

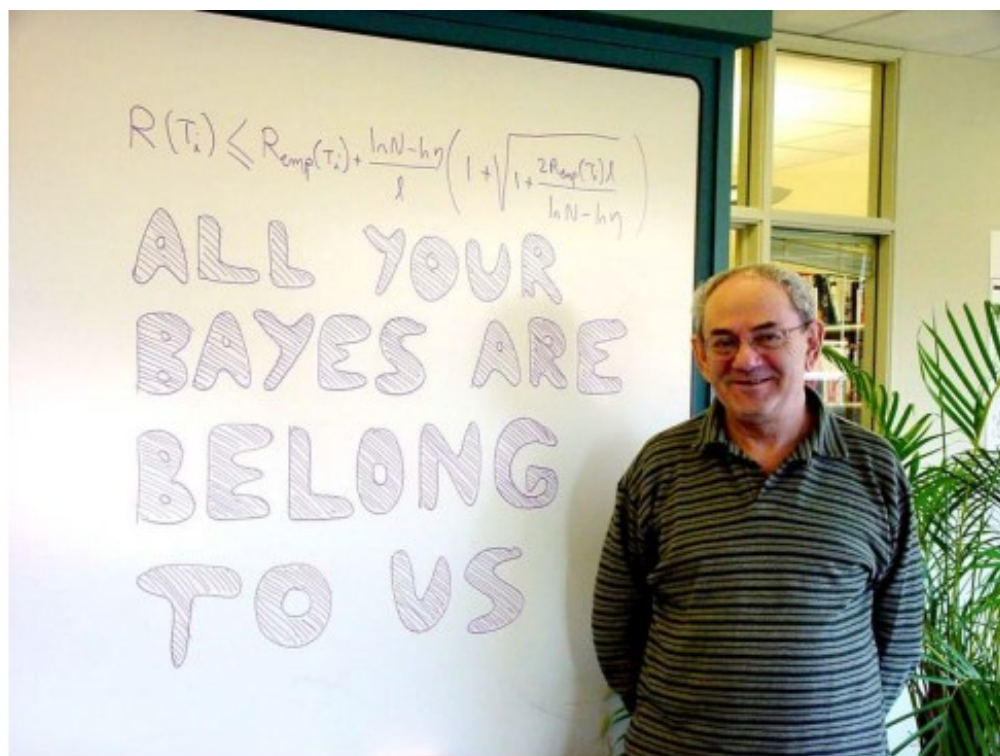
A decorative graphic on the left side of the slide consisting of a network of thin, light blue lines and small circles, resembling a circuit board or a neural network diagram. The lines are vertical and horizontal, with some diagonal segments, and the circles are small and white with blue outlines.

Lecture 9: Support Vector Machines

Agenda

- k-NN
- Decision Tree based Methods
- Support Vector Machines
- Neural Networks
- Deep learning
- Applications (Biophysical Modeling, Genomics, etc.)
- NLP

Support Vector Machines (SVM)



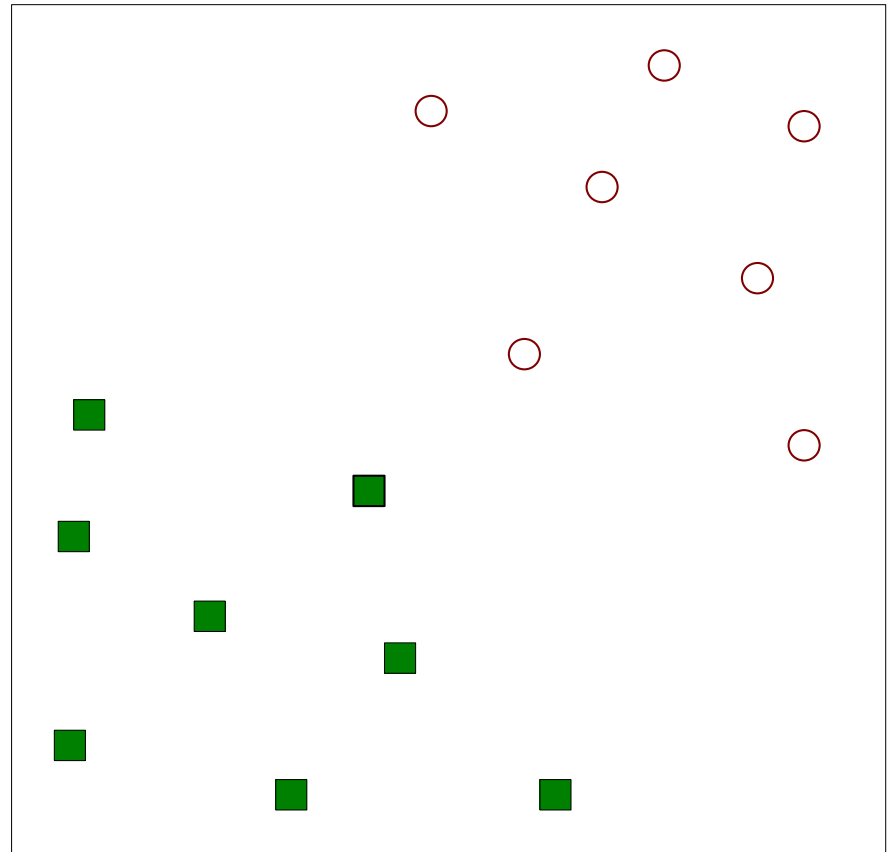
Vladimir Vapnik

Support Vector Machine

- Represents the decision boundary using only a **subset** of training examples known as **support vectors**
- Convex optimization problems with a **unique solution**
- Normally, a linear classifier
 - But, we can also use non-linear kernels
 - Kernels: application-specific measures of similarity
 - Leading to “**kernel machines**”

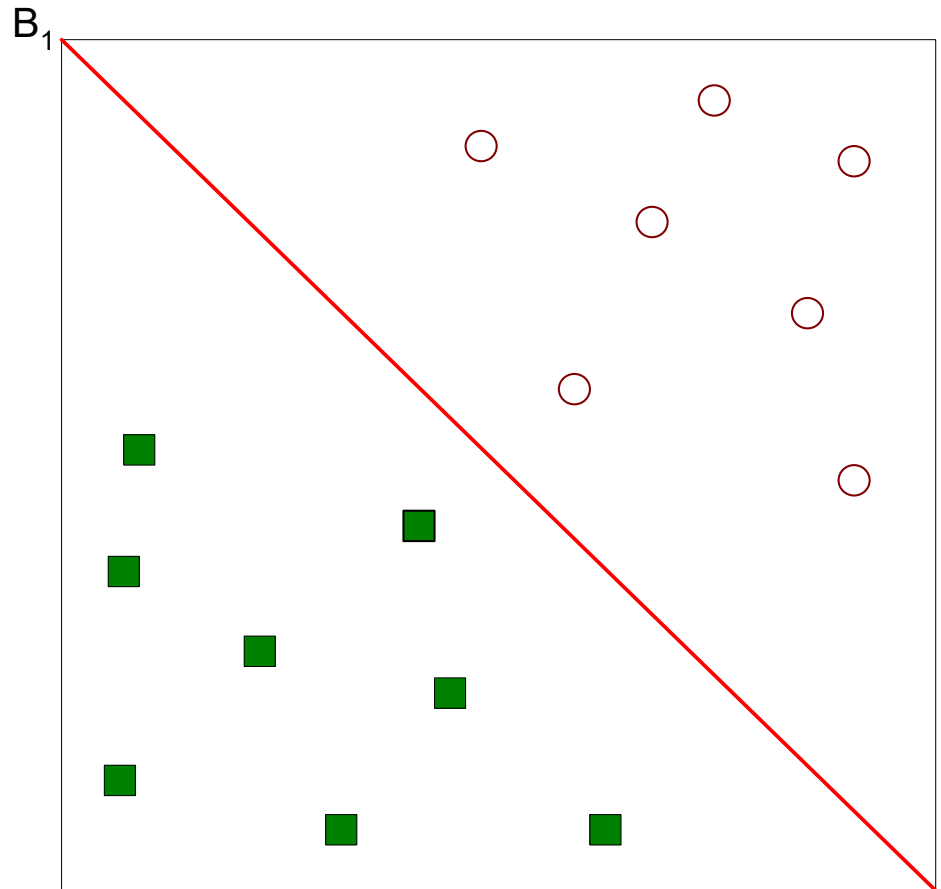
SVM Decision Boundary

- Question
 - Find a linear hyperplane to separate the data



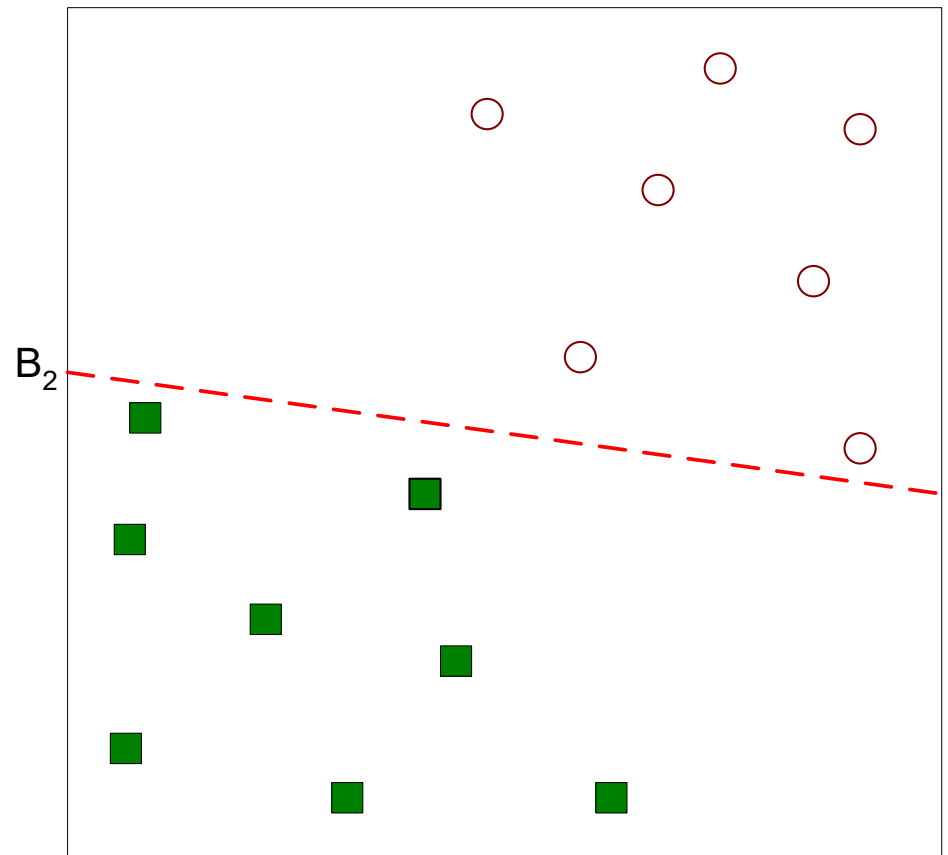
SVM Decision Boundary

- One possible solution



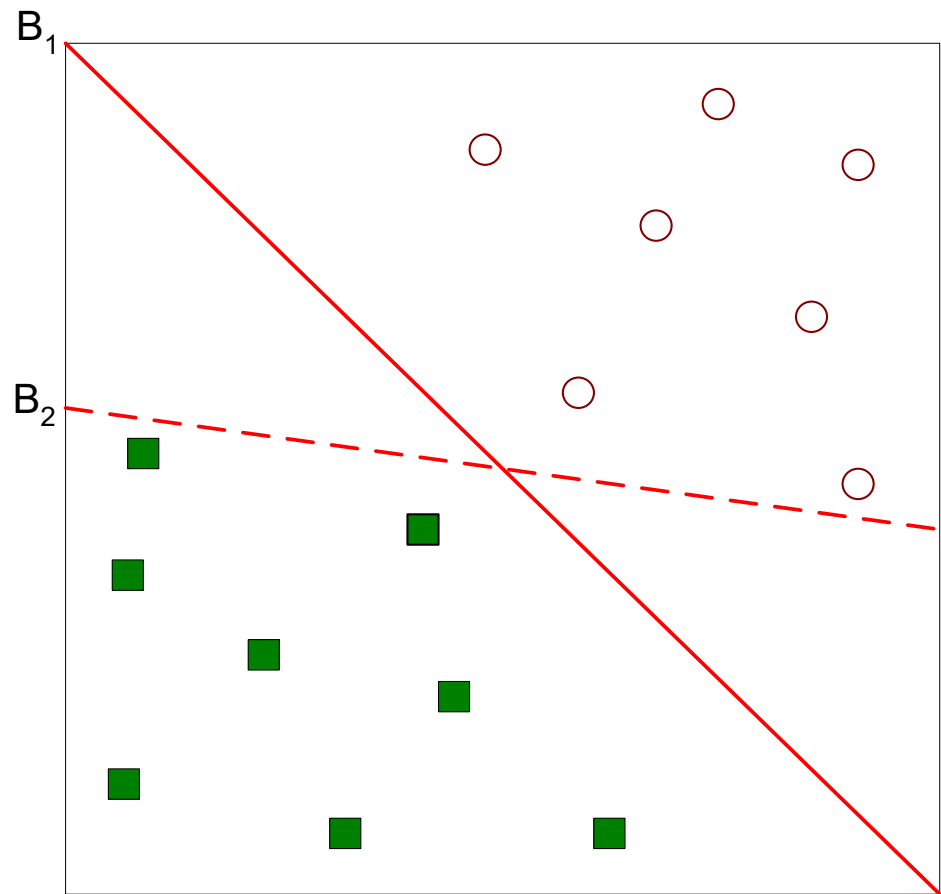
SVM Decision Boundary

- Another possible solution

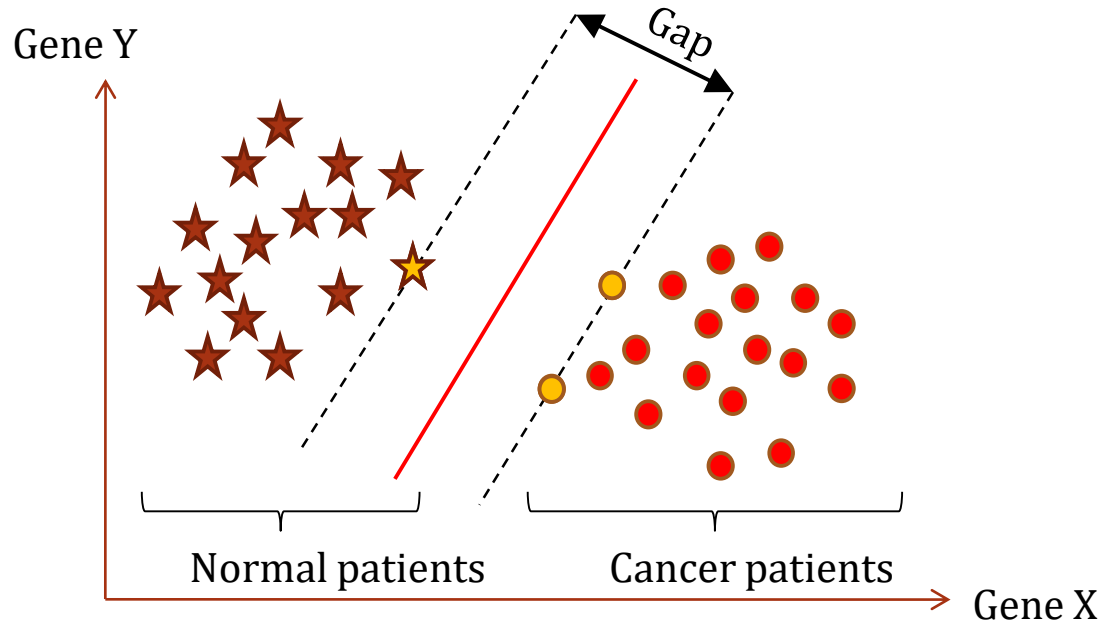


SVM Decision Boundary

- Which one is better?



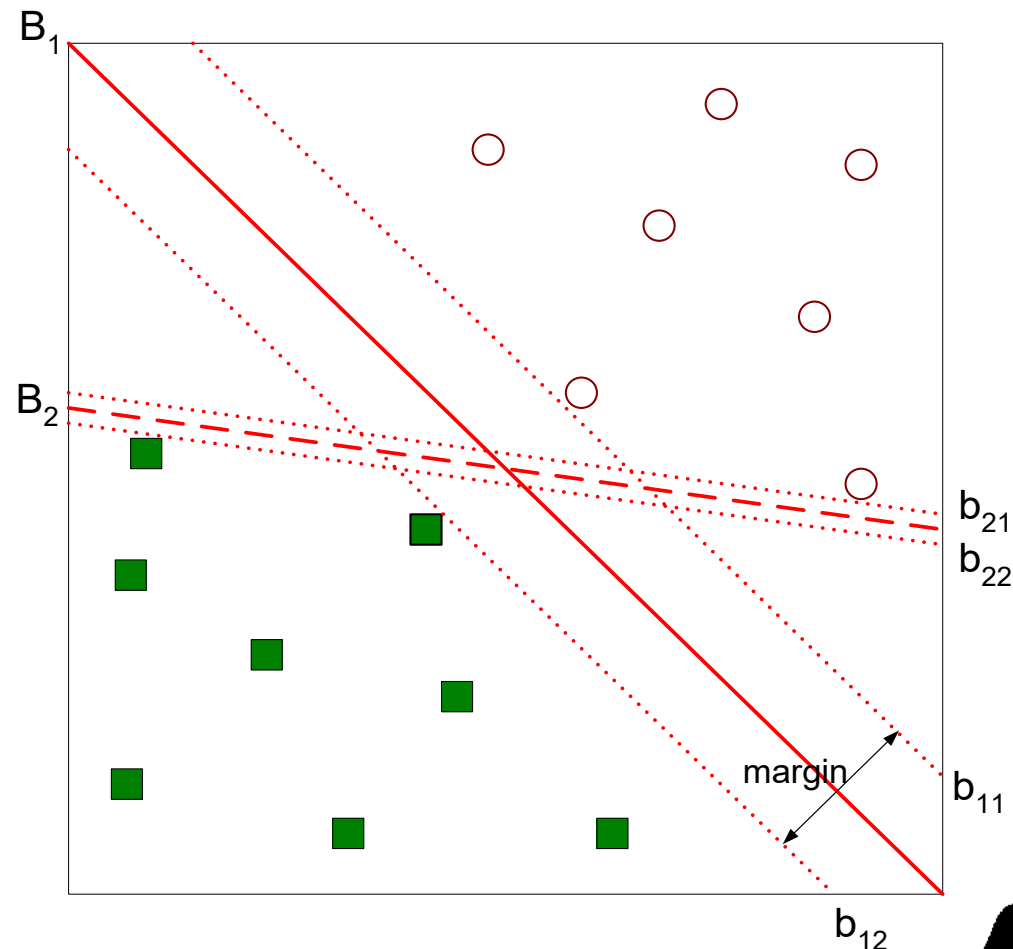
Main ideas of SVMs



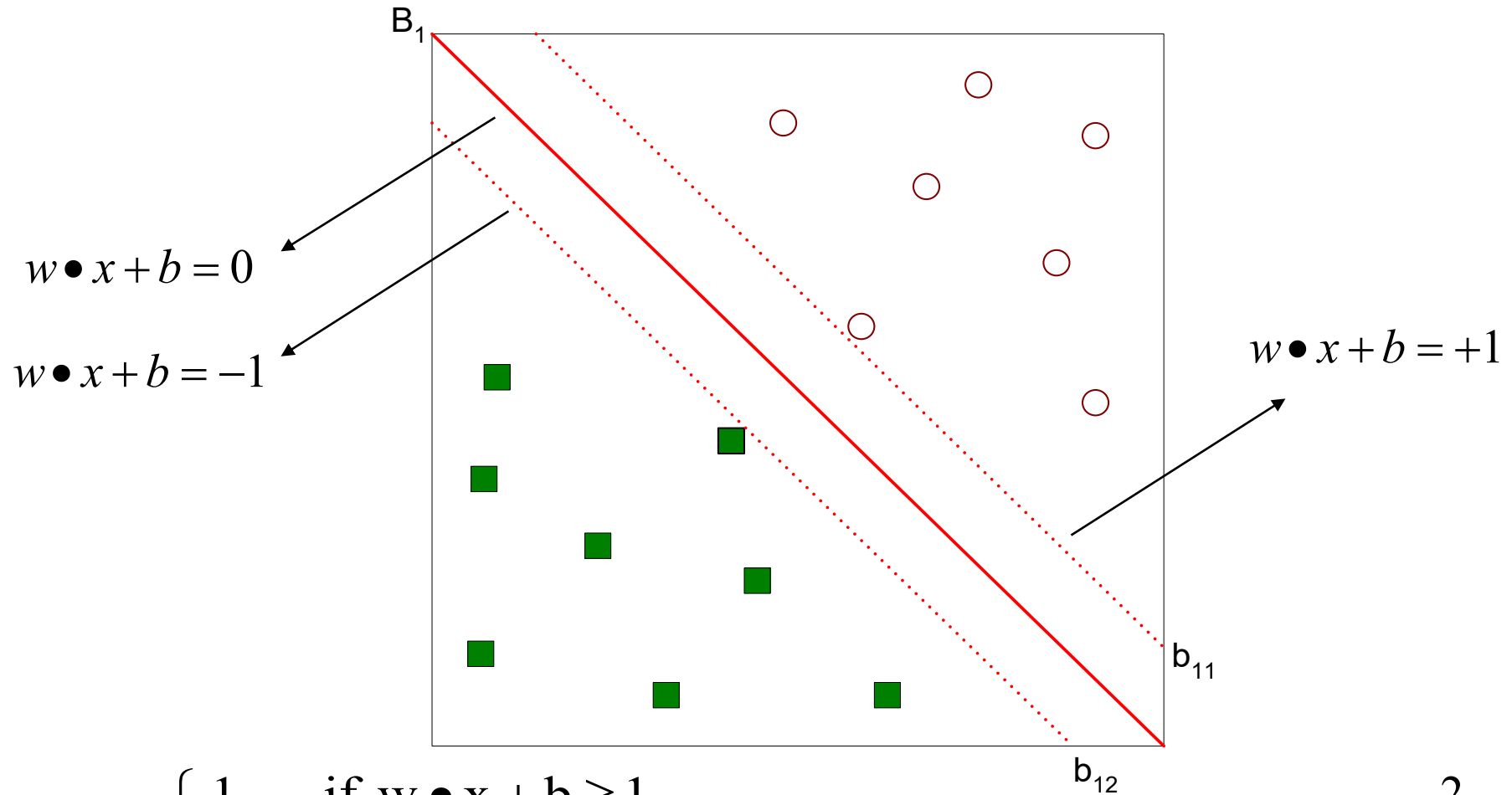
- Represent patients geometrically (by “vectors”);
- Find a linear decision surface (“hyperplane”) that can separate patient classes and has the largest distance (i.e., largest “gap” or “margin”) between border-line patients (i.e., “support vectors”);

Better Decision Boundary

- Margin definition:
 - Distance from the hyperplane to the closest instances on either side
- Find hyperplane that **maximizes** the **margin**
- B1 is better than B2
 - Wider margin, and better generalization



Support Vector Machine



$$f(x) = \begin{cases} 1 & \text{if } w \bullet x + b \geq 1 \\ -1 & \text{if } w \bullet x + b \leq -1 \end{cases}$$

$$\text{Margin} = \frac{2}{\|w\|^2}$$

Let's Formulate this ...

- Our simple rule

It can be re-written

$$f(x) = \begin{cases} 1 & \text{if } \mathbf{w} \bullet \mathbf{x} + b \geq 1 \\ -1 & \text{if } \mathbf{w} \bullet \mathbf{x} + b \leq -1 \end{cases} \quad \xrightarrow{\text{as}} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

- But, we also want to maximize the margin (or minimize the inverse of the margin), so we get:

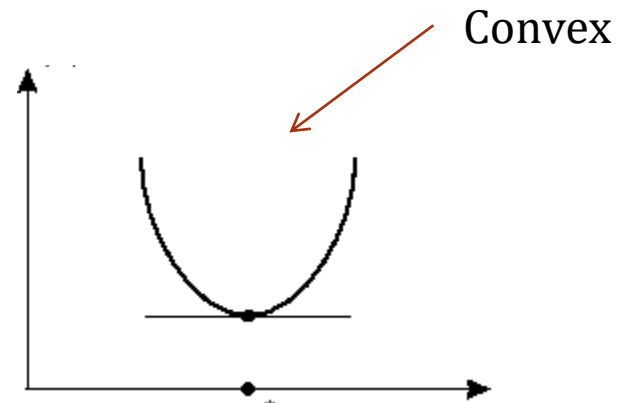
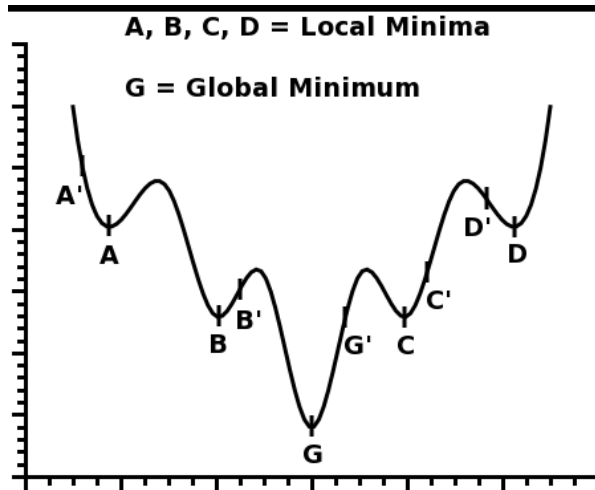
$$\min\left(\frac{\|\mathbf{w}\|^2}{2}\right)$$

$$\text{Subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad i = 1, 2, \dots, N$$

A Convex Optimization Problem

- This is a convex optimization problem (quadratic programming or QP)
 - Therefore, it has only a single optimum (great!)
 - Therefore no worries about convergence, learning rates, etc.

Non-convex



SVM Training

- Involves estimating parameters \mathbf{w} and b
- This is a quadratic optimization programming (QP) problem
 - Can be solved using Lagrange multiplier technique

$$\min(\frac{\|\mathbf{w}\|^2}{2})$$

$$\text{Subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad i = 1, 2, \dots, N$$

Lagrange Multiplier

- Write its Lagrangian:

$$(0) \quad L_p = \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^N \lambda_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1) \quad \text{where } \lambda_i \geq 0$$

- Here λ_i are called the Lagrange multipliers
- We will minimize L_p with respect to \mathbf{w} and b :

$$(1) \quad \frac{\partial L_p}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i$$

$$(2) \quad \frac{\partial L_p}{\partial b} = 0 \quad \Rightarrow \quad \sum_{i=1}^N \lambda_i y_i = 0$$

KKT Conditions

- Additionally, the following should hold (KKT condition, Karush-Kuhn-Tucker):

$$(3) \quad \lambda_i \geq 0$$

$$(4) \quad \lambda_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1] = 0$$

Support Vectors

- In equation (4) we had,

$$(4) \quad \lambda_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1] = 0$$

To be zero, either λ_i must be zero, or

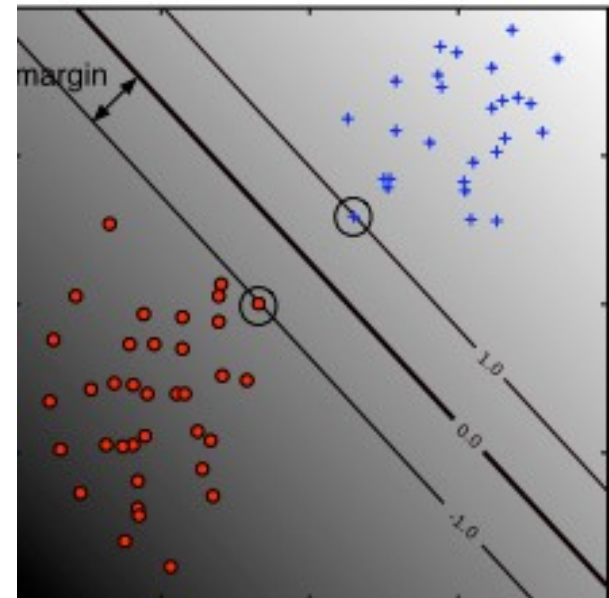
$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$$

which means that such a training instance \mathbf{x}_i lies along the hyperplanes b_1 or b_2

- Such a training instance is known as a **support vector**

Support Vectors

- A subset of training set, which are on the decision boundary
 - i.e. those for which $\lambda_i > 0$
- For the other training instances (i.e. those far away from the boundary), $\lambda_i = 0$
 - So, they have no effect on the hyperplane.



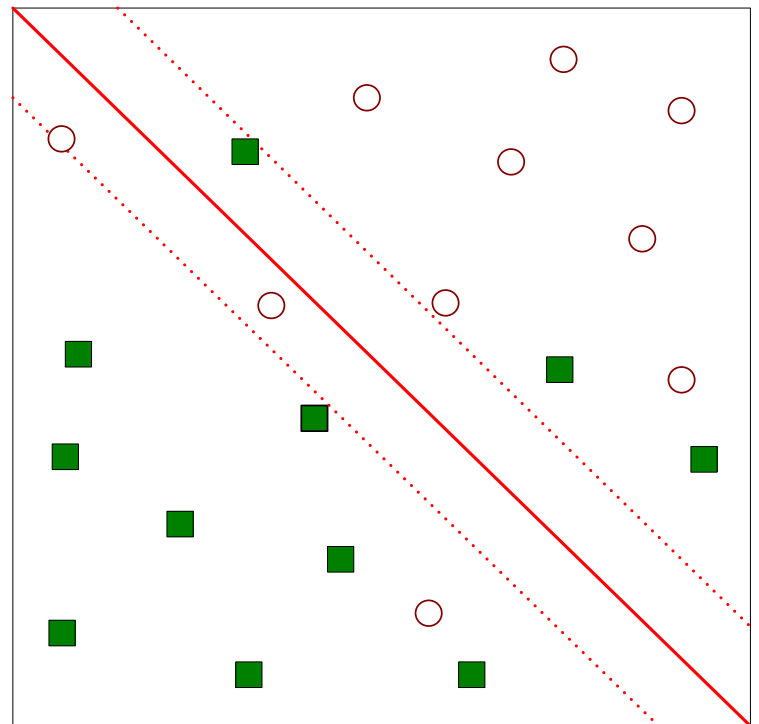
Final Solution

- We can obtain b and λ using numerical solutions from equations (0)-(4) (we won't go into details)
- The final solution will be:

We use this equation in order to classify a new data point z . $\longrightarrow f(z) = \text{sign}(\mathbf{w} \cdot \mathbf{z} + b) = \text{sign}\left(\sum_{i=1}^n \lambda_i y_i \mathbf{x}_i \cdot \mathbf{z} + b\right)$

Non-separable case

- How decision boundary can be modified to tolerate small training error?
 1. Instance may lie on the wrong side
 2. Instance may be inside the margin



Soft Margin

- We should consider a decision boundary that is **tolerable** to small training errors
 - This approach is known as **soft margin**
 - But SVM must consider the trade-off between the width of the margin and the number of training errors committed by the boundary

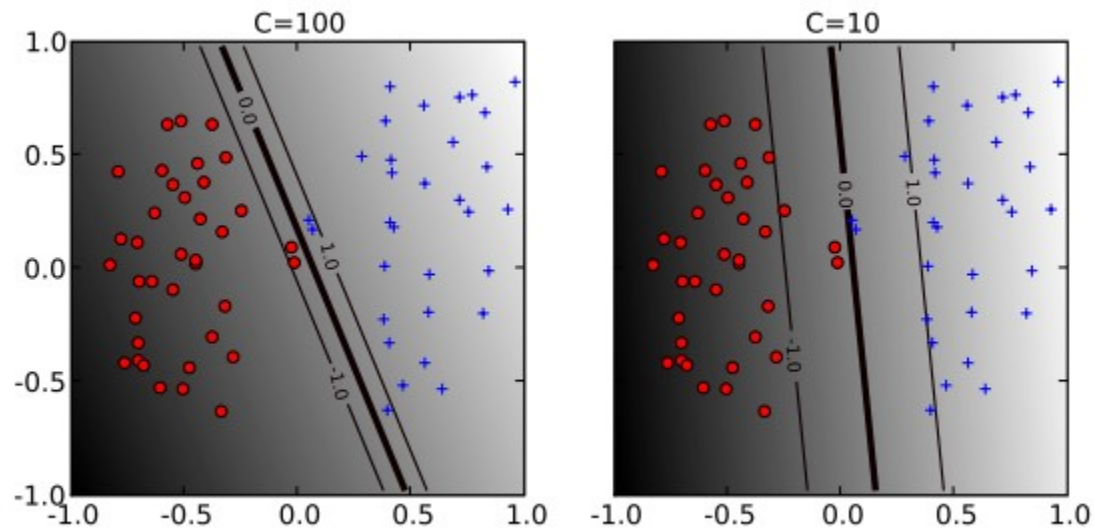
Soft Margin

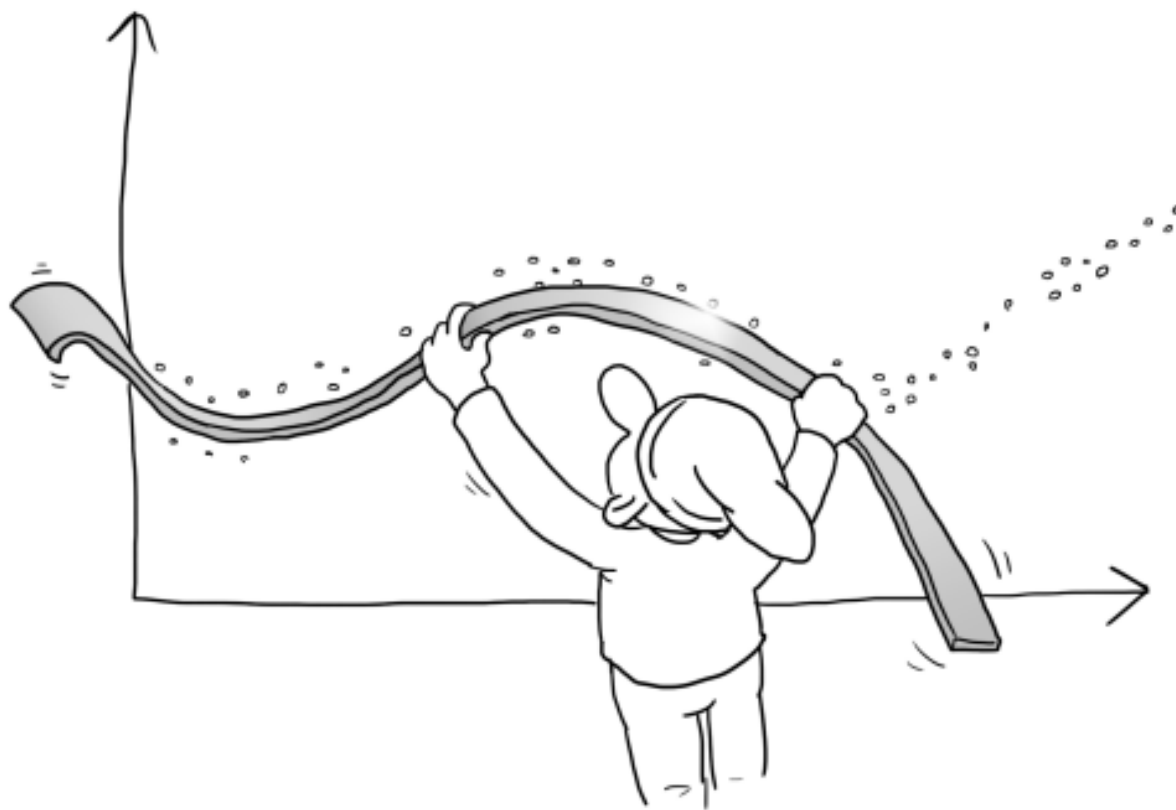
- The inequality constrained should be relaxed
 - Introduce the **slack variables** ξ
- The objective function must be modified to **penalize** a decision boundary with large values of slack variables

$$\min\left(\frac{\|w\|^2}{2} + C\left(\sum_{i=1}^N \xi_i\right)\right)$$

$$f(x_i) = \begin{cases} 1 & \text{if } w \bullet x_i + b \geq 1 - \xi_i \\ -1 & \text{if } w \bullet x_i + b \leq -1 + \xi_i \end{cases}$$

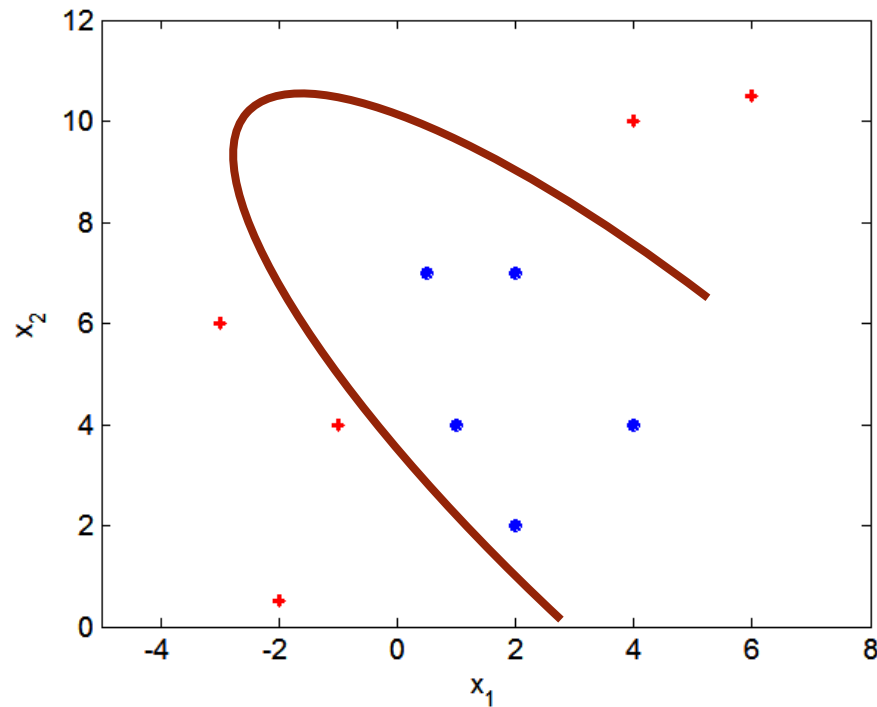
Effect of C



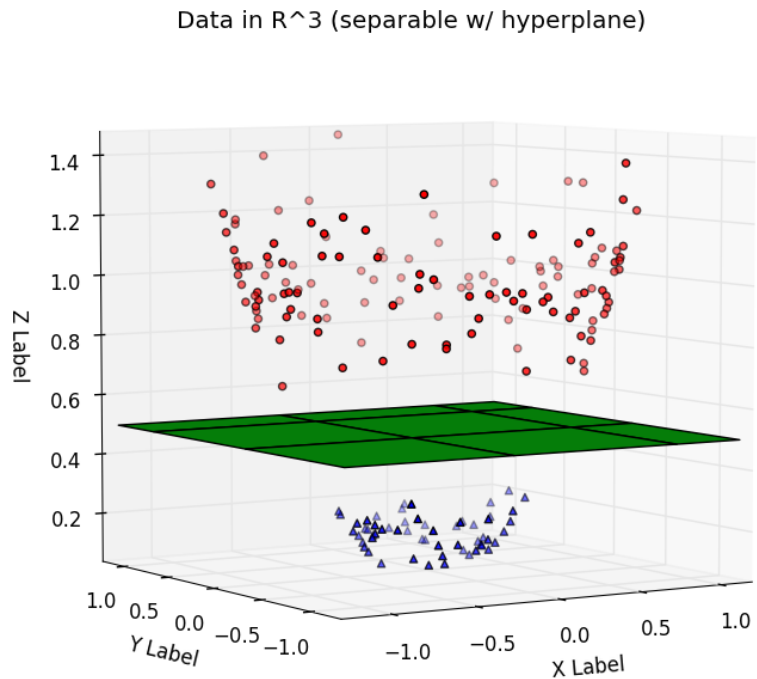
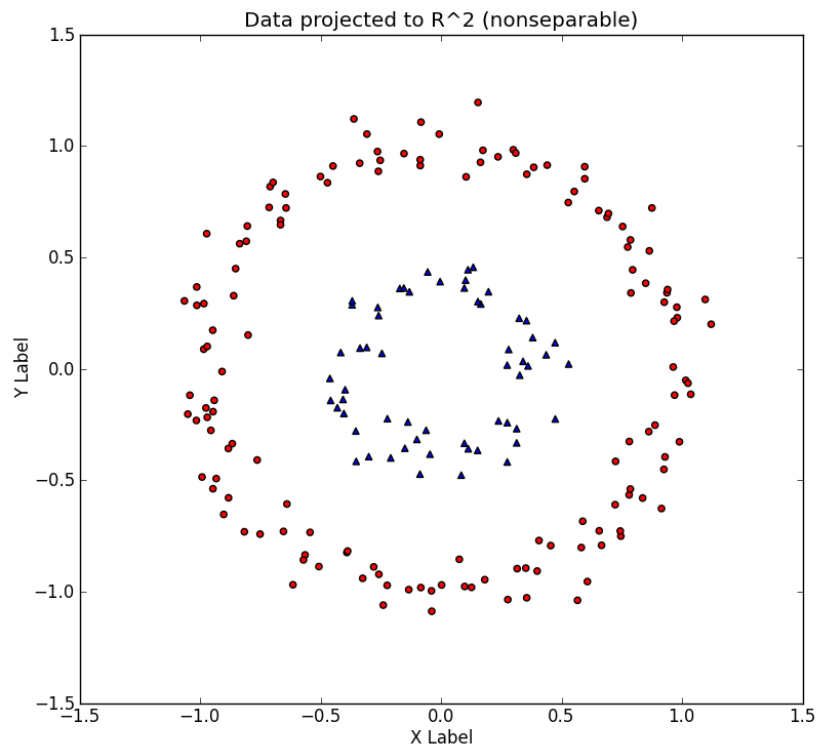


Nonlinear Decision Boundary

- What if the decision boundary is not linear?



Separable in higher-dimension



Kernels in Higher Dimensions

Animated

1. <https://www.youtube.com/watch?t=42&v=3liCbRZPrZA>
2. <https://www.youtube.com/watch?v=9NrALgHFwTo>

Nonlinear Decision Boundary

- The trick is to transform data from its original coordinates space x into a new space $(\phi(x))$

$$\min\left(\frac{\|w\|^2}{2}\right)$$

$$\text{Subject to } y_i(w \cdot \phi(x_i) + b) \geq 1 \quad i = 1, 2, \dots, N$$

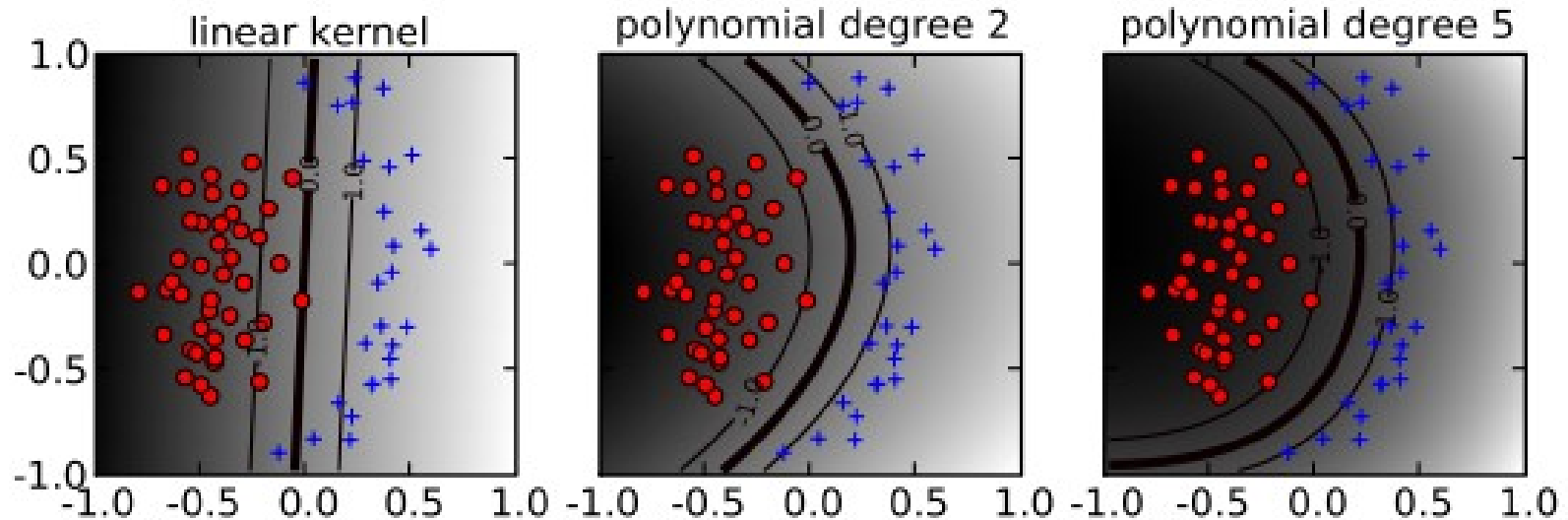
Solution



We use this equation in order to classify a new data point z .

$$\rightarrow f(z) = \text{sign}\left(\sum_{i=1}^n \lambda_i y_i \phi(x_i) \cdot \phi(z) + b\right)$$

Example



Popular kernels

A kernel is a dot product in *some* feature space:

$$K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$$

Examples:

$$K(x_i, x_j) = x_i \cdot x_j$$

Linear kernel

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

Gaussian kernel

$$K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|)$$

Exponential kernel

$$K(x_i, x_j) = (p + x_i \cdot x_j)^q$$

Polynomial kernel

$$K(x_i, x_j) = (p + x_i \cdot x_j)^q \exp(-\gamma \|x_i - x_j\|^2)$$

Hybrid kernel

$$K(x_i, x_j) = \tanh(kx_i \cdot x_j - \delta)$$

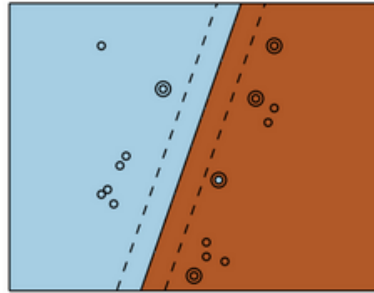
Sigmoidal

SVM Characteristics

- Convex optimization problem
 - Finds global minimum
 - Other methods such as neural network find local minimum
- Parameters should be tuned
 - Kernel type, C, ..
- High computational demands (both test and training)
 - Training stage
 - Naïve QP implementation
 - $O(n^3)$
 - More efficient techniques
 - Between $O(n)$ and $O(n^2)$

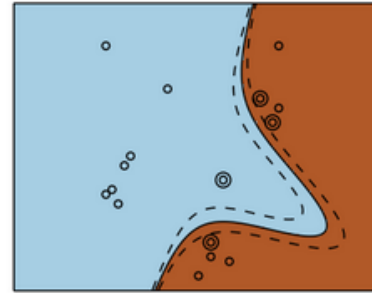
SVM in Scikit-learn (Python)

Linear kernel



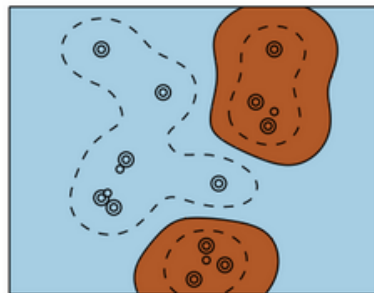
```
>>> svc = svm.SVC(kernel='linear')>>>
```

Polynomial kernel



```
>>> svc = svm.SVC(kernel='poly',>>>  
...               degree=3)  
>>> # degree: polynomial degree
```

RBF kernel (Radial Basis Function)

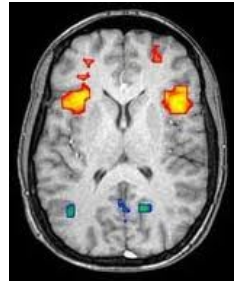


```
>>> svc = svm.SVC(kernel='rbf')>>>  
>>> # gamma: inverse of size of  
>>> # radial kernel
```


Neuroimaging Data (fMRI)

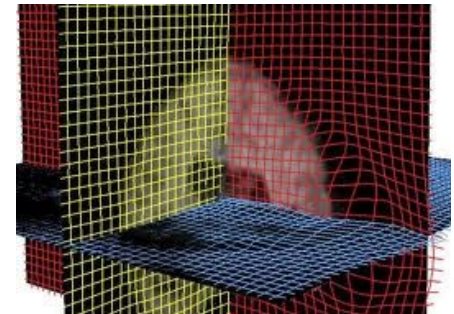
Neuroimaging Data Analysis

- Classify structural MR images
 - identify disease precursors as early as possible
 - e.g. MCI; Alzheimer's disorder (AD)
 - identify structural trajectories in neurodevelopment disorders
- Classify functional MR images
 - identify brain state associated with a stimulus



Neuroimaging Data

- Each image has lots and lots of voxels.
 - Of those, $\sim 30,000$ are actually brain voxels
 - grey matter $\approx 23,000$ voxels
 - white matter $\approx 8,000$ voxels
- To reduce the feature-space, we can:
 - resample data
 - do some latent variable extraction
 - e.g. using Principal Components Analysis (PCA) or some other automated feature extraction
- Then run SVM using a smaller number of “features”
 - e.g. 3-8 PCs rather than thousands of voxels

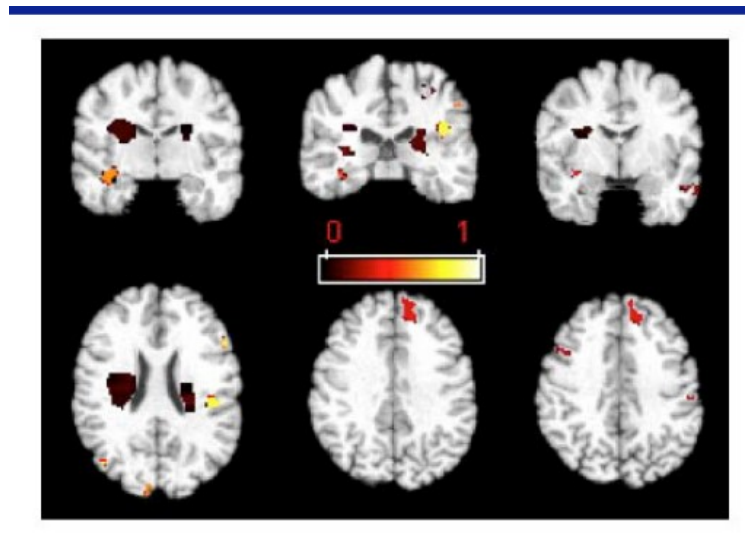


Example Features

- **Average.** For each ROI, calculate the mean activity over all voxels in the ROI. Use these ROI means as the input features.
- **ActiveAvg(n).** For each ROI, select the most active voxels, then calculate the mean of their values. Again, use these ROI means as the input features.
- **“most active”** voxels are those whose activity while performing the task varies the most from their activity when the subject is at rest.
- **Active(n).** Select the most active voxels over the entire brain. Use only these voxels as input features.

Example: Schizophrenia

- Motivation: Which brain areas best distinguish brains of female schizophrenic patients from typical?



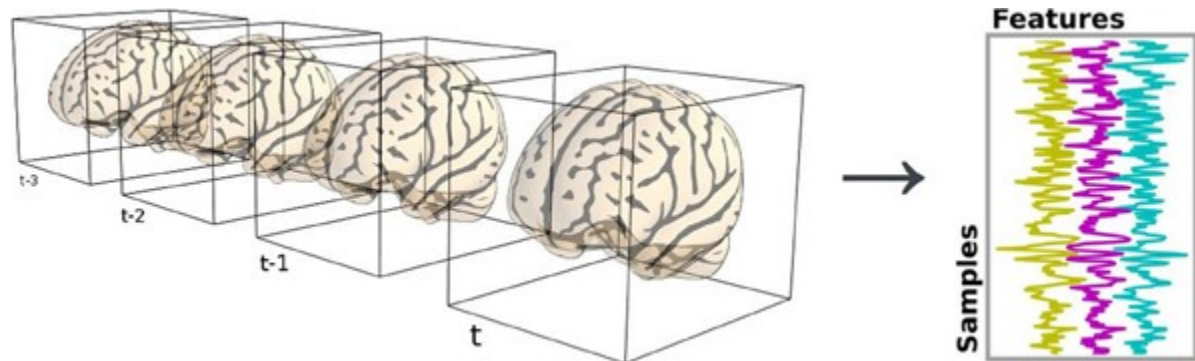
Analysis with scikit-learn

#First we need to read data (and its mask)

#nibabel is a reader of common neuroimaging formats

```
>>import nibabel as ni
```

```
>>> X = ni.load("bold.nii").get_data()
```



Preprocessing

- **De-trending**: remove linear trend
(`scipy.signal.detrend`)
- **Normalization**: all features will have the same range
(`sklearn.preprocessing.normalize`)
- **Frequency filtering**: e.g. low frequencies due to physiological mechanisms, use Fourier Transform
(`scipy.fftpack.fft`)

Prediction

```
from sklearn.linear_model import
    LogisticRegression as LR
from sklearn.cross_validation import
    cross_val_score

pipeline_LR = Pipeline([('selection',
    SelectKBest(f_classif, 500)),
    ('clf', LR(penalty='l1', C=0.05))])

scores_lr = []

for pixel in y_train.T:
    score = cross_val_score(pipeline_LR,
        X_train, pixel, cv=5)
    scores_lr.append(score)
```