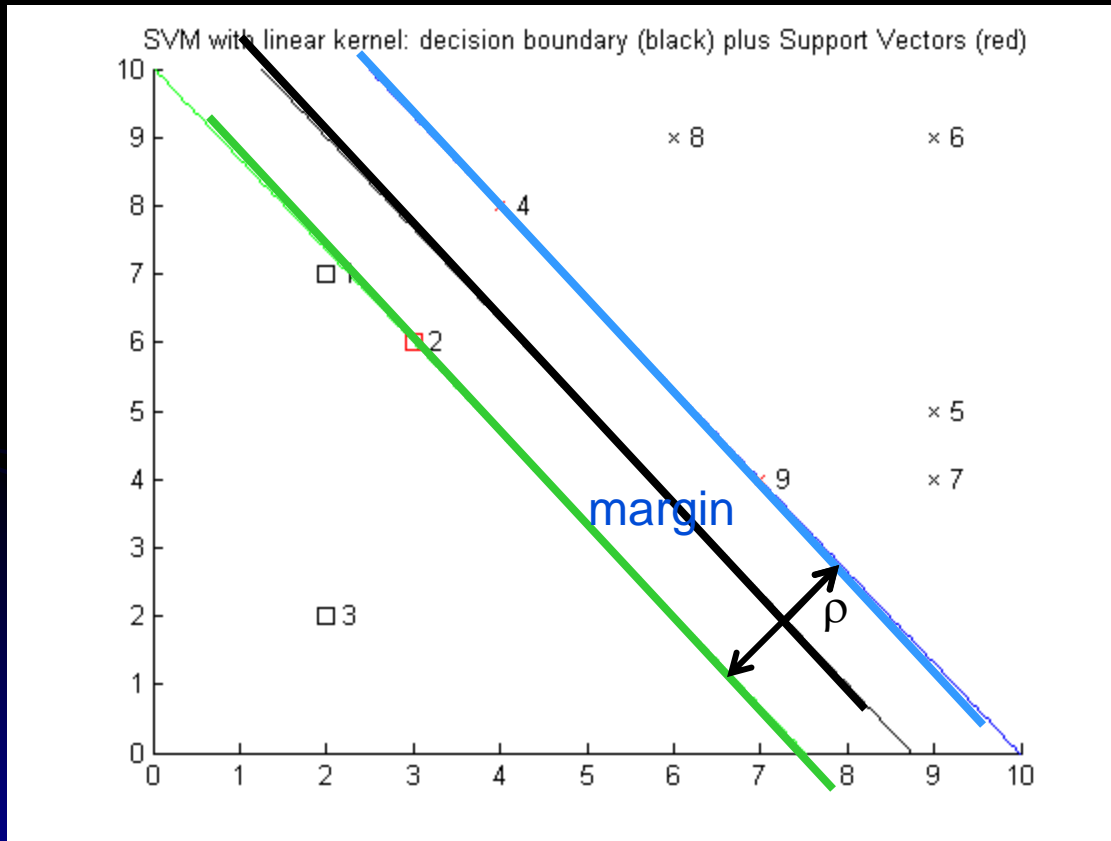# CSSE463: Image Recognition        Day 15

- Announcements:
    - Lab 5 posted, due Weds, Jan 14.
    - Sunset detector posted, due Weds, Jan 21.
- Today:
    - Project intro
    - Wrap up SVM and do demo
- Friday: Lab 5 (SVM)
- Next week:
    - Monday: Bayesian classifiers
    - Tuesday: lightning talks
    - Thursday: Mid-term exam
    - Friday: sunset detector lab
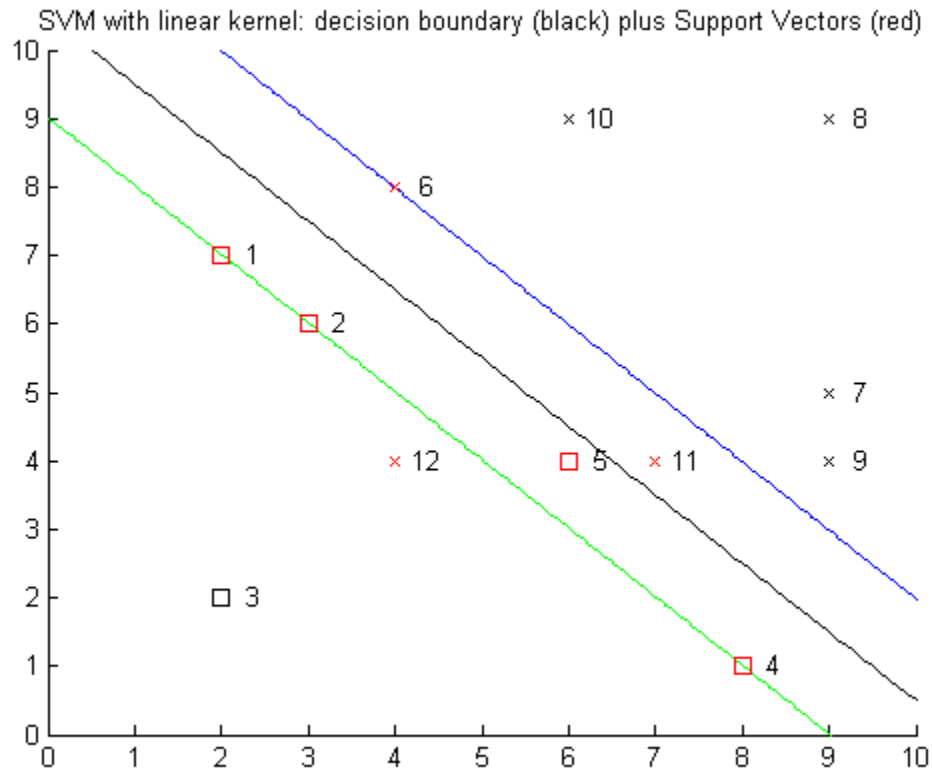
# Review: SVMs: "Best" decision boundary



SVM with linear kernel: decision boundary (black) plus Support Vectors (red)

- The "best" hyperplane is the one that *maximizes the margin, $\rho$,* between the classes. Equivalent to:

$$\min \; \phi(w) = \frac{1}{2} w^T w$$

$$d_i(w^T x_i + b) \geq 1 \; for \; i = 1, 2, \ldots N$$

- Solve using quadratic programming

# Non-separable data



SVM with linear kernel: decision boundary (black) plus Support Vectors (red)

- Allow data points to be misclassifed
- But assign a cost to each misclassified point.
- The cost is bounded by the parameter C (which you can set)
- You can set different bounds for each class. Why?
  - Can weigh false positives and false negatives differently

# Can we do better?

- Cover's Theorem from information theory says that we can map nonseparable data in the input space to a feature space where the data is separable, with high probability, if:
  - The mapping is nonlinear
  - The feature space has a higher dimension
- The mapping is called a *kernel function*.
- Lots of math would follow here

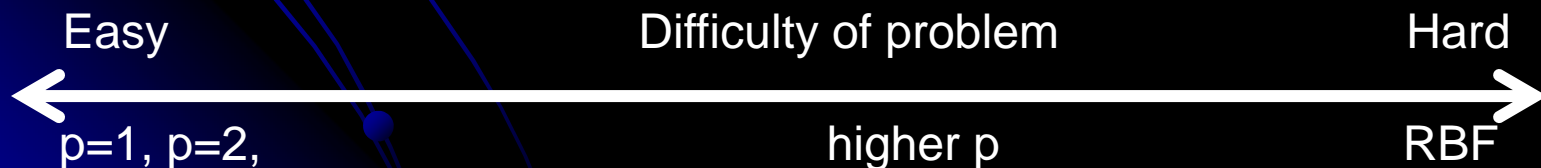# Most common kernel functions

- Polynomial
- Gaussian Radial-basis function (RBF)
- Two-layer perceptron

$$K(x, x_i) = (x^T x_i + 1)^p$$

$$K(x, x_i) = \exp\left(-\frac{1}{2\sigma^2}\|x - x_i\|^2\right)$$

$$K(x, x_i) = \tanh\left(\beta_0 x^T x_i + \beta_1\right)$$

- You choose p, $\sigma$, or $\beta_i$
- My experience with real data: **use Gaussian RBF!**

Easy           Difficulty of problem           Hard

p=1, p=2,           higher p           RBF

# Demo

- Software courtesy of http://ida.first.fraunhofer.de/~anton/software.html (GNU public license)
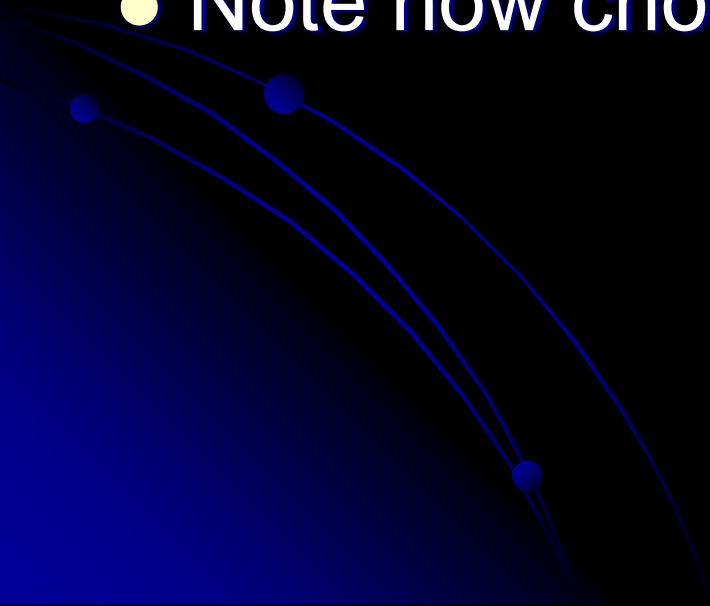
- Preview of Lab 4 (posted):
  - Download the Matlab functions that train and apply the SVM.
  - The demo script contains examples of how to call the system
  - Write a similar script to classify data in another toy problem

- Directly applicable to sunset detector

# Kernel functions

- Note that a hyperplane (which by definition is linear) in the feature space = a nonlinear boundary in the input space
  - Recall the RBFs
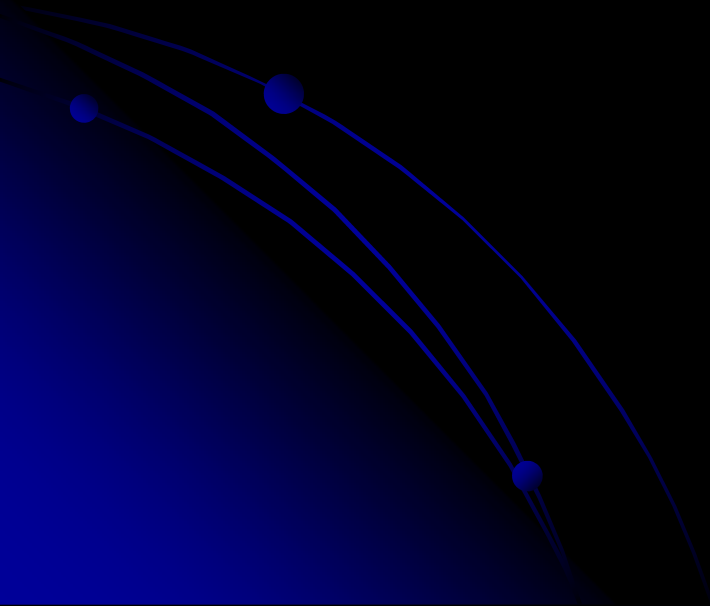- Note how choice of $\sigma$ affects the classifier

# Comparison with neural nets

- Expensive
  - Training can take a *long* time with large data sets. Consider that you'll want to experiment with parameters…
  - The classification runtime and space are *O(sd)*, where *s* is the number of support vectors, and d is the dimensionality of the feature vectors.
  - In the worst case, *s* = size of whole training set (like nearest neighbor)
  - But no worse than implementing a neural net with *s* perceptrons in the hidden layer.
- Empirically shown to have good generalizability even with relatively-small training sets and no domain knowledge.

# Speaking of neural nets:

- Demo matlabNeuralNetDemo.m

# Preview of demsvm2 (if time)

- Shows differing values of C for positive and negative examples.
- Part of tomorrow's lab