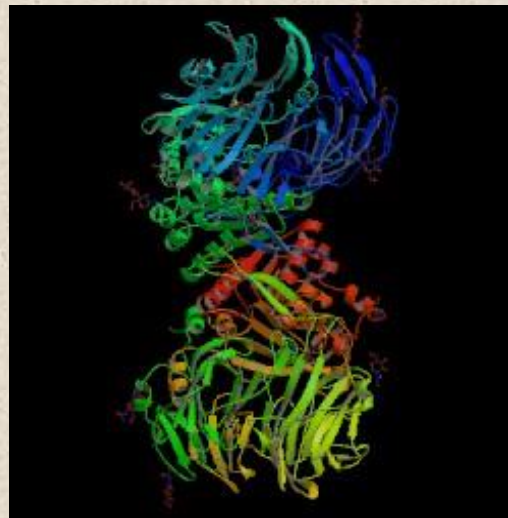
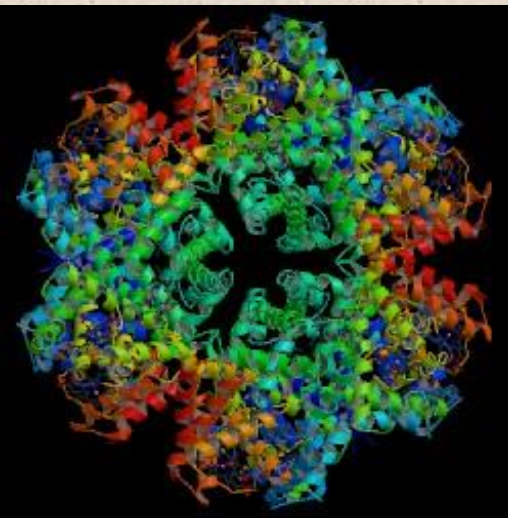
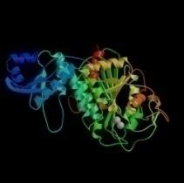


# Understanding SVMs

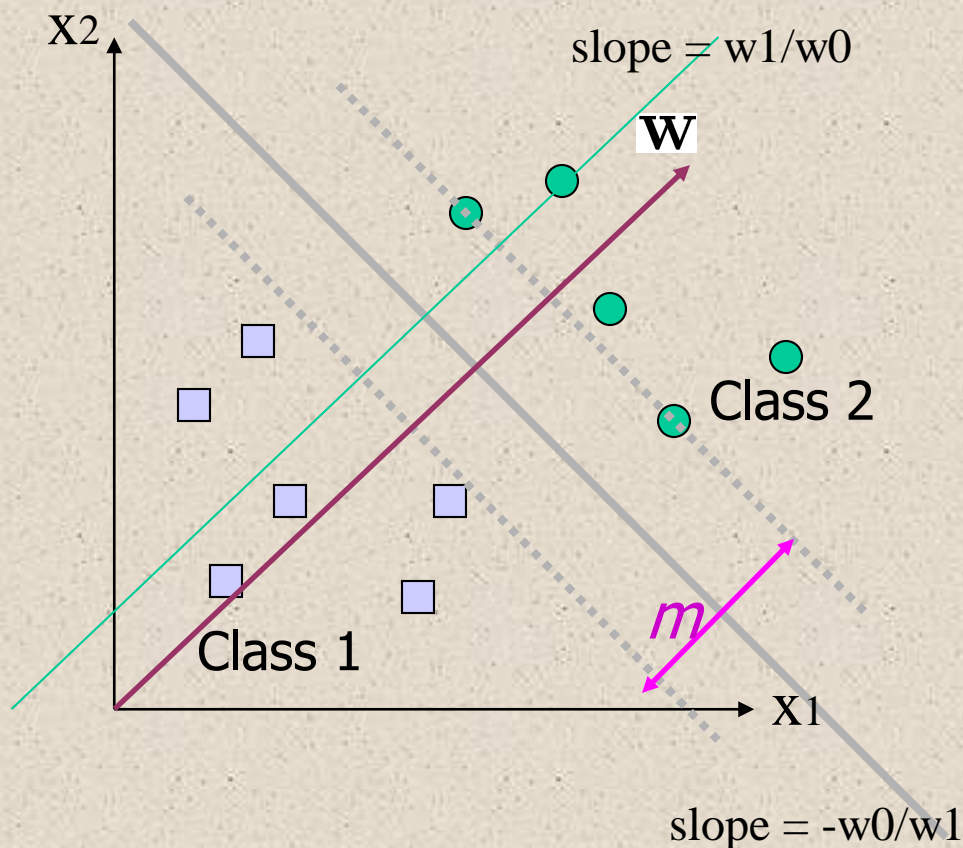
Professor Ami M. Gates  
Georgetown University





# A Review of Vector Math

- Any line in dimension D can be represented as  $\mathbf{w}^T \mathbf{x} + b = 0$ .
- $\mathbf{w}$  is a vector of coefficients,  $\mathbf{x}$  is the vector of variables,  $b$  is the translation (can be thought of in 2D as the y intercept).



Suppose we are in 2D. (the common Cartesian coordinate system).

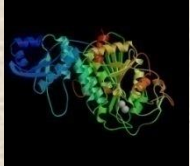
Then, a line can be written as

$$\mathbf{w}_0 x_0 + \mathbf{w}_1 x_1 + b = 0$$

$$x_1 = -b/w_1 - (w_0/w_1)(x_0)$$

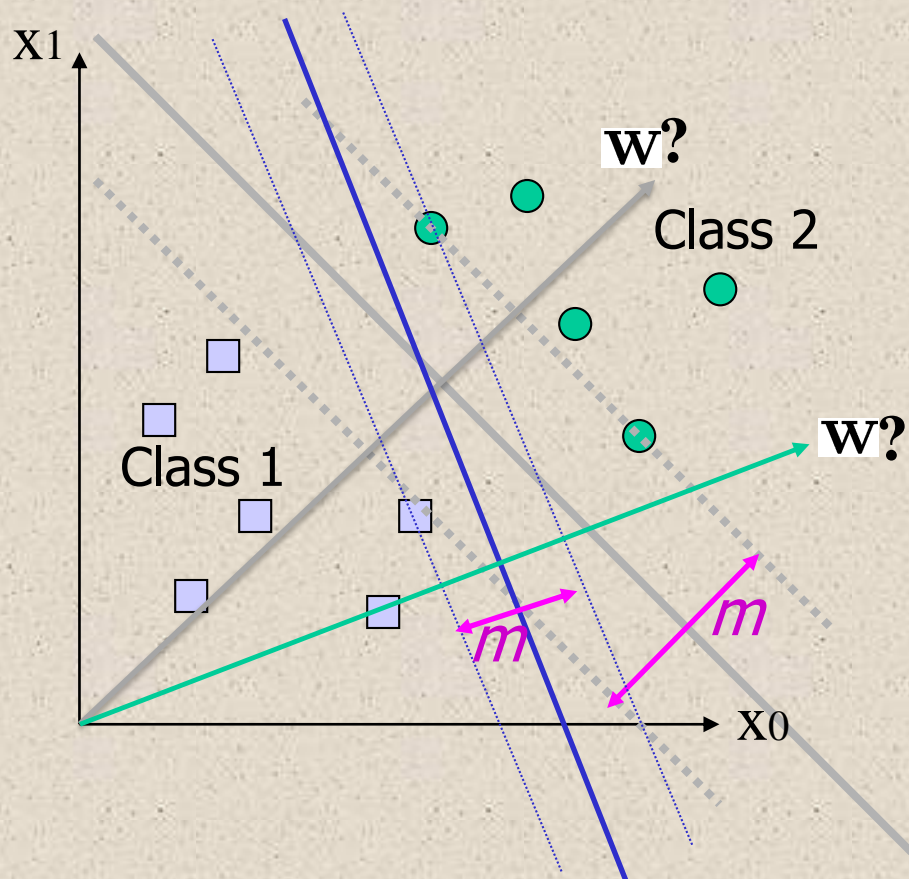
→ The “slope” is  $-w_0/w_1$

Given any vector  $\mathbf{w}$  (in 2D,  $[w_0, w_1]$ , from the origin, a line parallel to  $\mathbf{w}$  has slope  $w_1/w_0$   
The line **perpendicular** to that vector has slope  $-w_0/w_1$



# A Review of Vector Math, cont.

- The goal of an SVM is to determine (via training) the “best” line in 2D (or hyperplane in higher D) that **separates** two classes.
- **There are an infinite number of possible lines (hyperplanes)** that may work.



A general linear equation in two variables is:

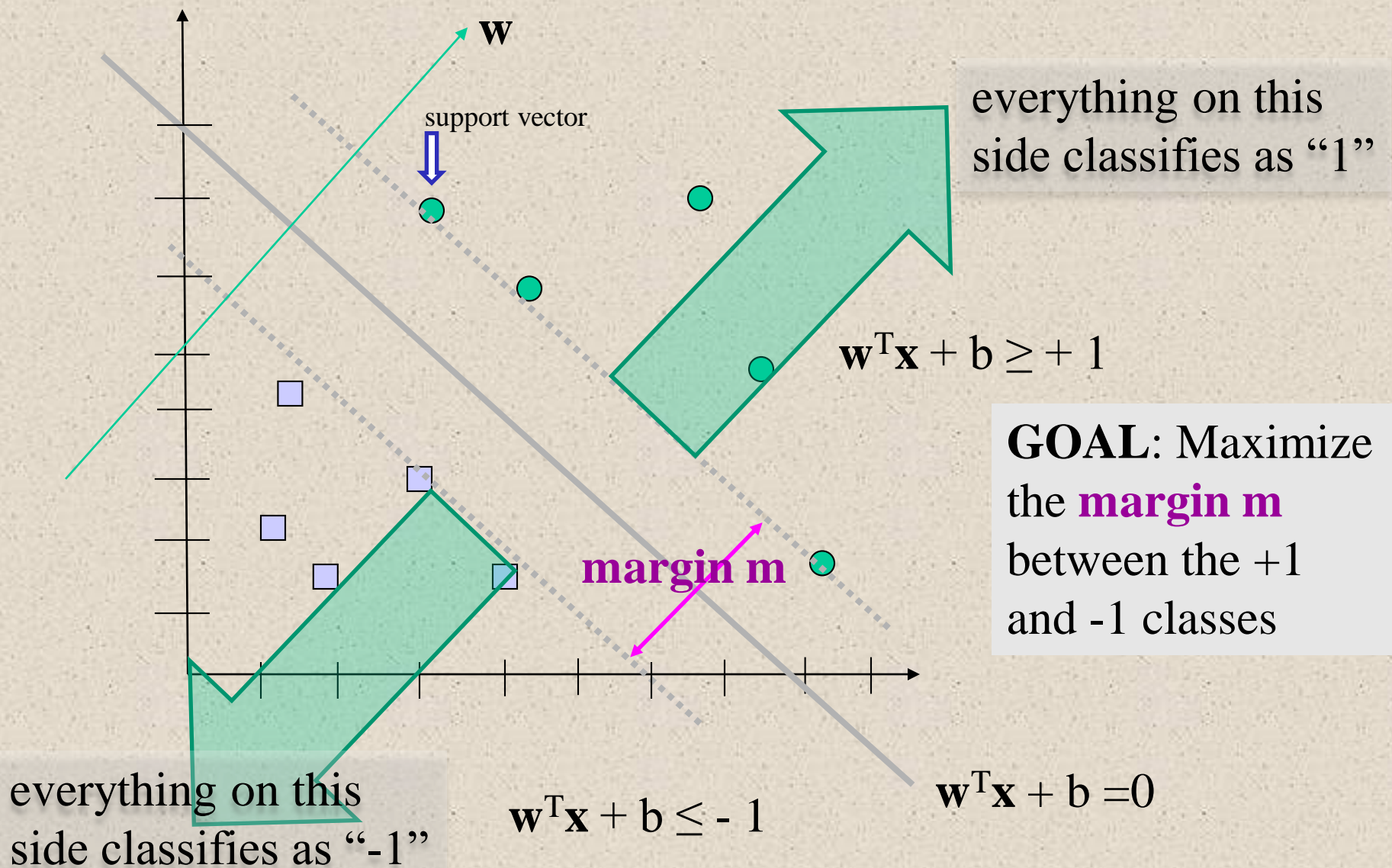
$$w_0x_0 + w_1x_1 + b = 0$$

$$\rightarrow w^T x + b = 0$$

An SVM algorithm must determine the weight vector  $w$  that will result in a line (hyperplane) that best separates the two classes.



# Setting Up the SVM Problem



# Important Concepts

- The **margin m** can be created such that

$$w_0x_0 + w_1x_1 + b \geq 1$$

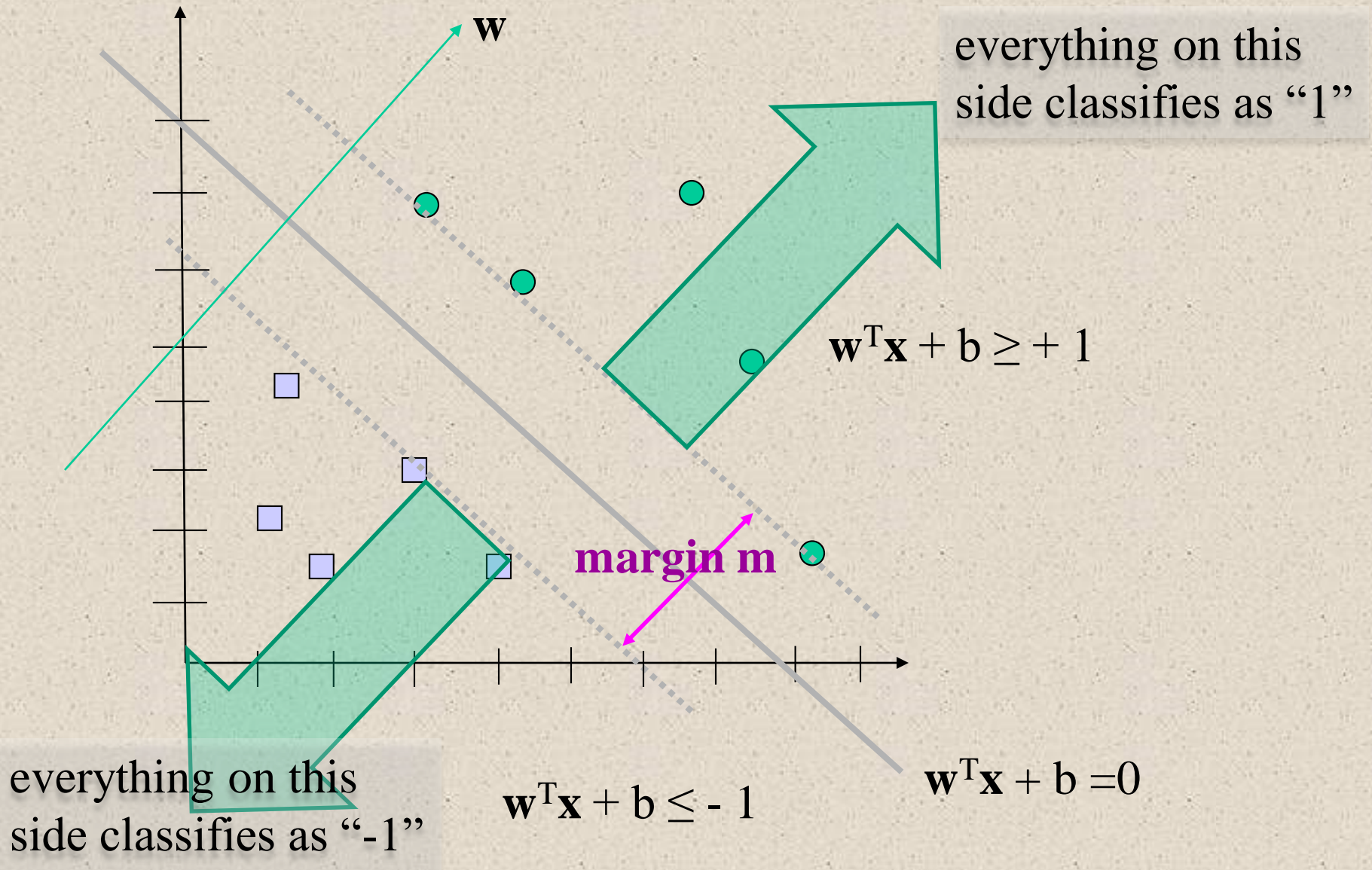
$$\rightarrow \mathbf{w}^T \mathbf{x} + b \geq 1$$

$$w_0x_0 + w_1x_1 + b \leq -1$$

$$\rightarrow \mathbf{w}^T \mathbf{x} + b \leq -1$$

- Remember that **w will have to be determined (using the training data)** so that the classes are “best” separated.
- The value of “b” is a constant that affects the location (but not slope) of the line (in 2D).
- These concepts are the same in higher D.

The **margin  $m$**  between the two classes should be maximized.



# Key Math Concepts

1) Any line in 2D or hyperplane in higher D can be represented as:

$$w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n + b = 0$$

→ same as:  $\mathbf{w}^T \mathbf{x} + b = 0$

2) This linear equation can be represented with vector  $\mathbf{w}$ , with vector  $\mathbf{x}$ , and with  $b$ .

If vector  $\mathbf{w}$  is plotted, it will be perpendicular to the linear equation above.

The distance between any hyperplane  $\mathbf{w}^T \mathbf{x} = k$  and the origin is  $k/\|\mathbf{w}\|$

If vector  $\mathbf{w}$  is altered, the “slope” (in 2D) or “rotation” (in higher D) is affected. If “ $b$ ” is altered the “y intercept (in 2D) or the “translation” (in high D) is affected.

3) We can use  $\mathbf{w}^T \mathbf{x} + b = 0$  to try different lines. **To do this, we need  $\mathbf{w}$  and  $b$  and a measure of “best”.**



# The SVM Margin $m$

Assume that in some dimension that we have two classes that can be linearly separated.

- Then, there must be a line (hyperplane) that separates the classes.
- This line is defined by  $\mathbf{w}^T \mathbf{x} + b = 0$ .
- This line can be contained within two parallel lines that create a “**margin**” between the two classes.

The two parallel lines can initially be defined as:

$$\mathbf{w}^T \mathbf{x} + b = 1 \quad \text{and} \quad \mathbf{w}^T \mathbf{x} + b = -1$$

(Note that  $b$  is a constant that can be altered as needed)

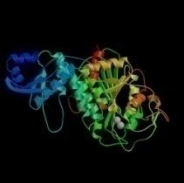
**margin  $m$ :**

Distance between two parallel lines (hyperplanes)

**$\mathbf{w}^T \mathbf{x} + b = 1$  and  $\mathbf{w}^T \mathbf{x} + b = -1$  is...**

$$|b - 1| - |b + 1| / \|\mathbf{w}\| = 2 / \|\mathbf{w}\|$$

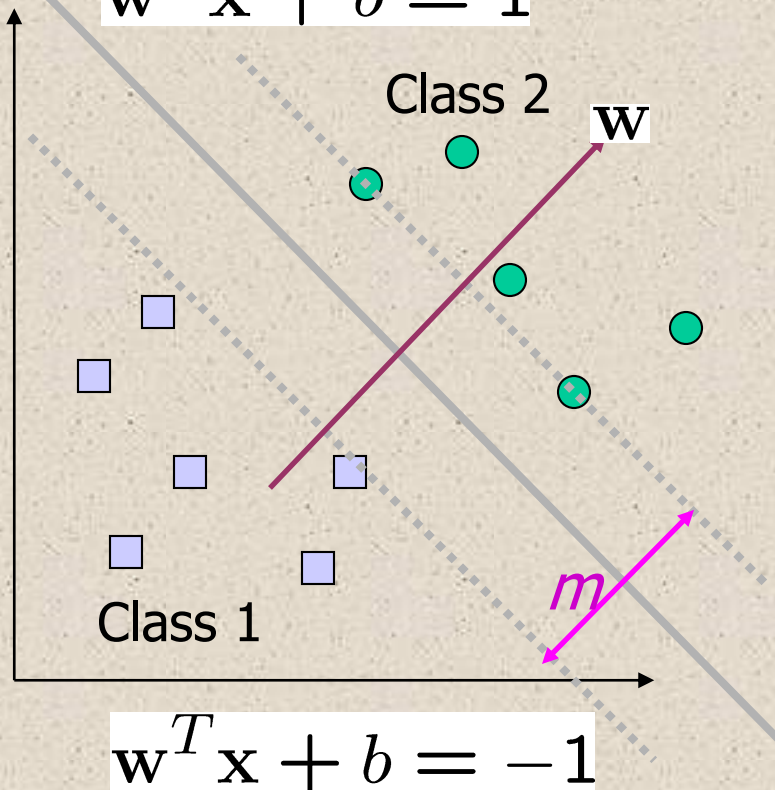




Maximizing margin  $m$  will create the “best” separation between the Classes.

$$\mathbf{w}^T \mathbf{x} + b = 0$$

$$\mathbf{w}^T \mathbf{x} + b = 1$$



**margin  $m$**  is the **distance** between the two parallel boundary lines (or two parallel hyperplanes)

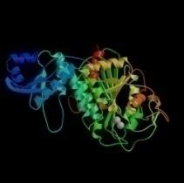
Distance between a line  $\mathbf{w}^T \mathbf{x} + b = 0$  (hyperplane) and the origin is  $b/\|\mathbf{w}\|$

**Margin  $m$ :**

Distance between two parallel lines (hyperplanes)

$\mathbf{w}^T \mathbf{x} + b = 1$  and  $\mathbf{w}^T \mathbf{x} + b = -1$

$$|b - 1| - |b + 1| / \|\mathbf{w}\| = 2/\|\mathbf{w}\|$$



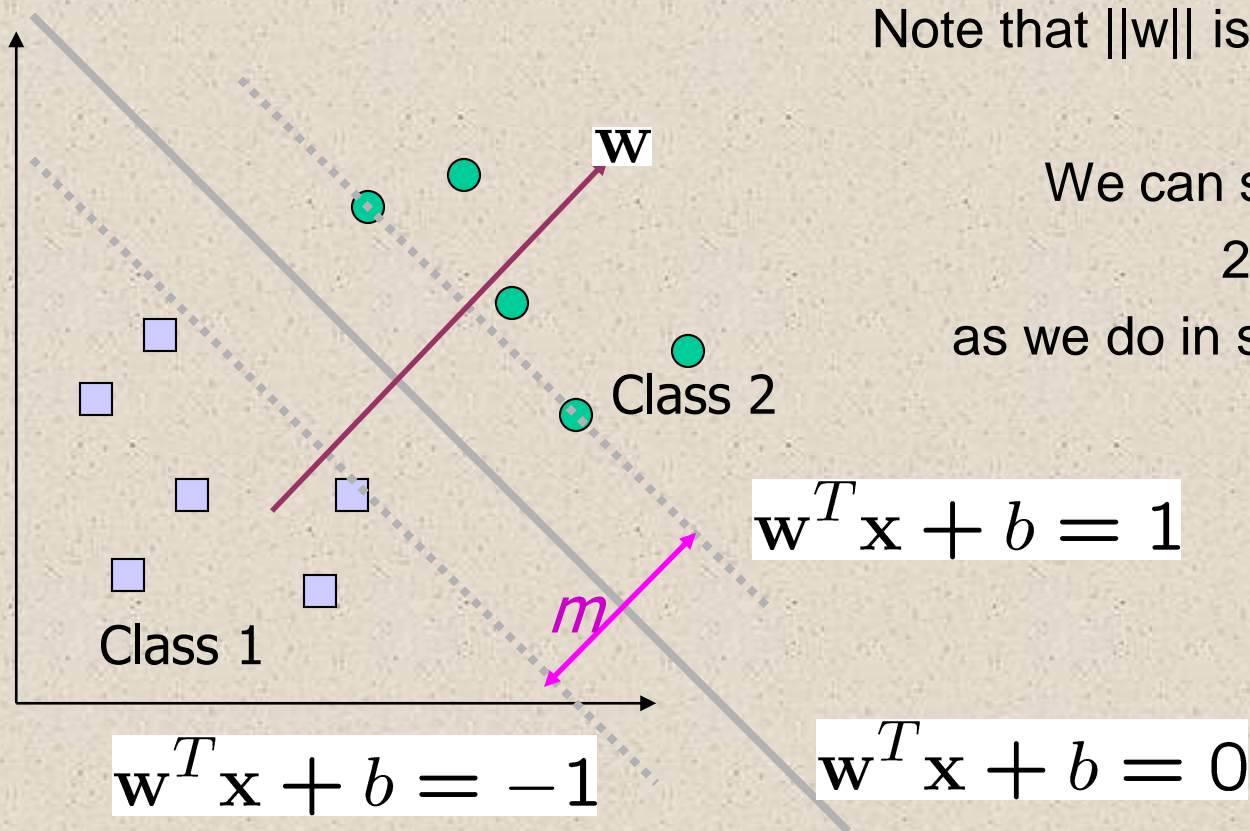
# Maximizing margin m

To **maximize** the margin  $m = 2/||w||$

→ We need to **minimize**  $||w||$  which is same as  $\min \frac{1}{2}||w||^2$

→ which makes the math easier.

Note that  $||w||$  is  $\sqrt{w_0^2 + w_1^2}$ .



We can say that margin  $m =$

$$\frac{2}{\sqrt{w_0^2 + w_1^2}}$$

as we do in scikit-learn example

# The Optimization Problem

1) Find the **min** of  $\|\mathbf{w}\|$  and  $\mathbf{b}$

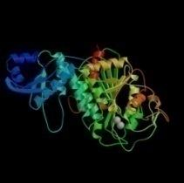
**This is the same min  $\frac{1}{2}\|\mathbf{w}\|^2$  and  $\mathbf{b}$**

2) Such that this **constraint** holds:

$y_i (\mathbf{w}^T \mathbf{x}_i + \mathbf{b}) \geq 1$  for all input vectors  $\mathbf{x}_i$

**Recall that  $y_i$  will be +1 or -1 (depending on the class)**

This is a **quadratic** optimization problem.



# A Soft Margin Hyperplane SVM

- If we minimize  $\sum_i \xi_i$ ,

$\xi_i$  can be computed by

$$\begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 1 - \xi_i & y_i = 1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq -1 + \xi_i & y_i = -1 \\ \xi_i \geq 0 & \forall i \end{cases}$$

$\xi_i$  are “**slack variables**” in optimization

– Note that  $\xi_i=0$  if there is no error for  $\mathbf{x}_i$

$\xi_i$  is an upper bound of the number of errors

- We want to **minimize**:  $\frac{1}{2} \|\mathbf{w}\|^2 + \mathbf{C} \sum_i \xi_i$ 
  - $C$  : tradeoff parameter between error and margin
- The optimization problem becomes

$$\frac{1}{2} \|\mathbf{w}\|^2 + \mathbf{C} \sum_i \xi_i$$

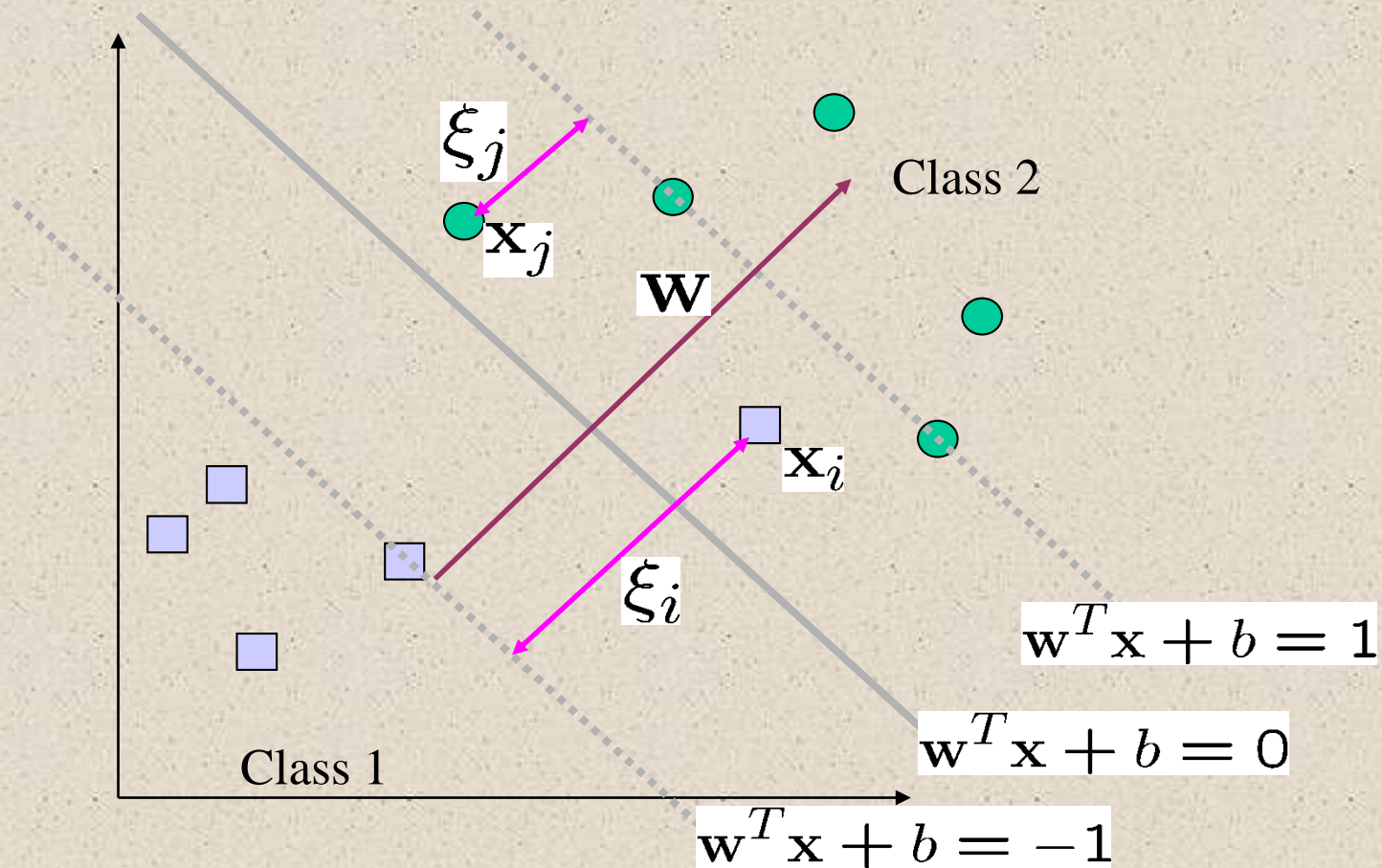
$$\text{subject to } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$



# Soft Margin Example

- Allow “error”  $\xi_i$  in classification

$\xi_i$  approximates the number of misclassified samples



# Constrained Optimization Quick Review

Suppose you want to minimize some function  $\mathbf{f}(\mathbf{x})$  subject to the constraint of  $\mathbf{g}(\mathbf{x})$ .

Then, for any  $\mathbf{x}$  to be a solution, the following must be true.

Note that the “ $\alpha$ ” is a Lagrange multiplier

$$\begin{cases} \frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x}) + \alpha g(\mathbf{x})) \Big|_{\mathbf{x}=\mathbf{x}_0} = 0 \\ g(\mathbf{x}) = 0 \end{cases}$$

# Using Lagrange in SVM Solving

Goal: Minimize  $\frac{1}{2} \|\mathbf{w}\|^2$

subject to constraint:

$$1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0, \quad i = 1, 2, \dots, n$$

The **Lagrangian**  $L$  is

$$L = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{(i=1 \text{ to } n)} \alpha_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

Them, setting the gradient of  $L$  w.r.t.  $\mathbf{w}$  and  $b$  to 0, gives:

$$\mathbf{w} + \sum_{(i=1 \text{ to } n)} \alpha_i (-y_i) \mathbf{x}_i = 0 \rightarrow$$

$$\mathbf{w} = \sum_{(i=1 \text{ to } n)} \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_{(i=1 \text{ to } n)} \alpha_i y_i = 0$$

Finally, if we substitute  $\mathbf{w}$  into  $L$  we get (after a bit of math):  $-\frac{1}{2} \sum_{(i=1 \text{ to } n)} \sum_{(j=1 \text{ to } n)} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{(i=1 \text{ to } n)} \alpha_i$

# The Dual Problem in SVM

The new objective function is called the **dual problem** in SVM math.

$$W(\alpha) = -1/2 \sum_{(i=1 \text{ to } n)} \sum_{(j=1 \text{ to } n)} \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{(i=1 \text{ to } n)} \alpha_i$$

It is a function of alpha.

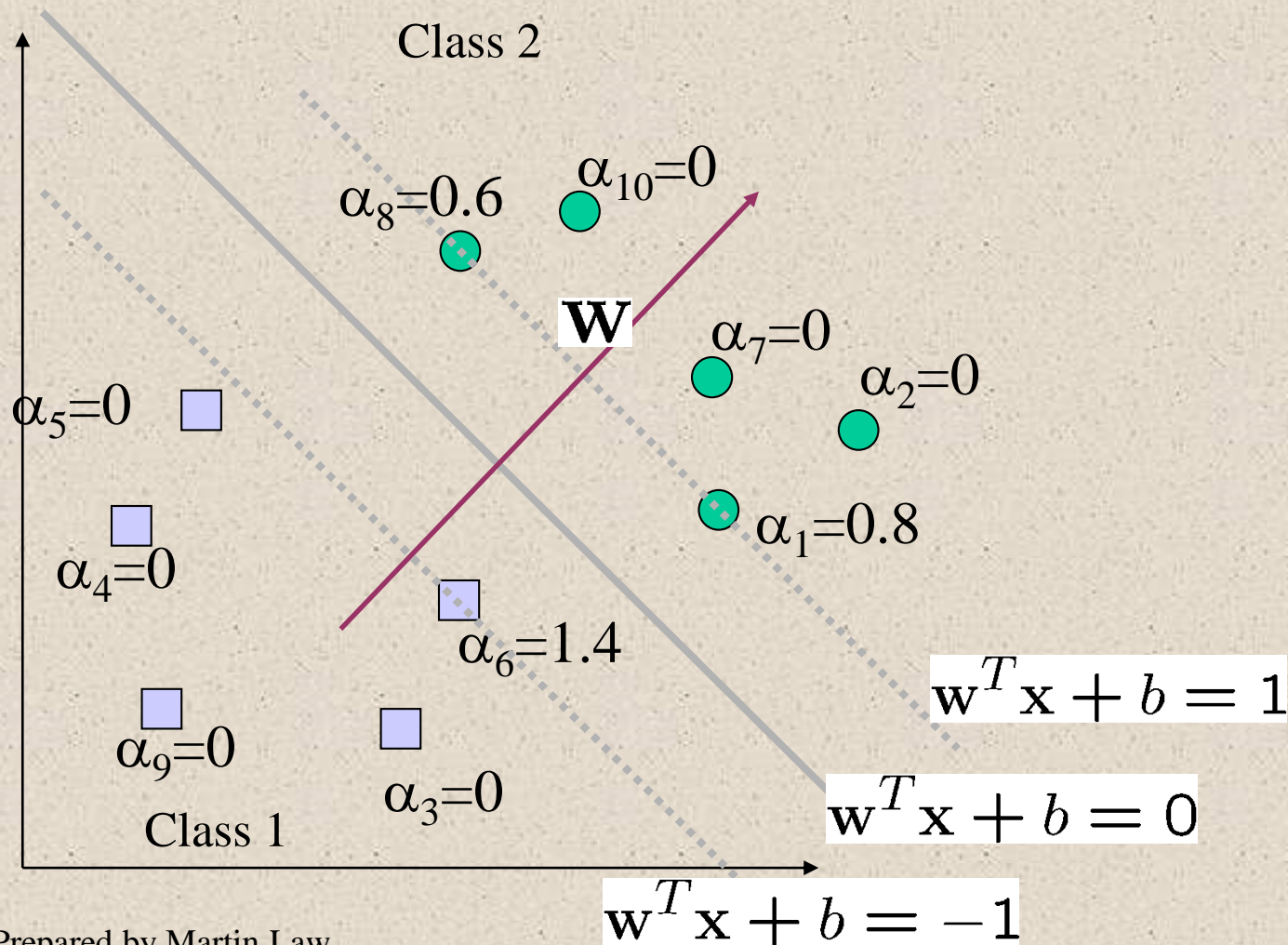
The goal is now to maximize  $W(\alpha)$  subject to alphas all being positive and  $\sum_{(i=1 \text{ to } n)} \alpha_i y_i = 0$

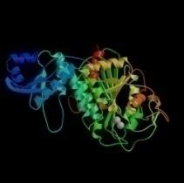
This is a **quadratic programming** problem with

$$\mathbf{w} = \sum_{(i=1 \text{ to } n)} \alpha_i y_i x_i$$

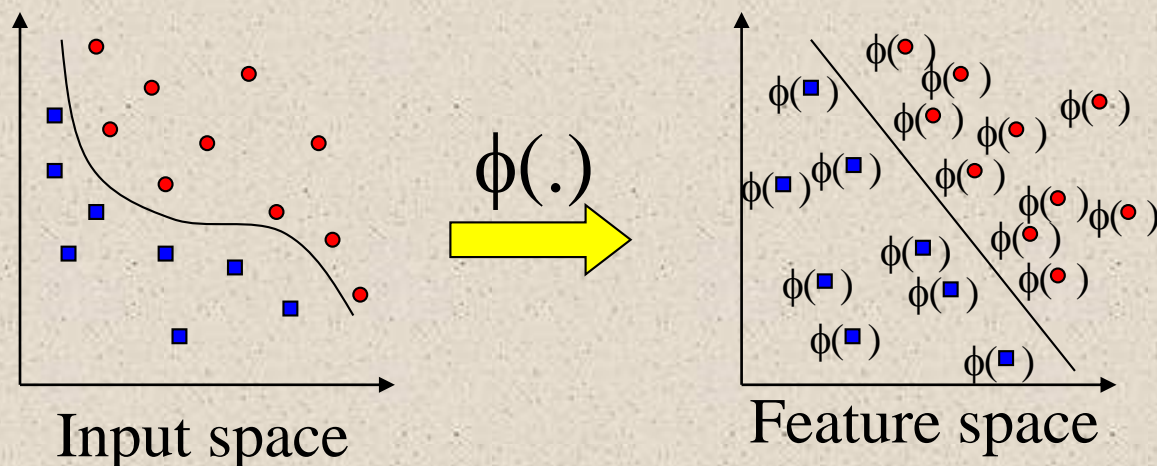


# A Visual Interpretation of $\alpha$





# SVM kernel Transformation: $\phi(\mathbf{x})$



Recall that the SVM optimization problem solved is:

$$W(\alpha) = -1/2 \sum_{(i=1 \text{ to } n)} \sum_{(j=1 \text{ to } n)} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{(i=1 \text{ to } n)} \alpha_i$$

Here, only the inner product of the  $\mathbf{x}$  values is needed.

Therefore, if a kernel transformation is used, the updated optimization problem becomes

$$W(\alpha) = -1/2 \sum_{(i=1 \text{ to } n)} \sum_{(j=1 \text{ to } n)} \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) + \sum_{(i=1 \text{ to } n)} \alpha_i$$

where kernel  $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$

We only need the inner product of the feature space – not an explicit mapping.

# Common Kernels

Polynomial with degree  $d$

$$K(x, y) = (x^T y + 1)^d$$

Radial

$$K(x, y) = \exp(-\|x - y\|^2 / 2\sigma^2)$$

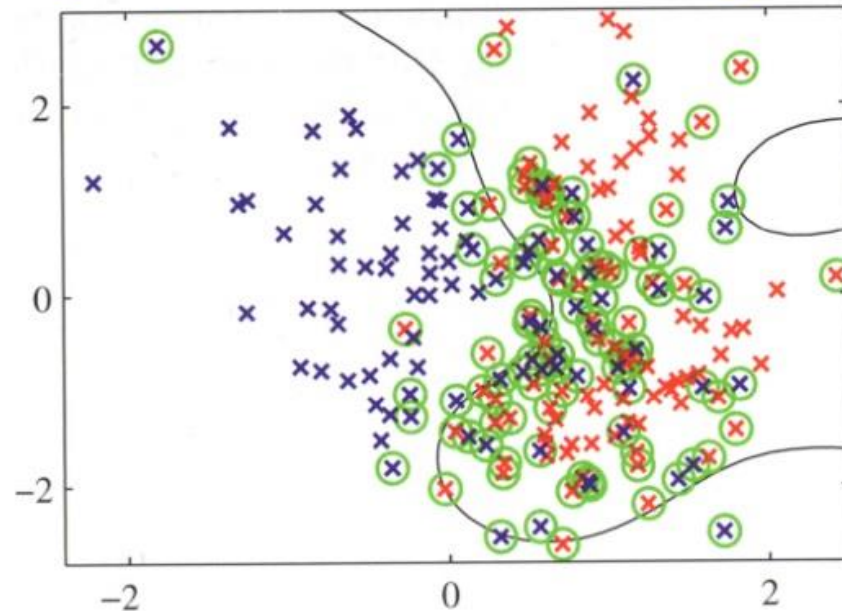
Sigmoid:

$$K(x, y) = \tanh(cx^T y + d)$$

Gaussian:

$$K(x, y) = \exp(-a\|x - y\|^2), a > 0$$

## SVM Soft Margin Decision Surface using Gaussian Kernel



[from Bishop, figure 7.4]

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) + b$$

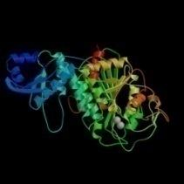
Circled points are the support vectors: training examples with non-zero  $\alpha_l$

Points plotted in original 2-D space.

Contour lines show constant  $\hat{f}(\mathbf{x})$

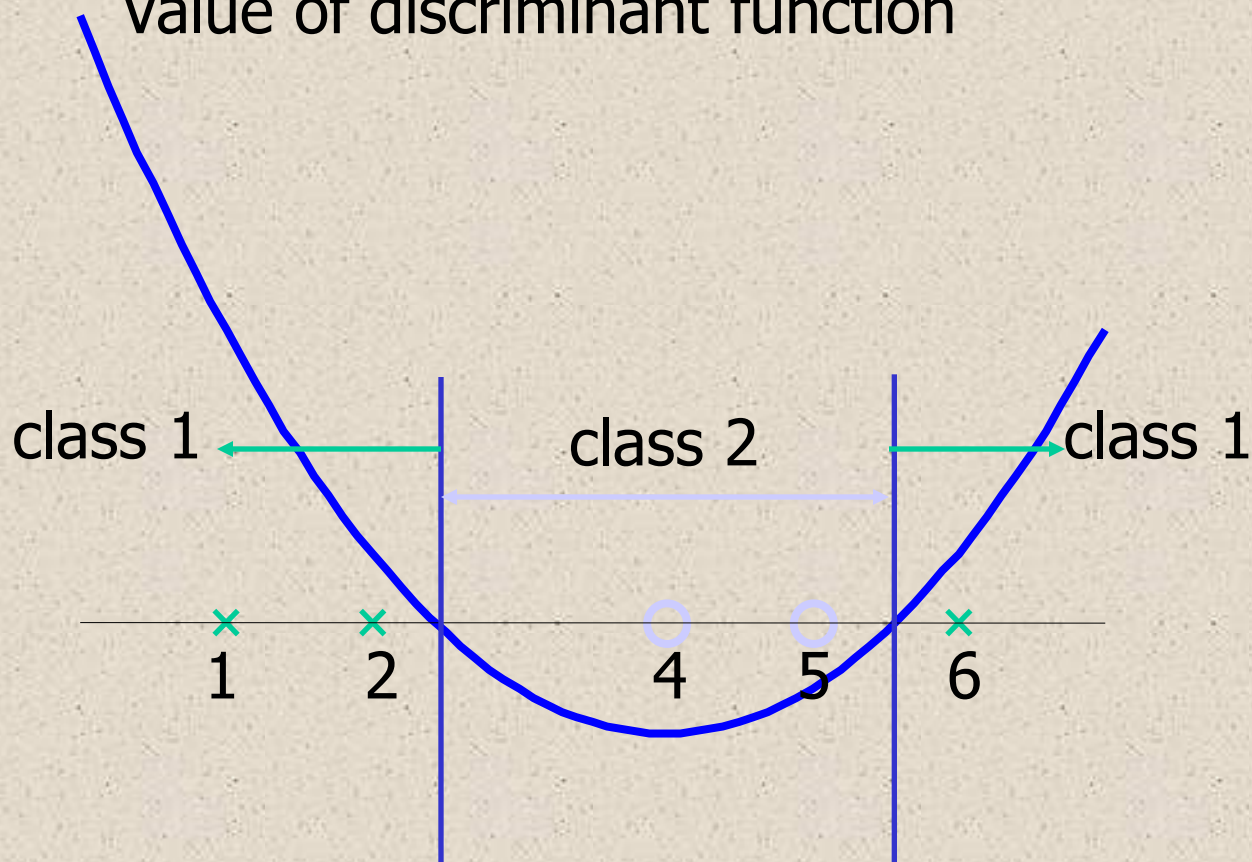
$$\hat{f}(\mathbf{x}) = b + \sum_{l=1}^M \alpha_l y_l \kappa(\mathbf{x}, \mathbf{x}_l) = b + \sum_{l=1}^M \alpha_l y_l \exp(-\|\mathbf{x} - \mathbf{x}_l\|^2 / 2\sigma^2)$$





# Example non-linear

Value of discriminant function



# Using Python 3 scikit-learn: SVM

<http://scikit-learn.org/stable/modules/svm.html>

## **SVM:**

- 1) supervised learning method
- 2) for classification, regression, outlier detection

## **Advantages:**

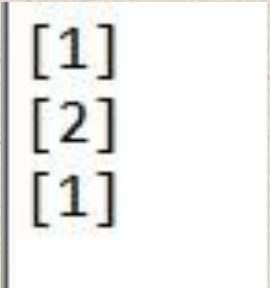
- 1) Effective in high D space
- 2) Effective when dimension  $>$  # samples
- 3) Uses “subset” of training data – just the “support vectors”
- 4) Allows or the use of kernel functions

## **Disadvantages:**

- 1) If features  $\gg$  samples  $\rightarrow$  poor performance
- 2) No prob estimates – these must be “done by hand” using at least 5-fold cross validation

# SVM Classification with scikit-learn: Example 1

```
from sklearn.svm import SVC
import numpy as np
#Train
#SVC takes in two arrays [n_samples, n_features]
#This holds the training samples
X = np.array([[-1, -1], [-2, -1], [-3, -4], [3,3], [1, 1], [2, 1]])
#The classes of the training data
y = np.array([1, 1, 1, 2, 2, 2])
#TRAIN
#Define the model with SVC
# Fit SVM with training data
clf = SVC(C=1, kernel="rbf", verbose=False)
clf.fit(X, y)
#PREDICT - perform classification
print(clf.predict([[-.8, -1]]))
print(clf.predict([[.8, 1]]))
print(clf.predict([[0, 0]]))
```



[1]  
[2]  
[1]

# More scikit-learn Examples

```
from sklearn.svm import SVC
from sklearn import datasets
import numpy as np
import matplotlib.pyplot as plt

iris = datasets.load_iris()
#Place first two columns into X
X=np.array(iris.data[:, :2])
#Place the targets/classes into y
y = np.array(iris.target)
#Set up the SVM - linear kernel
clf = SVC(C=1, kernel="linear")
clf.fit(X, y)
#Margin is  $2 / ||w||$ 
margin = 2 /
np.sqrt(np.sum(clf.coef_ ** 2))

#Predict a new data points
print(clf.predict([[6, 3]]))
```

```
# get the separating hyperplane
#The weights vector w
w = clf.coef_[0]
#The slope of the SVM sep line
a = -w[0] / w[1]
#Create a variable xx that are values
between 4 and 8
xx = np.linspace(4, 8)
## This is the y values for the main sep
line
yy = a * xx - (clf.intercept_[0]) / w[1]
# adding or subtracting  $\frac{1}{2}$  of the margin
yy_down = yy + .5*margin
yy_up = yy - .5*margin
```



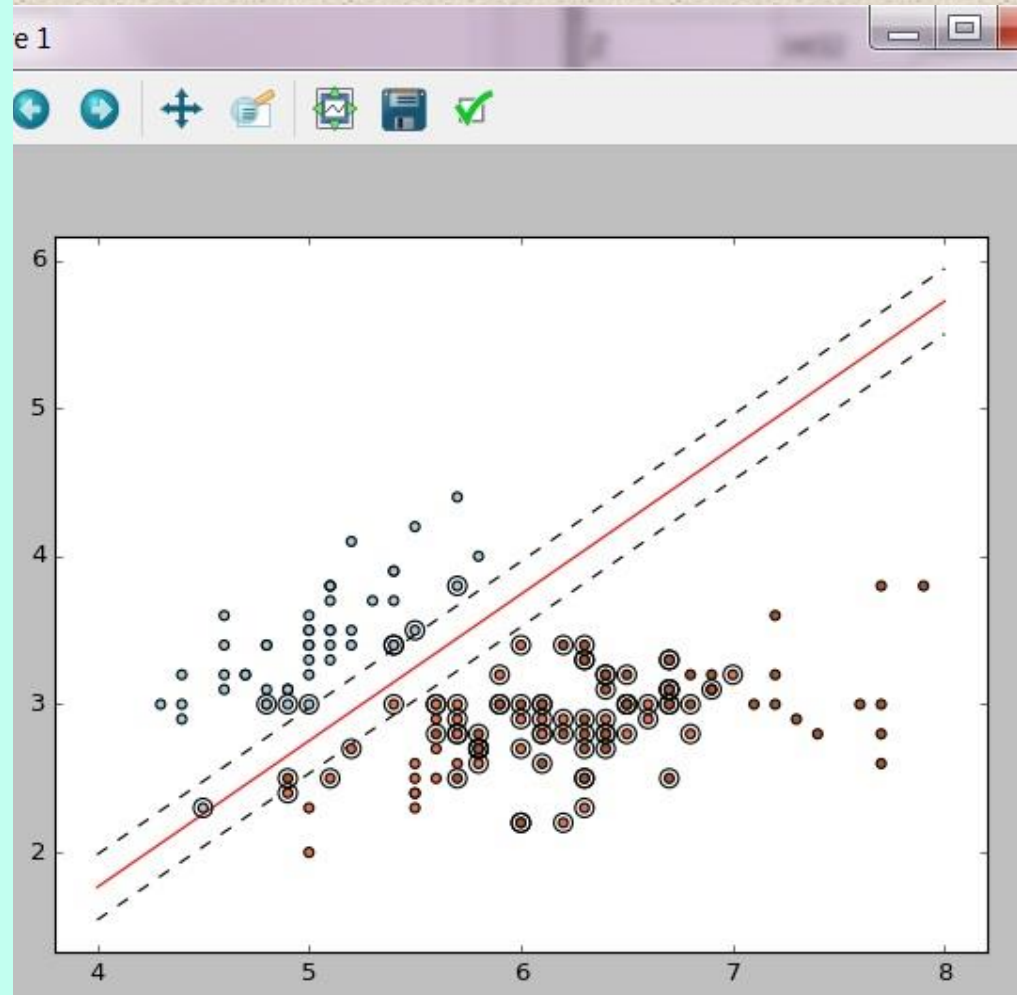
# continued...

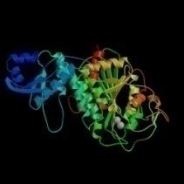
```
# plot the line, the points, and the
nearest vectors to the plane
plt.clf()
plt.plot(xx, yy, 'r-')
plt.plot(xx, yy_down, 'k--')
plt.plot(xx, yy_up, 'k--')

plt.scatter(clf.support_vectors_[0],
            clf.support_vectors_[1], s=80,
            facecolors='none', zorder=5)

#cmap is the color map
plt.scatter(X[:, 0], X[:, 1], c=y,
            zorder=5, cmap=plt.cm.Paired)

plt.axis('tight')
```





# The SVM Math - FYI

The **perpendicular distance** from a point to hyperplane  $\mathbf{y}(\mathbf{x})=0$ , where  $\mathbf{y}(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + \mathbf{b}$ , is defined as  $t_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + \mathbf{b}) / \|\mathbf{w}\|$ , since only correctly predicted points are of interest.

The **margin** is given by the perpendicular distance to the closest point  $\mathbf{x}_n$ , and  $\mathbf{w}$  and  $\mathbf{b}$  can be **optimized to maximize** the distance.

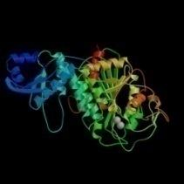
The max margin solution can be found by solving:

$$\arg \max_{\mathbf{w}, \mathbf{b}} \{ 1/\|\mathbf{w}\| \min_n [t_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + \mathbf{b})] \}$$

Note: rescaling  $\mathbf{w}$  and  $\mathbf{b}$  by the same constant does not affect the distance between  $\mathbf{x}_n$  and the decision boundary.

Therefore, set  $t_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + \mathbf{b}) = 1$  for the closest point and  $t_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + \mathbf{b}) \geq 1$  for all other points.

- Extra slides...



# The SVM Math - FYI

$t_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) \geq 1$  is the canonical representation of the hyperplane.

To maximize the margin, **minimize**  $\mathbf{w}^T \mathbf{w}$

Therefore, the **Quadratic Programming** problem can be stated as:

$$\arg \min \mathbf{w}, b : \mathbf{w}^T \mathbf{w}$$

$$\text{such that } t_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) \geq 1$$

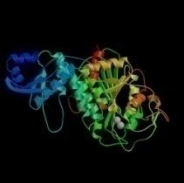
Lagrangian to solve:

$$L(\mathbf{w}, b, \mathbf{a}) = \mathbf{w}^T \mathbf{w} - \sum_{n=1 \text{ to } N} a_n \{t_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) - 1\}$$

Take derivatives of  $L(\mathbf{w}, b, \mathbf{a})$  w.r.t.  $\mathbf{w}$  and  $b$  to get:

$$\mathbf{w} = \sum_{n=1 \text{ to } N} a_n t_n \boldsymbol{\phi}(\mathbf{x}_n) \quad \text{and} \quad \sum_{n=1 \text{ to } N} a_n t_n = 0$$





# The SVM: The kernel

Substituting:

$$\mathbf{w} = \sum_{n=1 \text{ to } N} a_n \mathbf{t}_n \phi(\mathbf{x}_n) \quad \text{and} \quad \sum_{n=1 \text{ to } N} a_n \mathbf{t}_n = 0$$

into

$$L(\mathbf{w}, \mathbf{b}, \mathbf{a}) = \mathbf{w}^T \mathbf{w} - \sum_{n=1 \text{ to } N} a_n \{ \mathbf{t}_n (\mathbf{w}^T \phi(\mathbf{x}_n) + \mathbf{b}) - 1 \}$$

**The Dual Representation:**

$$L'(\mathbf{a}) = \sum_{n=1 \text{ to } N} a_n - \frac{1}{2} \sum_{n=1 \text{ to } N} \sum_{m=1 \text{ to } N} a_n a_m \mathbf{t}_n \mathbf{t}_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

**Kernel:**

$$K(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

**NOTE:** This portion of the math allows for the use of a kernel to represent the inner product of data points.

# Kernels

Because the solution to the SVM optimization problem (the Dual Form) involves only **inner products**, the data values can be “transformed” to a different dimension using a kernel.

- **The Dual Representation:**

$$L'(\mathbf{a}) = \sum_{n=1 \text{ to } N} \mathbf{a}_n - \frac{1}{2} \sum_{n=1 \text{ to } N} \sum_{m=1 \text{ to } N} \mathbf{a}_n \mathbf{a}_m \mathbf{t}_n \mathbf{t}_m \boldsymbol{\varphi}(\mathbf{x}_n)^T \boldsymbol{\varphi}(\mathbf{x}_m)$$

**Kernel:**  $K(\mathbf{x}_n, \mathbf{x}_m) = \boldsymbol{\varphi}(\mathbf{x}_n)^T \boldsymbol{\varphi}(\mathbf{x}_m)$

This can be thought of as a virtual transformation as the points themselves do not have to be projected.

This is the Gaussian Kernel for example.

$$K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2), \gamma > 0,$$