# Support Vector Machines

Auburn University

BUAL 5610/6610

# Content

- Maximal Margin Classifier

- Support Vector Classifier

- Support Vector Machine
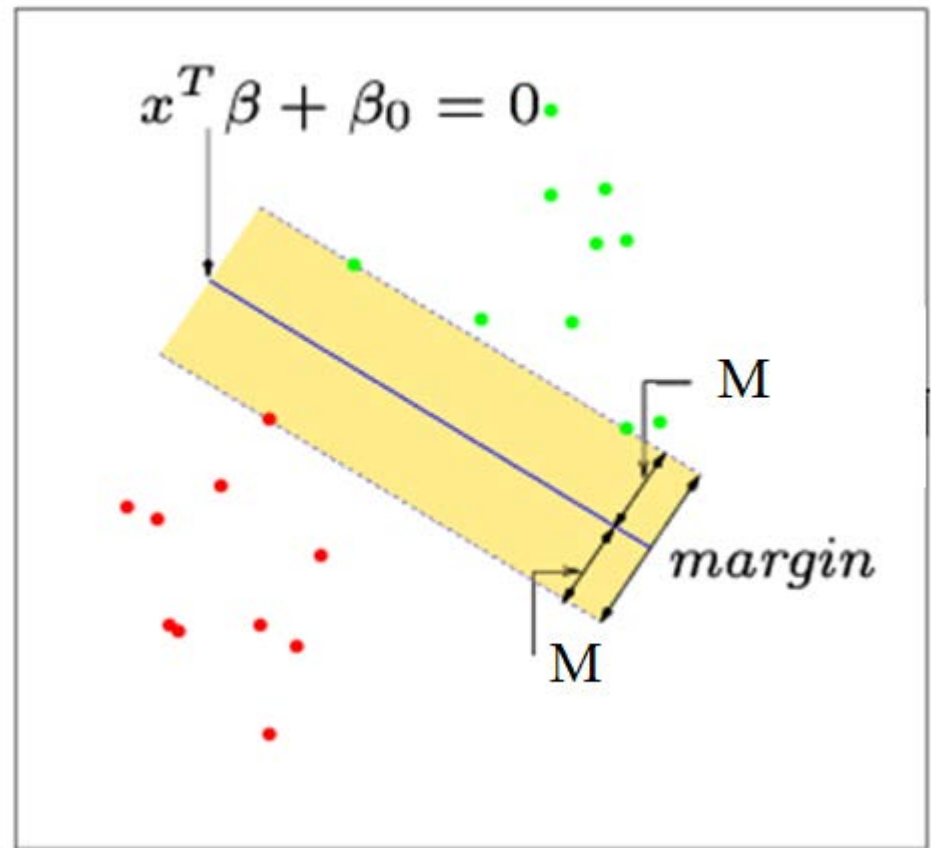
# Maximal Margin Classifier

# Separable Hyperplanes

- Imagine a situation where you have a two-class classification problem with two predictors $X_1$ and $X_2$.

- Suppose that the two classes are "linearly separable" i.e. one can draw a straight line in which all points on one side belong to the first class and points on the other side to the second class.

- Then a natural approach is to find the straight line that gives the biggest separation between the classes i.e. the points are as far from the line as possible

- This is the basic idea of a support vector classifier.

# Its Easiest To See With A Picture

- M is the minimum perpendicular distance between each point and the separating line.

- We find the line which maximizes M.

- This line is called the "optimal separating hyperplane"

- The classification of a point depends on which side of the line it falls on.
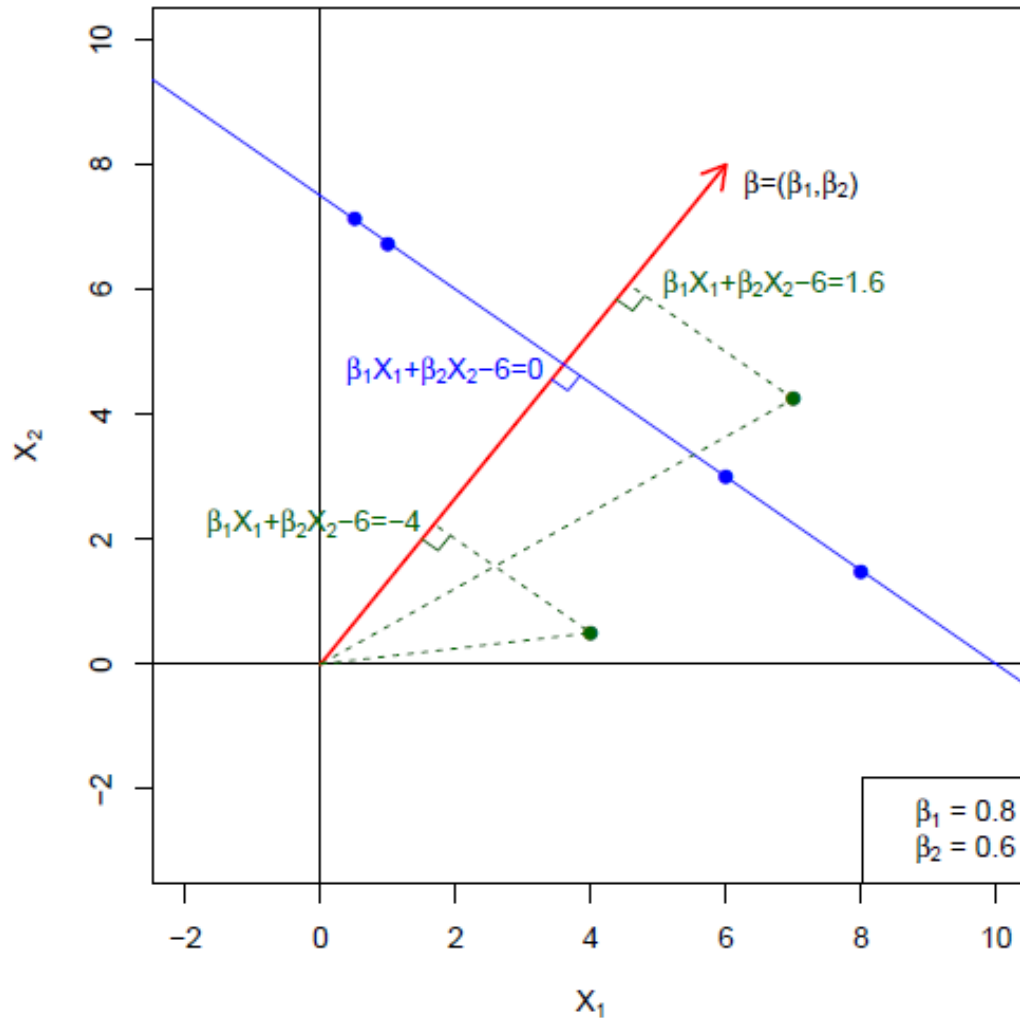
$$x^T \beta + \beta_0 = 0$$

M

margin

M

# More Than Two Predictors

- This idea works just as well with more than two predictors.

- For example, with three predictors you want to find the plane that produces the largest separation between the classes.

- With more than three dimensions it becomes hard to visualize a plane but it still exists. In general they are caller hyper-planes.

# Hyperplane in 2 Dimensions
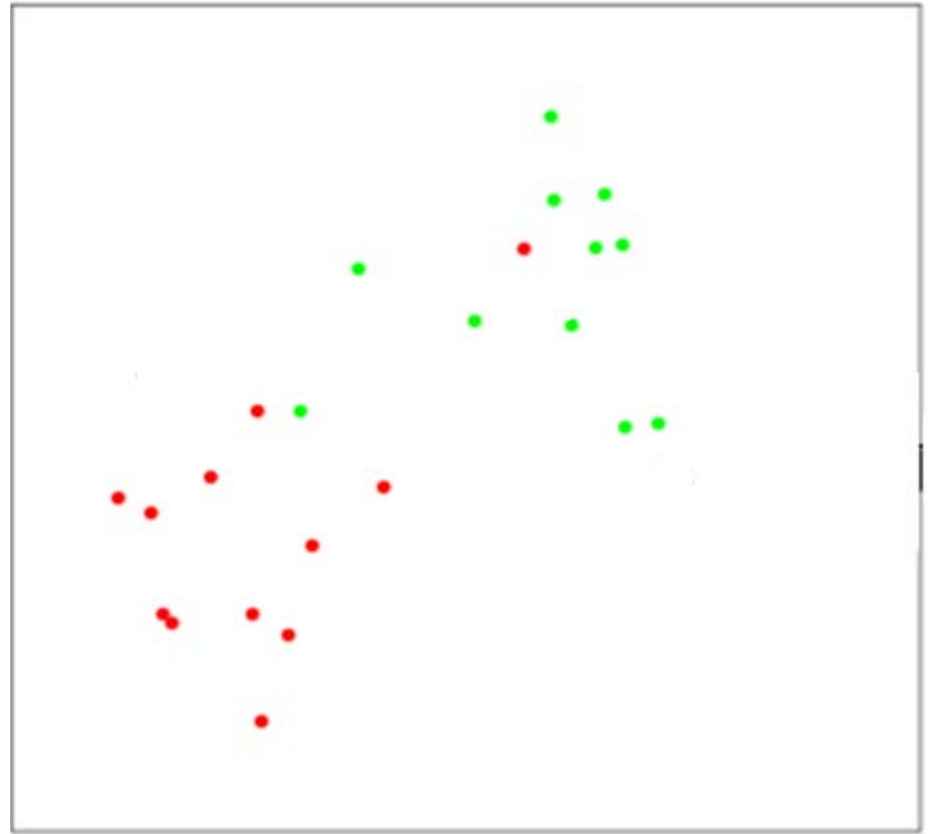
# Support Vector Classifier

# Non-Separating Classes

- In practice it is not usually possible to find a hyper-plane that perfectly separates two classes.

- In other words, for any straight line or plane that I draw there will always be at least some points on the wrong side of the line.

- In this situation we try to find the plane that gives the best separation between the points that are correctly classified <u>subject to the points on the wrong side</u> of the line not being off by too much.

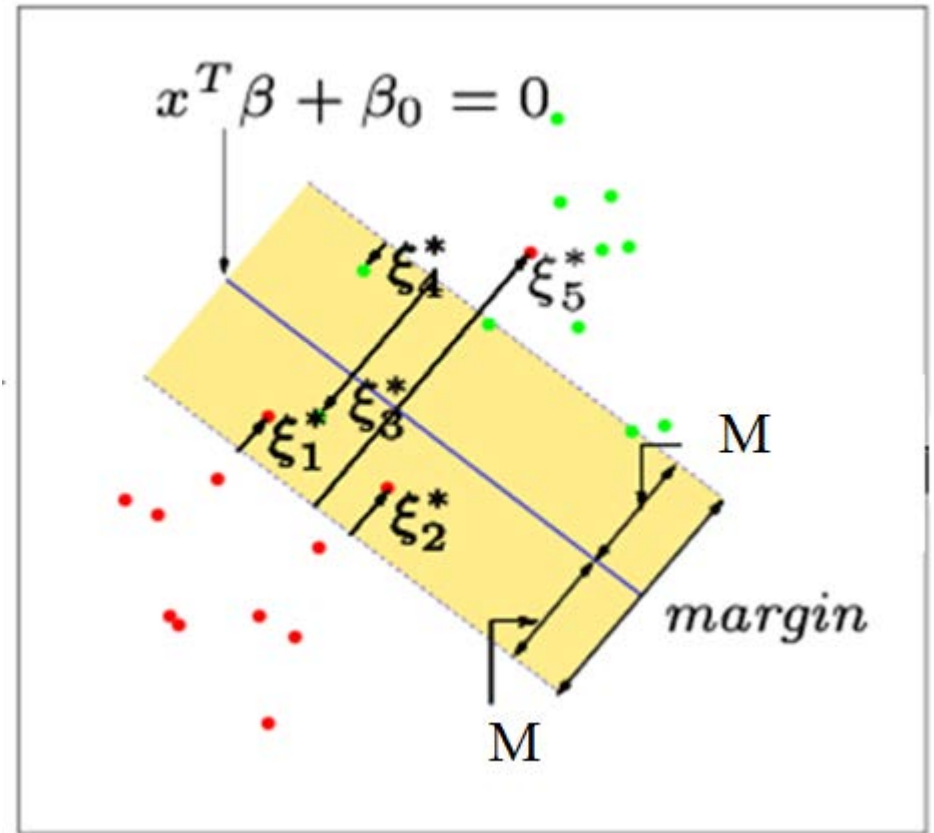- The support vector classier maximizes a <span style="color:red">soft margin</span>.

# Non-Separating Example

- Let $\xi^*_i$ represent the amount that the ith point is on the wrong side of the margin (the dashed line).

- Then we want to maximize M subject to

$$\frac{1}{M}\sum_{i=1}^{n}\xi^*_i \leq \text{Constant}$$

- The constant is a *tuning parameter* that we choose.

# Non-Separating Example

- Let $\xi^*_i$ represent the amount that the ith point is on the wrong side of the margin (the dashed line).

- Then we want to maximize M subject to

$$\frac{1}{M}\sum_{i=1}^{n}\xi^*_i \leq \text{Constant}$$

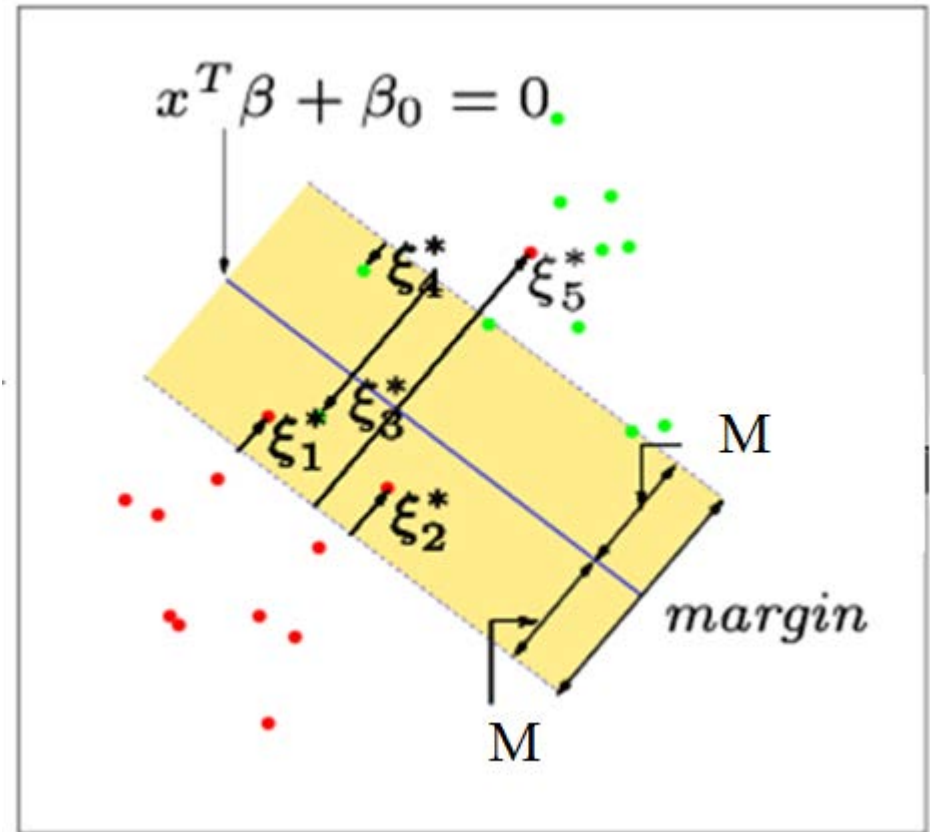- The constant is a *tuning parameter* that we choose.
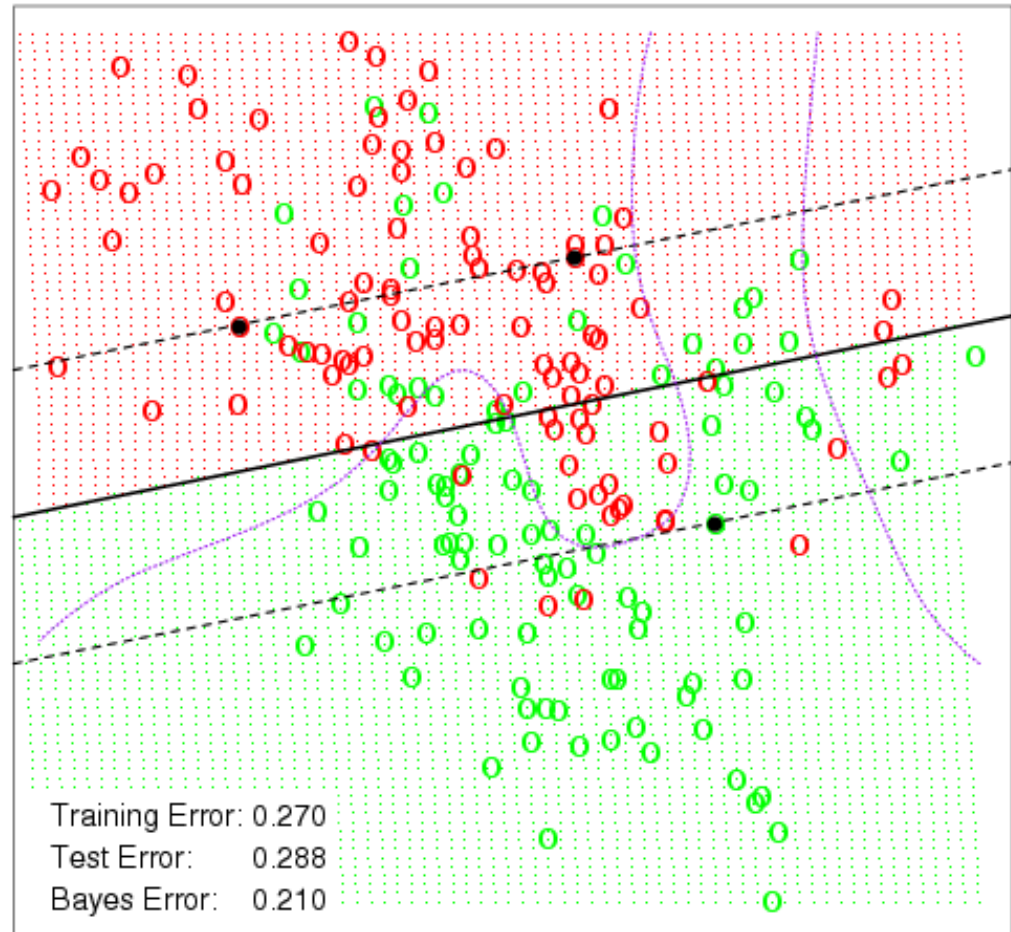


$$x^T\beta + \beta_0 = 0$$

# Non-Separating Example

- Observations that lie directly on the margin, or on the wrong side of the margin for their class, are known as *support vectors*. These observations do affect the support vector classifier.

- There are seven support vectors in this plot.



$$x^T \beta + \beta_0 = 0$$

$\xi_4^*$  $\xi_5^*$
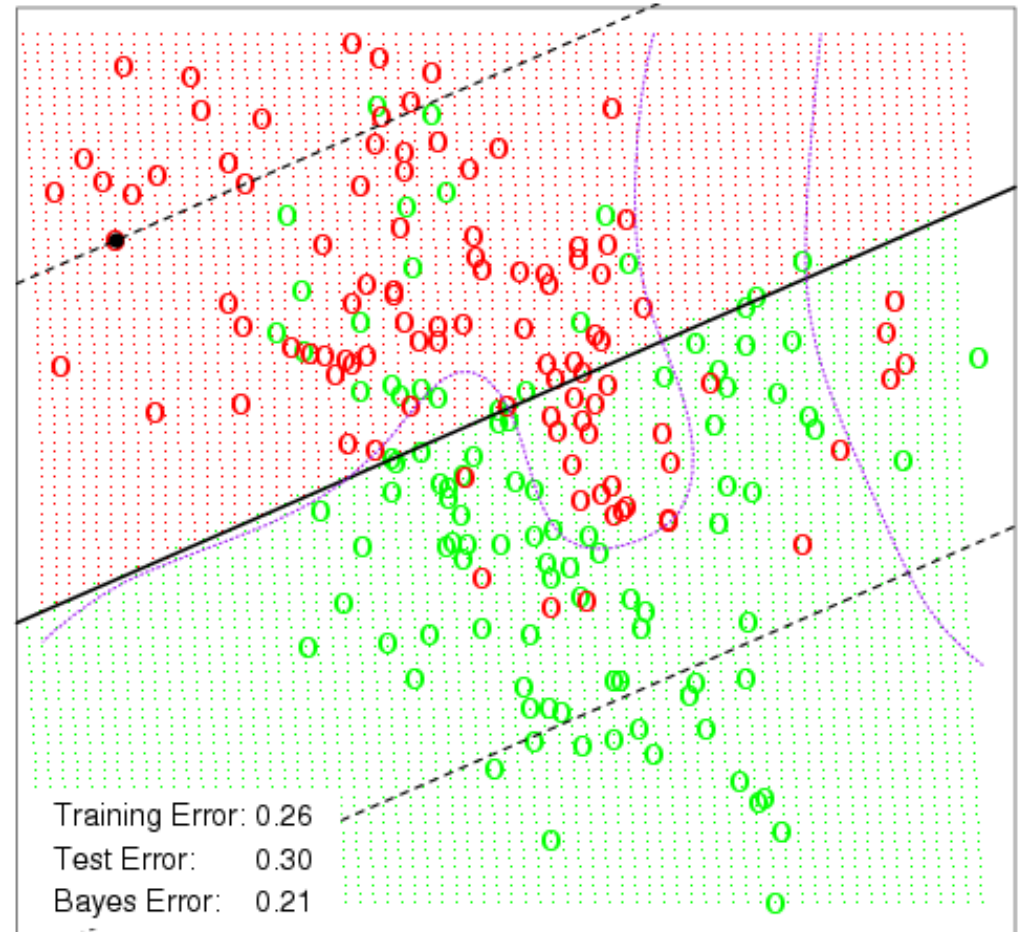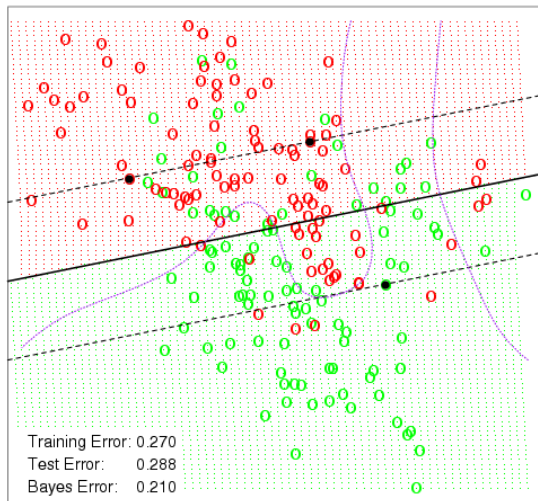
$\xi_1^*$  $\xi_3^*$

$\xi_2^*$

M

margin

M

# A Simulation Example With A Small Constant

- This is a plot generated from simulation.

- The distance between the dashed lines represents the margin or 2M.

- The purple lines represent the Bayes decision boundaries



Training Error: 0.270
Test Error:     0.288
Bayes Error:    0.210

# The Same Example With A Larger Constant

- Using a larger constant allows for a greater margin and creates a slightly different classifier.

- Notice, however, that the decision boundary must always be linear.



Training Error: 0.270
Test Error:    0.288
Bayes Error:   0.210



Training Error: 0.26
Test Error:    0.30
Bayes Error:   0.21

# Support Vector Machine

# Non-Linear Classifier

- The support vector classifier is fairly easy to think about. However, because it only allows for a linear decision boundary it may not be all that powerful.

- Recall that in chapter 3 we extended linear regression to non-linear regression using a basis function i.e.

$$Y_i = \beta_0 + \beta_1 b_1(X_i) + \beta_2 b_2(X_i) + \cdots + \beta_p b_p(X_i) + \varepsilon_i$$

# A Basis Approach

- Conceptually, we can take a similar approach with the support vector classifier.

- The support vector classifier finds the optimal hyper-plane in the space spanned by $X_1, X_2,\ldots, X_p$.

- Instead we can create transformations (or a basis) $b_1(x), b_2(x), \ldots, b_M(x)$ and find the optimal hyper-plane in the space spanned by $b_1(\mathbf{X}), b_2(\mathbf{X}), \ldots, b_M(\mathbf{X})$.

- This approach produces a linear plane in the transformed space but a non-linear decision boundary in the original space.

- This is called the Support Vector <span style="color:red">Machine</span> Classifier.
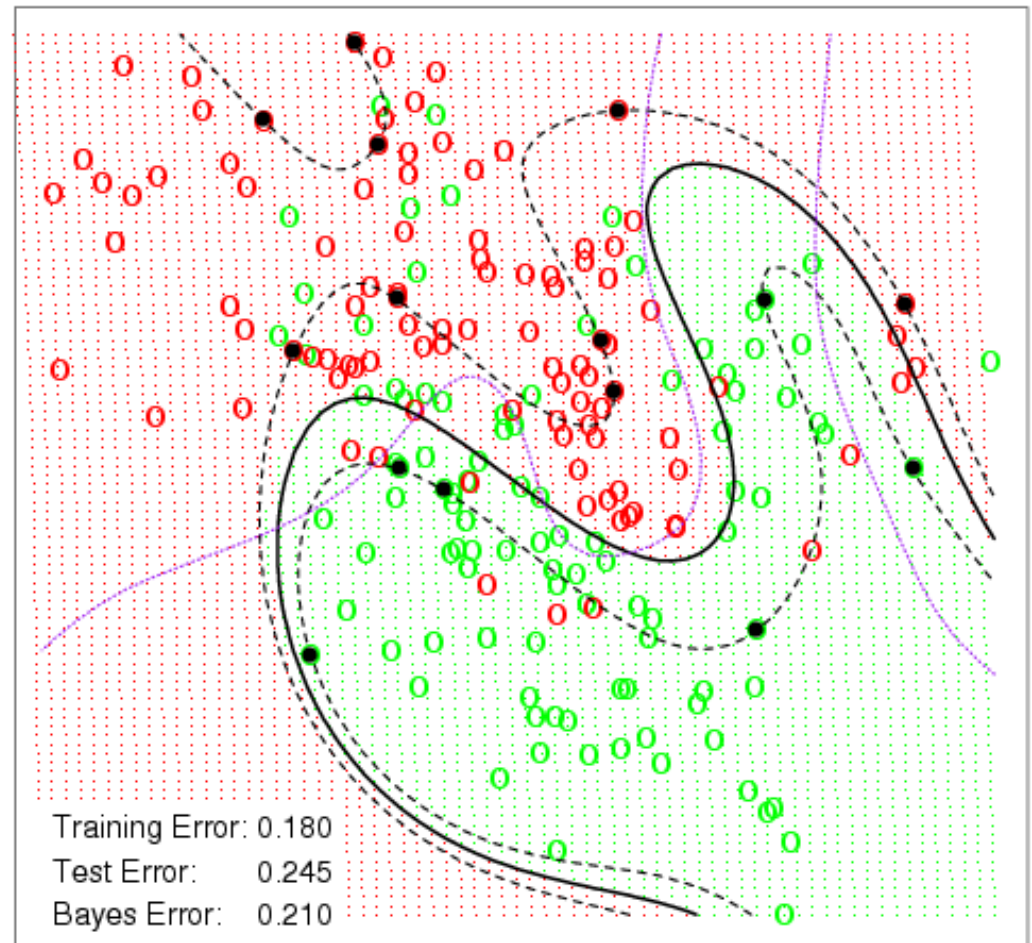
# In Reality

- While conceptually the basis approach is how the support vector machine works, there is some complicated math (which I will spare you) which means that we don't actually choose $b_1(x)$, $b_2(x)$, …, $b_M(x)$.

- Instead we choose something called a Kernel function which takes the place of the basis.

- Common kernel functions include
    - Linear
    - Polynomial
    - Radial Basis
    - Sigmoid

# Polynomial Kernel On Sim Data

- Using a polynomial kernel we now allow SVM to produce a non-linear decision boundary.

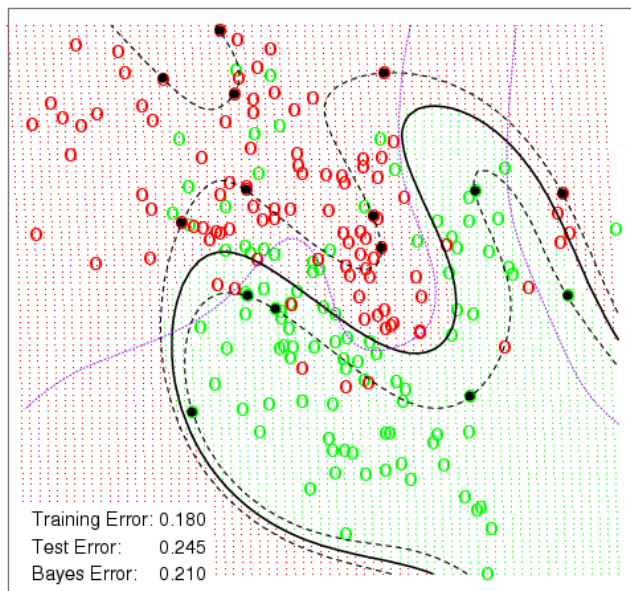- Notice that the test error rate is a lot lower.

SVM - Degree-4 Polynomial in Feature Space
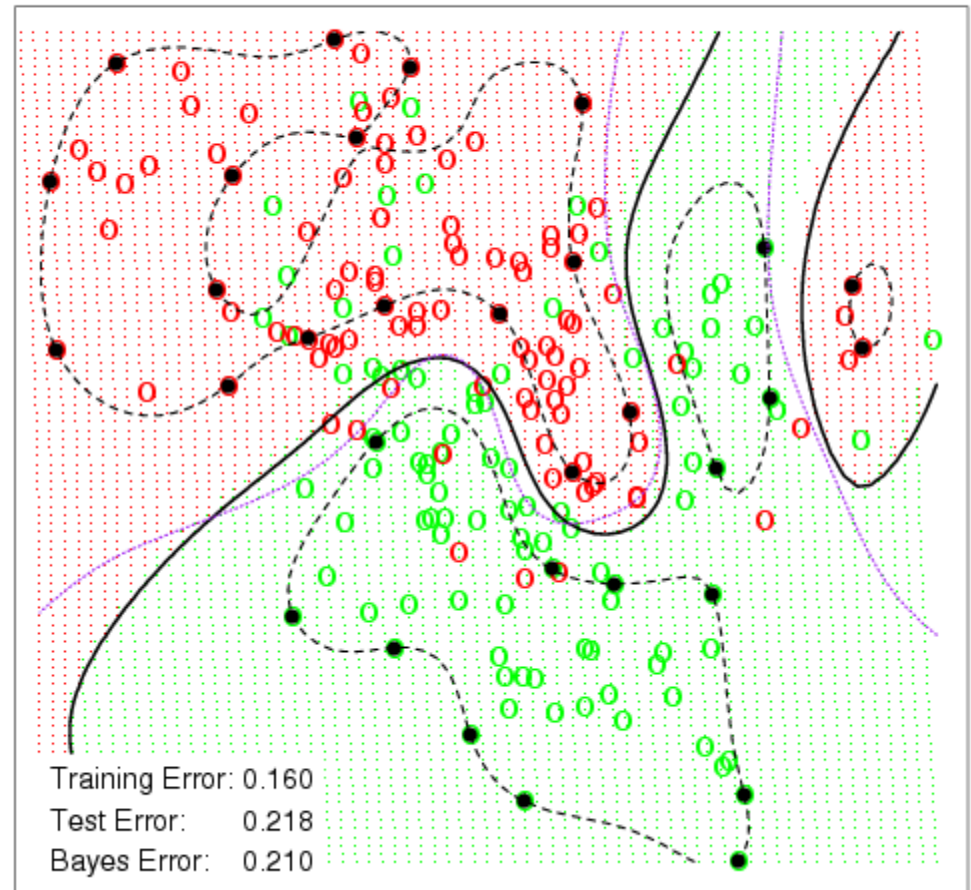


Training Error: 0.180
Test Error:     0.245
Bayes Error:   0.210

# Radial Basis Kernel

- Using a Radial Basis Kernel you get an even lower error rate.

SVM - Radial Kernel in Feature Space



```
Training Error: 0.160
Test Error:     0.218
Bayes Error:    0.210
```

SVM - Degree-4 Polynomial in Feature Space



```
Training Error: 0.180
Test Error:     0.245
Bayes Error:    0.210
```

# Python Code for SVM

```
from sklearn.svm import SVC, LinearSVC

svc = SVC(C= 1.0, kernel='linear')
svc.fit(X_train, y_train)

y_pred1 = svc.predict(X_test)
```

# Review Questions

1. Multiple Choice: Which of the following is the most different?

A) A Ph.D. in Mathematical Biology

B) A Ph.D. in Theoretical Mathematics

C) A Ph.D. in Statistics

D) A large pepperoni pizza

# Review Questions

2. We are given $n = 7$ observations in $p = 2$ dimensions. For each observation, there is an associated class label.

1) Sketch the observations.

2) Sketch the optimal separating hyperplane, and provide the equation for this hyperplane.

3) Describe the classification rule for the maximal margin classifier.

4) On your sketch, indicate the margin for the maximal margin hyperplane.

5) Indicate the support vectors for the maximal margin classifier.

| Obs. | X1 | X2 | Y |
|------|----|----|------|
| 1 | 3 | 4 | Red |
| 2 | 2 | 2 | Red |
| 3 | 4 | 4 | Red |
| 4 | 1 | 4 | Red |
| 5 | 2 | 1 | Blue |
| 6 | 4 | 3 | Blue |
| 7 | 4 | 1 | Blue |

# A few extensions

- Multi-class SVC takes care of y variables with more than two categories

- The method of Support Vector Classification can be extended to solve regression problems. This method is called Support Vector Regression.

Find out more at: http://scikit-learn.org/stable/modules/svm.html

# Disadvantage

- Besides the advantages of SVMs - from a practical point of view - they have some drawbacks. An important practical question that is not entirely solved, is the selection of the <u>parameters.</u>

- A most serious problem with SVMs is the <u>high algorithmic complexity</u> and <u>extensive memory requirements</u> of the required quadratic programming in large-scale tasks.

# References

- *"An Introduction to Statistical Learning."* James, Witten, Hastie, and Tibshirani.

- *"Applied Modern Statistical Learning Methods."* Abbass Al Sharif