

# An introduction to Support Vector Machine

Liyu Gong

Department of Statistics  
University of Kentucky

02/15/2018

# Outline

- 1 Concepts
- 2 Perceptron
- 3 Support Vector Machine for Classification
- 4 Support Vector Machine for Regression
- 5 Summary

# Machine Learning is An Optimization Problem

- Structural model: candidate decision function set
- Error model: criterion to be optimized
- Algorithm: computable way to do the optimization

# Linear Regression: An Example

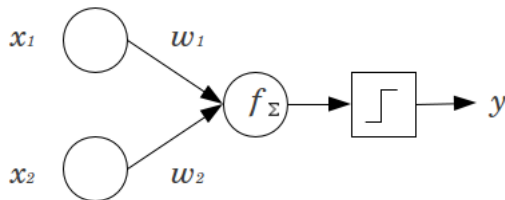
- Structural model:  $y(x) = \beta^T x$
- Error model:  $\sum (t_i - y(x_i))^2$
- Algorithm: closed form solution or gradient descent

# Perceptron: structural model

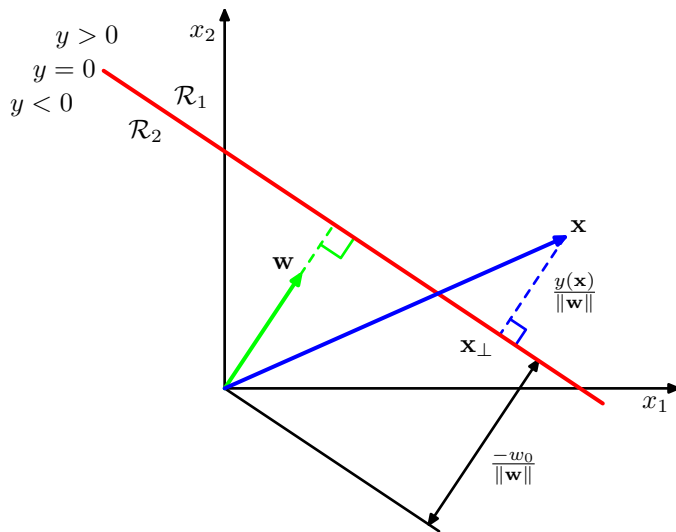
$$y(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b)$$

where

$$\text{sgn}(x) = \begin{cases} 1 & : x \geq 0 \\ -1 & : x < 0 \end{cases}$$



# Linear Discriminant Plane



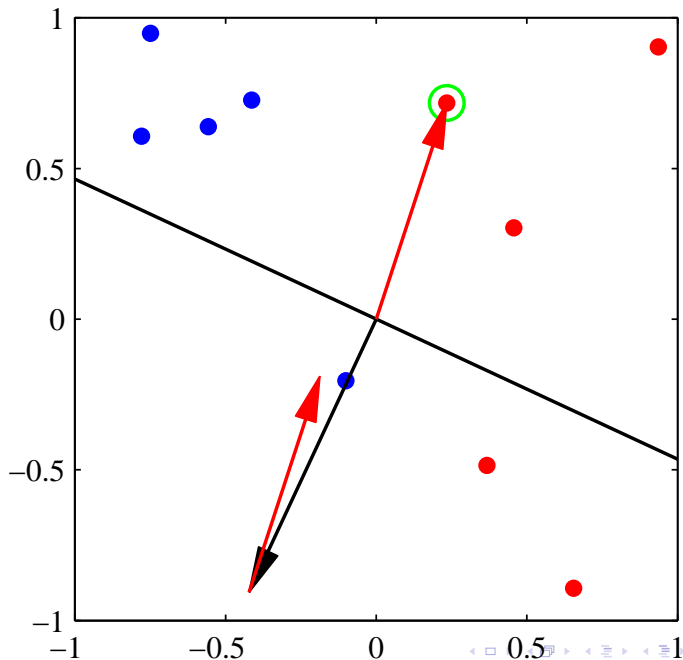
# Error Model and Algorithm

- perceptron criterion

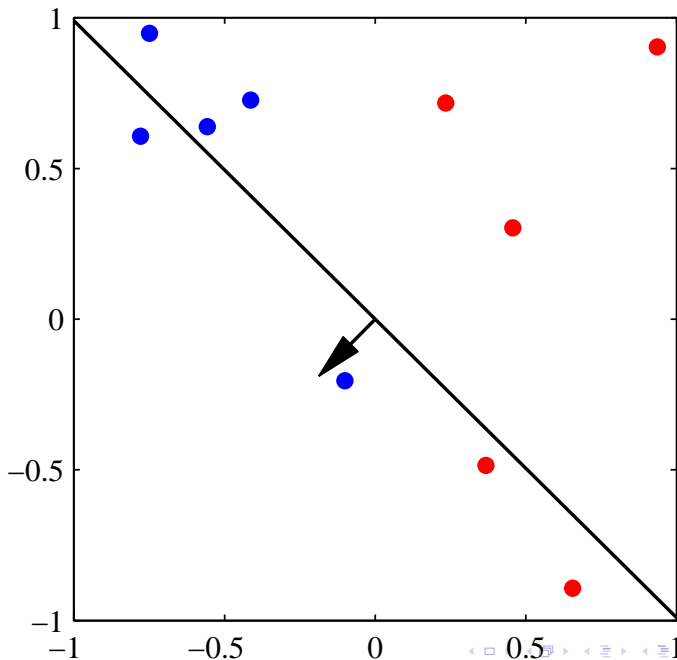
$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n t_n$$

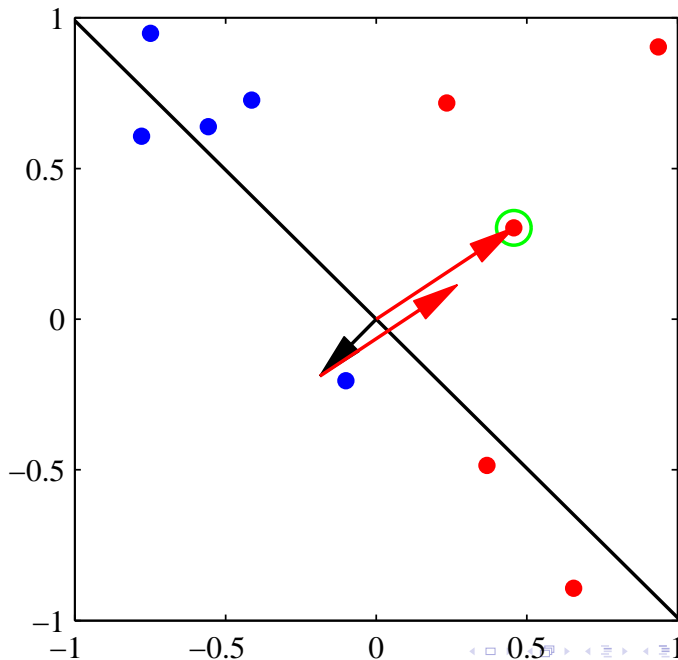
- stochastic gradient descent

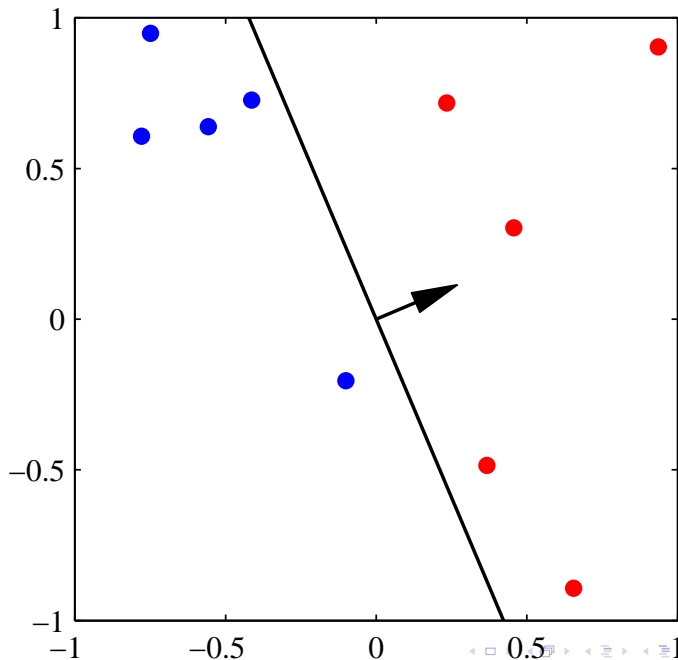
$$\begin{aligned} \tilde{\mathbf{w}}^{\tau+1} &= \mathbf{w}^{\tau} - \eta \nabla E_P(\mathbf{w}) \\ &= \mathbf{w}^{\tau} + \eta \mathbf{x}_n t_n \end{aligned}$$









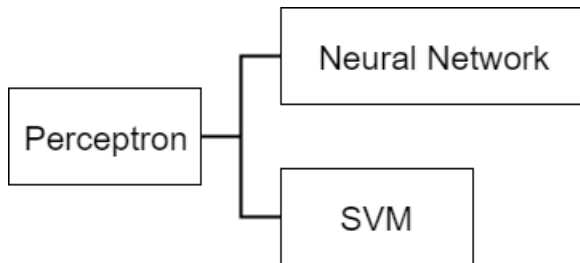


# Drawbacks and Insights

- Converges in finite steps only if there exist an exact solution
- It is a linear model
- There are multiple solutions, which one is selected depends on the initial value of and the order of training set
- If  $\tilde{w}$  is initialized as a linear combination of  $\tilde{x}_n$ , then the resultant  $\tilde{w}$  is also a linear combination of  $\tilde{x}$ , which means the final solution depends on  $x_n$  through their inner product only.
- Only a subset of the training data contributes to the decision function



# Overcome the Problems of Perceptron



- Problem 1: Non-Linearity
  - Neural Network: multiple layer perceptron
  - SVM: Kernels
- Problem 2: No unique solution
  - Neural Network: I don't care ...
  - SVM: maximum margin classification hyperplane

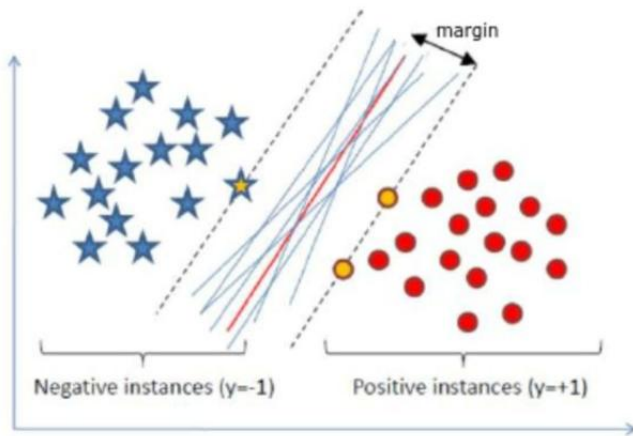


Figure: Multiple possible solutions, which one is better?

# Maximum Margin Decision Hyperplane

- Recall the distance of a point to the hyperplane is

$$\frac{t_n y(x_n)}{\|w\|} = \frac{t_n(w^T x_n + b)}{\|w\|}$$

- To find maximum margin hyperplane, we just need to find the solution for

$$\arg \max_{w,b} \left\{ \frac{1}{\|w\|} \min_n \left[ t_n(w^T x_n + b) \right] \right\}$$



# MMD: quadratic programming

Note that when we scaling  $\mathbf{w}$  and  $b$  simultaneously, the distance is unchanged. Therefore, the previous problem is equivalent to

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

w.r.t

$$t_n(\mathbf{w}^T \mathbf{x} + b) \geq 1, \quad n = 1, \dots, N.$$

## Equivalent Error Model

$$\sum_{n=1}^N E_{\infty}(y(\mathbf{x}_n)t_n - 1) + \lambda \|\mathbf{w}\|^2$$

where

$$E_{\infty}(z) = \begin{cases} 0 & z \geq 0 \\ \infty & 0 \end{cases}$$

# MMDH: Lagrange Multipliers

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \left[ t_n (\mathbf{w}^T \mathbf{x}_n + b) - 1 \right]$$

take the derivative with respect to  $\mathbf{w}$  and  $b$ , and set them to zero, we have

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$

$$0 = \sum_{n=1}^N a_n t_n$$

## KKT conditions

$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0$$

$$a_n \{ t_n y(\mathbf{x}_n) - 1 \} = 0$$

# MMDH: Dual Representation

Eliminating  $\mathbf{w}$  and  $b$  from  $L(\mathbf{w}, b, \mathbf{a})$ , the problem becomes to maximize

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \mathbf{x}_n^T \mathbf{x}_m$$

with respect to  $\mathbf{a}$  subject to the constraints

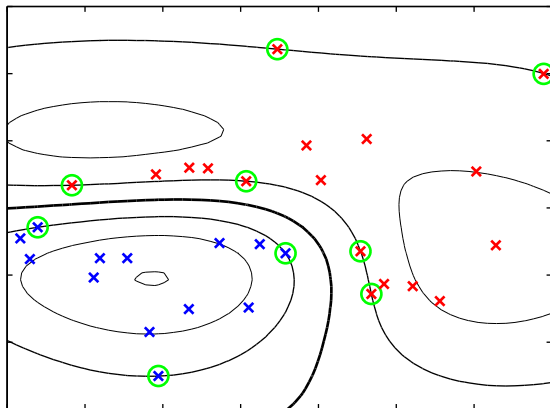
$$\begin{aligned} a_n &\geq 0, & n &= 1, \dots, N, \\ \sum_{n=1}^N a_n t_n &= 0. \end{aligned}$$

After solve the problem, the predict decision function could be written as

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n \mathbf{x}^T \mathbf{x}_n + b$$

- The problem depends on  $x_n$  through the inner product  $x_n^T x_m$
- If we can calculate the inner product of another feature space, we do not need to transform the feature actually
- Kernel function: compute the inner product of some transformed feature space using the current feature space
- The decision function will become  $y(x) = \sum_{n=1}^N a_n k(x, x_n)$
- Requirements for kernel function
  - symmetry
  - positive semi-definite
- Methods to find kernel
  - Construct a function, prove its a kernel
  - Combine existing kernels

# Overfitting

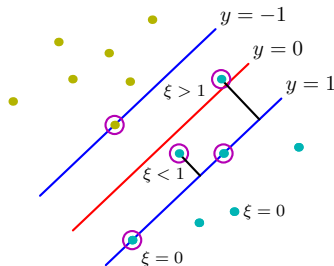


# hinge error function

$$\sum_{n=1}^N [1 - y_n t_n]_+ + \lambda \|w\|^2$$

- By not penalize misclassified points with  $\infty$ , we allow some points to be misclassified by the decision plane
- We also allow some points to fall into the decision margin

# Soft Margin



We minimize

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

with subject to

$$t_n y(x_n) \geq 1 - \xi_n, \quad n = 1, \dots, N$$

# Soft Margin: Dual Representation

We minimize

$$\tilde{L}(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(x_n, x_m)$$

with subject to

$$0 \leq a_n \leq C$$

$$\sum_{n=1}^N a_n t_n = 0$$



# Multiple Classes Classification

- No essential solution for multiple classes
- Needs high level strategies
- One-versus-the-Rest
  - Use  $K$  binary classifiers for  $K$  classes problem
- One-versus-One (pairwise)
  - Use  $\frac{K(K-1)}{2}$  binary classifiers for  $K$  classes problem

# SVR: Error Model

For regularized least square, we minimize

$$\frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

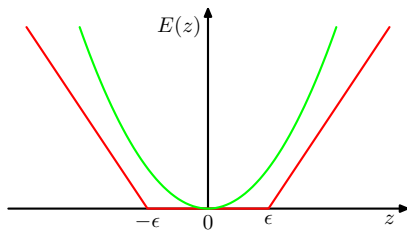
We replace the quadratic term with an  $\epsilon$ -insensitive error function

$$E_{\epsilon}(y(\mathbf{x}) - t) = \begin{cases} 0 & \text{if } |y(\mathbf{x}) - t| < \epsilon \\ |y(\mathbf{x}) - t| - \epsilon & \text{otherwise} \end{cases}$$

Then we minimize

$$E_{\epsilon}(y(\mathbf{x}_n) - t_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

# SVR: Error Model Illustration



# SVR: Optimization Problem

By introducing  $\xi_n \geq 0$  and  $\hat{\xi} \geq 0$ , we can reformulate the problem as minimize

$$C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|w\|$$

with subject to

$$t_n \leq y(\mathbf{x}_n) + \epsilon + \xi_n$$

$$t_n \leq y(\mathbf{x}_n) - \epsilon - \hat{\xi}_n$$

# SVR: Dual Representation

By manipulation using Lagrange multipliers, we can get the dual representation as to minimize

$$\begin{aligned}\tilde{L}(\mathbf{a}, \hat{\mathbf{a}}) = & -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) \\ & - \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n\end{aligned}$$

with subject to

$$0 \leq a_n \leq C$$

$$0 \leq \hat{a}_n \leq C$$

- Error model of SVC: hinge function
- Error model of SVR:  $\epsilon$ -insensitive error function
- Kernel tricks to extend a linear model to nonlinear one
- SVM is a kind of sparse kernel methods
- Well formulated optimization problem, thus can guarantee to find the global optimal
- practical tips
  - Always try to normalize  $x$
  - The parameter  $C$ ,  $\gamma$  and  $\epsilon$  should be sampled in log-space if we need to do a grid search

# References

- Bishop: Pattern recognition and Machine Learning
- Vapnik: Statistical Learning Theory
- Vapnik: The Nature of Statistical Learning Theory



Figure: Vladimir Vapnik