X(13) * X]
- 0.6275818453891167
```

Number of support vectors: 10

Number of kernel evaluations: 104

=== Stratified cross-validation ===

| | | |
|----------------------------------|-----------|-----------|
| Correctly Classified Instances | 11 | 78.5714 % |
| Incorrectly Classified Instances | 3 | 21.4286 % |
| Kappa statistic | 0.5116 | |
| Mean absolute error | 0.2143 | |
| Root mean squared error | 0.4629 | |
| Relative absolute error | 45 % | |
| Root relative squared error | 93.8273 % | |
| Total Number of Instances | 14 | |

=== Confusion Matrix ===

```
a b  <-- classified as
8 1 | a = yes
2 3 | b = no
```

Kernal: Polynomial (Cubic Function)

Predictive
Accuracy
= 85.7%

Support
Vectors = 12

SMD

Classifier for classes: yes, no

BinarySMD

```
-1 * 0.058598146492685896 * K[X(1) * X]  
+ -1 * 0.032159637558211406 * K[X(2) * X]  
+ -1 * 0.36378917190737614 * K[X(3) * X]  
+ -1 * 0.07019514690769074 * K[X(4) * X]  
+ -1 * 0.07107098581642621 * K[X(5) * X]  
+ -1 * 0.8766783011511141 * K[X(6) * X]  
+ -1 * 0.06751016568226335 * K[X(7) * X]  
+ 1 * 0.052713460422677855 * K[X(9) * X]  
+ 1 * 0.3664976239753861 * K[X(10) * X]  
+ 1 * 1.0 * K[X(11) * X]  
+ 1 * 0.027564792859058898 * K[X(12) * X]  
+ 1 * 0.0932256782586451 * K[X(13) * X]  
- 0.5679604722895839
```

Number of support vectors: 12

Number of kernel evaluations: 105

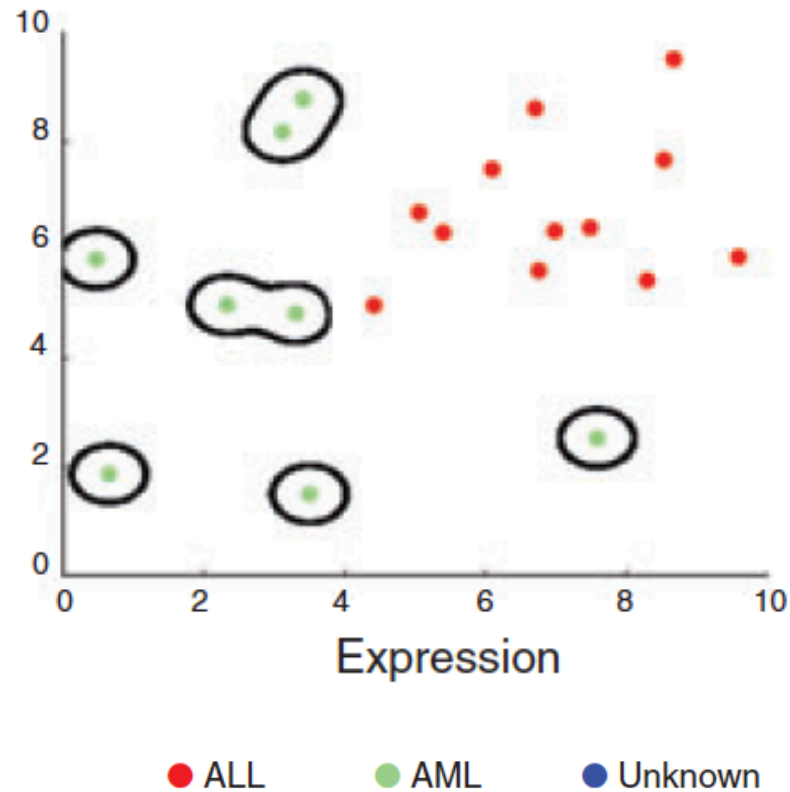
=== Stratified cross-validation ===

| | | |
|----------------------------------|---------|-----------|
| Correctly Classified Instances | 12 | 85.7143 % |
| Incorrectly Classified Instances | 2 | 14.2857 % |
| Kappa statistic | 0.6585 | |
| Mean absolute error | 0.1429 | |
| Root mean squared error | 0.378 | |
| Relative absolute error | 30 | % |
| Root relative squared error | 76.6097 | % |
| Total Number of Instances | 14 | |

=== Confusion Matrix ===

```
a b    <-- classified as  
9 0 | a = yes  
2 3 | b = no
```

Overfitting in SVM



(Adapted from Noble 2006)

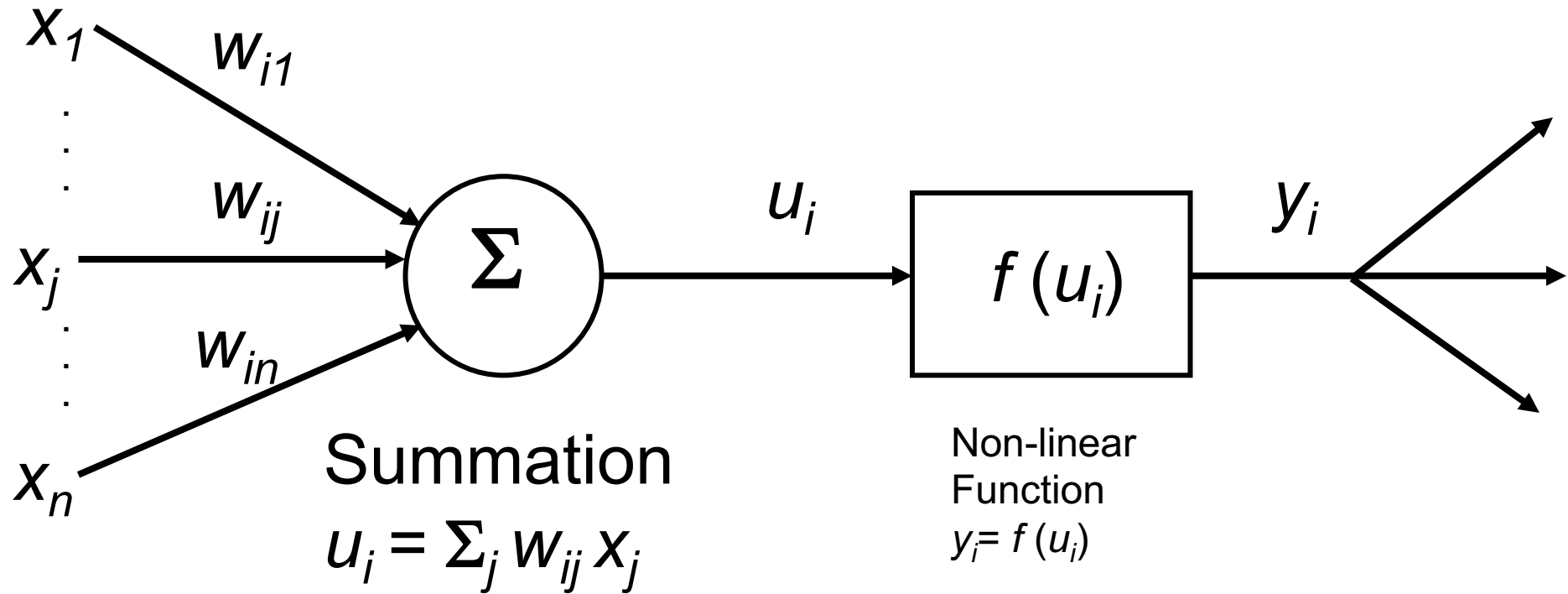
Artificial Neural Networks

- Based on the biological neural networks (brain)
- famous & widely used in bioinformatics

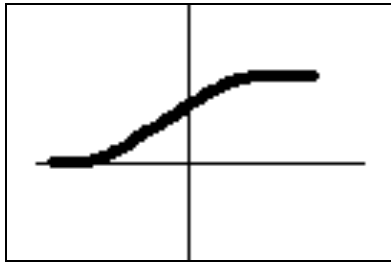
| ANN | Brain |
|-------------------|--|
| Feed Forward | Recurrent |
| Fully connected | Mostly local connections |
| Uniform structure | Functional modules |
| A few nodes type | > 100 types |
| 10 – 1000 nodes | Human brain: $O(10^{11})$ neurons, $O(10^{15})$ synapses |
| Mostly static | Dynamic |

(Reed & Marks, 1999)

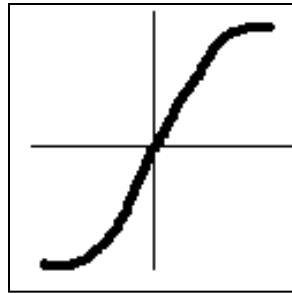
Artificial Neuron



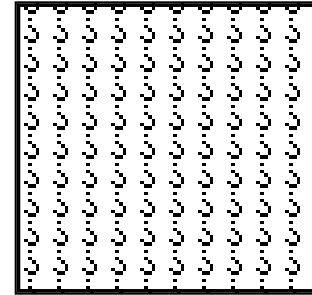
Common Node Nonlinearities



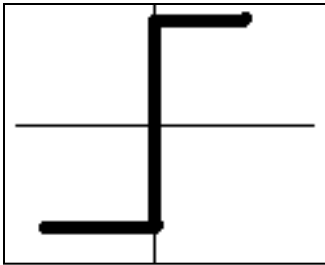
Sigmoid



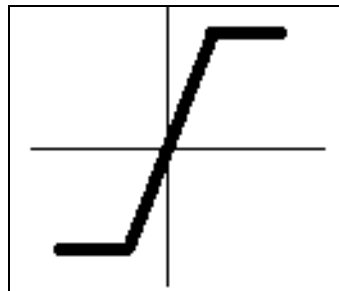
Tanh



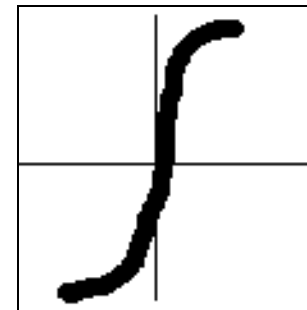
Step



Sign



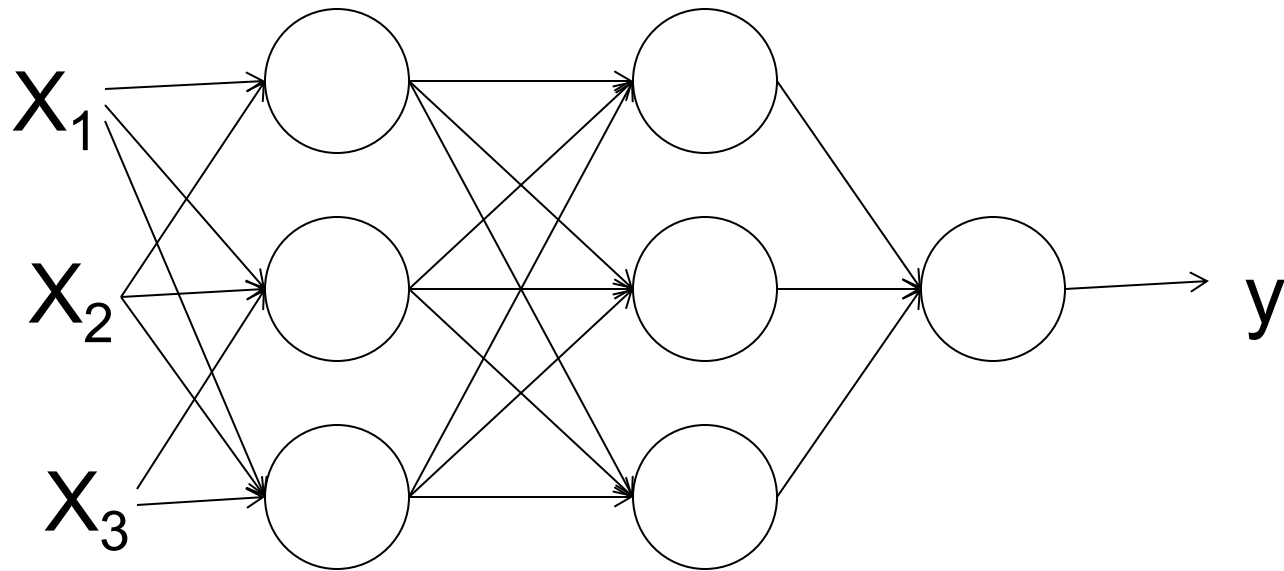
Clipped linear



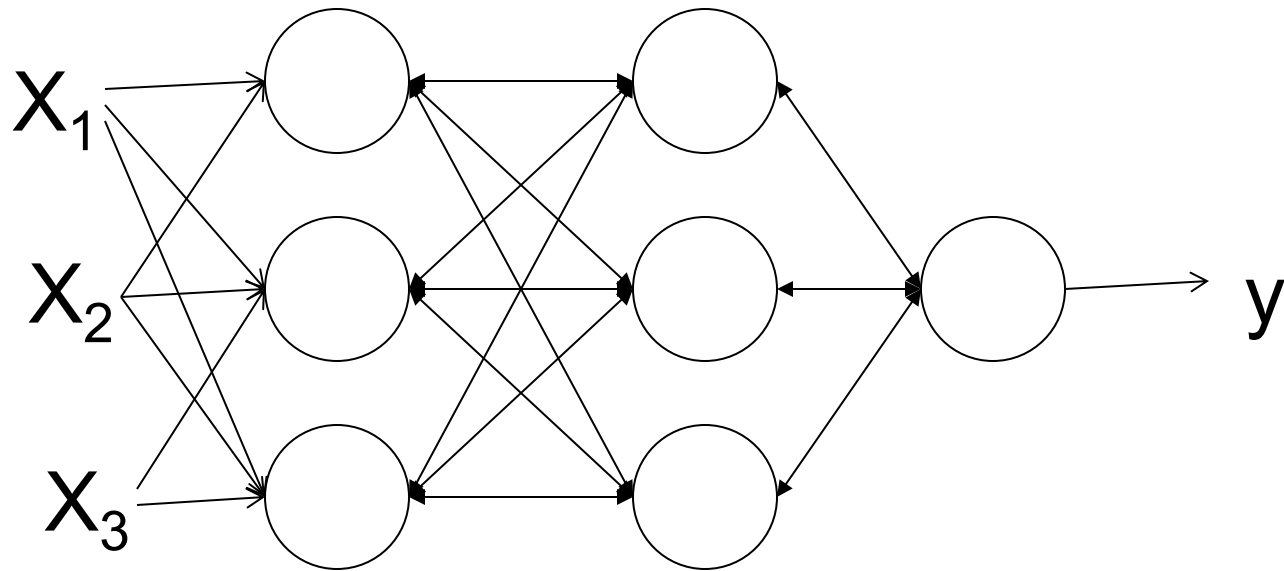
Inverse Abs

(Reed & Marks, 1999)

Feed Forward ANN Model



Back Propagation

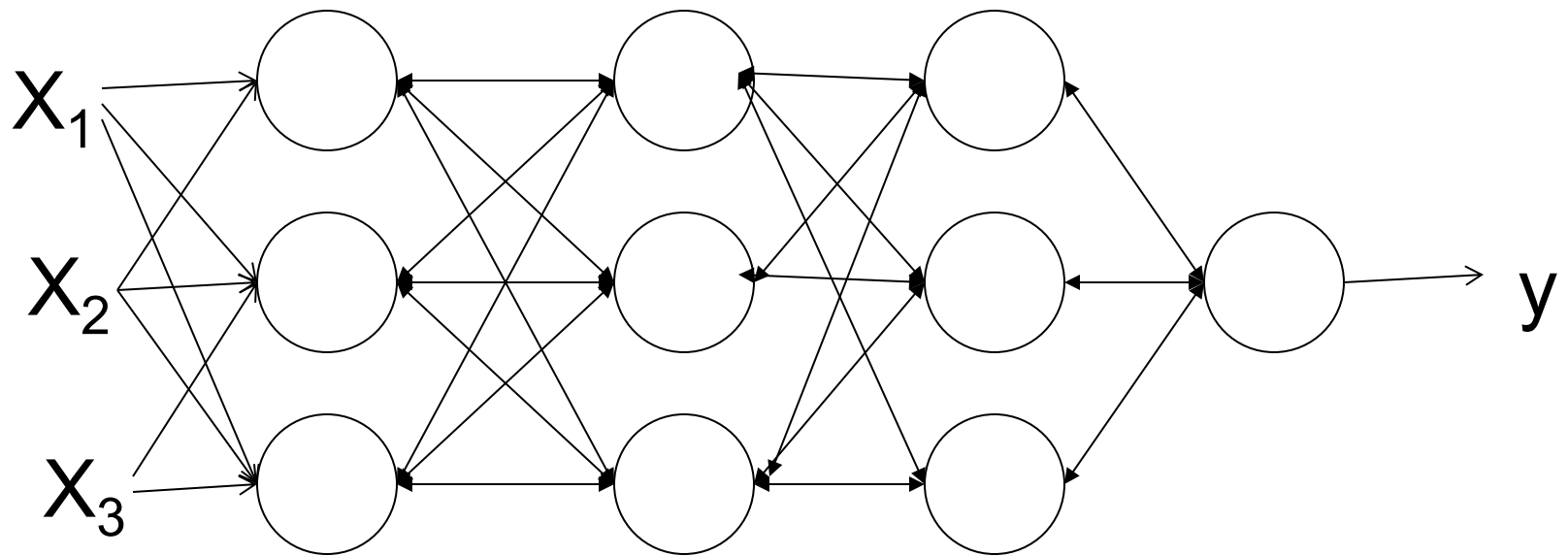


Deep Learning

Deep learning (also known as deep structured learning, hierarchical learning or deep machine learning) is a branch of machine learning based on a set of algorithms that attempt to model high-level abstractions in data by using a deep graph with multiple processing layers, composed of multiple linear and non-linear transformations.

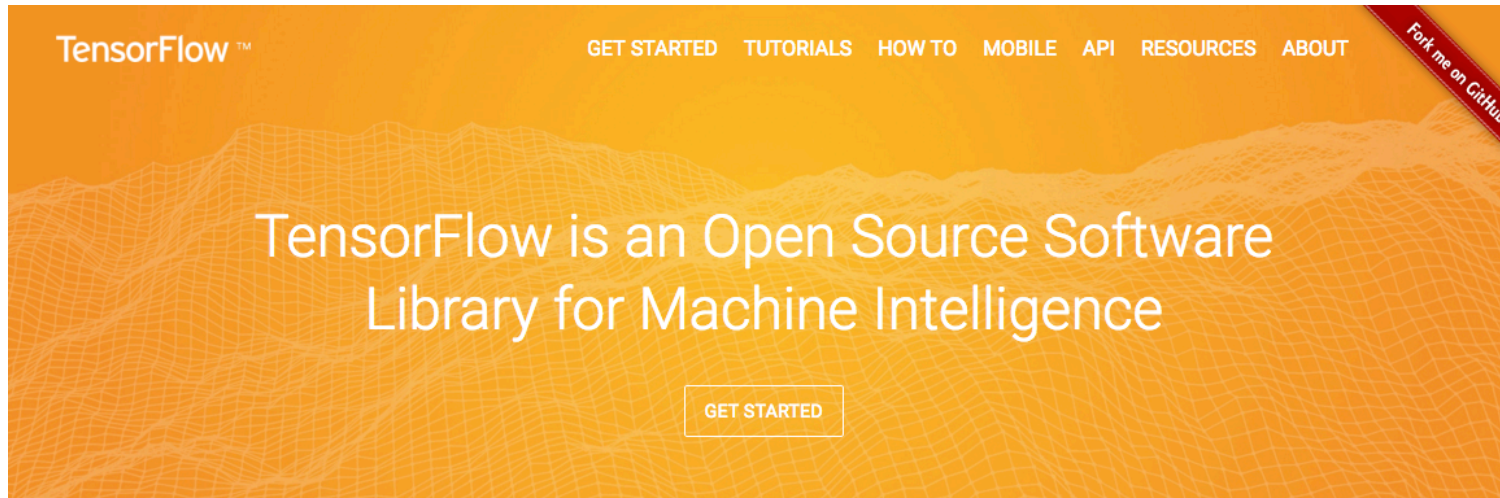
(From Wikipedia)

An illustration



Deep = more “nodes” and “hidden” layers

TensorFlow



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

About TensorFlow

TensorFlow™ is an open source software library for numerical computation using data flow graphs. Nodes in the graph represent mathematical operations, while the graph edges represent the multidimensional data arrays (tensors) communicated between them. The flexible architecture allows you to deploy computation to one or more CPUs or GPUs in a desktop, server, or mobile device with a single API. TensorFlow was originally developed by researchers and engineers working on the Google Brain Team within Google's Machine Intelligence research organization for the purposes of conducting machine learning and deep neural networks research, but the system is general enough to be applicable in a wide variety of other domains as well.



Example

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

<https://www.youtube.com/watch?v=lv0o9Lw3nz0>

https://www.ted.com/talks/fei_fei_li_how_we_re_teaching_computers_to_understand_pictures?language=en#t-118437

References

