Support Vector Machines

Emmanuel Jakobowicz

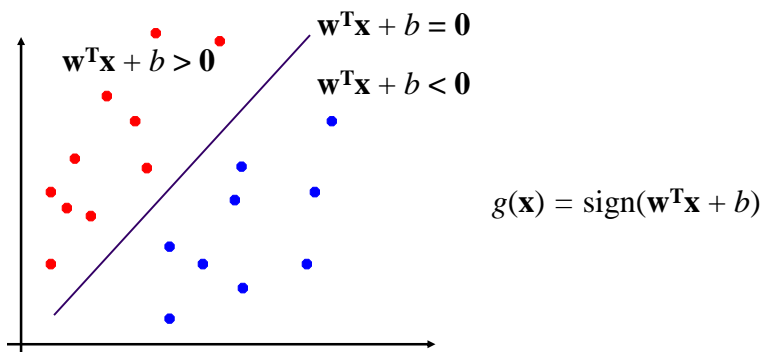1

## Goal of SVM

- SVM is a binary classifier
- Example:
  - Classify patients
  - Classify customers
- SVM is a supervised learning algorithm
- The dependent variable is binary

2

## Linear Separators

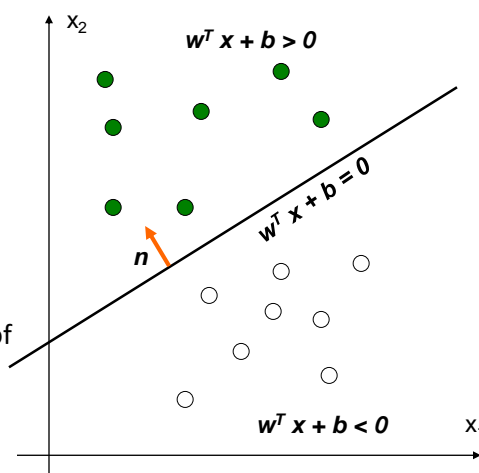- Binary classification can be viewed as the task of separating classes in feature space:

$$w^T x + b = 0$$
$$w^T x + b > 0$$
$$w^T x + b < 0$$

$$g(\mathbf{x}) = \text{sign}(\mathbf{w^T x} + b)$$

3

## Linear Discriminant Function

- g(x) is a linear function:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

- A hyper-plane in the feature space

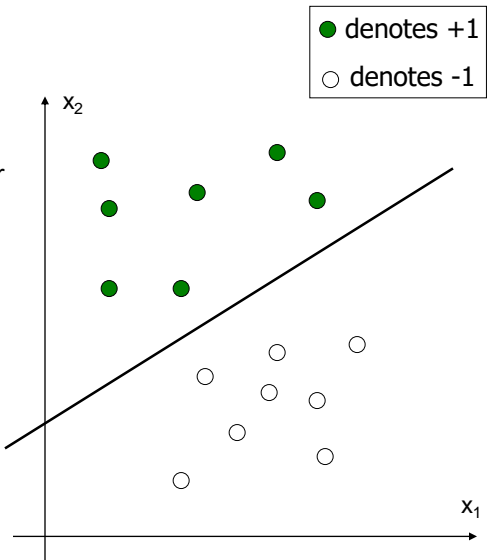- (Unit-length) normal vector of the hyper-plane:

$$\mathbf{n} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$x_2$

$w^T x + b > 0$

$w^T x + b = 0$

$n$

$w^T x + b < 0$    $x_1$

4

# Linear Discriminant Function



- How would you classify these points using a linear discriminant function in order to minimize the error rate?

- Infinite number of answers!

# Linear Discriminant Function
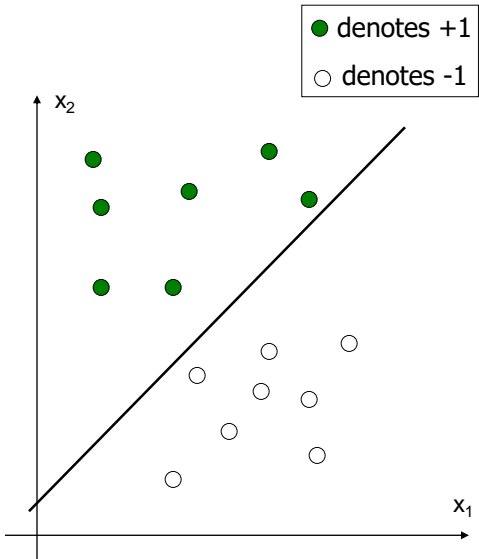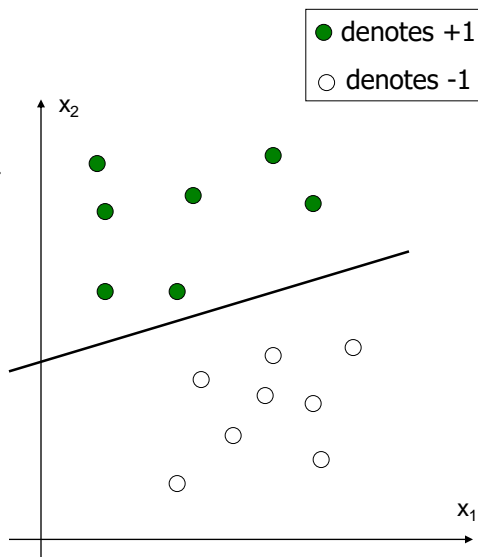


- How would you classify these points using a linear discriminant function in order to minimize the error rate?

- Infinite number of answers!

## Linear Discriminant Function

denotes +1
denotes -1

- How would you classify these points using a linear discriminant function in order to minimize the error rate?

- Infinite number of answers!

$x_2$

$x_1$

7

## Linear Discriminant Function

denotes +1
denotes -1

- How would you classify these points using a linear discriminant function in order to minimize the error rate?

- Infinite number of answers!

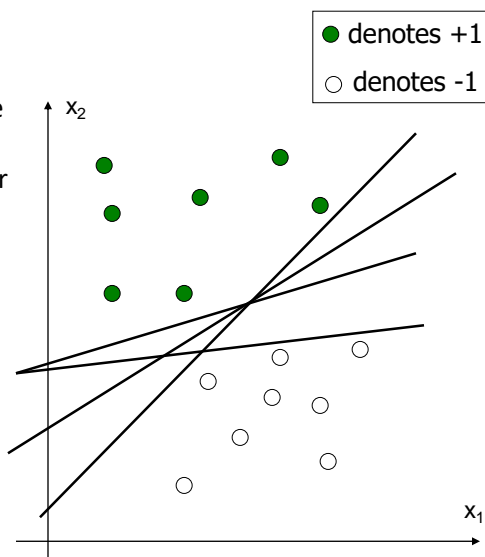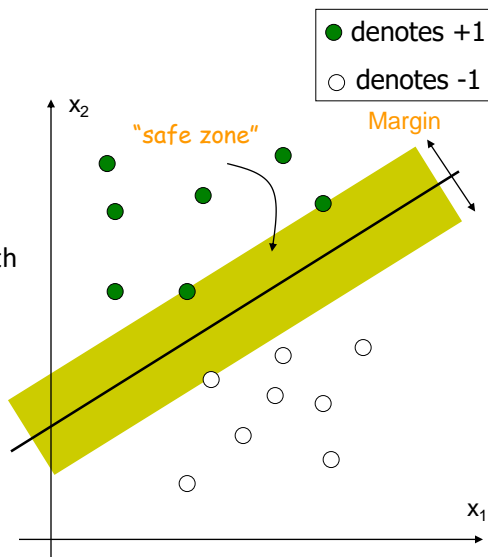- Which one is the best?

$x_2$

$x_1$

8

# Large Margin Linear Classifier

- The linear discriminant function (classifier) with the maximum margin is the best

- Margin is defined as the width that the boundary could be increased by before hitting a data point

- Why it is the best?
  - Robust to outliners and thus strong generalization ability



# Classification Margin

- Distance from example $\mathbf{x}_i$ to the separator is $r = \dfrac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|}$
- Examples closest to the hyperplane are **support vectors**.
- **Margin** $\rho$ of the separator is the distance between support vectors.

# Maximum Margin Classification

- Maximizing the margin is good according to intuition.
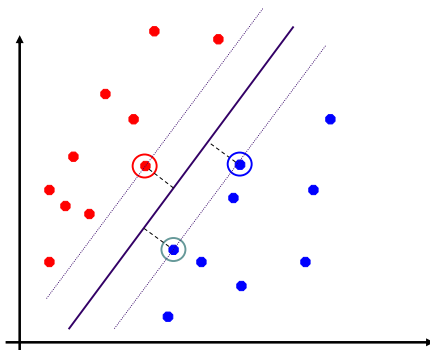- Implies that only support vectors matter; other training examples are ignorable.



11

# Large Margin Linear Classifier

- We know that

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

- The margin width is:

$$M = (\mathbf{x}^+ - \mathbf{x}^-) \cdot \mathbf{n}$$
$$= (\mathbf{x}^+ - \mathbf{x}^-) \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$



- ● denotes +1
- ○ denotes -1

$x_2$

Margin

$\mathbf{w}^T x + b = 1$
$\mathbf{w}^T x + b = 0$
$\mathbf{w}^T x + b = -1$

$x^+$
$n$
$x^-$

Support Vectors

$x_1$

12

6

## Linear SVMs Mathematically

- Then we can formulate the *quadratic optimization problem:*

Find $\mathbf{w}$ and $b$ such that
$\rho = \dfrac{2}{\|\mathbf{w}\|}$ is maximized
and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ : $\quad y_i(\mathbf{w}^{\mathbf{T}}\mathbf{x}_i + b) \geq 1$

Which can be reformulated as:

Find $\mathbf{w}$ and $b$ such that
$\mathbf{\Phi}(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^{\mathrm{T}}\mathbf{w}$ is minimized
and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ : $\quad y_i(\mathbf{w}^{\mathbf{T}}\mathbf{x}_i + b) \geq 1$

13

## Solving the Optimization Problem

Find $\mathbf{w}$ and b such that
$\mathbf{\Phi}(\mathbf{w}) = \mathbf{w}^{\mathrm{T}}\mathbf{w}$ is minimized
and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ : $\quad y_i(\mathbf{w}^{\mathbf{T}}\mathbf{x}_i + b) \geq 1$

- Need to optimize a *quadratic* function subject to *linear* constraints.
- Quadratic optimization problems are a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist.
- The solution involves constructing a *dual problem* where a *Lagrange multiplier* $\alpha_i$ is associated with every inequality constraint in the primal (original) problem:

Find $\alpha_1...\alpha_n$ such that
$\mathbf{Q}(\boldsymbol{\alpha}) = \Sigma\alpha_i - \tfrac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j \mathbf{x}_i^{\mathbf{T}}\mathbf{x}_j$ is maximized and
(1) $\Sigma\alpha_i y_i = 0$
(2) $\alpha_i \geq 0$ for all $\alpha_i$

14

## The Optimization Problem Solution

- Given a solution $\alpha_1 ... \alpha_n$ to the dual problem, solution to the primal is:

$$\mathbf{w} = \Sigma \alpha_i y_i \mathbf{x}_i \qquad b = y_k - \Sigma \alpha_i y_i \mathbf{x}_i{}^{\mathbf{T}} \mathbf{x}_k \quad \text{for any } \alpha_k > 0$$

- Each non-zero $\alpha_i$ indicates that corresponding $\mathbf{x}_i$ is a support vector.
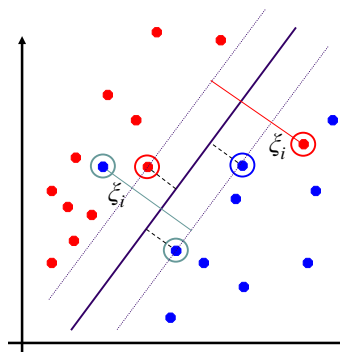- Then the classifying function is (note that we don't need $\mathbf{w}$ explicitly):

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \mathbf{x}_i{}^{\mathbf{T}} \mathbf{x} + b$$

- Notice that it relies on an *inner product* between the test point $\mathbf{x}$ and the support vectors $\mathbf{x}_i$
- Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x}_i{}^{\mathbf{T}} \mathbf{x}_j$ between all training points.

15

## Soft Margin Classification

- What if the training set is not linearly separable?
- *Slack variables $\xi_i$* can be added to allow misclassification of difficult or noisy examples, resulting margin called *soft*.



What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k$$

*C is an hyper-parameter*

16

### Soft Margin Classification Mathematically

- The old formulation:

> Find $\mathbf{w}$ and b such that
> $\Phi(\mathbf{w}) = \mathbf{w}^T\mathbf{w}$ is minimized
> and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ : $\quad y_i(\mathbf{w^T x}_i + b) \geq 1$

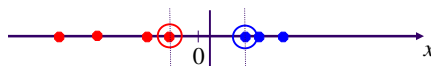- Modified formulation incorporates slack variables:

> Find $\mathbf{w}$ and b such that
> $\Phi(\mathbf{w}) = \mathbf{w}^T\mathbf{w} + C\Sigma\xi_i$ is minimized
> and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ : $\quad y_i(\mathbf{w^T x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$

- Parameter *C* can be viewed as a way to control overfitting: it "trades off" the relative importance of maximizing the margin and fitting the training data.
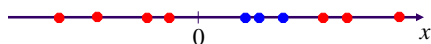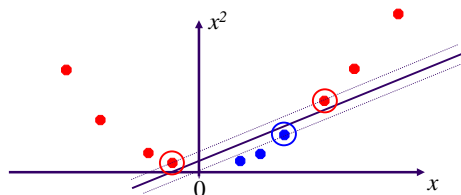
17

# Non-linear SVMs

- Datasets that are linearly separable with some noise work out great:



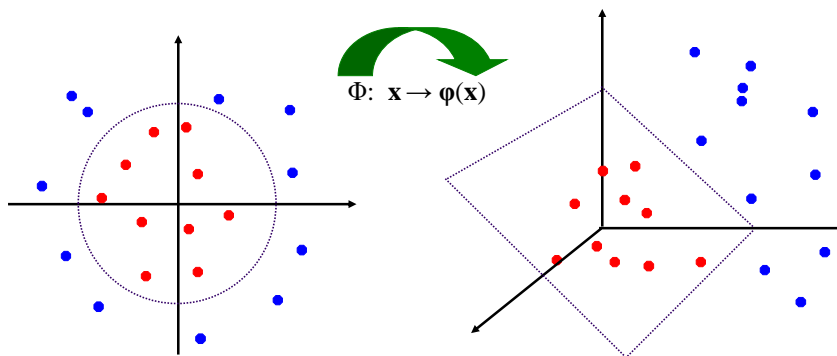- But what are we going to do if the dataset is just too hard?



- How about… mapping data to a higher-dimensional space:



18

# Non-linear SVMs:  Feature spaces

- General idea:  the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:

$$\Phi: \ \mathbf{x} \to \varphi(\mathbf{x})$$

19

# The "Kernel Trick"

- The linear classifier relies on inner product between vectors $K(\mathbf{x}_i,\mathbf{x}_j)=\mathbf{x}_i^{\mathsf{T}}\mathbf{x}_j$
- If every datapoint is mapped into high-dimensional space via some transformation $\Phi: \ \mathbf{x} \to \phi(\mathbf{x})$, the inner product becomes:

$$K(\mathbf{x}_i,\mathbf{x}_j)= \boldsymbol{\phi}(\mathbf{x}_i)^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_j)$$

- A *kernel function* is a function that is equivalent to an inner product in some feature space.
- Example:

  2-dimensional vectors $\mathbf{x}=[x_1 \ x_2]$; let $K(\mathbf{x}_i,\mathbf{x}_j)=(1 + \mathbf{x}_i^{\mathsf{T}}\mathbf{x}_j)^2$,

  Need to show that $K(\mathbf{x}_i,\mathbf{x_j})= \boldsymbol{\phi}(\mathbf{x}_i)^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_j)$:

  $K(\mathbf{x}_i,\mathbf{x}_j)=(1 + \mathbf{x}_i^{\mathsf{T}}\mathbf{x}_j)^2, = 1+ x_{i1}^2 x_{j1}^2 + 2\ x_{i1}x_{j1}\ x_{i2}x_{j2}+ x_{i2}^2 x_{j2}^2 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2}=$
  $= [1 \ \ x_{i1}^2 \ \sqrt{2}\ x_{i1}x_{i2} \ \ x_{i2}^2 \ \sqrt{2}x_{i1} \ \sqrt{2}x_{i2}]^{\mathsf{T}} [1 \ \ x_{j1}^2 \ \sqrt{2}\ x_{j1}x_{j2} \ \ x_{j2}^2 \ \sqrt{2}x_{j1} \ \sqrt{2}x_{j2}] =$
  $= \boldsymbol{\phi}(\mathbf{x}_i)^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_j), \quad \text{where } \boldsymbol{\phi}(\mathbf{x}) = [1 \ \ x_1^2 \ \sqrt{2}\ x_1x_2 \ \ x_2^2 \ \sqrt{2}x_1 \ \sqrt{2}x_2]$

- Thus, a kernel function *implicitly* maps data to a high-dimensional space (without the need to compute each $\boldsymbol{\phi}(\mathbf{x})$ explicitly).

20

## What Functions are Kernels?

- For some functions $K(\mathbf{x}_i,\mathbf{x}_j)$ checking that $K(\mathbf{x}_i,\mathbf{x}_j)= \phi(\mathbf{x}_i)^\mathsf{T}\phi(\mathbf{x}_j)$ can be cumbersome.
- Mercer's theorem:

    ***Every semi-positive definite symmetric function is a kernel***
- Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:

$$K=$$

| $K(\mathbf{x}_1,\mathbf{x}_1)$ | $K(\mathbf{x}_1,\mathbf{x}_2)$ | $K(\mathbf{x}_1,\mathbf{x}_3)$ | ... | $K(\mathbf{x}_1,\mathbf{x}_n)$ |
|---|---|---|---|---|
| $K(\mathbf{x}_2,\mathbf{x}_1)$ | $K(\mathbf{x}_2,\mathbf{x}_2)$ | $K(\mathbf{x}_2,\mathbf{x}_3)$ | | $K(\mathbf{x}_2,\mathbf{x}_n)$ |
| | | | | |
| ... | ... | ... | ... | ... |
| $K(\mathbf{x}_n,\mathbf{x}_1)$ | $K(\mathbf{x}_n,\mathbf{x}_2)$ | $K(\mathbf{x}_n,\mathbf{x}_3)$ | ... | $K(\mathbf{x}_n,\mathbf{x}_n)$ |

For any non-zero vector x, $x^\mathsf{T}Kx>0$

21

## Examples of Kernel Functions

- Linear: $K(\mathbf{x_i},\mathbf{x_j})= \mathbf{x_i}^\mathsf{T}\mathbf{x_j}$

- Polynomial of power $p$: $K(\mathbf{x_i},\mathbf{x_j})= (1+ \mathbf{x_i}^\mathsf{T}\mathbf{x_j})^p$

- Gaussian (radial-basis function network):

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

- Sigmoid: $K(\mathbf{x_i},\mathbf{x_j})= \tanh(\beta_0\mathbf{x_i}^\mathsf{T}\mathbf{x_j} + \beta_1)$

22

# Support Vector Machine: Algorithm

- 1. Choose a kernel function

- 2. Choose a value for *C (soft)*

- 3. Solve the quadratic programming problem (many software packages available)

- 4. Construct the discriminant function from the support vectors

23

# Some Issues

- Choice of kernel
  - Gaussian or polynomial kernel are the mostly used non-linear kernels
  - if ineffective, more elaborate kernels are needed
  - domain experts can give assistance in formulating appropriate similarity measures

- Choice of kernel parameters
  - e.g. $\sigma$ in Gaussian kernel
  - $\sigma$ is the distance between closest points with different classifications
  - In the absence of reliable criteria, applications rely on the use of a validation set or cross-validation to set such parameters.

- Optimization criterion – Hard margin v.s. Soft margin
  - a lengthy series of experiments in which various parameters are tested

24

### Why Is SVM Effective on High Dimensional Data?

- The complexity of trained classifier is characterized by the # of support vectors rather than the dimensionality of the data

- The support vectors are the essential or critical training examples — they lie closest to the decision boundary

- If all other training examples are removed and the training is repeated, the same separating hyperplane would be found

- The number of support vectors found can be used to compute an (upper) bound on the expected error rate of the SVM classifier, which is independent of the data dimensionality

- Thus, an SVM with a small number of support vectors can have good generalization, even when the dimensionality of the data is high

## Weakness of SVM

- **It is sensitive to noise**
  - **A relatively small number of mislabeled examples can dramatically decrease the performance**

- **It only considers two classes**
  - **how to do multi-class classification with SVM?**
  - **Answer:**
  
  **1) with output arity m, learn m SVM's**
    - **SVM 1 learns "Output==1" vs "Output != 1"**
    - **SVM 2 learns "Output==2" vs "Output != 2"**
    - **:**
    - **SVM m learns "Output==m" vs "Output != m"**
  
  **2)To predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the positive region.**

26

# SVM applications

- SVMs were originally proposed by Boser, Guyon and Vapnik in 1992 and gained increasing popularity in late 1990s.
- SVMs are currently among the best performers for a number of classification tasks ranging from text to genomic data.
- SVMs can be applied to complex data types beyond feature vectors (e.g. graphs, sequences, relational data) by designing kernel functions for such data.
- SVM techniques have been extended to a number of tasks such as regression [Vapnik *et al.* '97], principal component analysis [Schölkopf *et al.* '99], etc.
- Most popular optimization algorithms for SVMs use *decomposition* to hill-climb over a subset of $\alpha_i$'s at a time, e.g. SMO [Platt '99] and [Joachims '99]
- Tuning SVMs remains a black art: selecting a specific kernel and parameters is usually done in a try-and-see manner.

# Summary: Support Vector Machine

- 1. Large Margin Classifier
  - Better generalization ability & less over-fitting

- 2. The Kernel Trick
  - Map data points to higher dimensional space in order to make them linearly separable.
  - Since only dot product is used, we do not need to represent the mapping explicitly.

28

14

## SVM resources

- Many references here:
  http://www.kernel-machines.org

- One of the first SVM software:
  http://www.csie.ntu.edu.tw/~cjlin/libsvm/

- SVM are available in R and python
  - In R: package e1071
  - In python: library scikit-learn

29

## WORKSHOP

- We work on the titanic dataset
- Either with R or python, run a SVM model on the training set and test it with the testing set
- The dependent variable is the survival variable

30

## WORKSHOP – Parameters tuning

- Using python scikit-learn or R e1071, use a grid search to obtain the best parameters for your SVM model
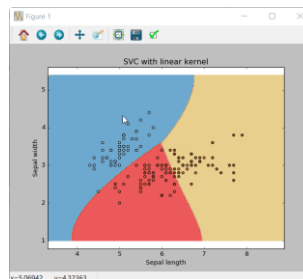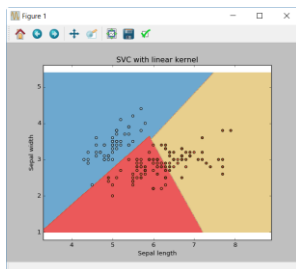- Choose an apropriate metric to fit parameters
- Obtain the ROC curve

31

## How to tune Parameters of SVM?

- Here is the list of parameters
- sklearn.svm.SVC(C=1.0, kernel='rbf', degree=3, gamma=0.0, coef0=0.0, shrinking=True, probability=False,tol=0.001, cache_size=200, class_weight=None, verbose=False, max_iter=-1, random_state=None)

- Important parameters are C, kernel and gamma
- C is the regularization parameter
  - Penalty parameter C of the error term. It also controls the trade off between smooth decision boundary and classifying the training points correctly.
- Kernel is the type of kernel
- Gamma is kernel parameter
  - Kernel coefficient for 'rbf', 'poly' and 'sigmoid'. Higher the value of gamma, will try to exact fit the as per training data set i.e. generalization error and cause over-fitting problem.

16

Use  np.meshgrid to build a grid and project the predicted classes

Try to modify the 3 main parameters to understand what happens

Work with the iris dataset



# Pros and cons

- **Pros:**
    - It works really well with clear margin of separation
    - It is effective in high dimensional spaces.
    - It is effective in cases where number of dimensions is greater than the number of samples.
    - It uses a subset of training points in the decision function (called support vectors), so it is also memory efficient.
- **Cons:**
    - It doesn't perform well, when we have large data set because the required training time is higher
    - It also doesn't perform very well, when the data set has more noise i.e. target classes are overlapping
    - SVM doesn't directly provide probability estimates, these are calculated using an expensive five-fold cross-validation. It is related SVC method of Python scikit-learn library.