# What is Machine Learning?
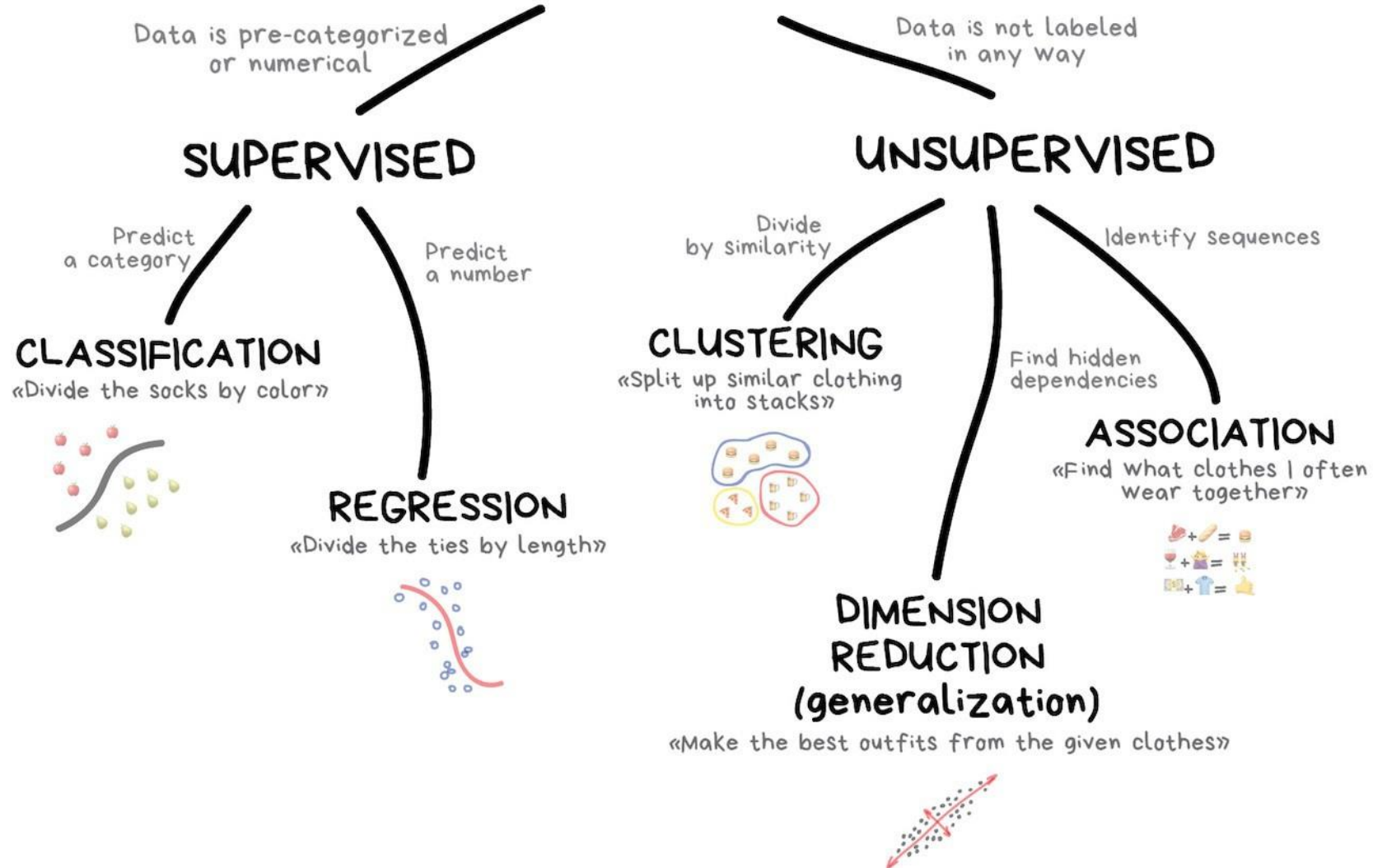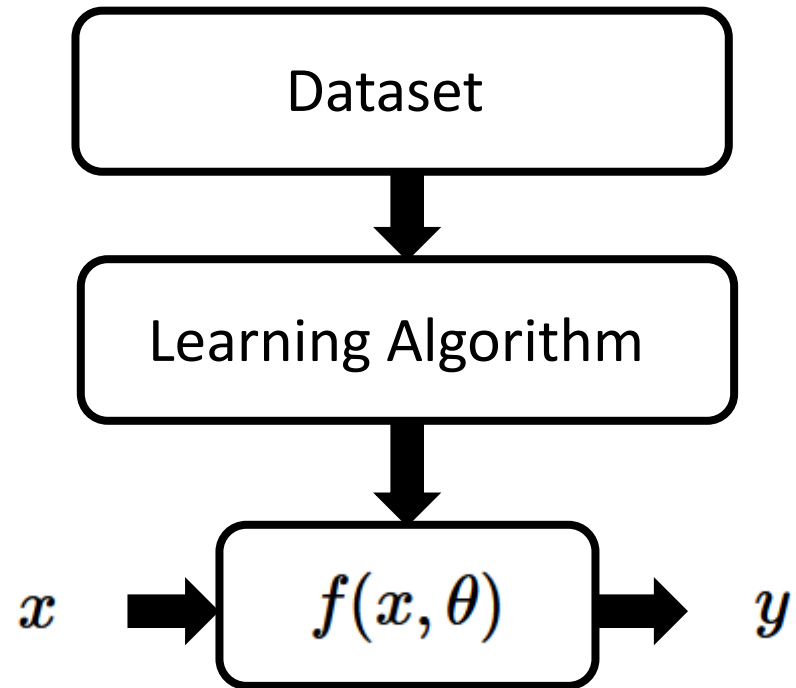
# CLASSICAL MACHINE LEARNING

Data is pre-categorized
or numerical

Data is not labeled
in any way

## SUPERVISED

## UNSUPERVISED

Predict
a category

Predict
a number

Divide
by similarity

Identify sequences

### CLASSIFICATION
«Divide the socks by color»

### CLUSTERING
«Split up similar clothing
into stacks»

Find hidden
dependencies

### ASSOCIATION
«Find what clothes I often
wear together»

### REGRESSION
«Divide the ties by length»

### DIMENSION
### REDUCTION
### (generalization)
«Make the best outfits from the given clothes»

## This is incredibly important - PLEASE READ

aaabeline@aol.com  to you    show details ⌄                                                May 2

IPad_Image.Jpg (36 KB)

Open this now. Hi, This is incredibly important, so PLEASE READ THIS EMAIL IN IT'S ENTIRETY

Check it here now <br/> A stay at home wor dpress programmer discoved a major Google loophole using a little know technique and built a plugin to automate the method. </p> This technique replaces the ori ginal search box with a Google "Adsense for Search" Custom Search box quickly and easily using a simple plugin he developed. As with most Wordpress themes, you'd have to change this manually, which would require some php programming... But now you can do it with the simple "push of a button" by adding the Google Adsense Accelerat or plugin to your wordpress site. NO PROGRAMMING...

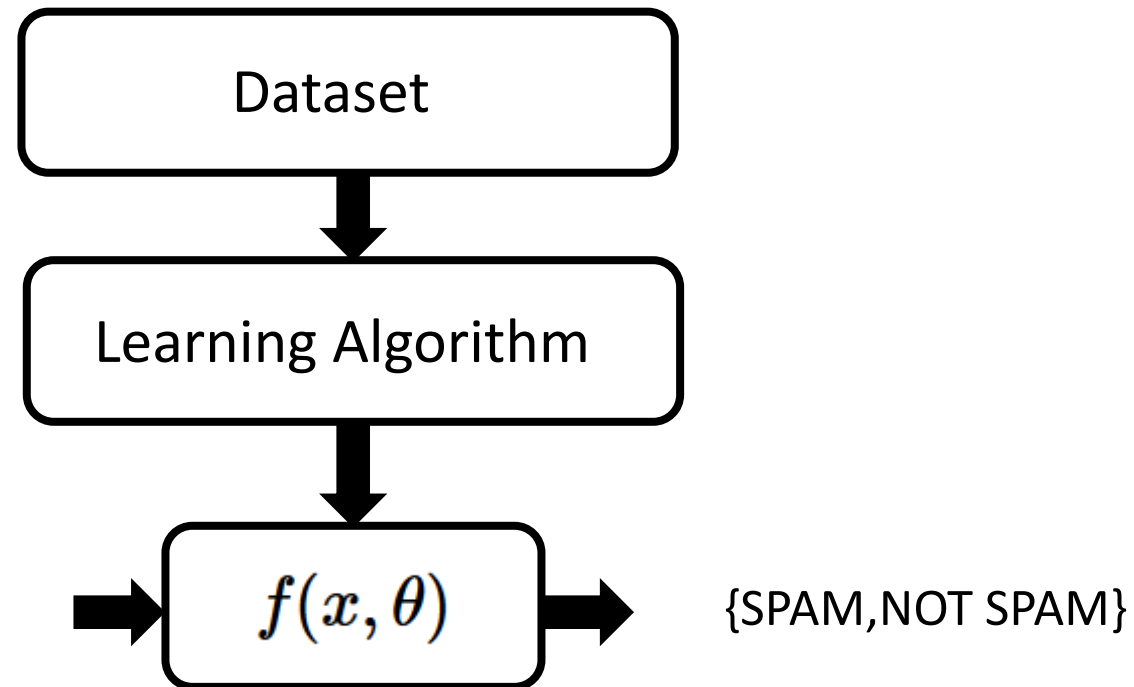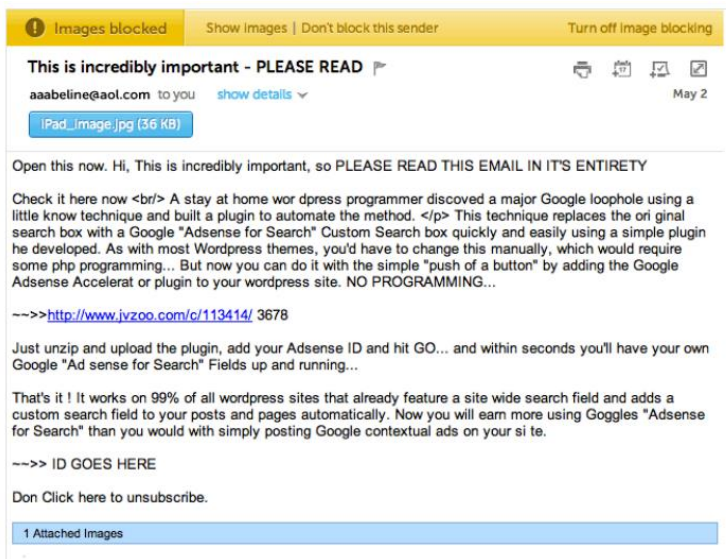~~>>http://www.jvzoo.com/c/113414/ 3678

Just unzip and upload the plugin, add your Adsense ID and hit GO... and within seconds you'll have your own Google "Ad sense for Search" Fields up and running...
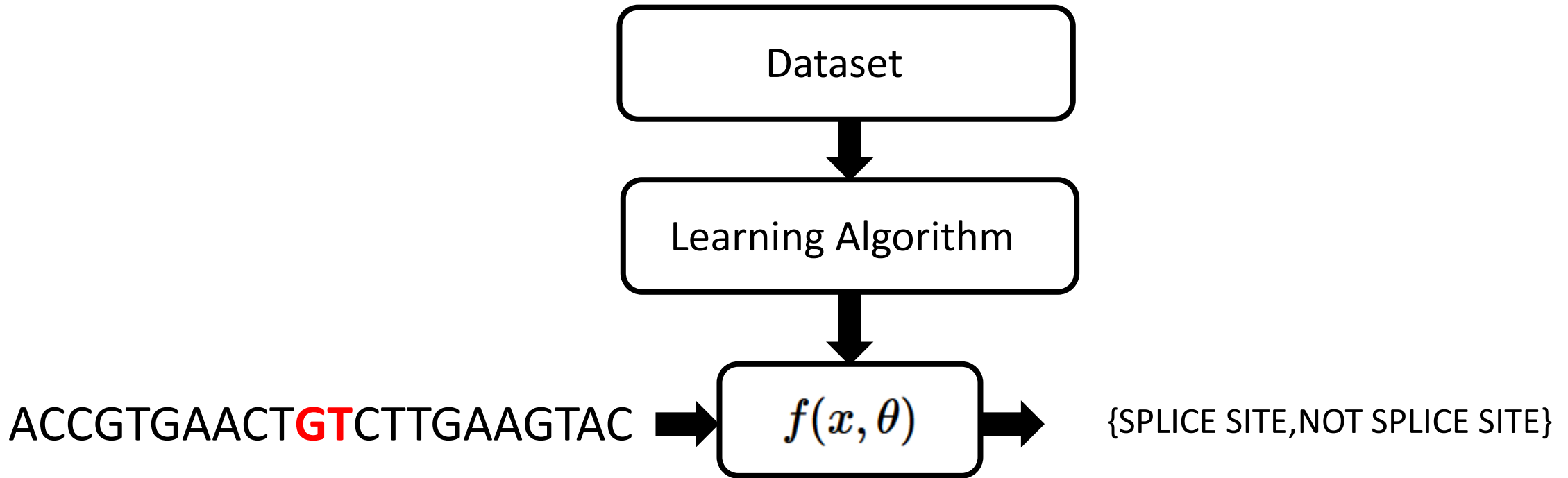
That's it ! It works on 99% of all wordpress sites that already feature a site wide search field and adds a custom search field to your posts and pages automatically. Now you will earn more using Goggles "Adsense for Search" than you would with simply posting Google contextual ads on your si te.
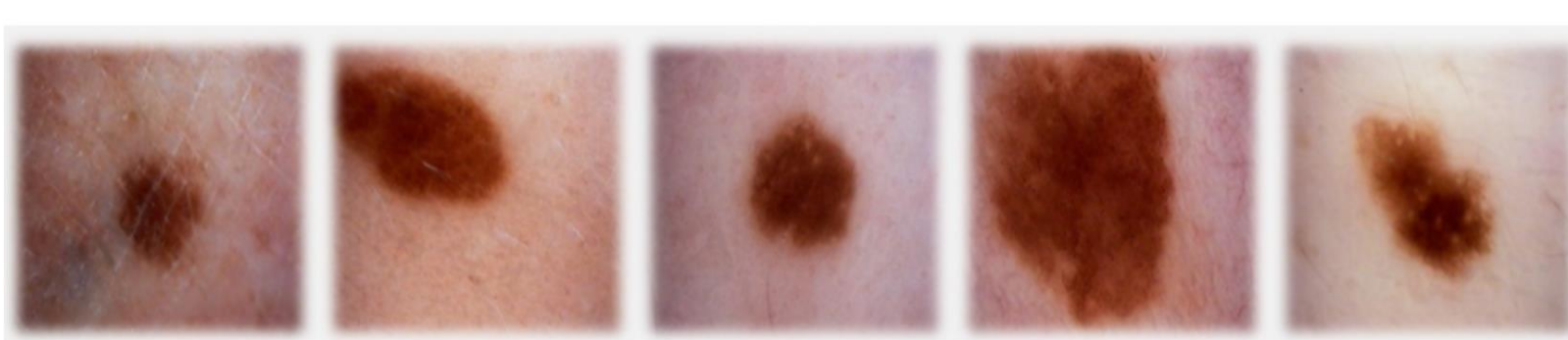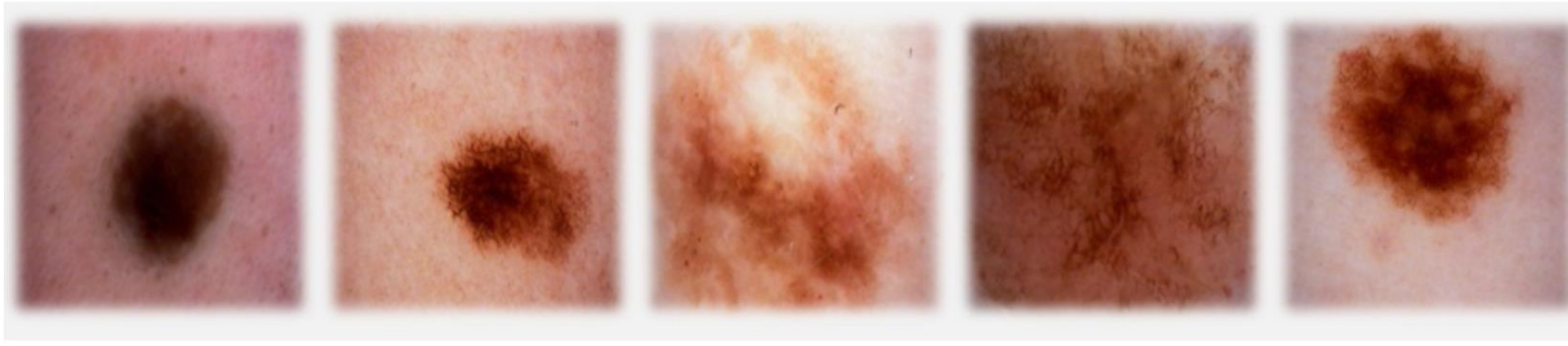
~~>> ID GOES HERE

Don Click here to unsubscribe.

1 Attached Images

Dataset
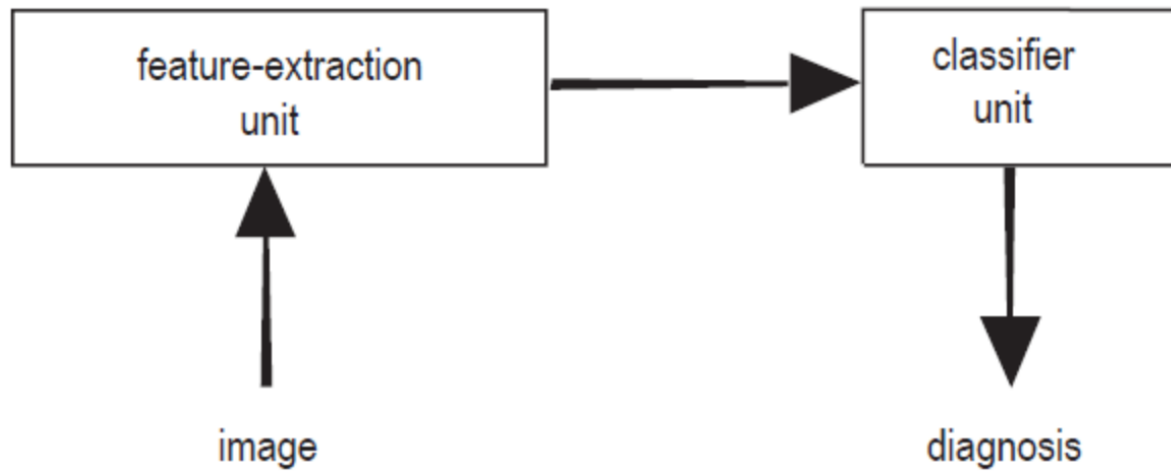
Learning Algorithm

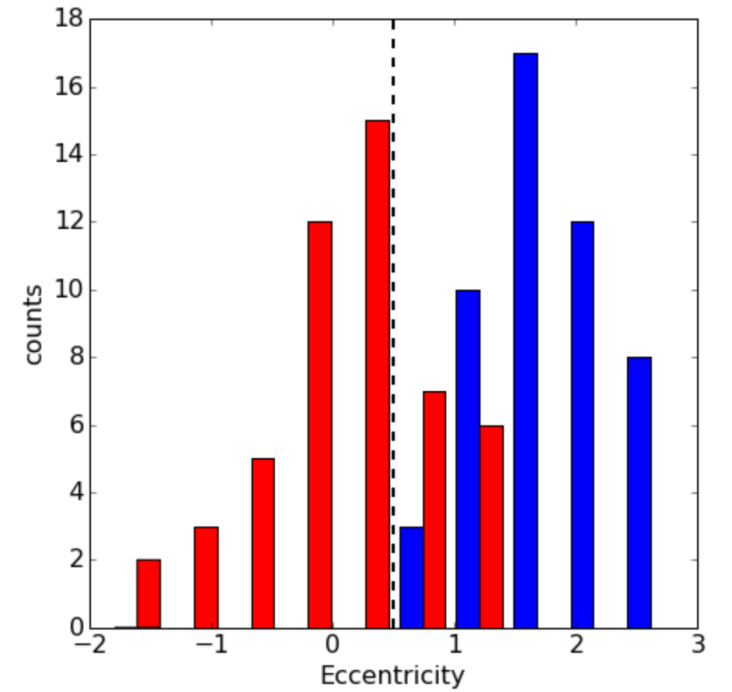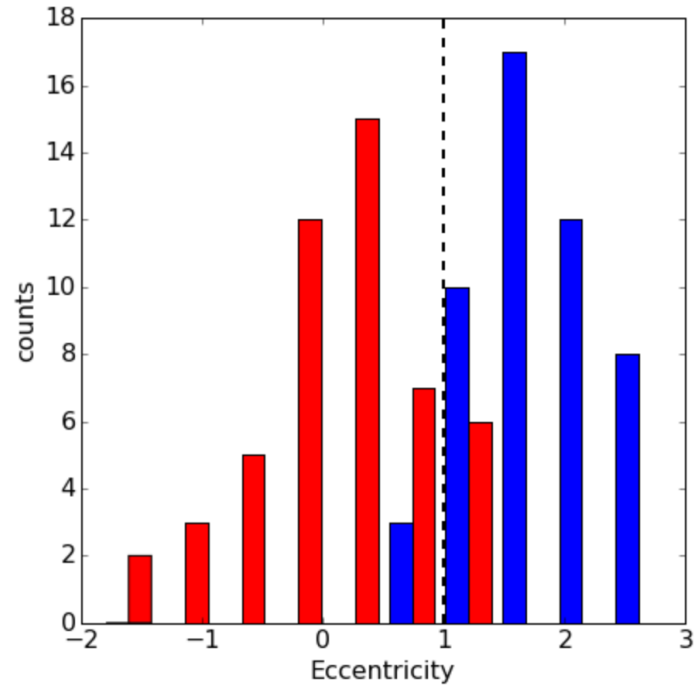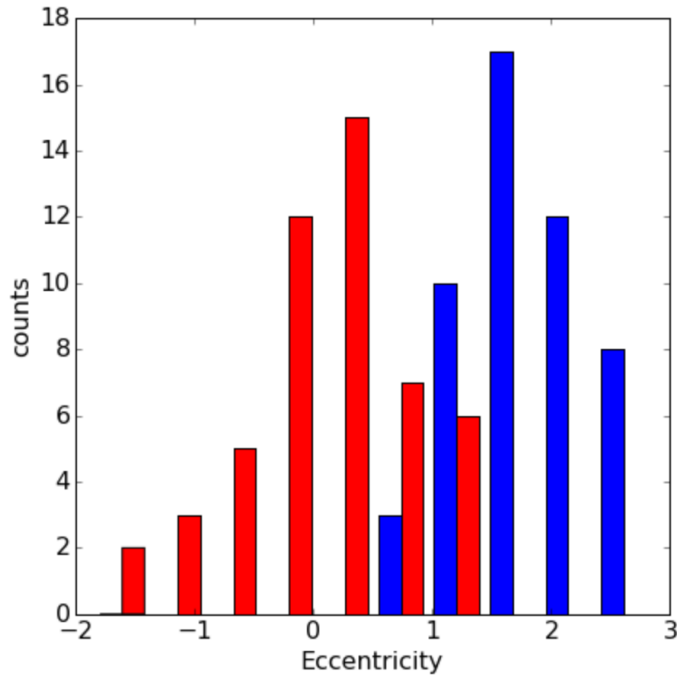$$f(x, \theta)$$

{SPAM,NOT SPAM}

concepts

- sign of cancer
- top row malignant
- bottom row benign

o feature extraction: features (a.k.a. properties or attributes)

o data set, sample (a.k.a. example, instance or data point), label (a.k.a. target)

- eccentricity of lesion (how nearly circular the lesion is)

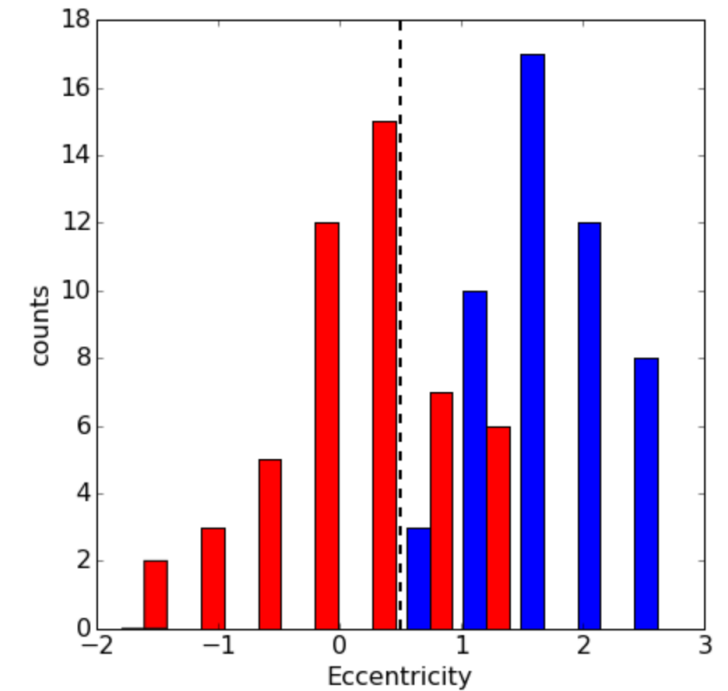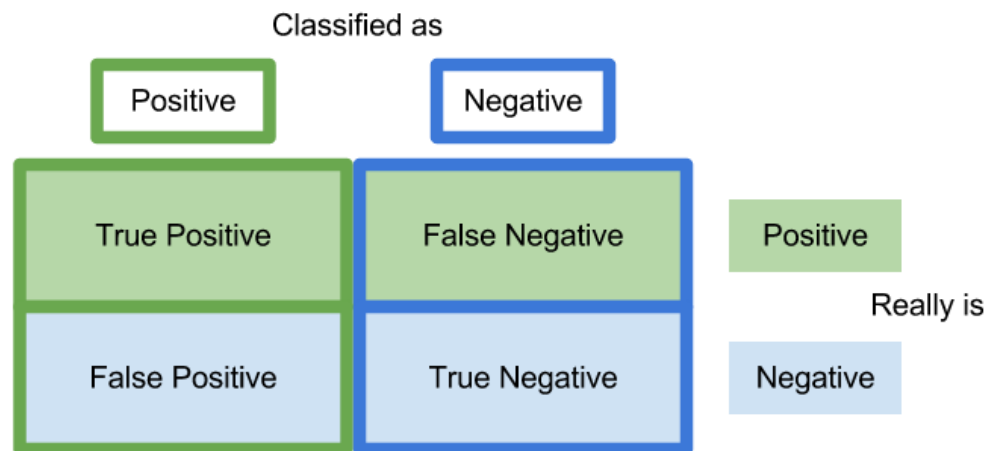- threshold *t*

- consequence of the predictions

o   malignant positive class, benign negative

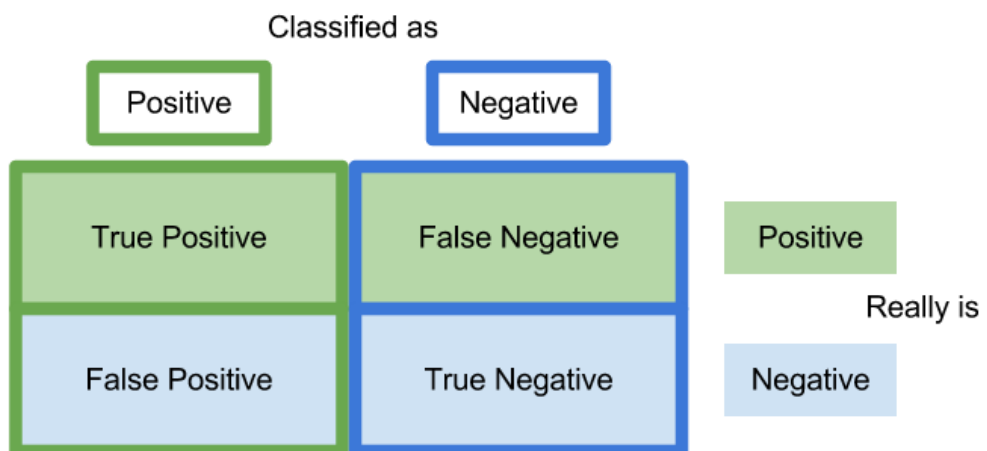o   consequence of the predictions



count the number of malignant images with eccentricity value $\geq t$: **true positive** predictions (TP)

count the number of malignant images with eccentricity value $< t$: **false negative** predictions (FN)

count the number of benign images with eccentricity value $\geq t$: **false positive** predictions (FP)

count the number of benign images with eccentricity value $< t$: **true negative** predictions (TN)

Classified as

| Positive | Negative | | |
|---|---|---|---|
| True Positive | False Negative | Positive | |
| False Positive | True Negative | Negative | |

Really is

## Classified as

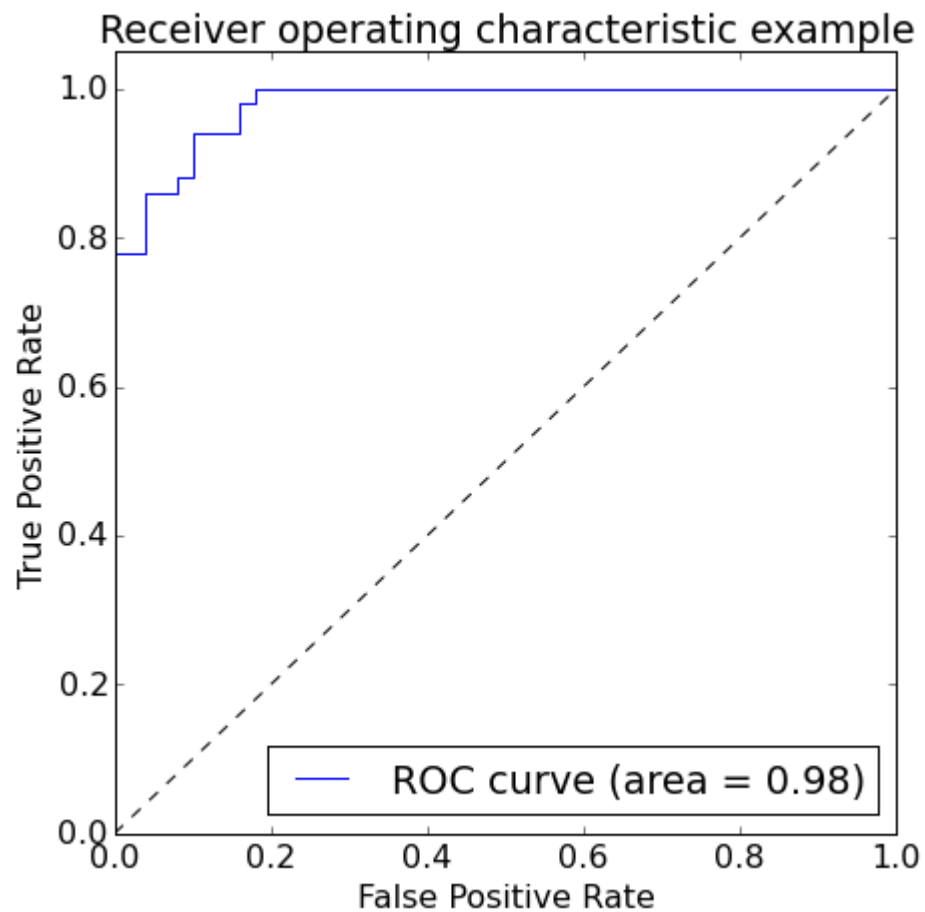|  | Positive | Negative |
|---|---|---|
| **Positive** (Really is) | True Positive | False Negative |
| **Negative** | False Positive | True Negative |

$$\text{accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

$$\text{TPR} = \frac{TP}{TP + FN}$$

$$\text{FPR} = \frac{FP}{FP + TN}$$
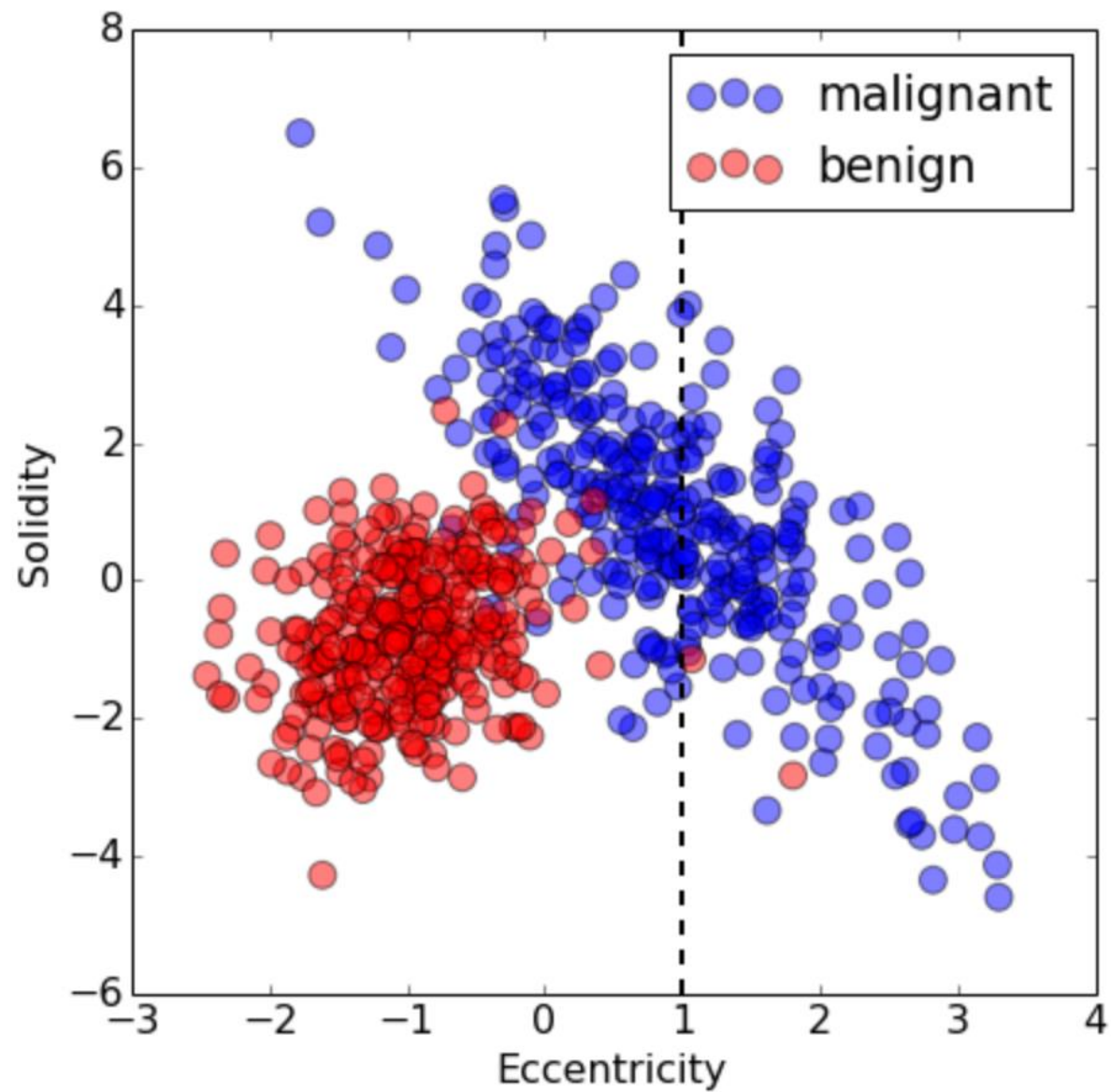
$$\text{TPR} = \frac{TP}{TP + FN}$$

## Receiver operating characteristic example

True Positive Rate

1.0

0.8

0.6

0.4

0.2

0.0

0.0    0.2    0.4    0.6    0.8    1.0

False Positive Rate

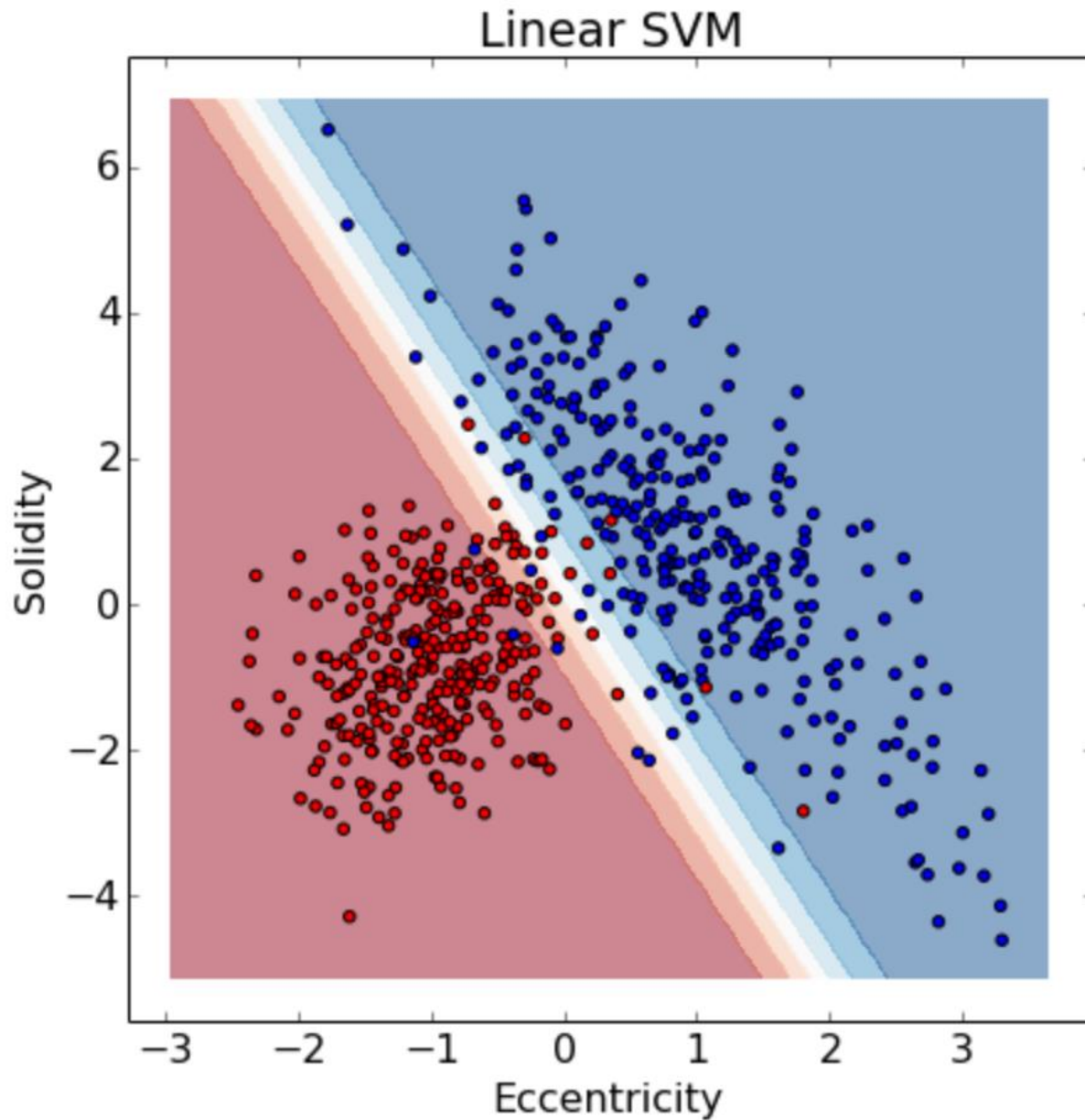— ROC curve (area = 0.98)

$$\text{FPR} = \frac{FP}{FP + TN}$$

- o  model that classifies all images as malignant: TPR=1 and FPR=1

- o  model that classifies all images a benign: TPR=0 and FPR=0

- o  vary threshold *t*

- o  *AUC*

o add another feature?

o feature vector X
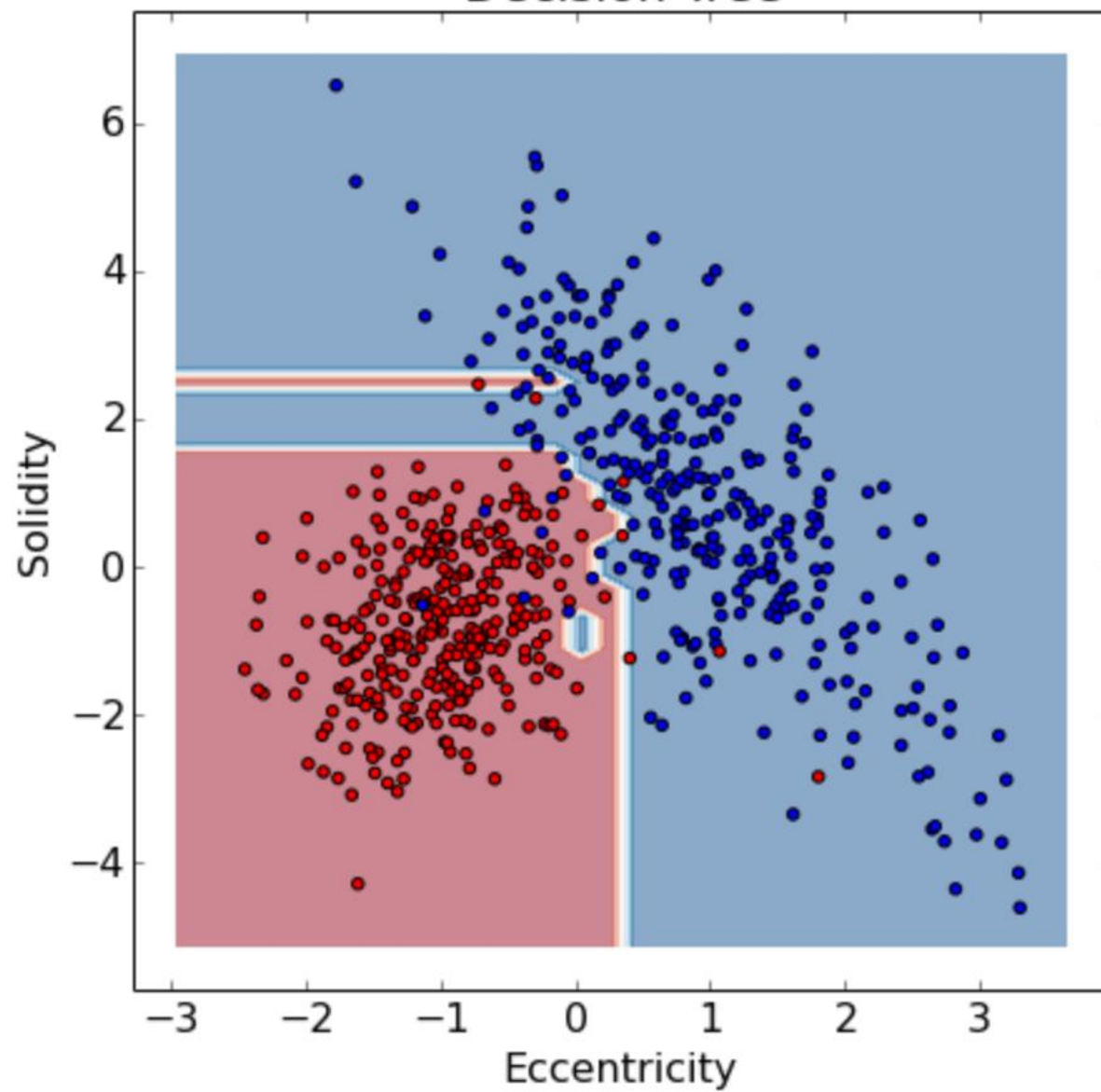
o Euclidean vector space

- feature vector X

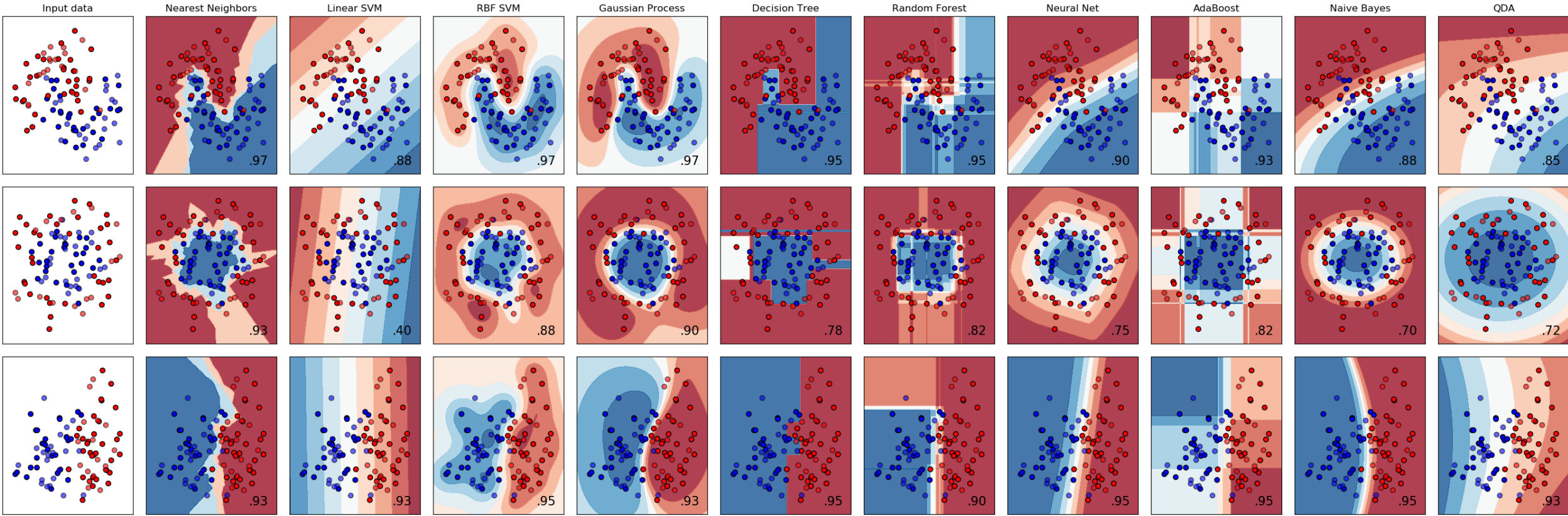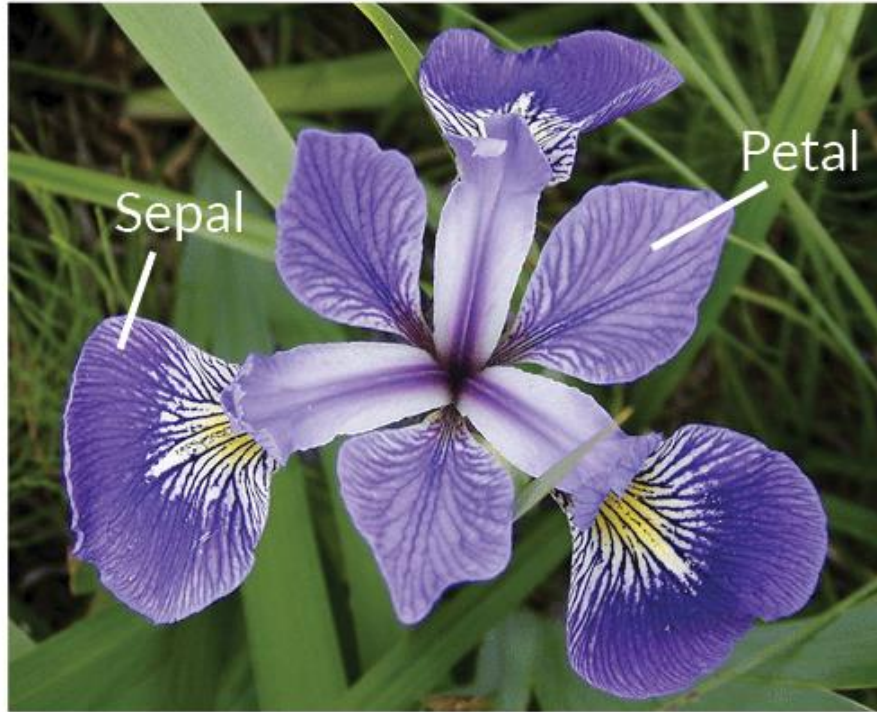- Euclidean vector space

Linear SVM

- o linear decision boundary

- o blue region malignant, red region benign

- o yet more features

- o can't look at the decision boundary

- o more complex

- unseen external images

- generalization

- overfitting

| Input data | Nearest Neighbors | Linear SVM | RBF SVM | Gaussian Process | Decision Tree | Random Forest | Neural Net | AdaBoost | Naive Bayes | QDA |
|---|---|---|---|---|---|---|---|---|---|---|
| | .97 | .88 | .97 | .97 | .95 | .95 | .90 | .93 | .88 | .85 |
| | .93 | .40 | .88 | .90 | .78 | .82 | .75 | .82 | .70 | .72 |
| | .93 | .93 | .95 | .93 | .95 | .90 | .95 | .95 | .95 | .93 |

**Iris Versicolor**  **Iris Setosa**  **Iris Virginica**

## Table I

| Iris setosa | | | | Iris versicolor | | | | Iris virginica | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Sepal length | Sepal width | Petal length | Petal width | Sepal length | Sepal width | Petal length | Petal width | Sepal length | Sepal width | Petal length | Petal width |
| 5·1 | 3·5 | 1·4 | 0·2 | 7·0 | 3·2 | 4·7 | 1·4 | 6·3 | 3·3 | 6·0 | 2·5 |
| 4·9 | 3·0 | 1·4 | 0·2 | 6·4 | 3·2 | 4·5 | 1·5 | 5·8 | 2·7 | 5·1 | 1·9 |
| 4·7 | 3·2 | 1·3 | 0·2 | 6·9 | 3·1 | 4·9 | 1·5 | 7·1 | 3·0 | 5·9 | 2·1 |
| 4·6 | 3·1 | 1·5 | 0·2 | 5·5 | 2·3 | 4·0 | 1·3 | 6·3 | 2·9 | 5·6 | 1·8 |
| 5·0 | 3·6 | 1·4 | 0·2 | 6·5 | 2·8 | 4·6 | 1·5 | 6·5 | 3·0 | 5·8 | 2·2 |
| 5·4 | 3·9 | 1·7 | 0·4 | 5·7 | 2·8 | 4·5 | 1·3 | 7·6 | 3·0 | 6·6 | 2·1 |
| 4·6 | 3·4 | 1·4 | 0·3 | 6·3 | 3·3 | 4·7 | 1·6 | 4·9 | 2·5 | 4·5 | 1·7 |
| 5·0 | 3·4 | 1·5 | 0·2 | 4·9 | 2·4 | 3·3 | 1·0 | 7·3 | 2·9 | 6·3 | 1·8 |
| 4·4 | 2·9 | 1·4 | 0·2 | 6·6 | 2·9 | 4·6 | 1·3 | 6·7 | 2·5 | 5·8 | 1·8 |
| 4·9 | 3·1 | 1·5 | 0·1 | 5·2 | 2·7 | 3·9 | 1·4 | 7·2 | 3·6 | 6·1 | 2·5 |
| 5·4 | 3·7 | 1·5 | 0·2 | 5·0 | 2·0 | 3·5 | 1·0 | 6·5 | 3·2 | 5·1 | 2·0 |
| 4·8 | 3·4 | 1·6 | 0·2 | 5·9 | 3·0 | 4·2 | 1·5 | 6·4 | 2·7 | 5·3 | 1·9 |
| 4·8 | 3·0 | 1·4 | 0·1 | 6·0 | 2·2 | 4·0 | 1·0 | 6·8 | 3·0 | 5·5 | 2·1 |
| 4·3 | 3·0 | 1·1 | 0·1 | 6·1 | 2·9 | 4·7 | 1·4 | 5·7 | 2·5 | 5·0 | 2·0 |
| 5·8 | 4·0 | 1·2 | 0·2 | 5·6 | 2·9 | 3·6 | 1·3 | 5·8 | 2·8 | 5·1 | 2·4 |

| | sepal_length | sepal_width | petal_length | petal_width |
|---|---|---|---|---|
| count | 150.000000 | 150.000000 | 150.000000 | 150.000000 |
| mean | 5.843333 | 3.057333 | 3.758000 | 1.199333 |
| std | 0.828066 | 0.435866 | 1.765298 | 0.762238 |
| min | 4.300000 | 2.000000 | 1.000000 | 0.100000 |
| 25% | 5.100000 | 2.800000 | 1.600000 | 0.300000 |
| 50% | 5.800000 | 3.000000 | 4.350000 | 1.300000 |
| 75% | 6.400000 | 3.300000 | 5.100000 | 1.800000 |
| max | 7.900000 | 4.400000 | 6.900000 | 2.500000 |

Fisher R., *The use of multiple measurements in taxonomic problems* (1936)
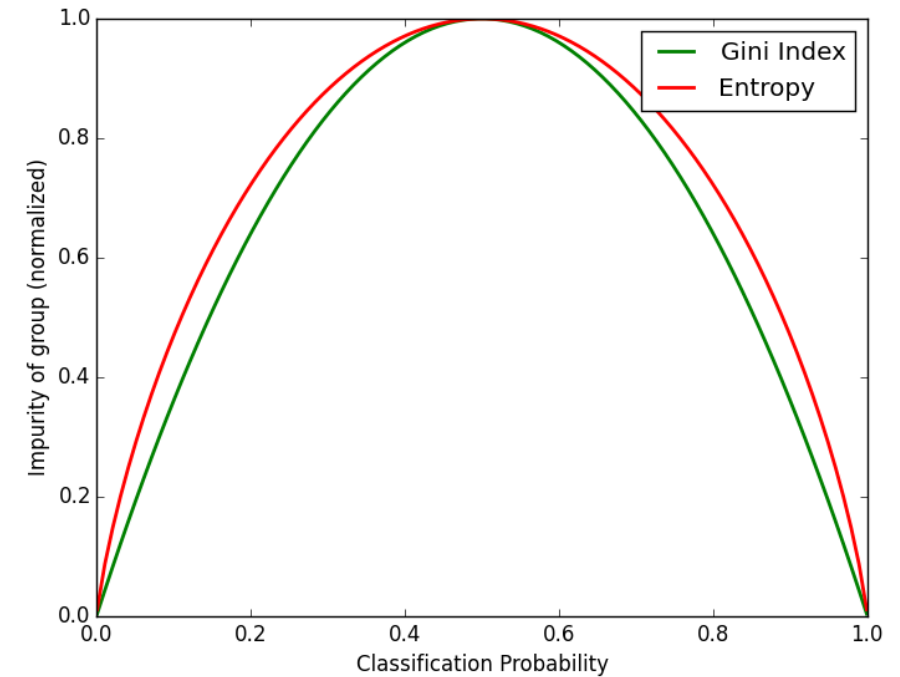
Decision surface of a decision tree using paired features

$$Gini: Gini(E) = 1 - \sum_{j=1}^{c} p_j^2$$

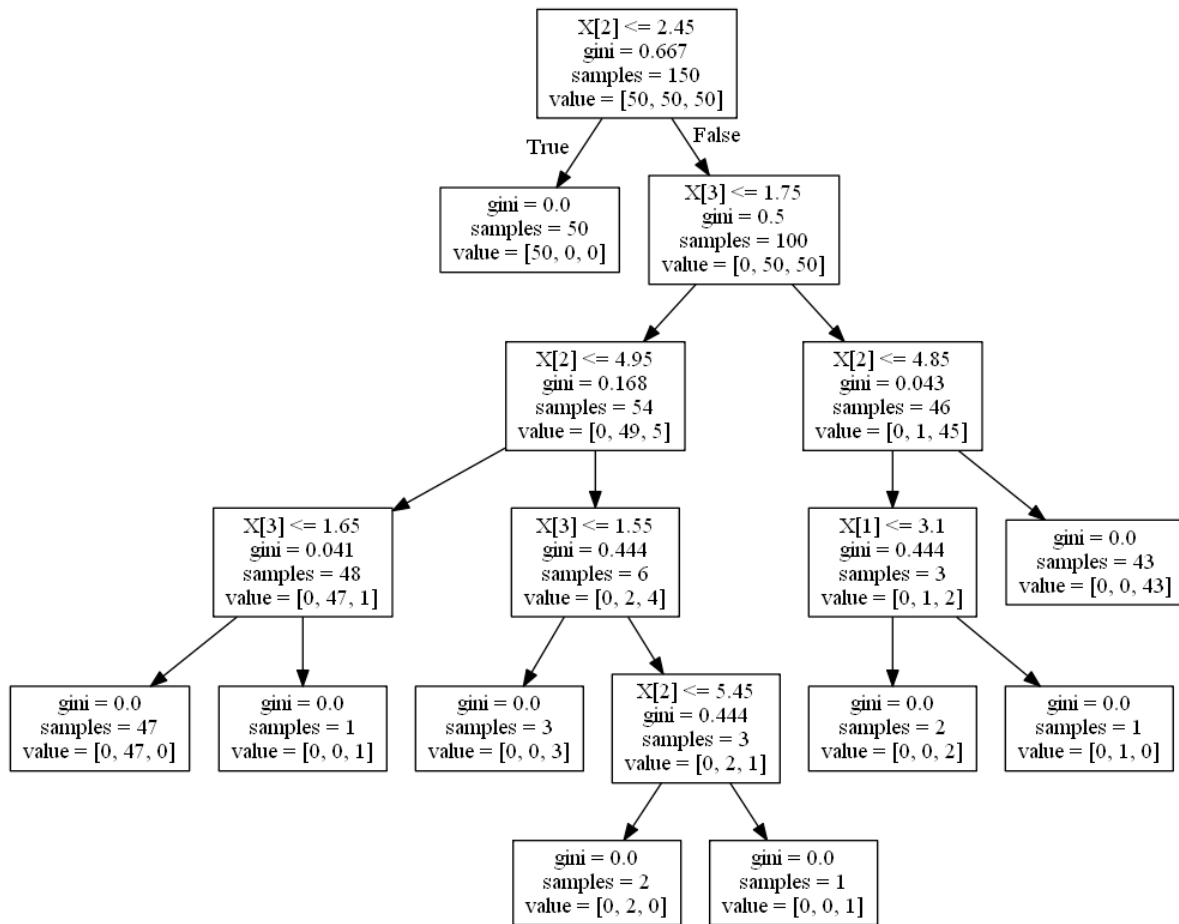$$Entropy: H(E) = - \sum_{j=1}^{c} p_j \log p_j$$

# Let's get real!

1. Building a decision tree
2. Kaggle evaluation
3. The test set
4. Hyperparameters

advantages:

o ease of interpretation
o handles continuous and discrete features
o invariant to monotone transformation of features
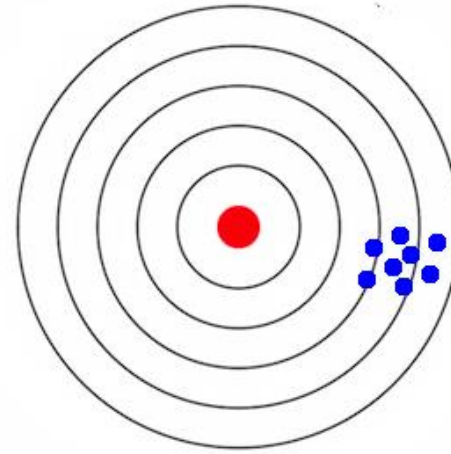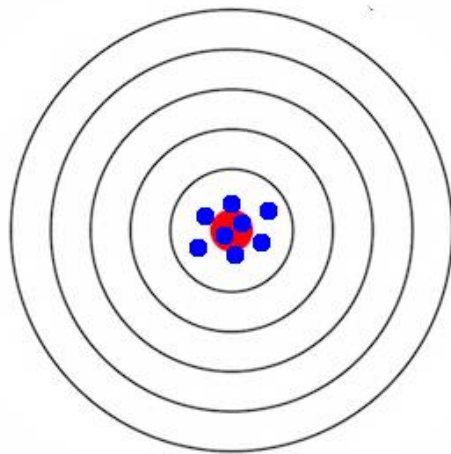o variable selection automated
o **low bias** (deep trees)

disadvantages:

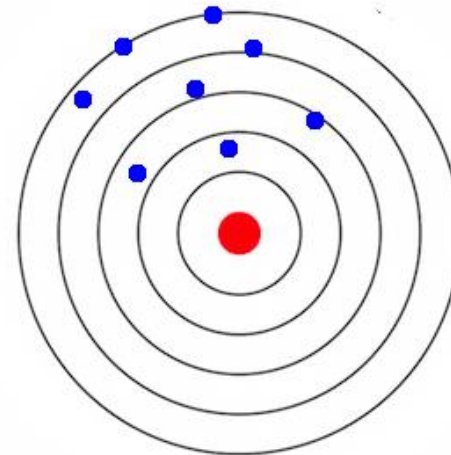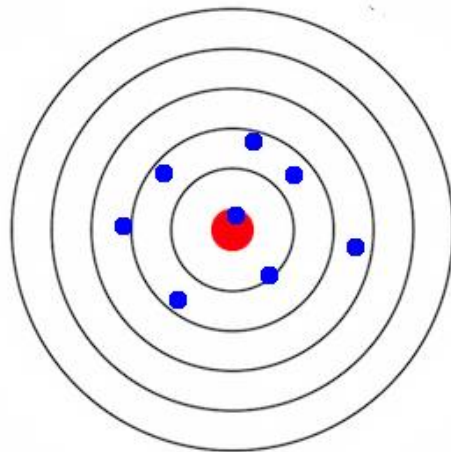o **high variance**
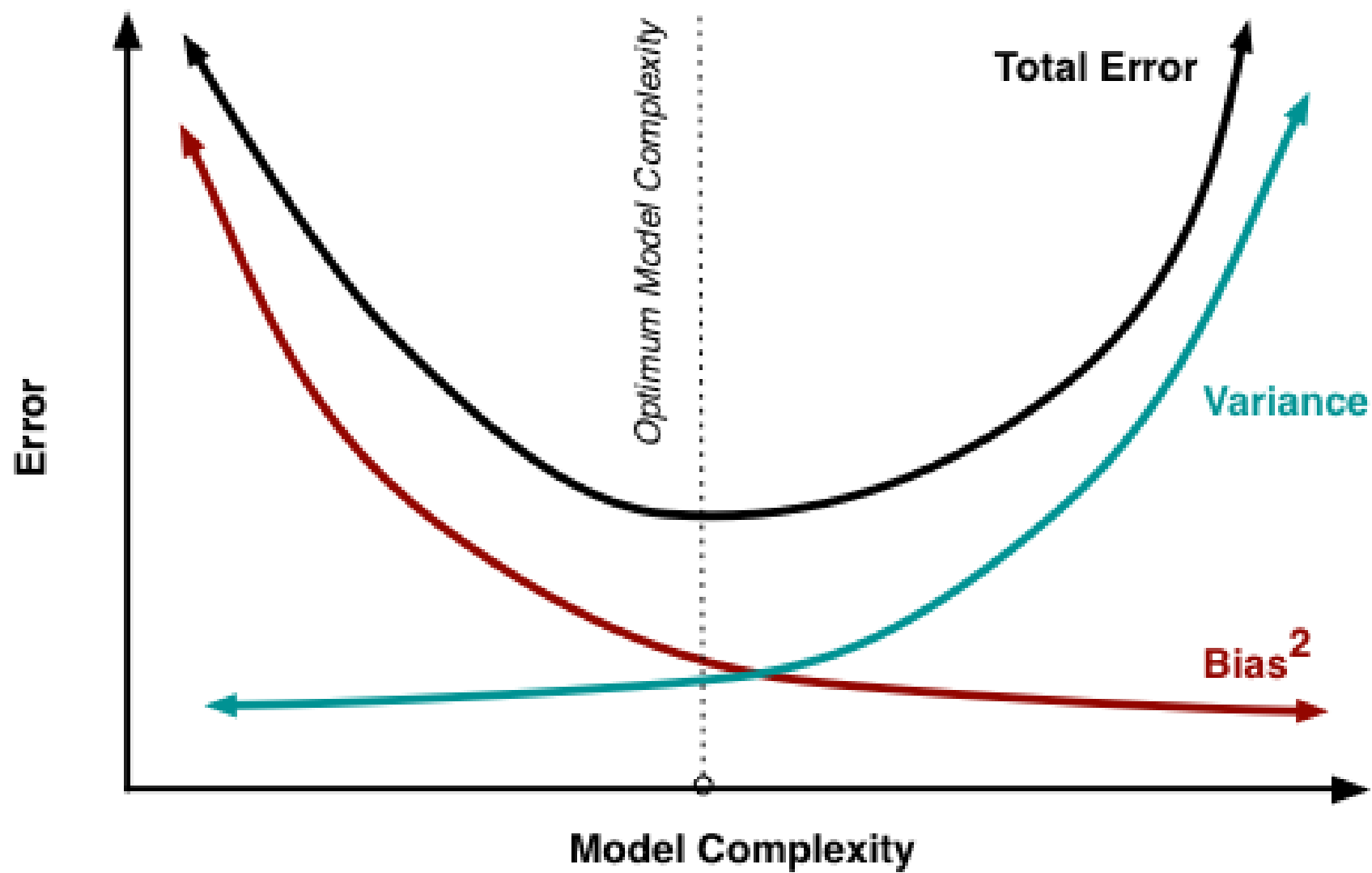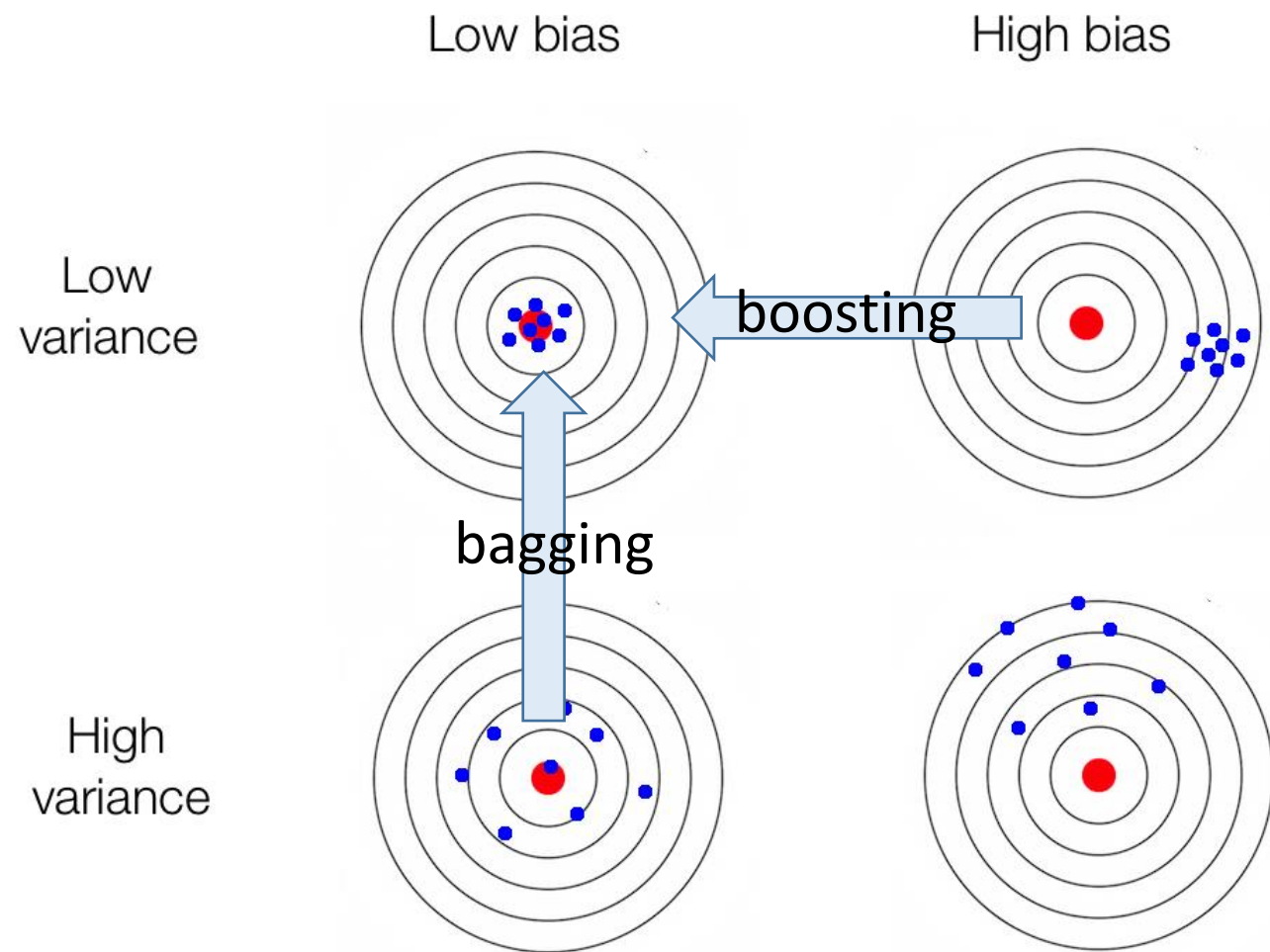o overfitting

Low bias

High bias
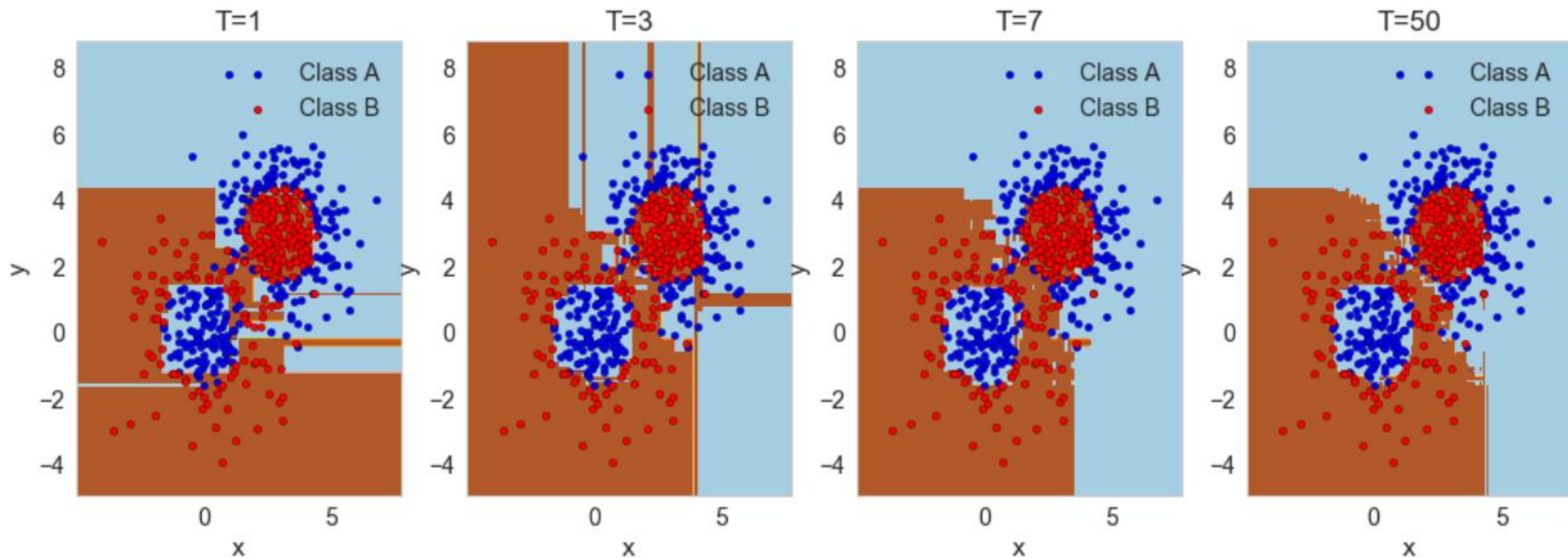
Low
variance

High
variance

# *Bagging*

o   Simulate the notion of different train set samples:

  1.   sample data points from train set to create new train set
  2.   fit low bias high variance model on new train set
  3.   repeat steps (1) and (2) *T* times
  4.   average the predictions of the *T* models

$$\hat{f}(x, \theta) = \frac{1}{T} \sum_{t=1}^{T} f_t(x, \theta)$$

# Bagging: Random Forests

o Train set contains $n$ data points with $m$ features.

o Construct  $T$ low bias high variance decision trees by following these steps:

1. Sample $n$ data points at random **with replacement** from the train set.

2. At each node,  select $h<<m$  features at random and compute the best split using only these $h$ features.

3. Each tree is grown to the largest extent possible. There is no pruning or early stopping.

o Step 3 ensures that the bagged models are low bias by learning deep complex decision trees.

# Bagging: Random Forests

# Let's get real!

**5.   Ensemble learning: bagging**

# Boosting

o   reduce the bias of a high bias low variance model

o   turning an ensemble of weak learners into a strong learner

o   the meta-model is additive, i.e. adaboost:

$$\hat{f}(x, \theta) = \sum_{t=1}^{T} \alpha_t f_t(x, \theta)$$

# Boosting: adaboost

Initialize the weights $D_1(i) = 1/n, i = 1, 2, \ldots, n.$

$$\boxed{y_i \in \{-1, +1\}}$$

For $t = 1$ to $T$:

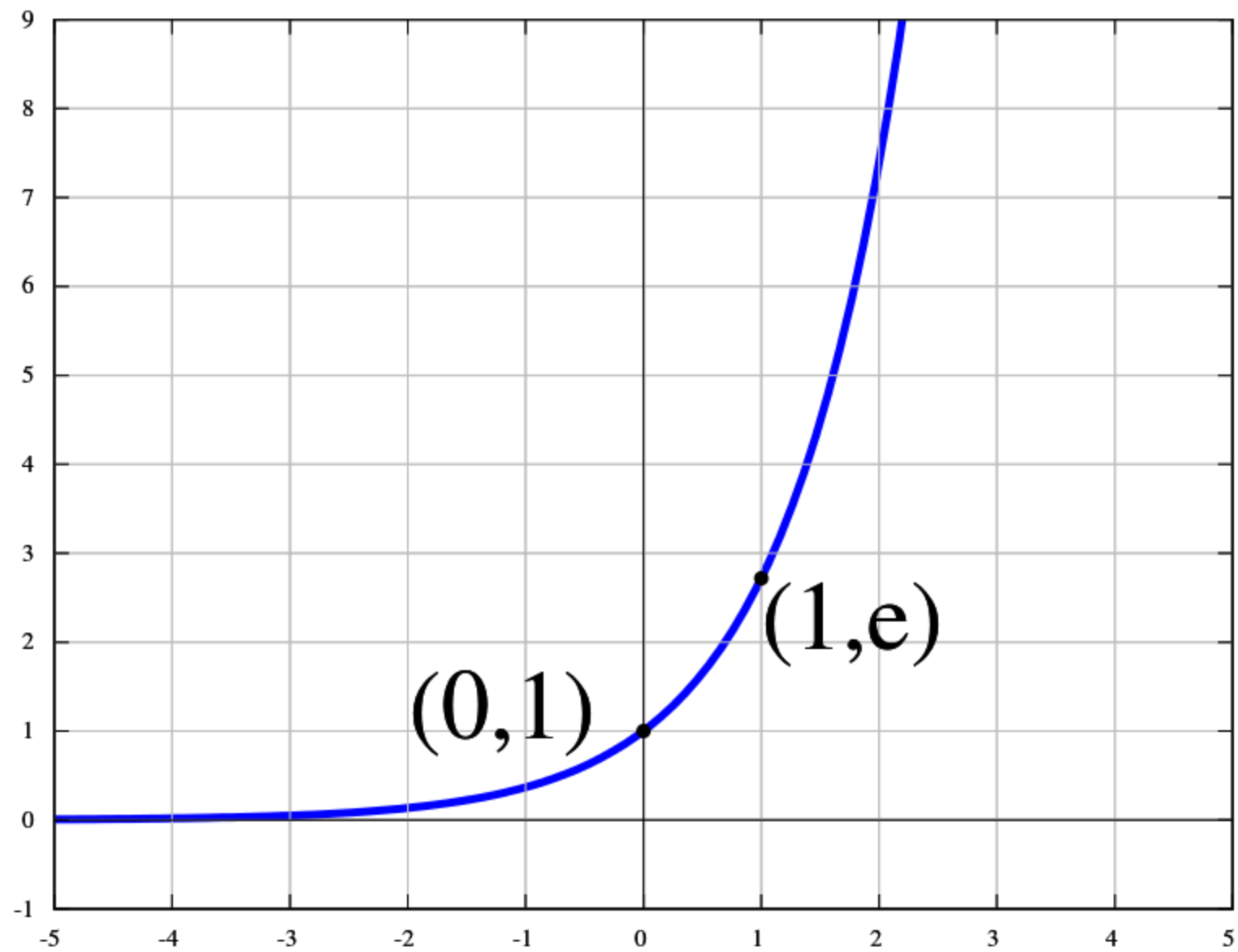    1. Fit a weak classifier $f_t(x, \theta)$ to the trainset data using weights $D_1(i)$.

    2. Set $\alpha_t = \frac{1}{2} ln(\frac{1-error}{error})$.
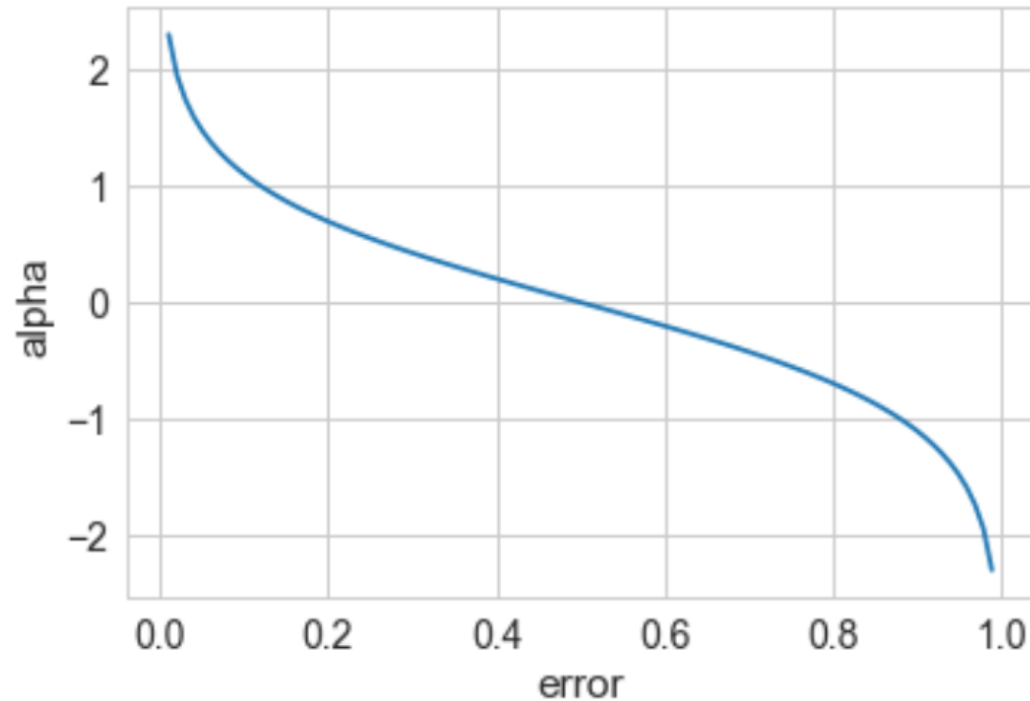
    3. Update weights:

$$D_{t+1}(i) = \frac{D_t(i)exp(-\alpha_t y_i f_t(x_i, \theta))}{Z_t},$$

    where $Z_t$ is a normalizing factor that makes sure that $\sum D_{t+1}(i) = 1.$

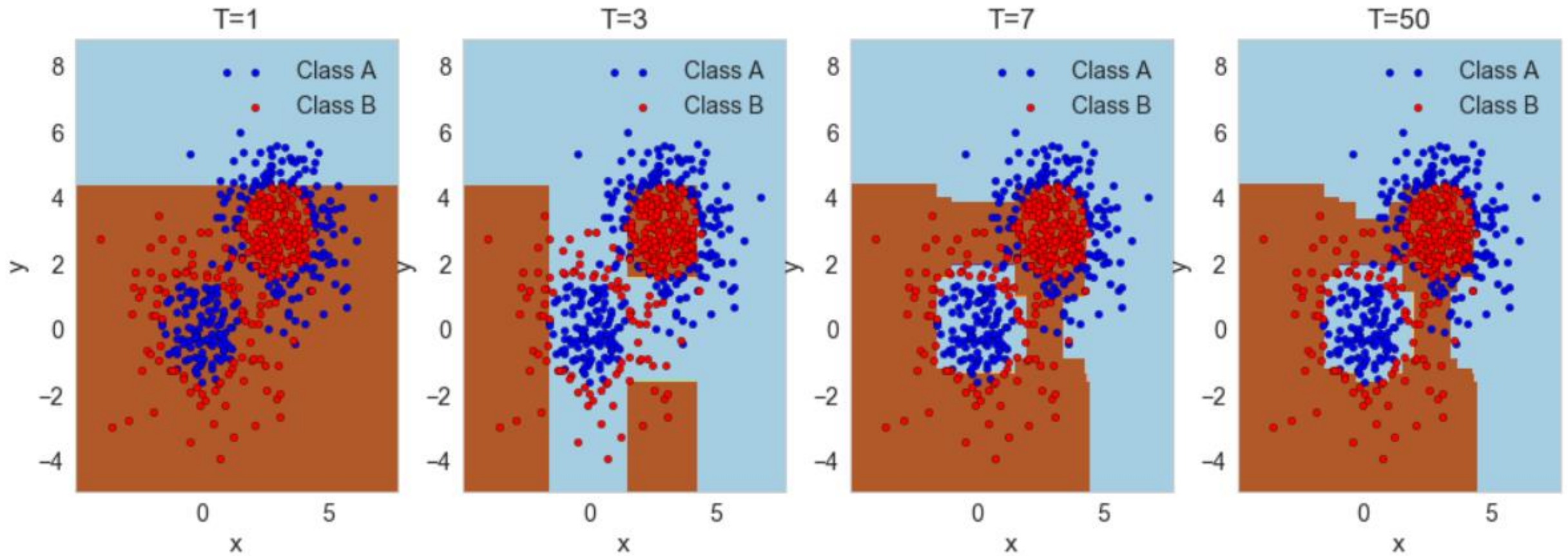$$\hat{f}(x, \theta) = \sum_{t=1}^{T} \alpha_t f_t(x, \theta)$$

# Boosting: adaboost



o The weight of a weak model in the boosted meta-model increases exponentially as the error approaches 0. Better models are given exponentially more weight.

o The weight is zero if the error rate is 0.5. A model with 50% accuracy is no better than random guessing, so it is ignored.

o The weight decreases exponentially as the error approaches 1. A negative weight is given to classifiers with worse than 50% accuracy. "Whatever that classifier says, do the opposite!".

$$\hat{f}(x, \theta) = \sum_{t=1}^{T} \alpha_t f_t(x, \theta)$$
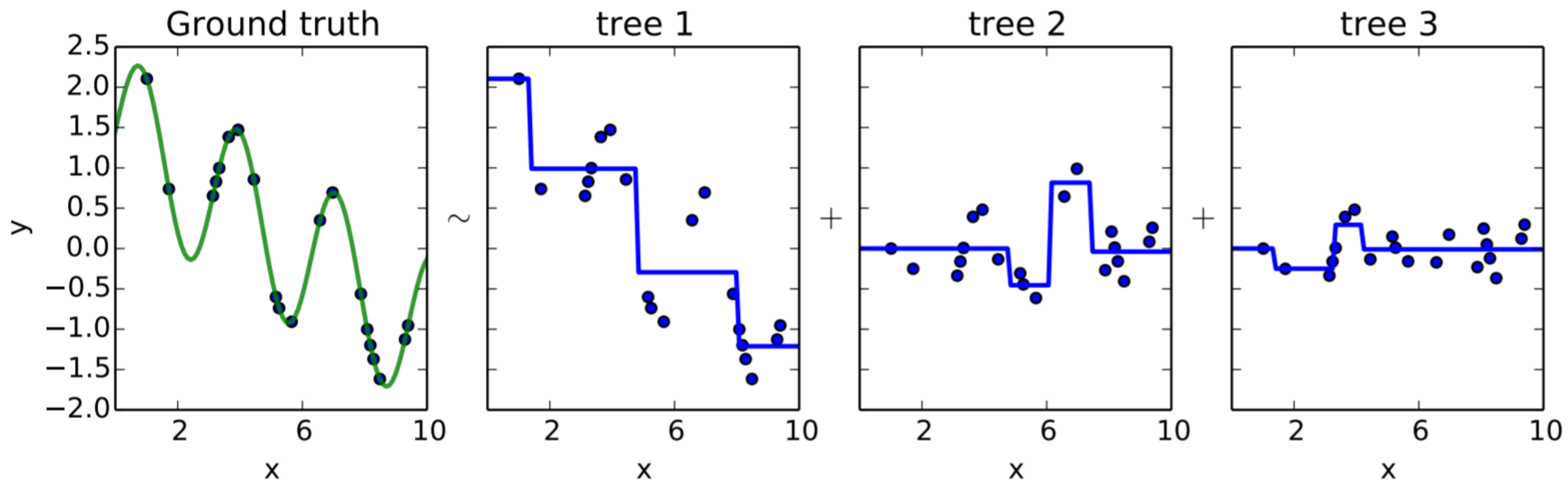
# Boosting: adaboost

# Boosting: Gradient Boosting

$$\hat{f}(x, \theta) = \sum_{t=1}^{T} f_t(x, \theta)$$

1. Fit a model $f_1(x, \theta) = y$

2. Fit a model to the **residuals** $h_1(x) = y - f_1(x, \theta)$

3. Create a new model $f_2(x, \theta) = f_1(x, \theta) + h_1(x)$

$$f_0(x, \theta) = \frac{1}{n} \sum_{i=1}^{n} y_i$$

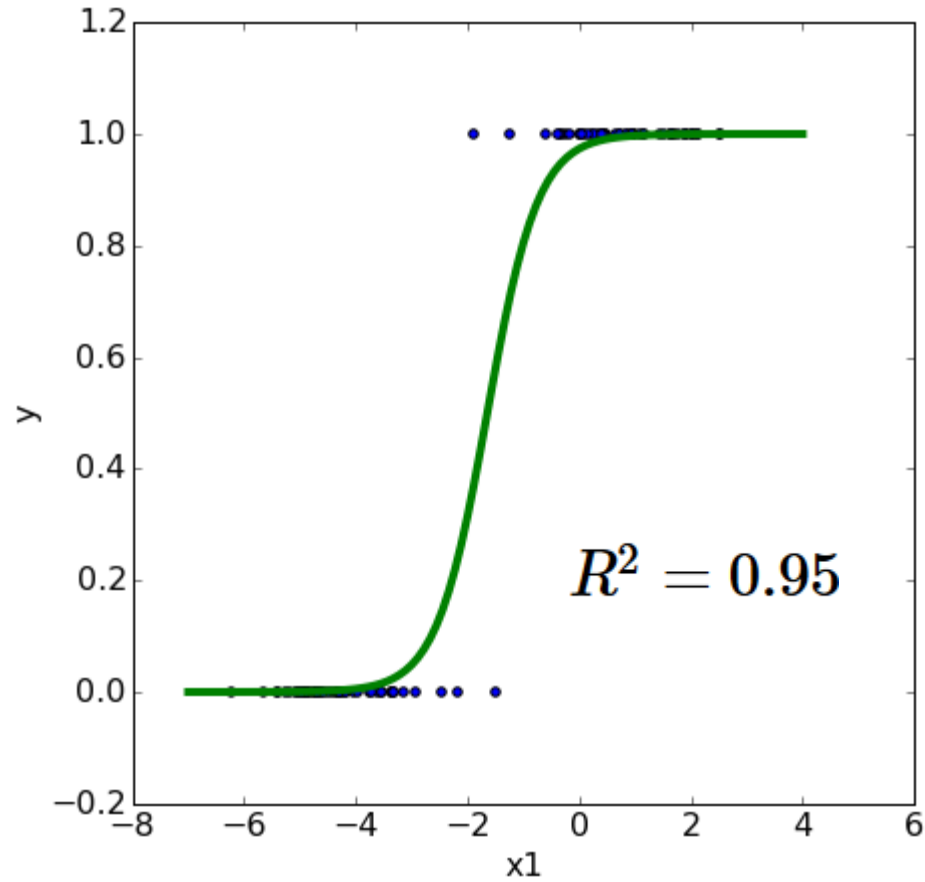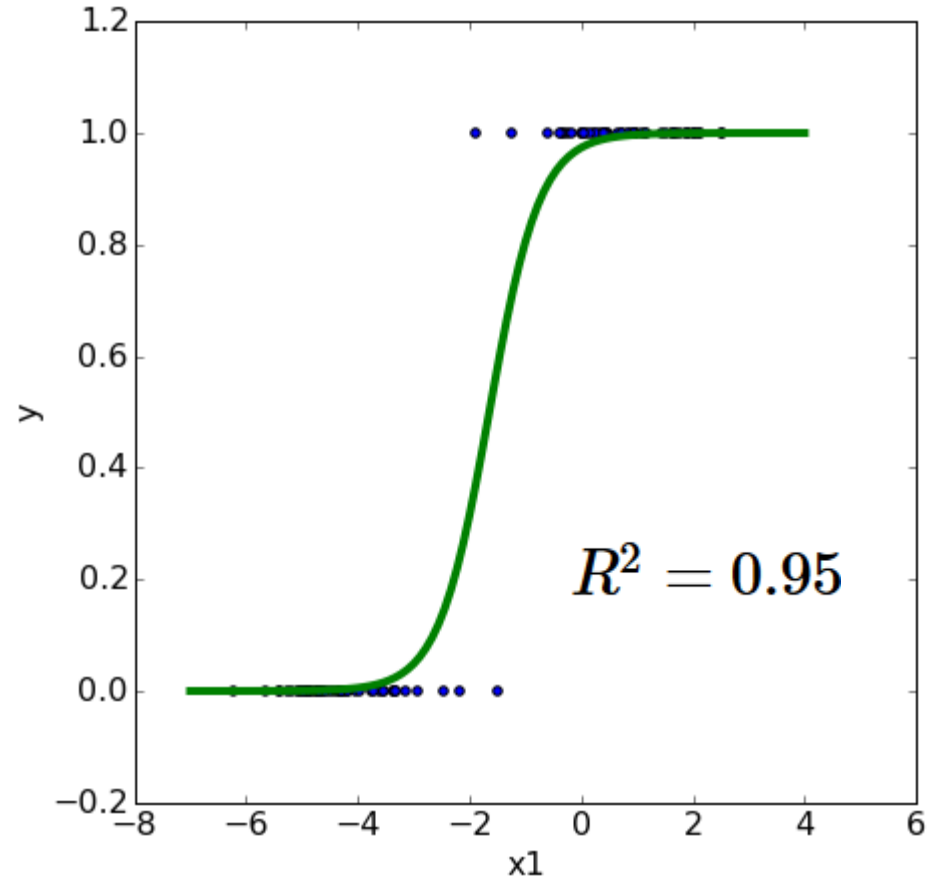$$f_{t+1}(x, \theta) = f_t(x, \theta) + h_t(x)$$

Ground truth $\sim$ tree 1 $+$ tree 2 $+$ tree 3

# Let's get real!

**6.   Ensemble learning: boosting**

# *logistic regression*



We need to make

**assumptions** *linearly separable*

about the

**model** *logistic model*

that generated the data.

# *logistic regression: logistic model*



$$f(x, \theta) = g(\theta_0 + \theta_1 x_1)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$
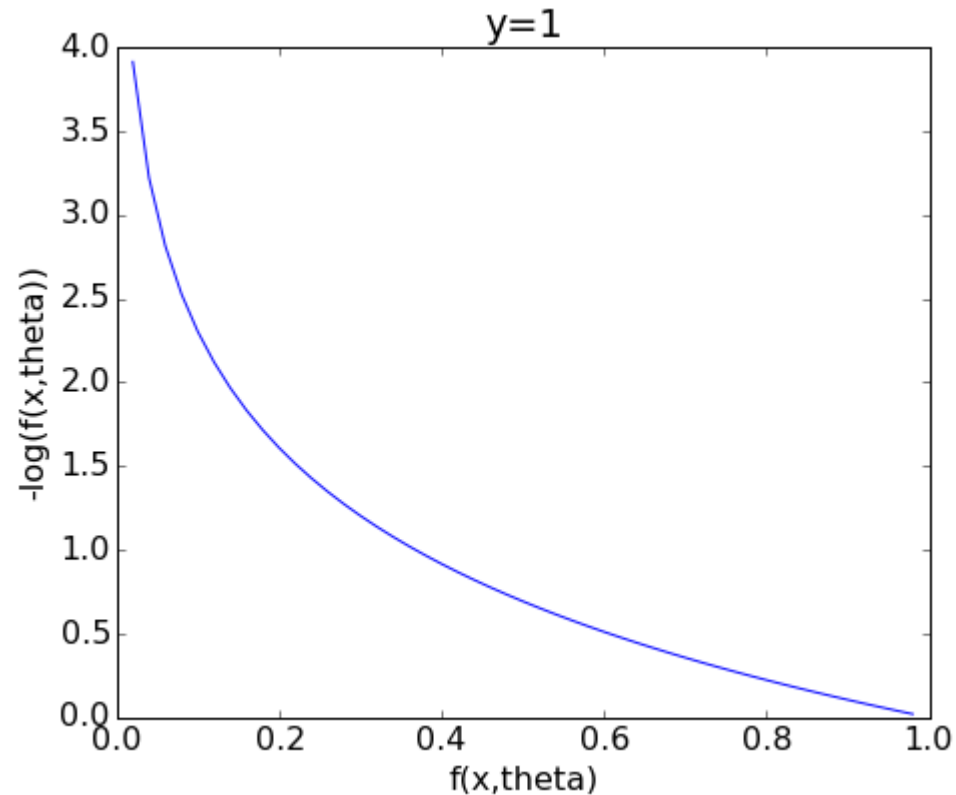
# logistic regression: cost function

$$J(\theta) = -[\frac{1}{n} \sum_{i=1}^{n} y^{(i)} log(f(x^{(i)}, \theta)) + (1 - y^{(i)})log(1 - f(x^{(i)}, \theta))]$$

We know that $y^{(i)}$ is either 0 or 1. If $y^{(i)} = 1$ then the cost function $J(\theta)$ is incremented by

$-log(f(x^{(i)}, \theta))$.

Similarly, if $y^{(i)} = 0$ then the cost function $J(\theta)$ is incremented by
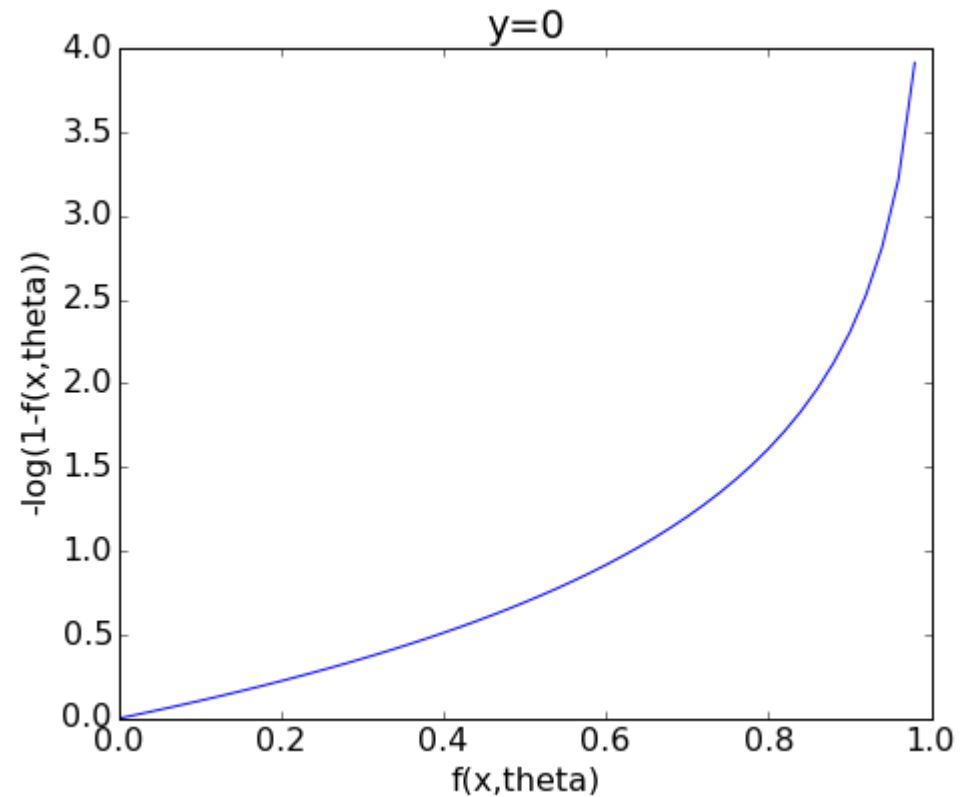
$-log(1 - f(x^{(i)}, \theta))$.

# logistic regression: cost function

We know that $y^{(i)}$ is either 0 or 1. If $y^{(i)} = 1$ then the cost function $J(\theta)$ is incremented by

$$-log(f(x^{(i)}, \theta)).$$

# *logistic regression: cost function*

Similarly, if $y^{(i)} = 0$ then the cost function $J(\theta)$ is incremented by

$$-log(1 - f(x^{(i)}, \theta)).$$

# logistic regression
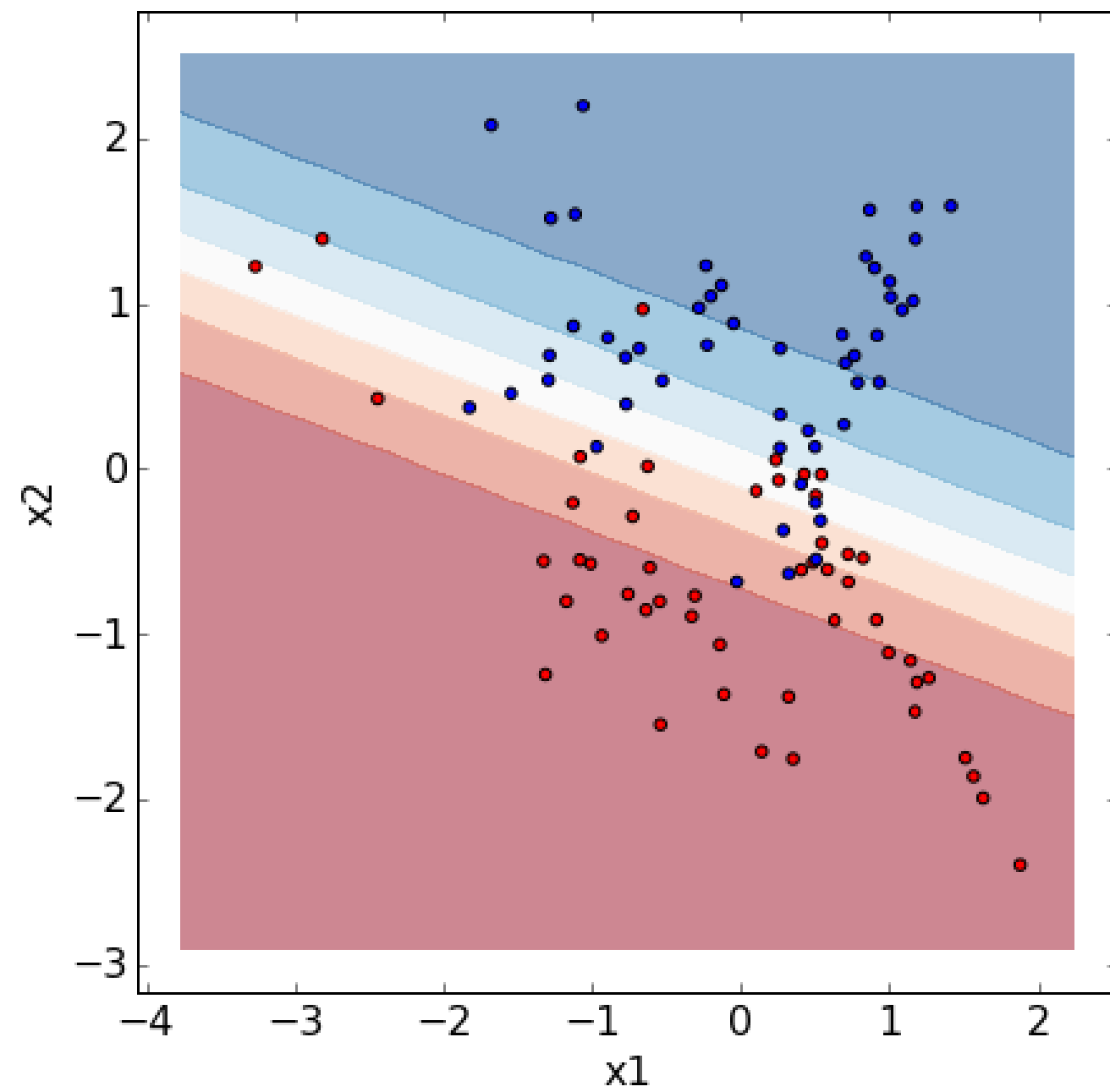
Fit a logistic model

$$f(x, \theta) = g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_m x_m) = g(\theta' x)$$

to the data set such that the cost function

$$J(\theta) = -[\frac{1}{n} \sum_{i=1}^{n} y^{(i)} log(f(x^{(i)}, \theta)) + (1 - y^{(i)})log(1 - f(x^{(i)}, \theta))]$$

is minimal using gradient descent
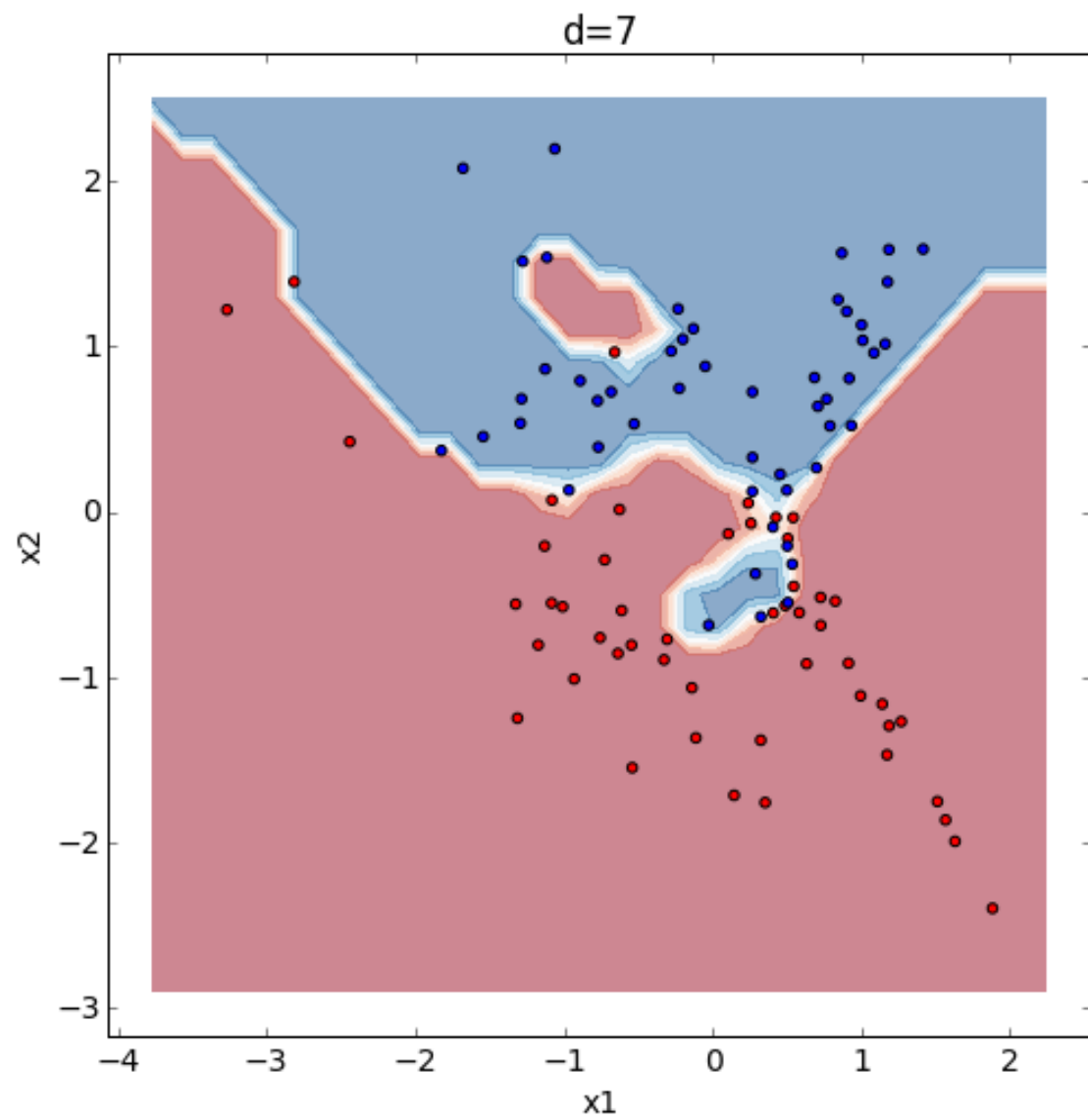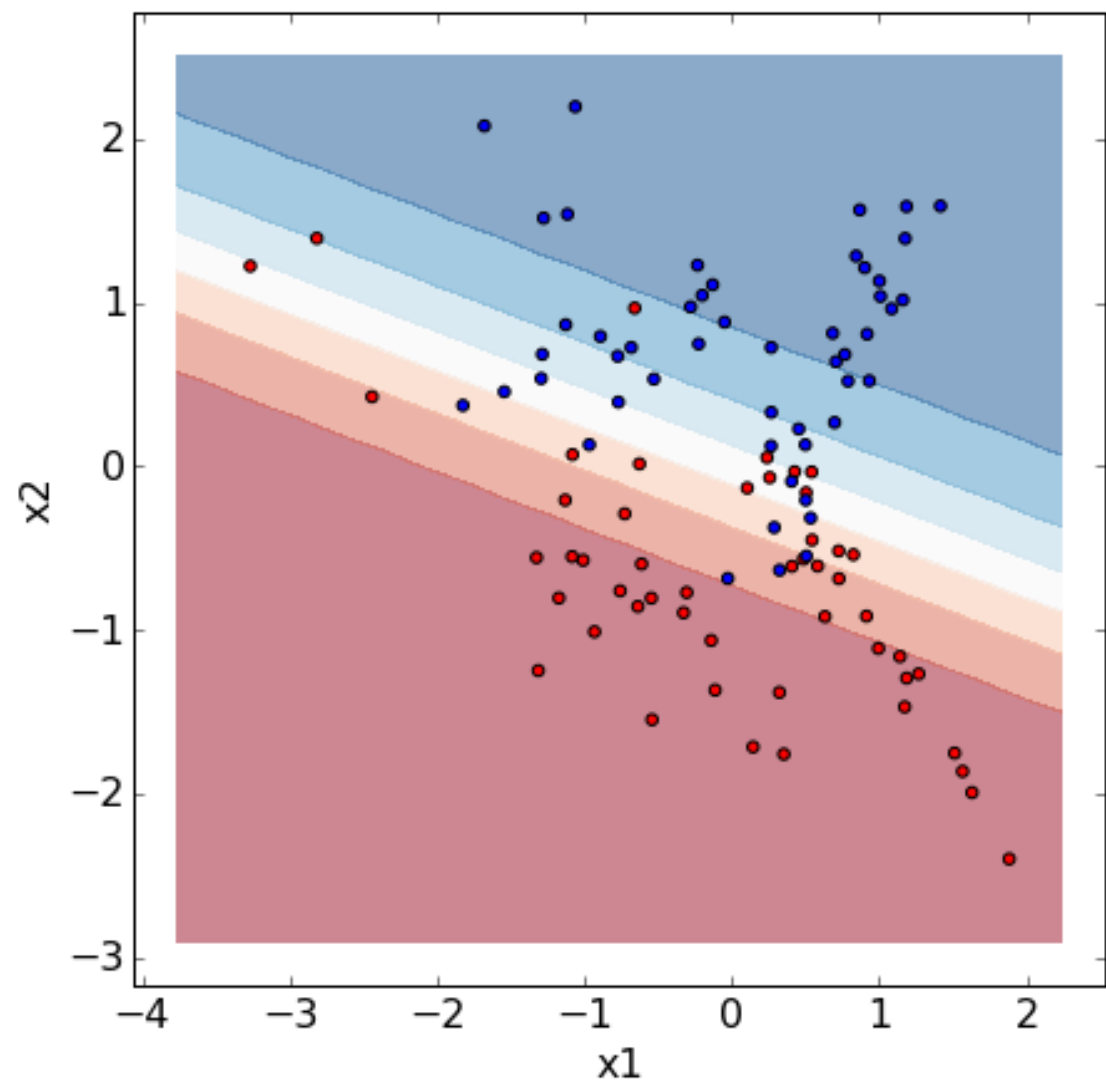
$$\theta_j := \theta_j - \alpha \frac{1}{n} \sum_{i=1}^{n} (f(x^{(i)}, \theta) - y^{(i)})x_j^{(i)}$$

# Let's get real!

7. **Logistic regression**
8. **Non-linear transformations**

## regularized logistic regression

$$f(x, \theta) = g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_m x_m)$$
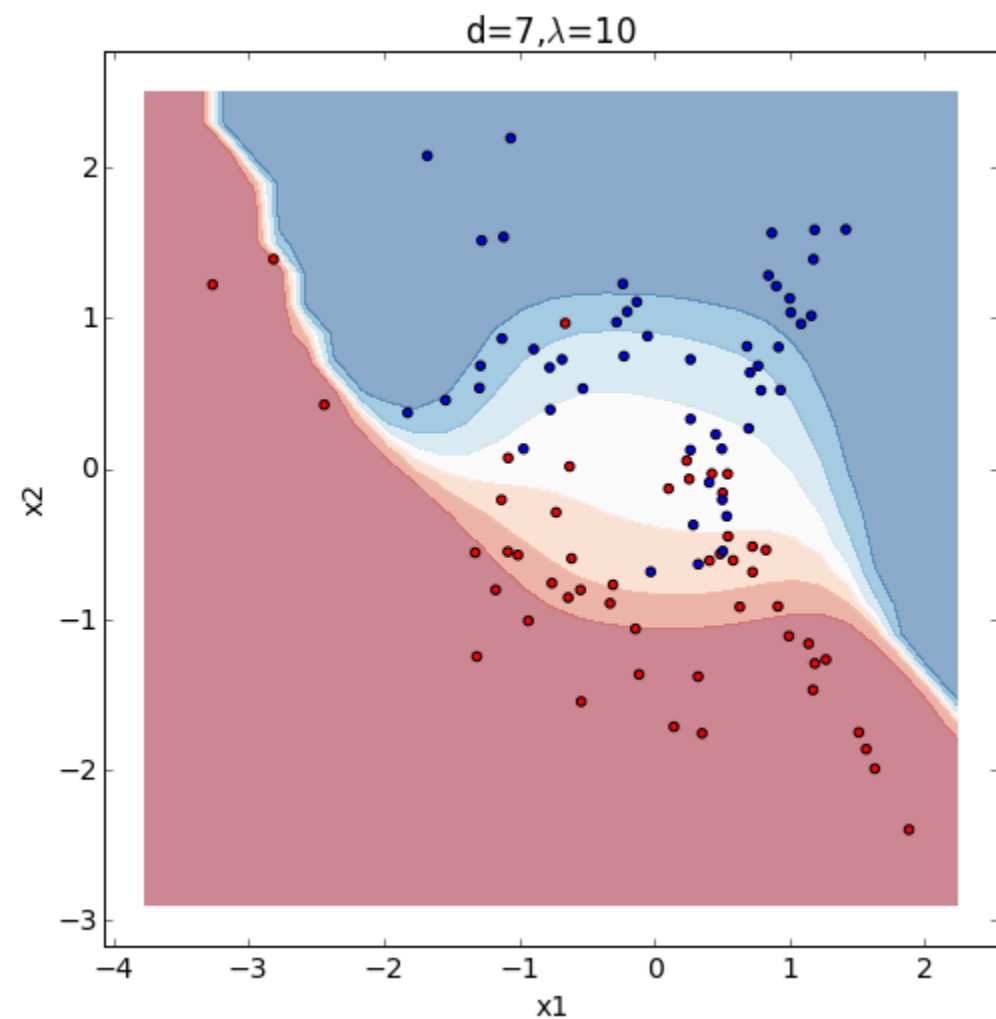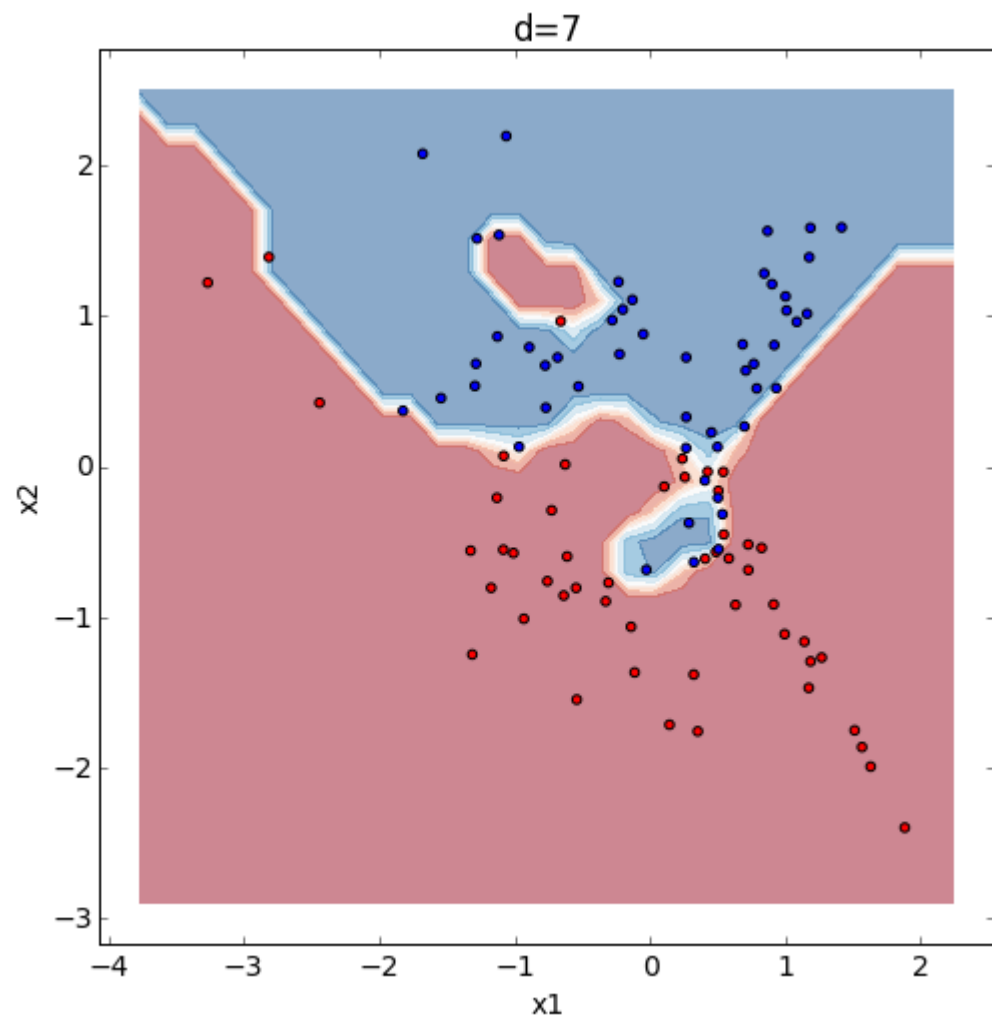
*cost function*

$$J(\theta) = -[\frac{1}{n} \sum_{i=1}^{n} y^{(i)} log(f(x^{(i)}, \theta)) + (1 - y^{(i)})log(1 - f(x^{(i)}, \theta))]$$

$$J(\theta) = -[\frac{1}{n} \sum_{i=1}^{n} y^{(i)} log(f(x^{(i)}, \theta)) + (1 - y^{(i)})log(1 - f(x^{(i)}, \theta))] + \frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2$$

*regularized cost function*

# regularized logistic regression

# support vector machines
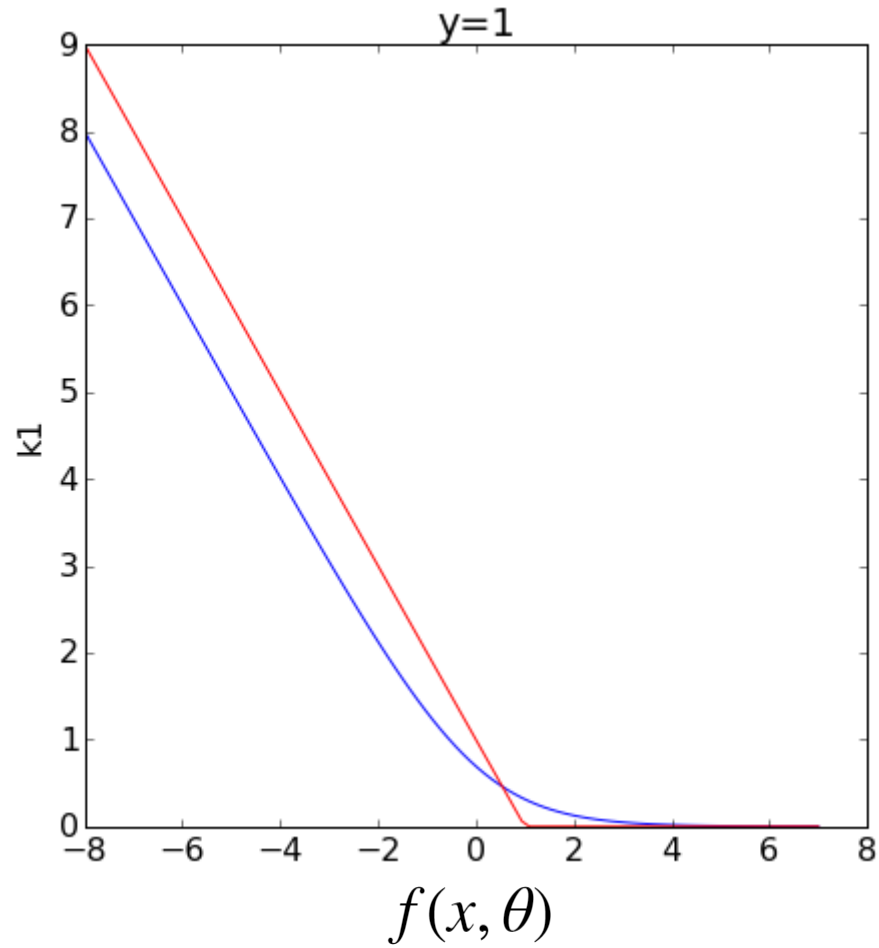
Fit a linear model

$$f(x, \theta) = \theta' x$$

such that

$$J(\theta) = \left[ C \sum_{i=1}^{n} y^{(i)} k_1(\theta' x^{(i)}) + (1 - y^{(i)}) k_0(\theta' x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^{m} \theta_j^2$$

with  $k_1(\theta' x) = max(0, 1 - \theta' x)$   and   $k_0(\theta' x) = max(0, 1 + \theta' x)$
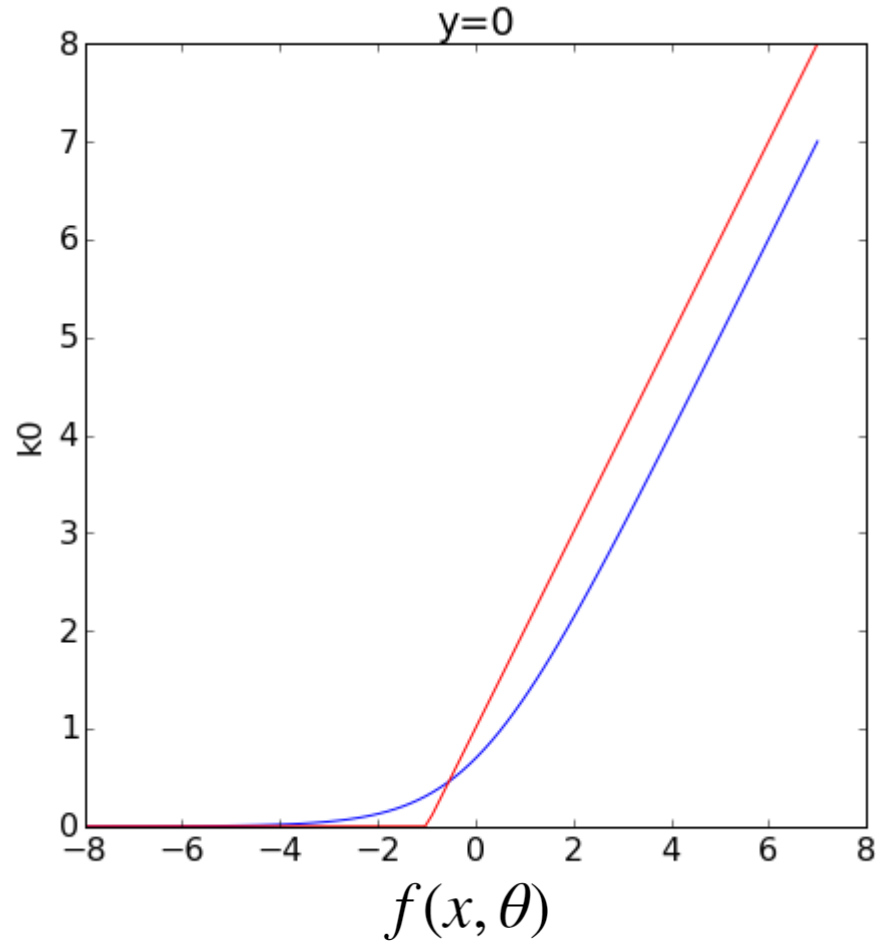
is minimized.

# support vector machines



- o replace cost function by piecewise linear function

- o if $y = 1$ then the contribution to the cost is

$$k_1(f(x, \theta)) = max(0, 1 - f(x, \theta))$$

# support vector machines



y=0

- replace cost function by piecewise linear function

- if *y* = 1 then the contribution to the cost is

$$k_1(f(x, \theta)) = max(0, 1 - f(x, \theta))$$

- if *y* = 0 then the contribution to the cost is

$$k_0(f(x, \theta)) = max(0, 1 + f(x, \theta))$$

# kernel support vector machines

SVMs can also be formulated as a linear function of the samples (dual form) instead of the features as

$$f(x, \theta) = \sum_{i=1}^{n} \theta_i (x \cdot x^{(i)}) + \theta_0$$

that can be reformulated as a non-linear function using what is know as a kernel function
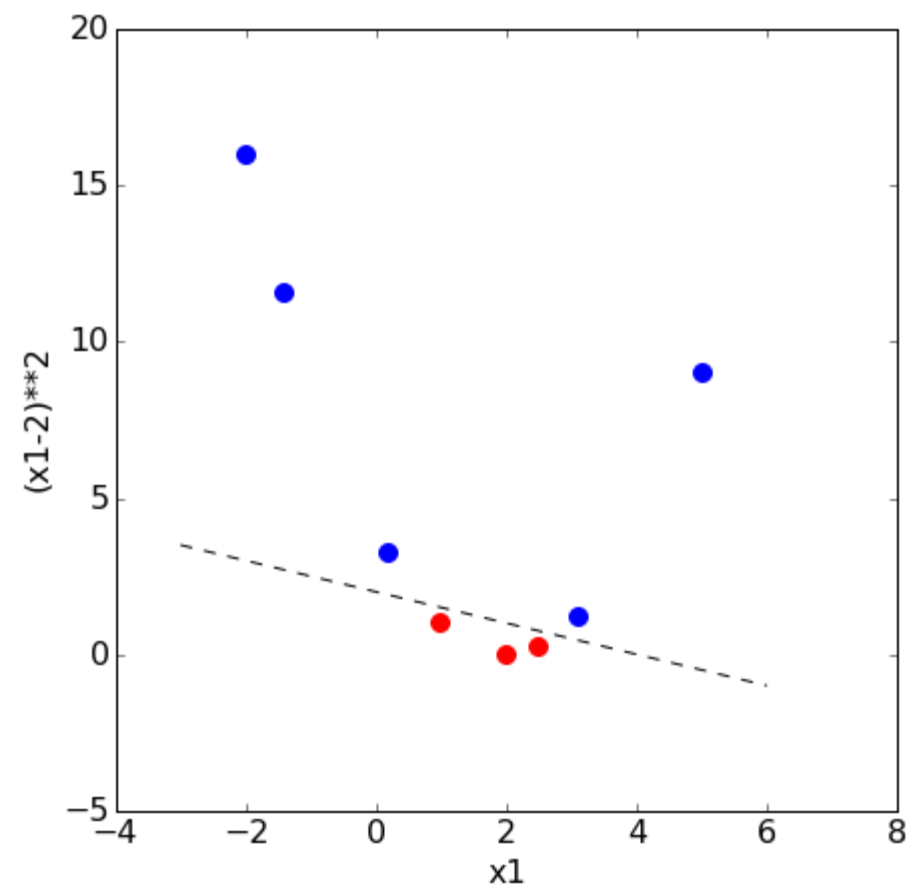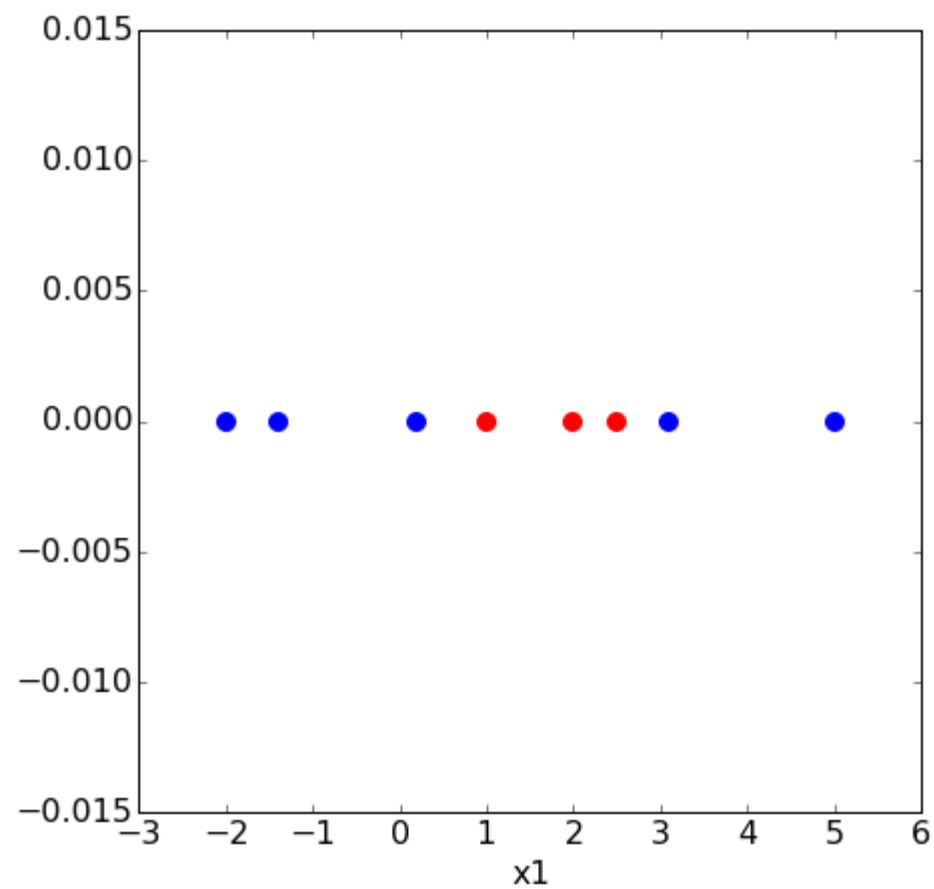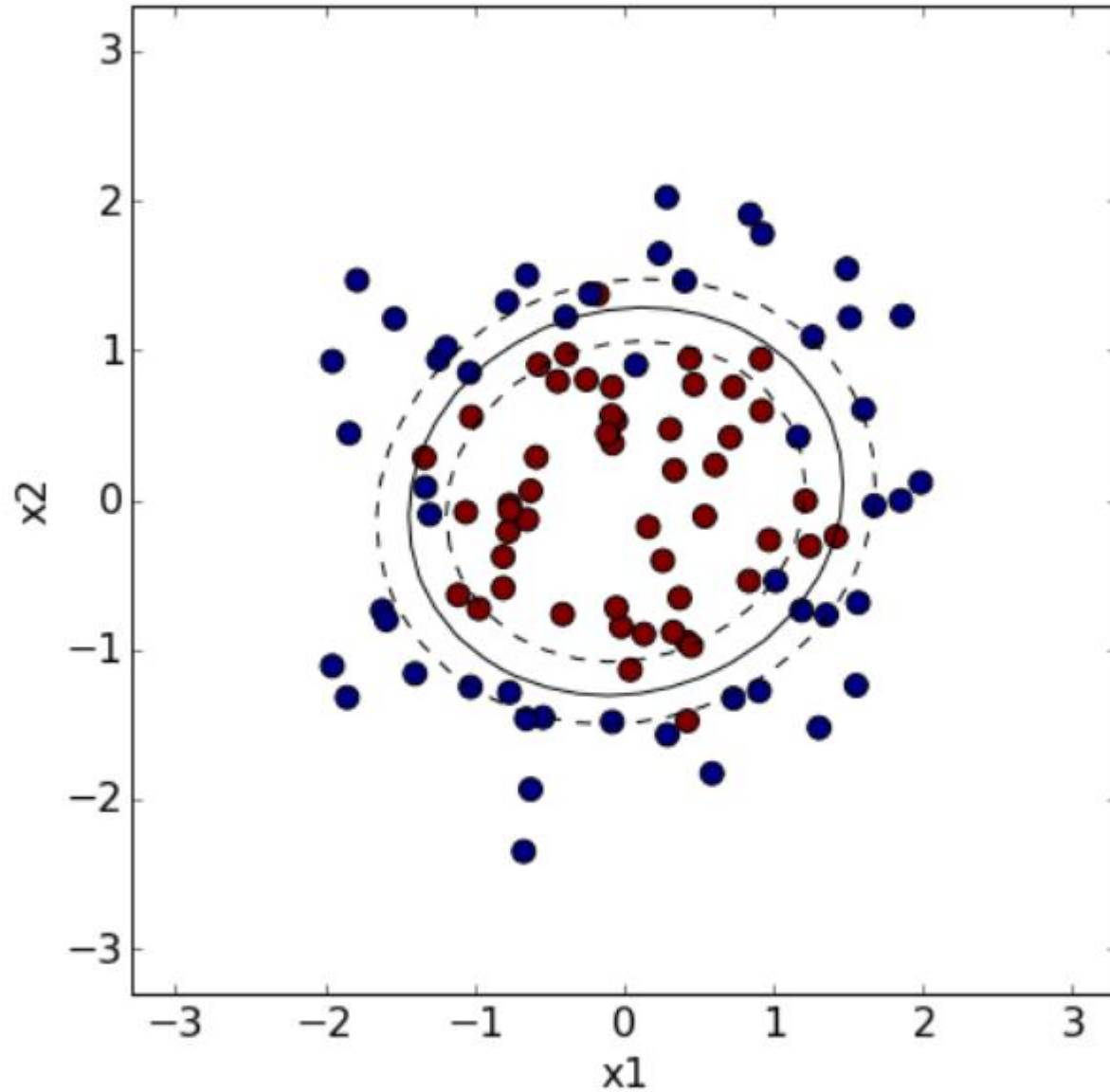
$$K(x^{(i)}, x^{(j)})$$

to become

$$f(x, \theta) = \sum_{i=1}^{n} \theta_i K(x, x^{(i)}) + \theta_0$$

The data points $x^{(i)}$ for which $\theta_i > 0$ are called the support vectors.

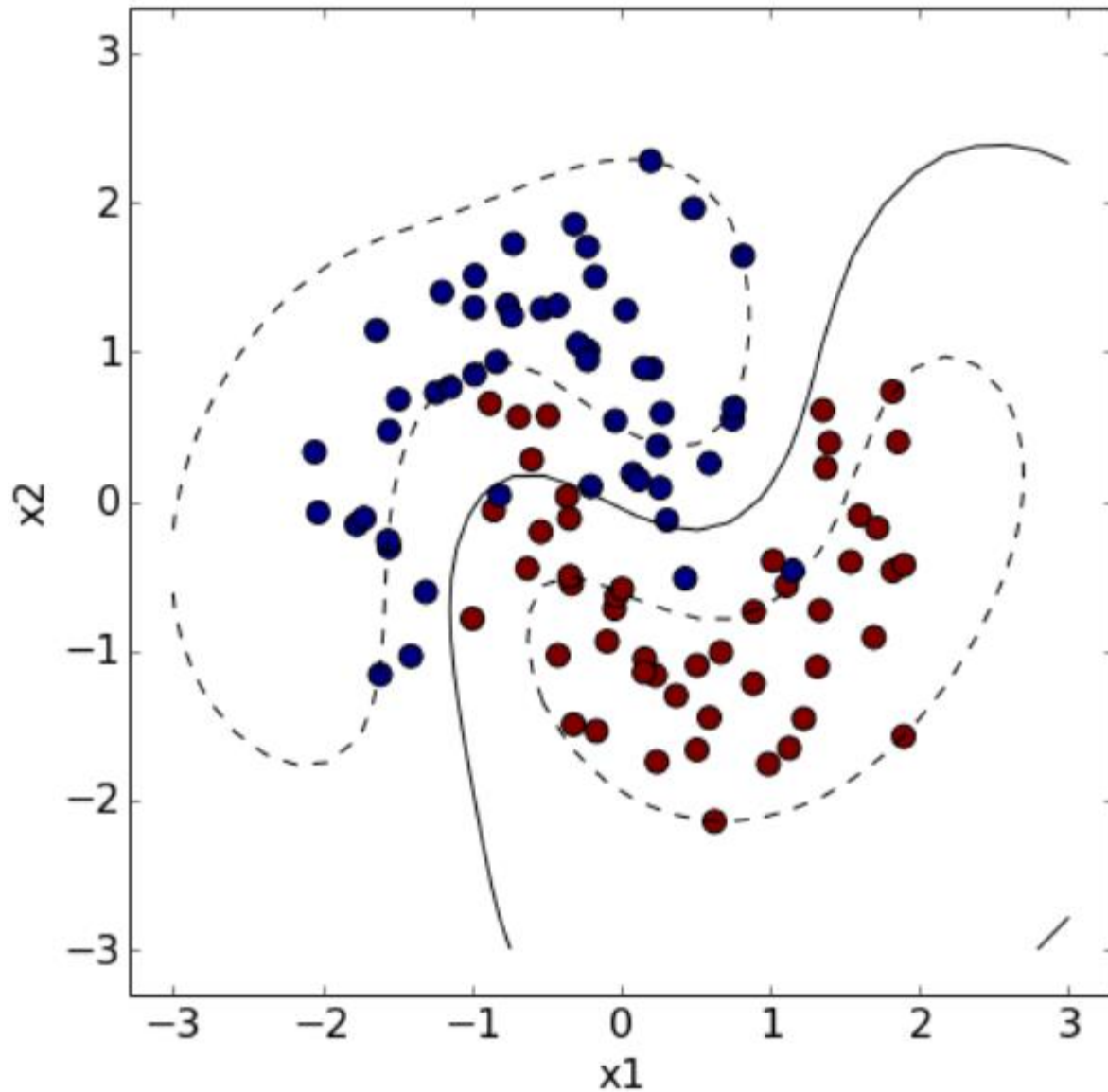# kernel support vector machines

# kernel support vector machines



$$f(x, \theta) = \sum_{i=1}^{n} \theta_i K(x, x^{(i)}) + \theta_0$$

$$K(x^{(i)}, x^{(j)}) = (x^{(i)} \cdot x^{(j)} + c)^d$$

# kernel support vector machines



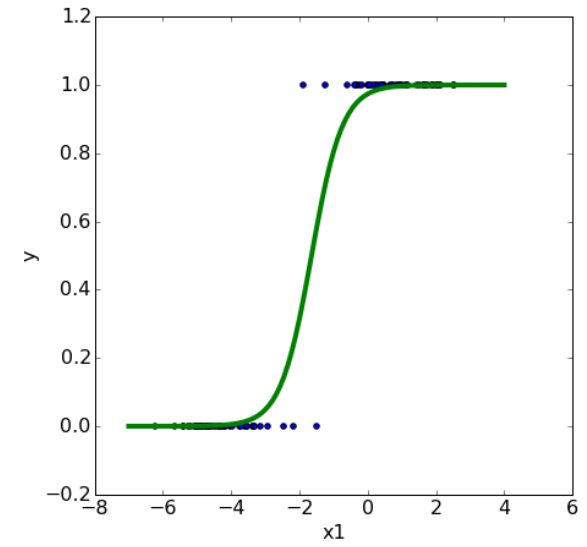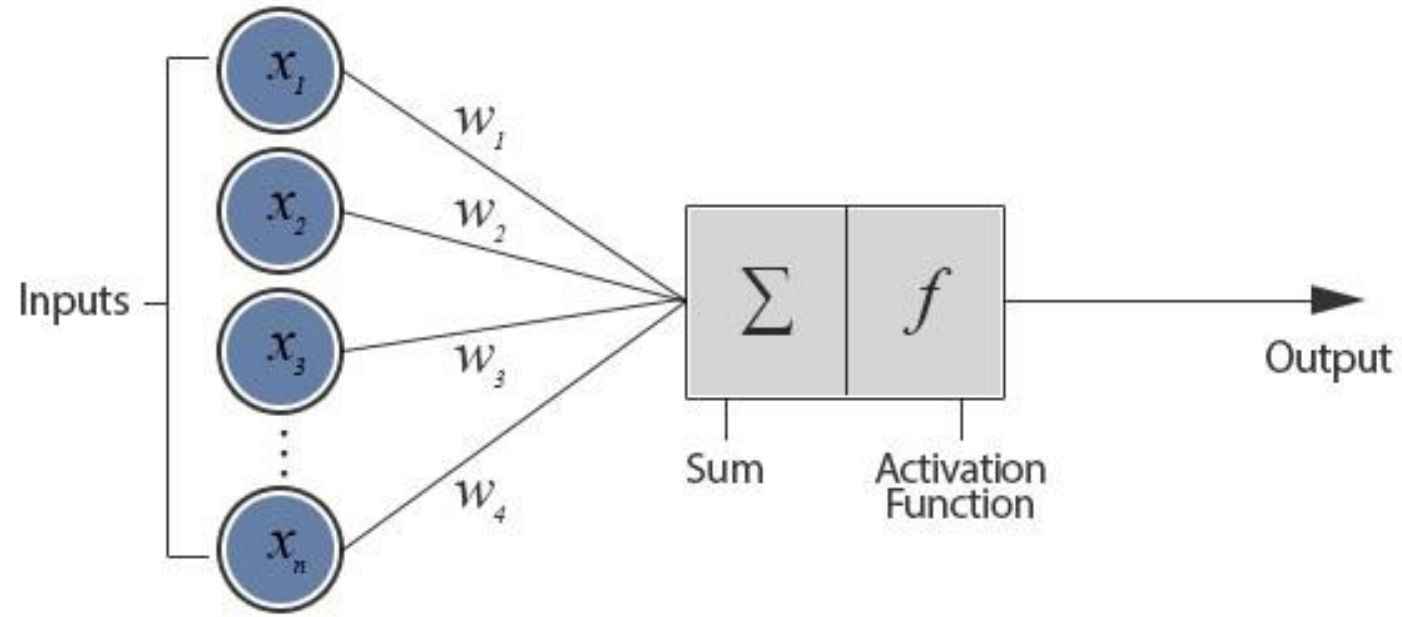$$f(x, \theta) = \sum_{i=1}^{n} \theta_i K(x, x^{(i)}) + \theta_0$$

$$K(x^{(i)}, x^{(j)}) = \exp\left[-\frac{\|x^{(i)} - x^{(j)}\|^2}{2\sigma^2}\right]$$

# Let's get real!

**9. Support Vector Machines**

$$f(x, \theta) = g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_m x_m) = g(\theta' x)$$



Inputs

$x_1$

$x_2$

$x_3$

$\vdots$

$x_n$

$w_1$

$w_2$

$w_3$

$w_4$

$\Sigma$ $\quad$ $f$
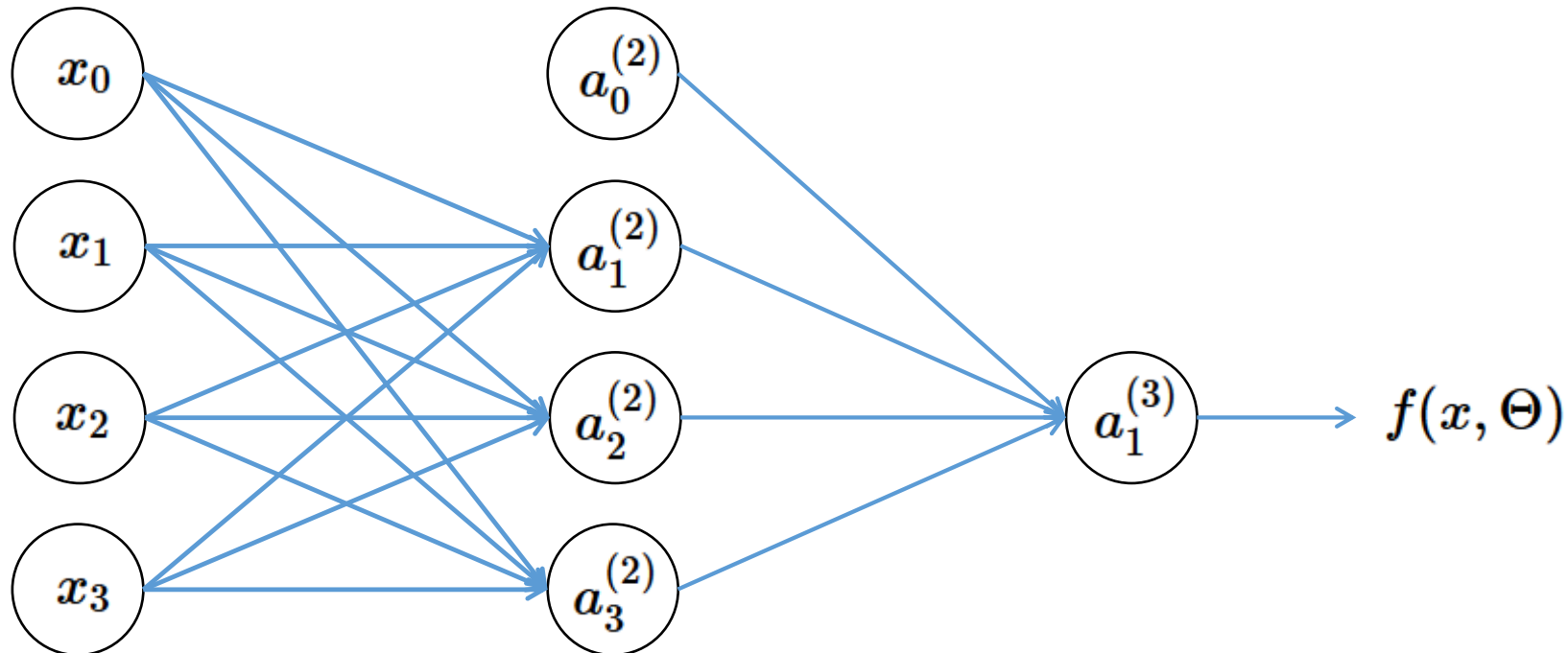
Sum

Activation Function

Output

Model: $f(x, \Theta) = g(\Theta_{10}^{(2)} a_0 + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$

$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

Cost function logistic regression:

$$J(\theta) = -[\frac{1}{m} \sum_{i=1}^{n} y^{(i)} log(f(x^{(i)}, \theta) + (1 - y^{(i)}) log(1 - f(x^{(i)}, \theta))] + \frac{\lambda}{2n} \sum_{j=1}^{n} \theta^2$$

Cost function feedforward neural network:

$$J(\theta) = -[\frac{1}{m} \sum_{i=1}^{n} \sum_{k=1}^{K} y_k^{(i)} log(f(x^{(i)}, \Theta)_k) + (1 - y_k^{(i)}) log(1 - f(x^{(i)}, \Theta)_k)] + \frac{\lambda}{2n} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

https://playground.tensorflow.org/

# Let's get real!

**10.  Neural Networks**

# Blending/Stacking