

CZ4041/SC4000: Machine Learning

Additional Notes: Nonseparable & Kernel SVMs

Li Boyang, Albert
School of Computer Science and Engineering,
NTU, Singapore

Acknowledgements: some figures are adapted from the lecture notes of the book “Introduction to Data Mining” (Chap. 5). Slides are modified from the version prepared by Dr. Sinno Pan.

Linear SVMs: Separable Case

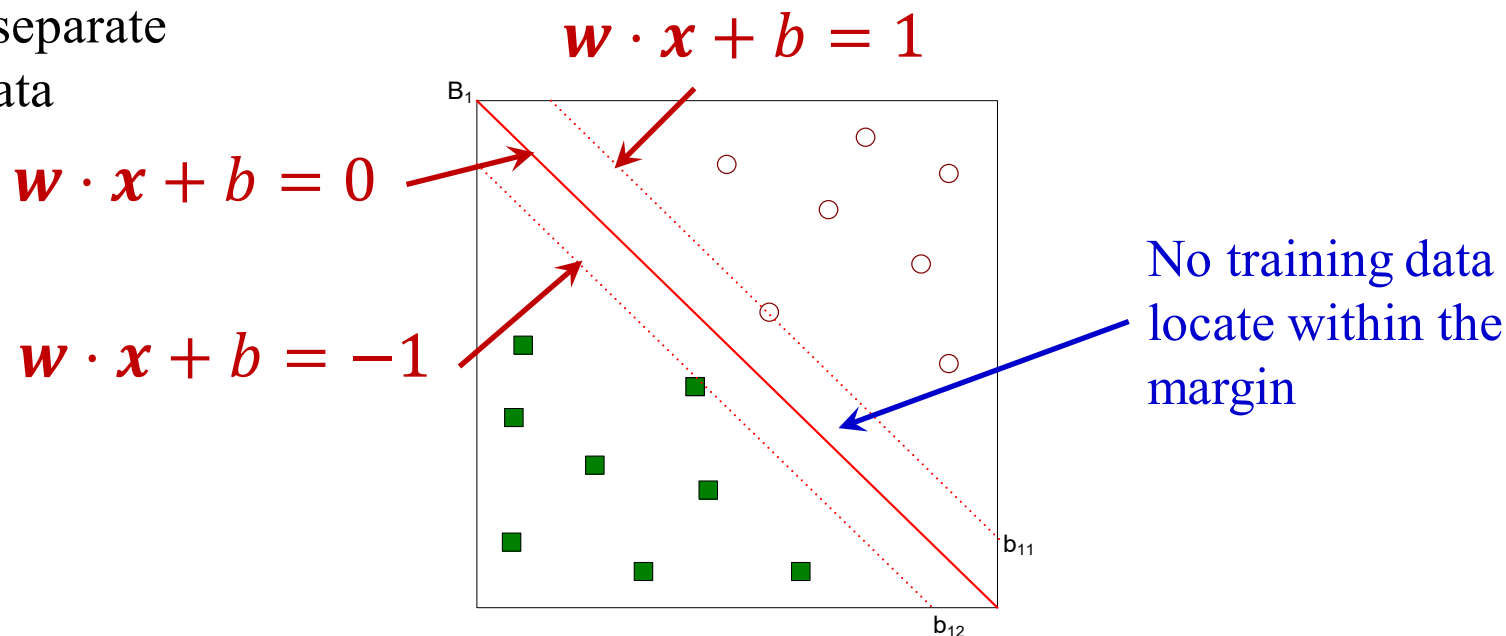
- Optimization problem of linear SVM

$$\min_{w,b} \frac{\|w\|_2^2}{2}$$

Maximize margin

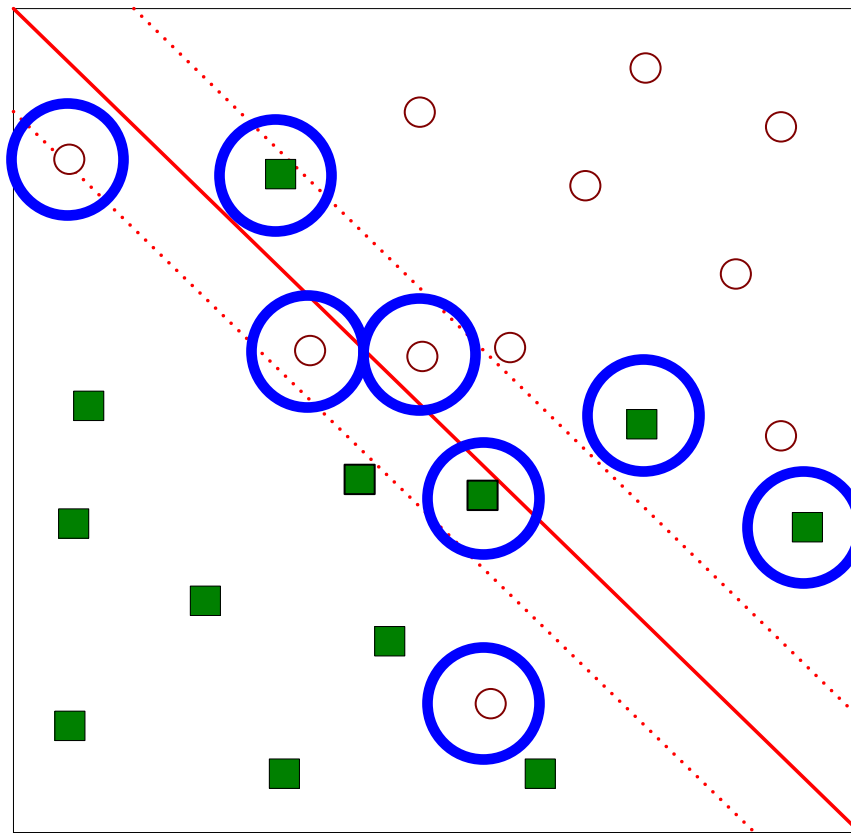
$$\text{s.t. } y_i \times (w \cdot x_i + b) \geq 1, i = 1, \dots, N$$

Perfectly separate
training data



Linear SVMs: Nonseparable Case

- What if the problem is not separable?



Slack variables
 $\xi_i \geq 0$ need to
be introduced to
absorb errors

Slack Variables

- For Separable Case:

$$\mathbf{w} \cdot \mathbf{x}_i + b \geq 1, \text{ if } y_i = 1$$

$$\mathbf{w} \cdot \mathbf{x}_i + b \leq -1, \text{ if } y_i = -1$$

OR $y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

- For Nonseparable Case:

$$\mathbf{w} \cdot \mathbf{x}_i + b \geq 1 - \xi_i, \text{ if } y_i = 1$$

$$\mathbf{w} \cdot \mathbf{x}_i + b \leq -1 + \xi_i, \text{ if } y_i = -1$$

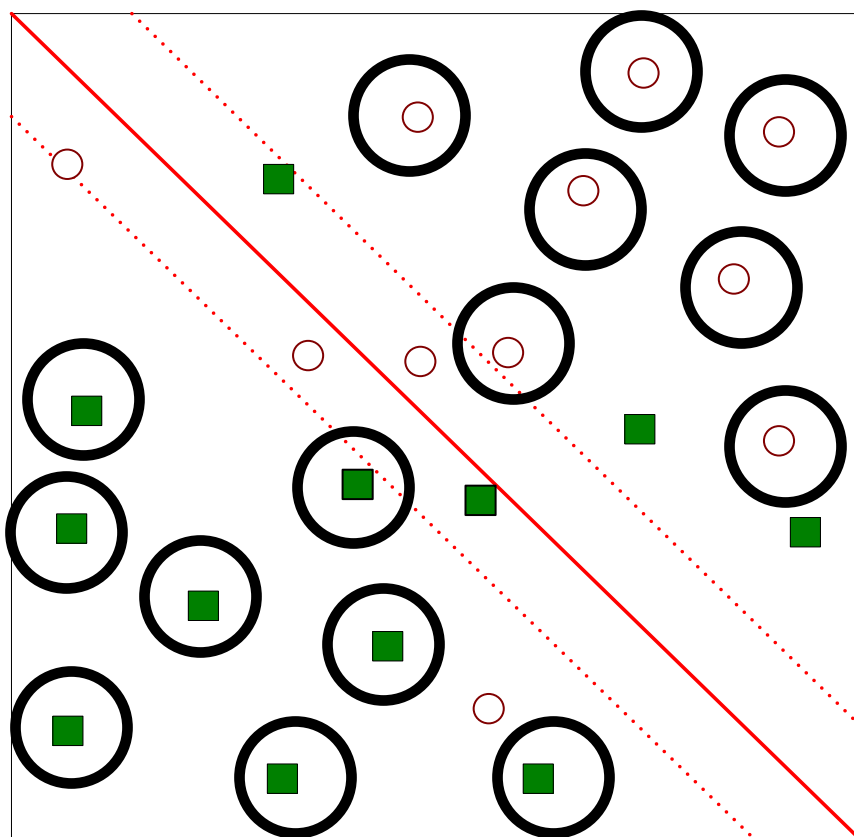
$$\xi_i \geq 0$$

Slack variables

OR

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

Slack Variables (cont.)



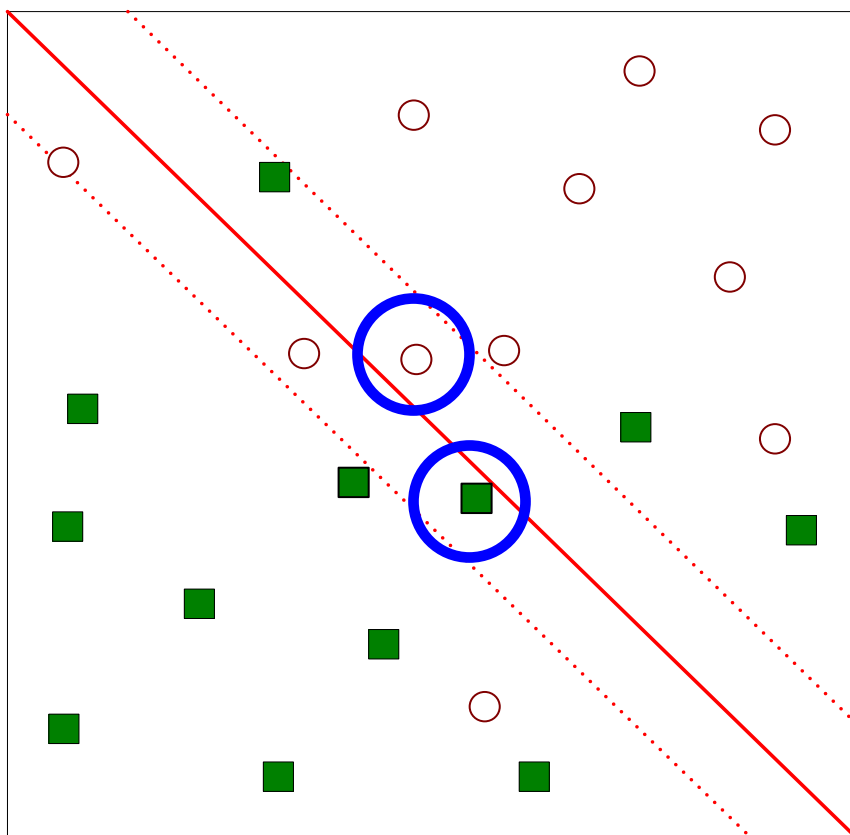
If $\xi_i = 0$, there is
no problem with \mathbf{x}_i

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$



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

Slack Variables (cont.)



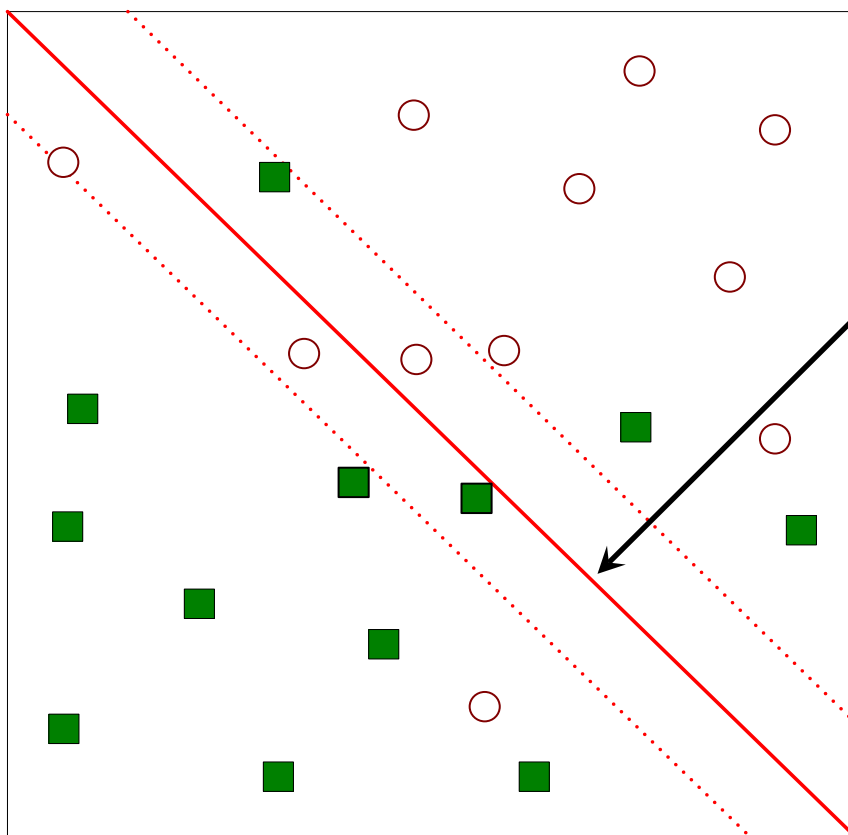
If $0 < \xi_i < 1$, \mathbf{x}_i is correctly classified but in the margin

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$



$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq k$$
$$0 < k < 1$$

Slack Variables (cont.)



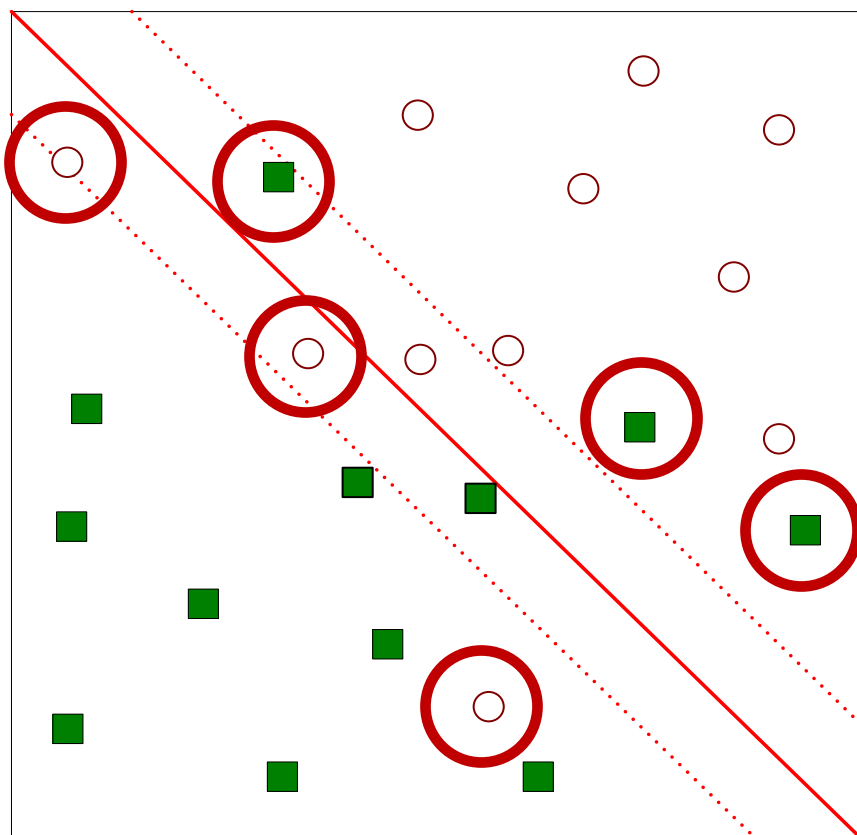
If $\xi_i = 1$, \mathbf{x}_i is on the decision boundary (random guess)

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$



$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 0$$

Slack Variables (cont.)



If $\xi_i > 1$, \mathbf{x}_i is misclassified

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$



$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq -k$$

$k > 0$



Can be negative

$$\underline{\mathbf{w} \cdot \mathbf{x}_i + b} \geq -k, \text{ if } y_i = 1$$

$$\underline{\mathbf{w} \cdot \mathbf{x}_i + b} \leq k, \text{ if } y_i = -1$$

Soft Error

- The number of misclassifications is $\#\{\xi_i > 1\}$
- The number of nonseparable points is $\#\{\xi_i > 0\}$
- Soft errors:

$$\sum_i \xi_i$$

Linear SVMs: Nonseparable Case

- Linear SVMs with soft errors:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i} \quad & \frac{\|\mathbf{w}\|_2^2}{2} + C \left(\sum_{i=1}^N \xi_i \right) \\ \text{s.t.} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, N, \\ & \xi_i \geq 0 \end{aligned}$$

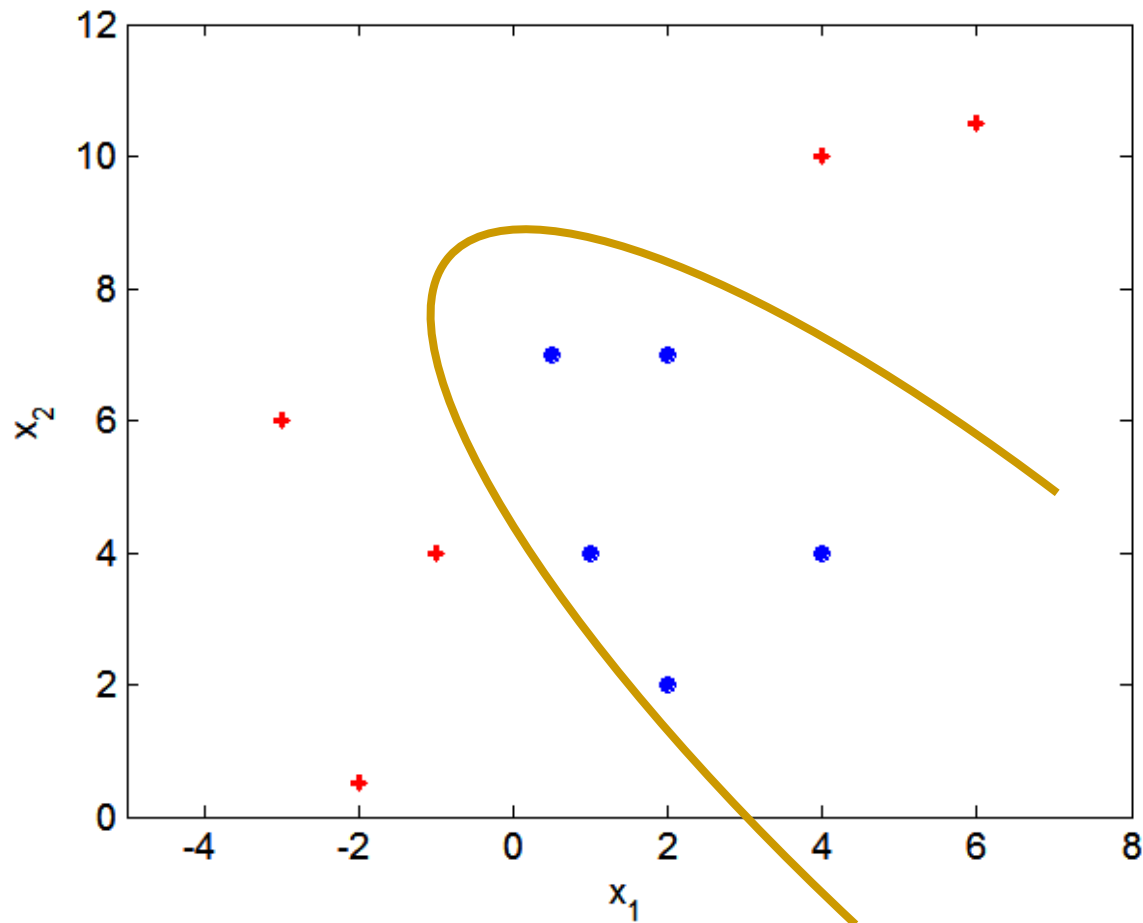
Penalize the decision boundary with large values of slack variables

$C \geq 0$ is a parameter to tradeoff the impact of margin maximization and tolerable errors

Nonnegative ξ_i provides an estimate of the error of the decision boundary on the training example \mathbf{x}_i

Nonlinear SVMs

- What if decision boundary is not linear?



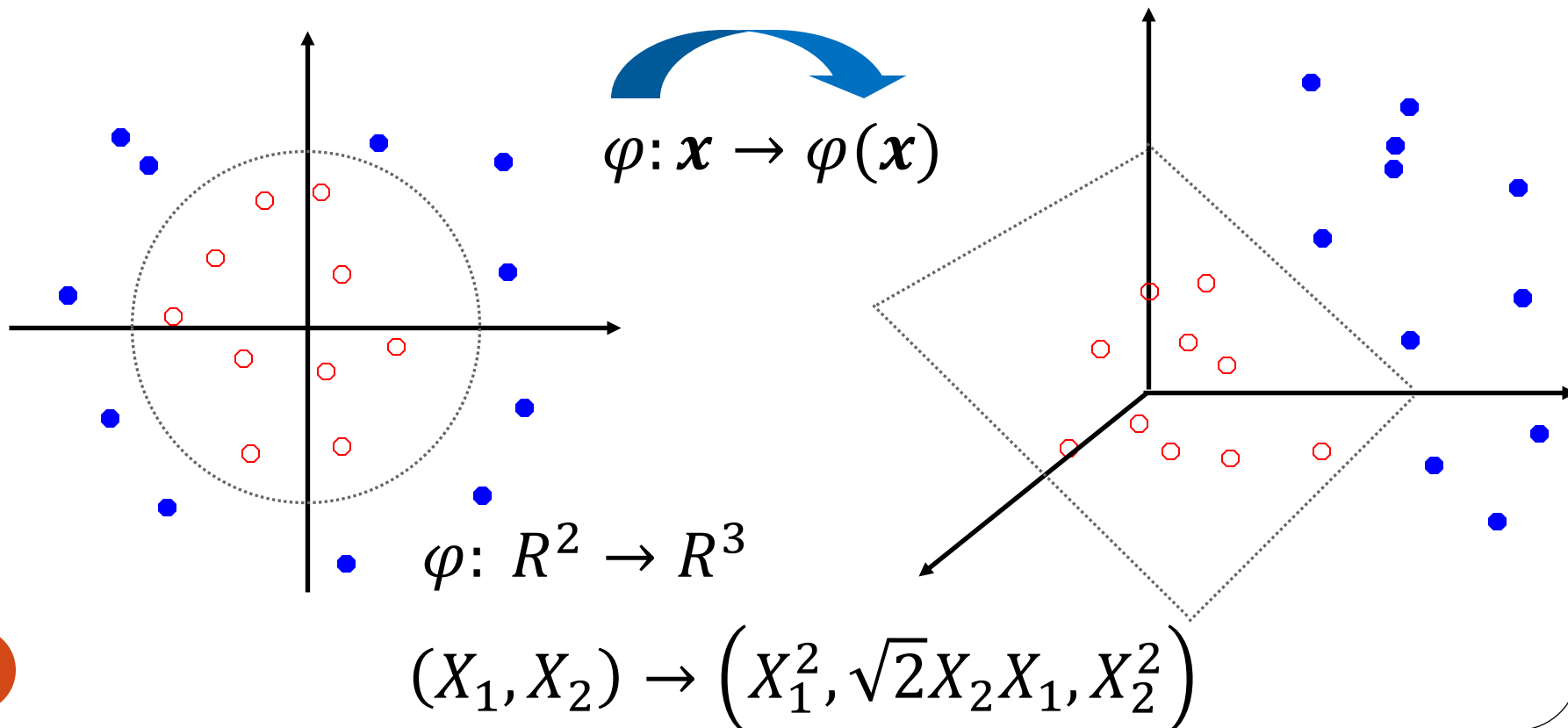
Kernel trick in
the dual form

Nonlinear SVMs (cont.)

- How to generalize linear decision boundary to become nonlinear?
- Key idea: transform \mathbf{x}_i to a higher dimensional space to “make life easier”
 - Input space: the space the point \mathbf{x}_i are located
 - Feature space: the space of $\varphi(\mathbf{x}_i)$ after transformation
- Assumption: in a higher dimensional space, it is easier to find a linear hyperplane to classify data

Feature Mapping

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



Nonlinear SVMs (cont.)

- Optimization problem of nonlinear SVMs

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{\|\mathbf{w}\|_2^2}{2} \\ \text{s.t.} \quad & y_i \times (\mathbf{w} \cdot \varphi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, \dots, N \\ & \mathbf{w} \cdot \varphi(\mathbf{x}_i) + b = 0 \quad \text{Hyperplane in feature space} \end{aligned}$$

- Computation in the feature space can be costly because it is high dimensional
 - The feature space is typically very high dimensional!
- The kernel trick comes to rescue

Nonlinear SVM: Kernel Trick

- Suppose $\varphi(\cdot)$ is given as follows, mapping an instance from 2-dimensional space to 6-dimensional space:

$$\varphi([X_1, X_2]) = [1, \sqrt{2}X_1, \sqrt{2}X_2, X_1^2, X_2^2, \sqrt{2}X_1X_2]$$

- Given two data instances: $\mathbf{a} = [A_1, A_2]$ and $\mathbf{b} = [B_1, B_2]$

$$\varphi(\mathbf{a}) = [1, \sqrt{2}A_1, \sqrt{2}A_2, A_1^2, A_2^2, \sqrt{2}A_1A_2]$$

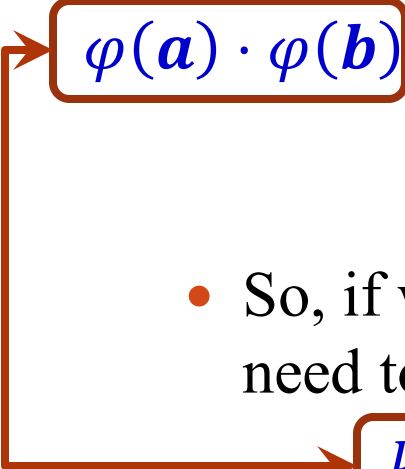
$$\varphi(\mathbf{b}) = [1, \sqrt{2}B_1, \sqrt{2}B_2, B_1^2, B_2^2, \sqrt{2}B_1B_2]$$

- Inner product of the two instances after feature mapping:

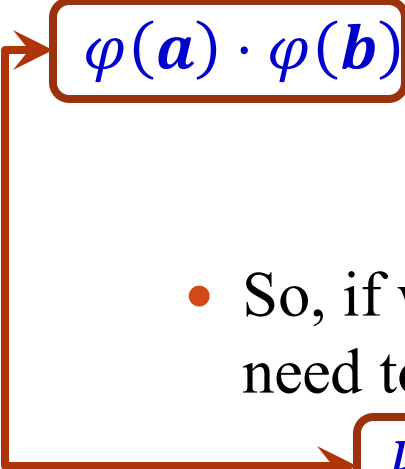
$$\begin{aligned}\varphi(\mathbf{a}) \cdot \varphi(\mathbf{b}) &= 1 + 2A_1B_1 + 2A_2B_2 + A_1^2B_1^2 + A_2^2B_2^2 + 2A_1A_2B_1B_2 \\ &= (1 + A_1B_1 + A_2B_2)^2\end{aligned}$$

Kernel Trick (cont.)

- Inner product of the two instances after feature mapping:


$$\begin{aligned}\varphi(\mathbf{a}) \cdot \varphi(\mathbf{b}) &= 1 + 2A_1B_1 + 2A_2B_2 + A_1^2B_1^2 + A_2^2B_2^2 + 2A_1A_2B_1B_2 \\ &= (1 + A_1B_1 + A_2B_2)^2\end{aligned}$$

- So, if we define the kernel function as follows, there is no need to carry out $\varphi(\cdot)$ explicitly

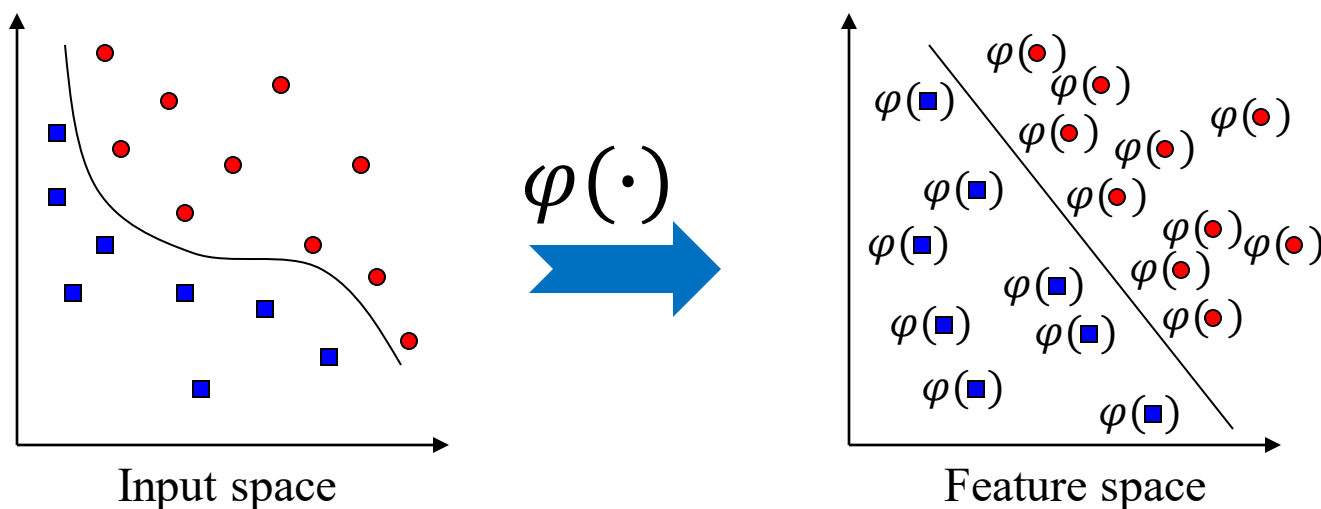

$$k(\mathbf{a}, \mathbf{b}) = (1 + A_1B_1 + A_2B_2)^2 = (1 + \mathbf{a} \cdot \mathbf{b})^2$$

- This use of kernel function to avoid carrying out $\varphi(\cdot)$ explicitly is known as the kernel trick

Kernel Trick: General Idea

- If $\varphi(\cdot)$ satisfies some conditions, then we can find a function $k(\cdot, \cdot)$ such that

Kernel function \rightarrow $k(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j)$




$$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}\right)$$

How to Apply Kernel Trick?

- Optimization problem for nonlinear SVMs (separable)

$$\begin{array}{ll} \min_{\mathbf{w}, b} & \frac{\|\mathbf{w}\|_2^2}{2} \\ \text{s.t.} & y_i \times (\mathbf{w} \cdot \boxed{\varphi(\mathbf{x}_i)} + b) \geq 1, \quad i = 1, \dots, N \end{array}$$



Primal Form

Instances in the feature space do not appear in the form of inner products

- The kernel trick is not applicable
- How about its dual form?

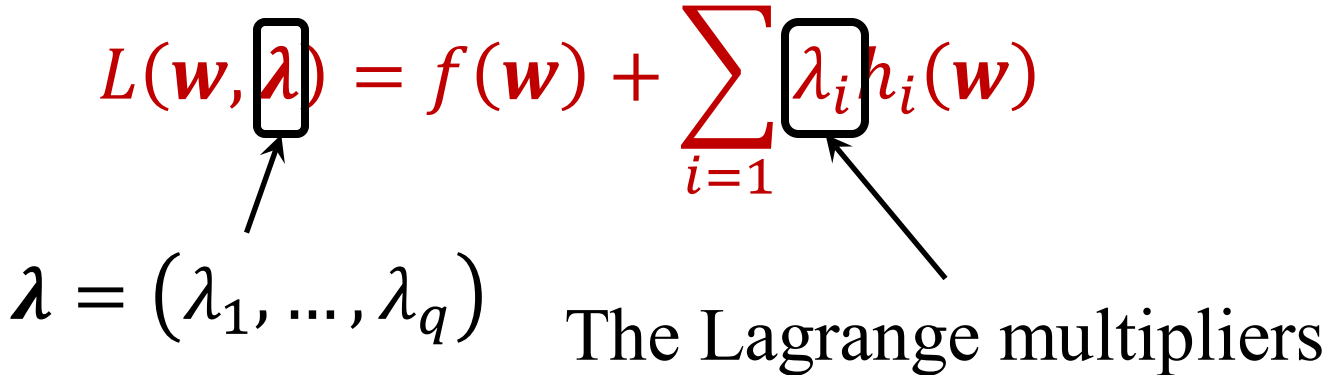
Lagrange multiplier method

Lagrange Multiplier Method: Idea

- Given: an objective $f(\mathbf{w})$ to be minimized, with a set of inequality constraints to be satisfied $h_i(\mathbf{w}) \leq 0, i = 1, 2, \dots, q$

$$\begin{array}{ll} \min_{\mathbf{w}} & f(\mathbf{w}) \\ \text{s.t.} & h_i(\mathbf{w}) \leq 0, i = 1, \dots, q \end{array}$$

- The Lagrangian for the optimization problem:

$$L(\mathbf{w}, \lambda) = f(\mathbf{w}) + \sum_{i=1}^q \lambda_i h_i(\mathbf{w})$$


$$\lambda = (\lambda_1, \dots, \lambda_q)$$

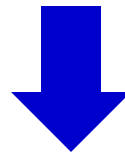
The Lagrange multipliers

The Dual Form (Separable)

- By using Lagrangian Multiplier method

$$\begin{array}{ll} \min_{w,b} & \frac{\|w\|_2^2}{2} \\ \text{s.t.} & y_i \times (w \cdot \varphi(x_i) + b) \geq 1, \quad i = 1, \dots, N \end{array}$$

Primal Form



$$\max_{\lambda} L_D(\lambda) = - \left(\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j (\varphi(x_i) \cdot \varphi(x_j)) - \sum_{i=1}^N \lambda_i \right)$$

Dual Form

Dual Optimization Problem

- The dual Lagrangian involves only the Lagrange multipliers and the training data
- The negative sign in the dual Lagrangian transforms a minimization problem of the primal form to a maximization problem of the dual form
- The objective is to maximize $L_D(\lambda)$
 - Can be solved using numerical techniques such as quadratic programming

Dual Optimization Problem (cont.)

- Once the λ_i 's are found, we can obtain the feasible solutions for \mathbf{w} and b from

$$\mathbf{w} = \sum_{i=1}^N \lambda_i y_i \varphi(\mathbf{x}_i) \quad \text{AND} \quad \lambda_i (y_i (\mathbf{w} \cdot \varphi(\mathbf{x}_i) + b) - 1) = 0$$

- The decision boundary can be expressed as

$$\mathbf{w} \cdot \varphi(\mathbf{x}) + b = \left(\sum_{i=1}^N \lambda_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}) \right) + b = 0$$

If \mathbf{x}_i is a support vector, then the corresponding $\lambda_i > 0$, otherwise, $\lambda_i = 0$

Dual Optimization Problem (cont.)

- For a test instance \mathbf{x}^* , it can be classified using

$$f(\mathbf{x}^*) = \text{sign} \left(\sum_{i=1}^N \lambda_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}^*) + b \right)$$

Nonlinear SVM via Kernel Trick

Training: $\max_{\lambda} \left(\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j (\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)) \right)$

Decision boundary: $\sum_{i=1}^N \lambda_i y_i (\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}^*)) + b = 0$

- The data points only appear as inner product
- As long as the inner product in the feature space can be calculated, no need for the explicit mapping

Nonlinear SVM via Kernel Trick (cont.)

- Replace inner product in feature space by kernel function

Training: $\max_{\lambda} \left(\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \right)$

Decision boundary: $\sum_{i=1}^N \lambda_i y_i k(\mathbf{x}_i, \mathbf{x}^*) + b = 0$

If \mathbf{x}_i is a support vector,
then the corresponding
 $\lambda_i > 0$, otherwise, $\lambda_i = 0$

$$k(\mathbf{x}_i, \mathbf{x}^*) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}^*)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$$

Kernel Functions: Examples

- Linear kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$$

- Radial basis function kernel with width σ

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

- Polynomial kernel with degree d

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j + 1)^d$$

Soft Margin Dual Form

- By using Lagrangian Multiplier method

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i} \quad & \frac{\|\mathbf{w}\|_2^2}{2} + C \left(\sum_{i=1}^N \xi_i \right) \\ \text{s.t.} \quad & y_i(\mathbf{w} \cdot \varphi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad i = 1, \dots, N, \\ & \xi_i \geq 0, \quad i = 1, \dots, N, \end{aligned}$$



$$\min_{\lambda} L_D(\lambda) = - \left(\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \boxed{\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)} - \sum_{i=1}^N \lambda_i \right)$$

$$\text{s.t., } 0 \leq \lambda_i \leq C$$

Kernel trick can be applied

Soft Margin Dual Form (cont.)

$$\mathbf{w} \cdot \varphi(\mathbf{x}) + b = \left(\sum_{i=1}^N \lambda_i y_i \underbrace{\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x})}_{k(\mathbf{x}_i, \mathbf{x})} \right) + b = 0$$

Kernel trick can be applied

- For a test instance \mathbf{x}^* , it can be classified using

$$f(\mathbf{x}^*) = \text{sign} \left(\sum_{i=1}^N \lambda_i y_i \underbrace{\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}^*)}_{k(\mathbf{x}_i, \mathbf{x}^*)} + b \right)$$

Popular Toolboxes of SVMs

- LIBSVM
 - <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- LIBLINEAR
 - <http://www.csie.ntu.edu.tw/~cjlin/liblinear/>
- SVM-light
 - <http://svmlight.joachims.org/>
- SVM-struct
 - http://www.cs.cornell.edu/People/tj/svm_light/svm_struct.html
- SVM-perf
 - http://www.cs.cornell.edu/People/tj/svm_light/svm_perf.html
- SVM-rank
 - http://www.cs.cornell.edu/People/tj/svm_light/svm_rank.html

Further Readings

- *A Tutorial on Support Vector Machines for Pattern Recognition*, by Christopher J. C. Burges, DMKD, 1998
- *Convex Optimization*, by Stephen Boyd and Lieven Vandenberghe, Cambridge University Press, 2004
- *Learning with Kernel*, by Bernhard Scholkopf and Alex Smola, The MIT Press, 2002
- *Statistical Learning Theory*, by Vladimir N. Vapnik, Wiley-Interscience, 1998

Thank you!