

01.112/50.007 Machine Learning

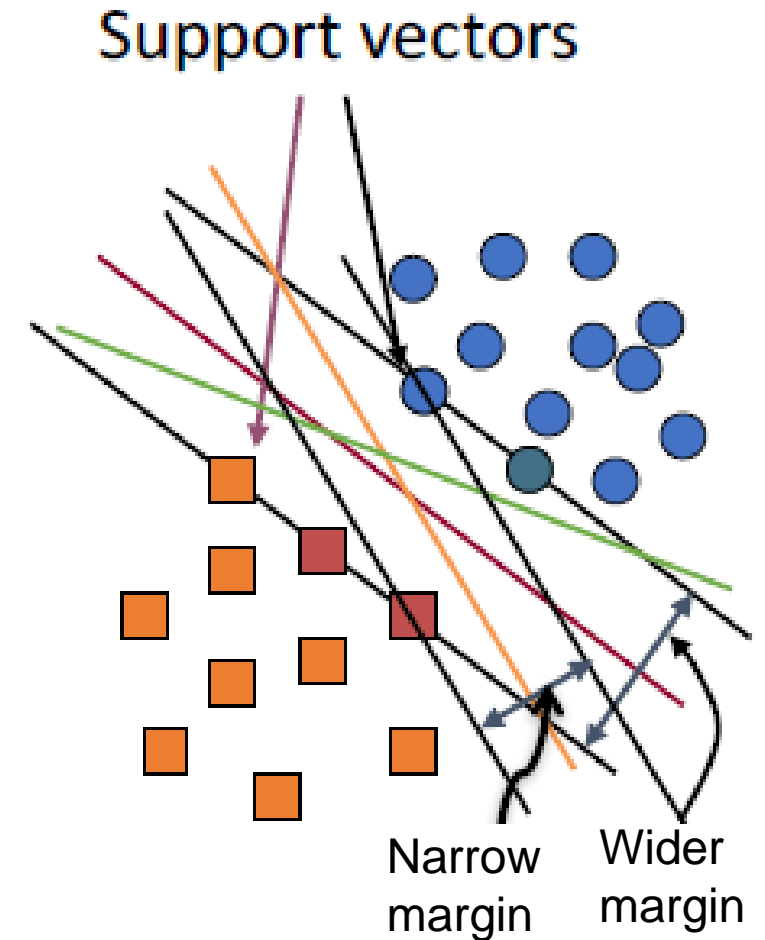
Lecture 8

Support Vector Machines (Part 2)

Recap

Support Vector Machine (SVM)

- SVMs **maximize the *margin*** around the separating hyperplane. A.k.a. **large margin classifiers**.
- The decision function is fully specified by a subset of training samples, ***the support vectors***.
- Solving SVMs is a ***quadratic programming problem***.
- Seen by many as the most successful current text classification method*



*but other discriminative methods often perform very similarly

Support Vector Machine (SVM)

- Distance of each point from decision boundary

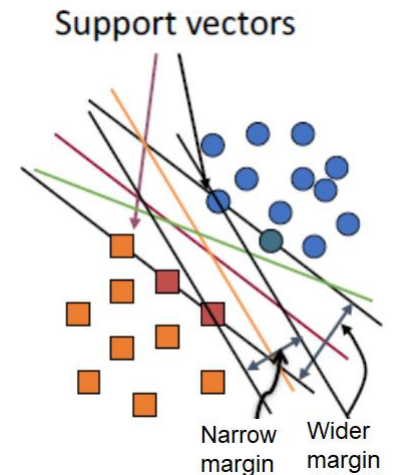
$$\gamma^{(t)}(\theta, \theta_0) = \frac{y^{(t)}(\theta \cdot x^{(t)} + \theta_0)}{\|\theta\|}$$

- Goal: Maximize minimum distance to the boundary

$$\min_{t=1, \dots, n} \gamma^{(t)}(\theta, \theta_0)$$

- Formulate the goal as **quadratic programming problem (SVM)**

$$\min \frac{1}{2} \|\theta\|^2 \text{ subject to } y^{(t)}(\theta \cdot x^{(t)} + \theta_0) \geq 1, t = 1, \dots, n$$



Constrained Optimization

Want to minimize some function $f(x)$, but there are some *constraints* on the values of x .

Method 1 (Dual Problem)

Solve a *dual optimization problem* where the constraints are nicer, and where it is easier to implement gradient descent.

Method 2 (Exact Solution)

Solve the *Lagrangian* system of equations.

Equality Constraints

Problem.

$$\begin{array}{ll}\text{minimize} & f(x) \\ \text{subject to} & h_1(x) = 0, \dots, h_l(x) = 0\end{array}$$

Lagrangian.

$$L(x, \lambda) = f(x) + \lambda_1 h_1(x) + \dots + \lambda_l h_l(x)$$

Example.

$$\begin{array}{ll}\text{minimize} & f(x) = n_1 \log x_1 + \dots + n_d \log x_d \\ \text{subject to} & h(x) = x_1 + \dots + x_d - 1 = 0\end{array}$$

$$L(x, \lambda) = n_1 \log x_1 + \dots + n_d \log x_d + \lambda(x_1 + \dots + x_d - 1)$$

Two-Player Game

$$L(x, \lambda) = f(x) + \lambda_1 h_1(x) + \cdots + \lambda_l h_l(x)$$

Rules.

- You get to choose the value of x .
Your goal is to minimize $L(x, \lambda)$.
- Your adversary gets to choose the value of λ .
His goal is to maximize $L(x, \lambda)$.

Primal Game

$$L(x, \lambda) = f(x) + \lambda_1 h_1(x) + \cdots + \lambda_l h_l(x)$$

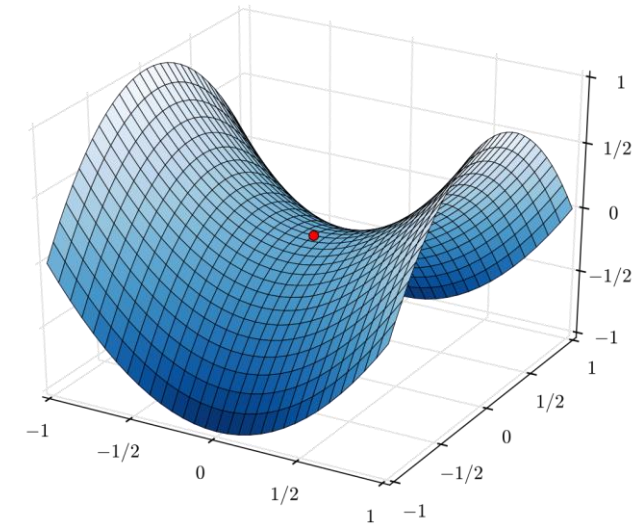
Primal Game. You go first.

Your Strategy.

- Ensure that $h_1(x) = 0, \dots, h_l(x) = 0$.
- Find x that minimizes $f(x)$.

Final Score. $p^* = \min_x \max_{\lambda} L(x, \lambda)$

The optimal x^*, λ^* are
saddle points of $L(x, \lambda)$.



Dual Game

$$L(x, \lambda) = f(x) + \lambda_1 h_1(x) + \cdots + \lambda_l h_l(x)$$

Dual Game. You go second.

Adversary's Strategy.

- For each λ , compute $\ell(\lambda) = \min_x L(x, \lambda)$
- Find λ that maximizes $\ell(\lambda)$.

Final Score. $d^* = \max_{\lambda} \min_x L(x, \lambda)$

Max-Min Inequality

Primal. $p^* = \min_x \max_{\lambda} L(x, \lambda)$

Dual. $d^* = \max_{\lambda} \min_x L(x, \lambda)$

“you do better if you
have the last say”

Weak Duality

$$\begin{aligned} p^* &= \min_x \max_{\lambda} L(x, \lambda) \\ &\geq \max_{\lambda} \min_x L(x, \lambda) = d^* \end{aligned}$$

If $p^* = d^*$, we can solve the primal by solving the dual.

**Strong
duality**

Max-Min Inequality

Example.

| | $x = 1$ | $x = 2$ |
|---------------|---------|---------|
| $\lambda = 1$ | ① | ④ |
| $\lambda = 2$ | ③ | ② |

Primal. $p^* = \min_x \max_{\lambda} L(x, \lambda) = \textcircled{3}$

Dual. $d^* = \max_{\lambda} \min_x L(x, \lambda) = \textcircled{2}$

Exact Solution

Problem.

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & h_1(x) = 0, \dots, h_l(x) = 0 \end{array}$$

Lagrange multipliers.

1. Write down the Lagrangian.

$$L(x, \lambda) = f(x) + \lambda_1 h_1(x) + \dots + \lambda_l h_l(x)$$

2. Solve for critical points x, λ .

$$\nabla_x L(x, \lambda) = 0, \quad h_1(x) = 0, \dots, h_l(x) = 0$$

3. Pick critical point which gives global minimum.

Example

$$\begin{array}{ll}\text{minimize} & f(x) = n_1 \log x_1 + \cdots + n_d \log x_d \\ \text{subject to} & h(x) = x_1 + \cdots + x_d - 1 = 0\end{array}$$

Lagrangian

$$L(x, \lambda) = n_1 \log x_1 + \cdots + n_d \log x_d + \lambda(x_1 + \cdots + x_d - 1)$$

Critical points

$$0 = n_i/x_i + \lambda$$

$$x_i = n_i/(-\lambda)$$

$$0 = x_1 + \cdots + x_d - 1$$

$$(-\lambda) = n_1 + \cdots + n_d$$

Inequality Constraints (Primal-Dual)

Primal Problem.

$$\begin{array}{ll}\text{minimize} & f(x) \\ \text{subject to} & g_1(x) \leq 0, \dots, g_m(x) \leq 0\end{array}$$

Lagrangian.

$$L(x, \alpha) = f(x) + \alpha_1 g_1(x) + \dots + \alpha_m g_m(x)$$

Dual Problem.

$$\begin{array}{ll}\text{maximize} & \ell(\alpha) \\ \text{subject to} & \alpha_1 \geq 0, \dots, \alpha_m \geq 0\end{array} \quad \text{where } \ell(\alpha) = \min_{x \in \mathbb{R}^d} L(x, \alpha)$$

Box constraints are
easier to work with!

Inequality Constraints (Exact Solution)

$$\begin{array}{ll}\text{minimize} & f(x) \\ \text{subject to} & g_1(x) \leq 0, \dots, g_m(x) \leq 0\end{array}$$

Lagrangian.

$$L(x, \alpha) = f(x) + \alpha_1 g_1(x) + \dots + \alpha_m g_m(x)$$

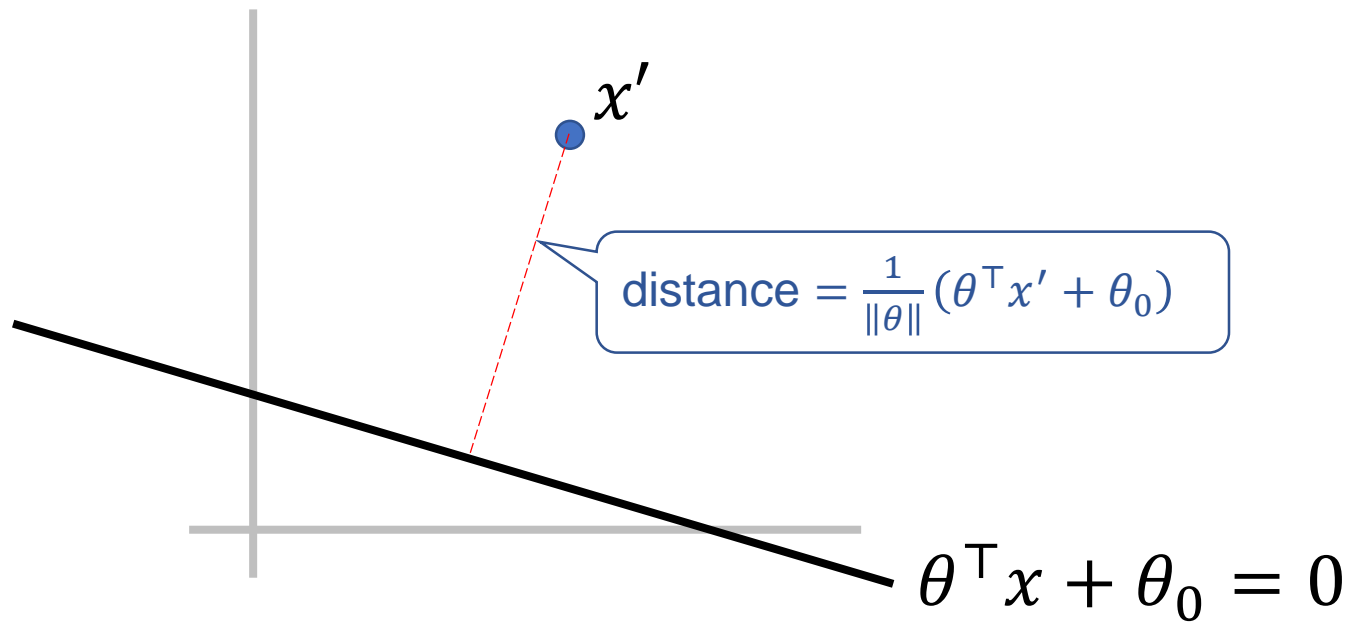
Solve for x, α satisfying

1. $\nabla_x L(x, \alpha) = 0$
2. $g_1(x) \leq 0, \dots, g_m(x) \leq 0$
3. $\alpha_1 \geq 0, \dots, \alpha_m \geq 0$
4. $\alpha_1 g_1(x) = 0, \dots, \alpha_m g_m(x) = 0$

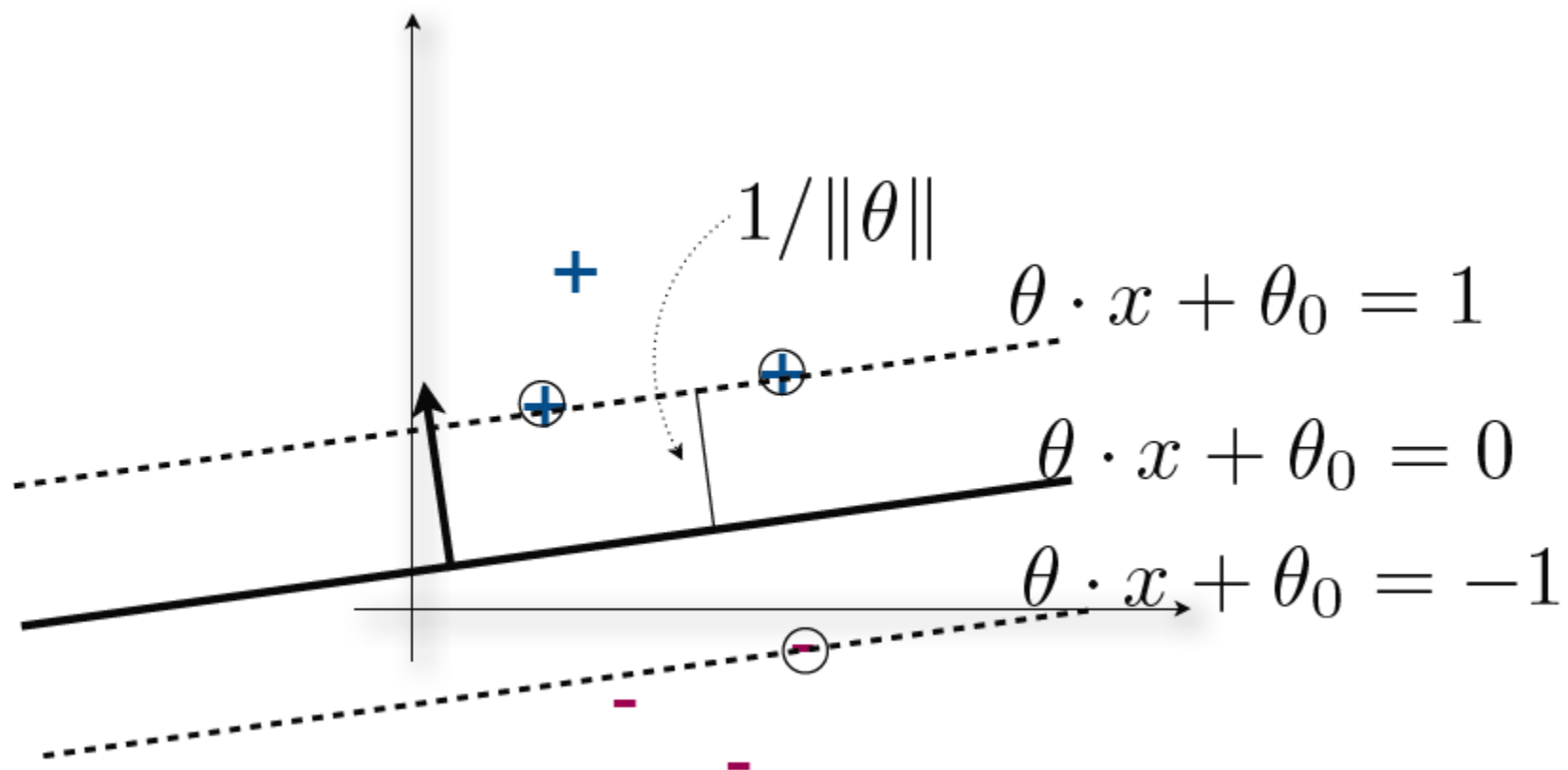
Complementary
Slackness

SVM: Maximum Margins

Computing the margin



Computing the margin



Maximum Margin

Our goal is to

$$\begin{array}{ll} \text{maximize} & 1/\|\theta\| \\ \text{subject to} & y(\theta^\top x + \theta_0) \geq 1 \text{ for all data } (x, y) \end{array}$$

Or equivalently,

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|\theta\|^2 \\ \text{subject to} & y(\theta^\top x + \theta_0) \geq 1 \text{ for all data } (x, y) \end{array}$$

Lagrangian

Primal. minimize $\frac{1}{2} \|\theta\|^2$
 subject to $y(\theta^\top x) \geq 1$ for all data (x, y)

Lagrangian. $L(\theta, \alpha) = \frac{1}{2} \|\theta\|^2 + \sum_{(x,y)} \alpha_{x,y} (1 - y(\theta^\top x))$

To find $\ell(\alpha) = \min_{\theta} L(\theta, \alpha)$, we solve

$$0 = \nabla_{\theta} L(\theta, \alpha) = \theta - \sum_{(x,y)} \alpha_{x,y} yx$$

to get $\theta = \sum_{(x,y)} \alpha_{x,y} yx$. Substituting into $L(\theta, \alpha)$ gives

$$\ell(\alpha) = \sum_{(x,y)} \alpha_{x,y} - \frac{1}{2} \sum_{(x,y)} \sum_{(x',y')} \alpha_{x,y} \alpha_{x',y'} y y' (x^\top x').$$

Primal-Dual

Primal.

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|\theta\|^2 \\ \text{subject to} & y(\theta^\top x) \geq 1 \text{ for all data } (x, y) \end{array}$$

It can be shown that the primal and dual problems are equivalent (*strong duality*).

Dual.

$$\begin{array}{ll} \text{maximize} & \sum_{(x,y)} \alpha_{x,y} - \frac{1}{2} \sum_{(x,y)} \sum_{(x',y')} \alpha_{x,y} \alpha_{x',y'} y y' (x^\top x') \\ \text{subject to} & \alpha_{x,y} \geq 0 \text{ for all } (x, y) \end{array}$$

After solving the dual to get the optimal $\alpha_{x,y}$'s, we obtain the optimal θ using $\theta = \sum_{(x,y)} \alpha_{x,y} yx$.

Support Vectors

Complementary Slackness.

$$\hat{\alpha}_{x,y} > 0: \quad y(\hat{\theta}^\top x) = 1$$

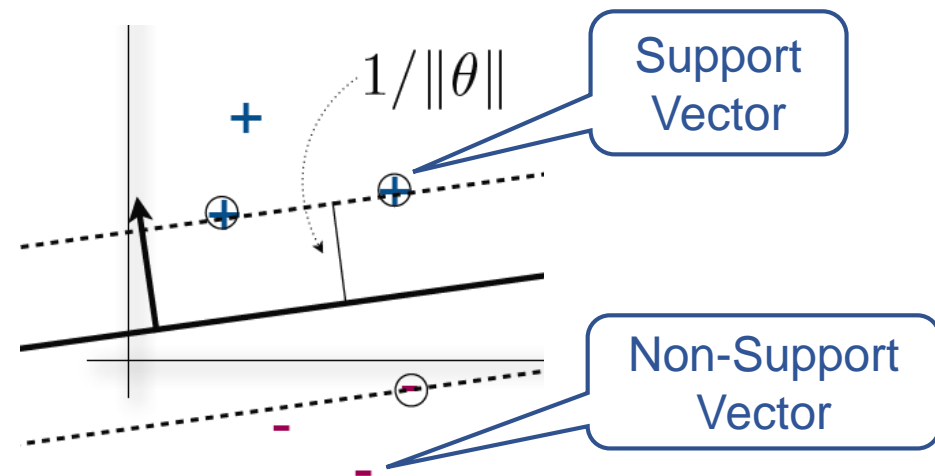
$$\hat{\alpha}_{x,y} = 0: \quad y(\hat{\theta}^\top x) > 1$$

Sparsity

Since very few data points are support vectors, **most of the $\hat{\alpha}_{x,y}$ will be zero.**

Support Vectors

Non-Support Vectors



Extensions

SVM with offset

Primal.

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|\theta\|^2 \\ \text{subject to} & y(\theta^\top x + \theta_0) \geq 1 \text{ for all data } (x, y) \end{array}$$

Dual.

$$\begin{array}{ll} \text{maximize} & \sum_{(x,y)} \alpha_{x,y} - \frac{1}{2} \sum_{(x,y)} \sum_{(x',y')} \alpha_{x,y} \alpha_{x',y'} y y' (x^\top x') \\ \text{subject to} & \alpha_{x,y} \geq 0 \text{ for all } (x, y) \\ & \sum_{(x,y)} \alpha_{x,y} y = 0 \end{array}$$

SVM with offset

Dual.

$$\begin{array}{ll}\text{maximize} & \sum_{(x,y)} \alpha_{x,y} - \frac{1}{2} \sum_{(x,y)} \sum_{(x',y')} \alpha_{x,y} \alpha_{x',y'} y y' (x^\top x') \\ \text{subject to} & \alpha_{x,y} \geq 0 \text{ for all } (x,y) \\ & \sum_{(x,y)} \alpha_{x,y} y = 0\end{array}$$

Parameters.

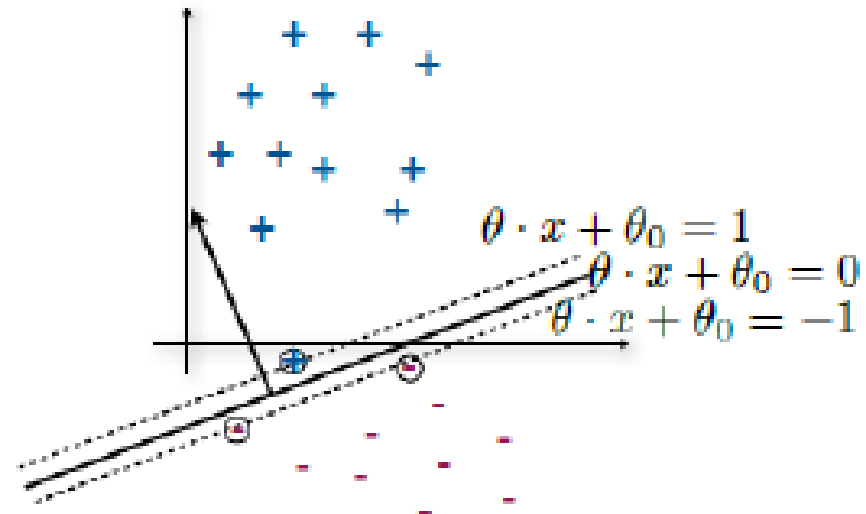
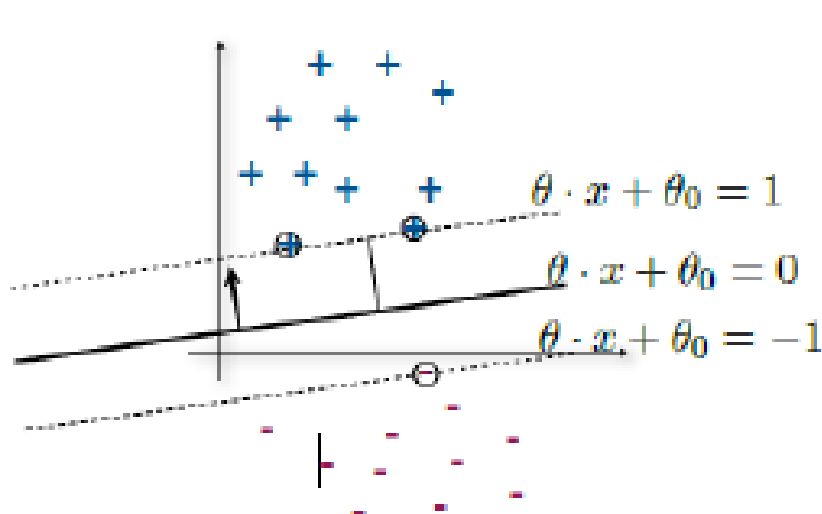
$$\begin{aligned}\hat{\theta} &= \sum_{(x,y)} \alpha_{x,y} y x \\ \hat{\theta}_0 &= y - \hat{\theta}^\top x\end{aligned}\quad \text{where } (x, y) \text{ is a support vector}$$

Derivation for $\hat{\theta}_0$

$$\begin{aligned}y^{(t)}(\hat{\theta} \cdot x^{(t)} + \hat{\theta}_0) &= y^{(t)}\left(\sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \theta_0\right) = 1 \\ \hat{\theta}_0 &= y^{(t)} - \left(\sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)})\right)\end{aligned}$$

SVM with errors

- Effect of errors in labelling training examples:



SVM with errors

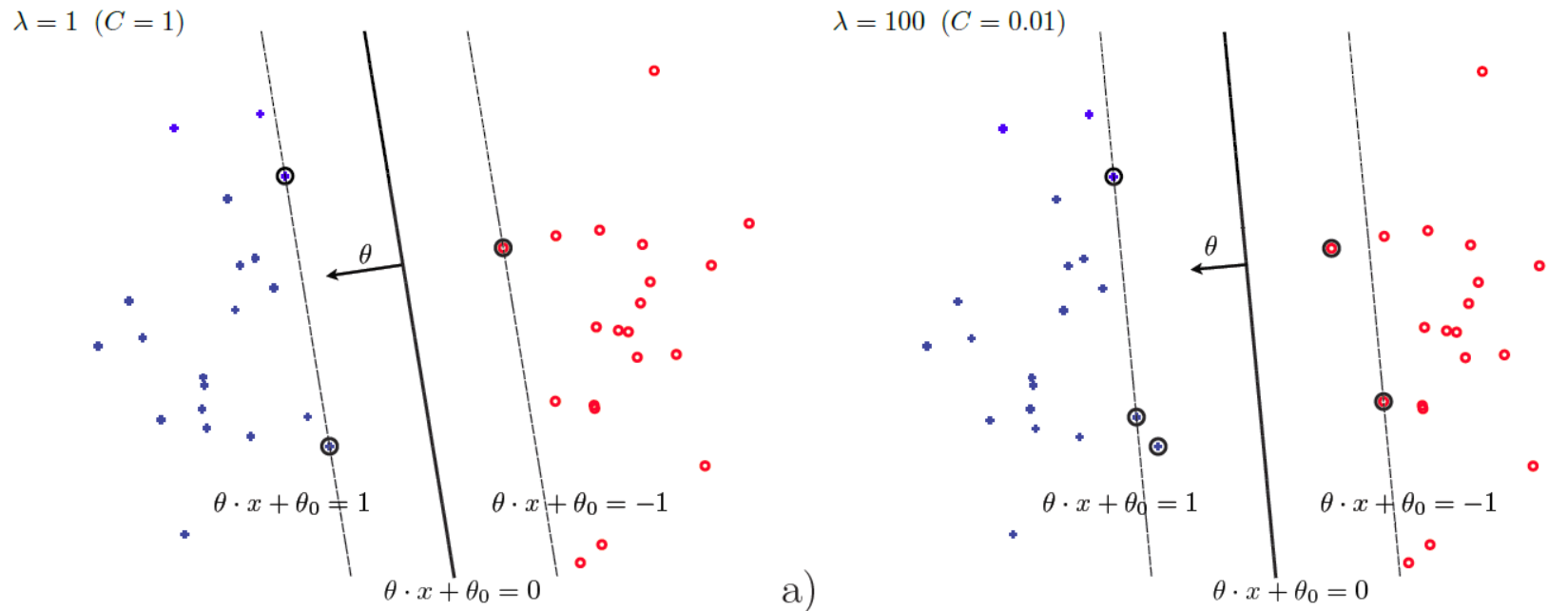
- Allowing misclassified points, yet maximize the margin
- Convert hard constraints to soft constraints
- Add a **slack variable** ξ and a **regularization parameter** λ
- **Slack variables** allow constraints to be violated for a cost.
- **Regularization** parameter balances favouring between increasing the margin and misclassifications.

Primal.

$$\begin{aligned} &\text{minimize} && \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{n} \sum_{(x,y)} \xi_{x,y} \\ &\text{subject to} && y(\theta^\top x + \theta_0) \geq 1 - \xi_{x,y} && \text{for all data } (x, y) \\ &&& \xi_{x,y} \geq 0 && \text{for all data } (x, y) \end{aligned}$$

SVM with errors

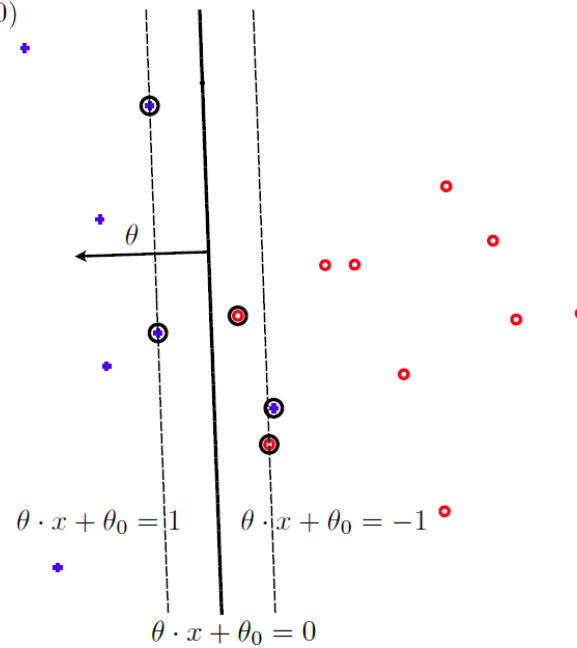
- Linearly Separable



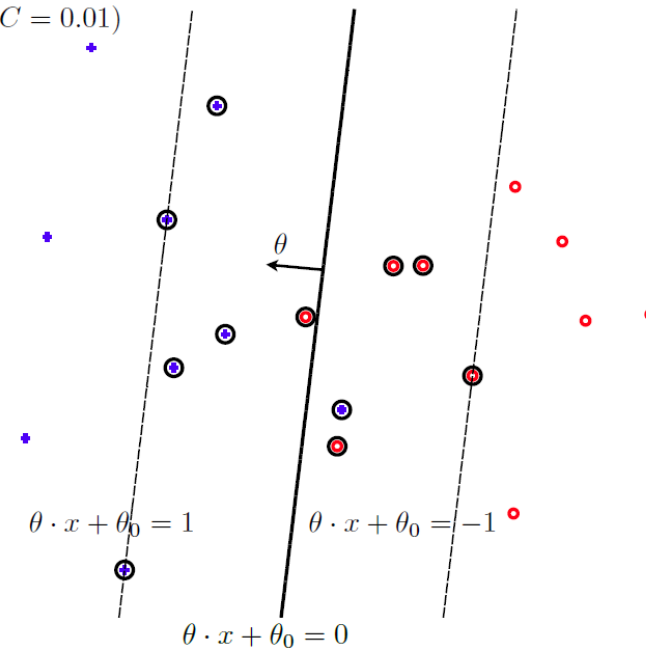
SVM with errors

- Not Linearly Separable

$\lambda = 0.1$ ($C = 10$)



$\lambda = 100$ ($C = 0.01$)



a)

SVM with errors

Primal.

$$\begin{aligned} &\text{minimize} && \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{n} \sum_{(x,y)} \xi_{x,y} \\ &\text{subject to} && y(\theta^\top x + \theta_0) \geq 1 - \xi_{x,y} \text{ for all data } (x, y) \\ &&& \xi_{x,y} \geq 0 \text{ for all data } (x, y) \end{aligned}$$

Equivalent Primal.

$$\text{minimize} \quad \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{n} \sum_{(x,y)} \text{Loss}_H(y(\theta^\top x + \theta_0))$$

Hinge-loss classifier with regularization!

SVM with errors

Primal.

$$\begin{array}{ll} \text{minimize} & \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{n} \sum_{(x,y)} \xi_{x,y} \\ \text{subject to} & y(\theta^\top x + \theta_0) \geq 1 - \xi_{x,y} \quad \text{for all data } (x, y) \\ & \xi_{x,y} \geq 0 \quad \text{for all data } (x, y) \end{array}$$

Dual.

$$\begin{array}{ll} \text{maximize} & \sum_{(x,y)} \alpha_{x,y} - \frac{1}{2} \sum_{(x,y)} \sum_{(x',y')} \alpha_{x,y} \alpha_{x',y'} y y' (x^\top x') \\ \text{subject to} & 1/\lambda \geq \alpha_{x,y} \geq 0 \text{ for all } (x, y) \\ & \sum_{(x,y)} \alpha_{x,y} y = 0 \end{array}$$

There are many efficient solvers for quadratic problems with box constraints.

Putting limits on what the adversary can do.

SVM with errors

- Complementary slackness

$$\hat{\alpha}_t = 0 \Rightarrow y^{(t)} \left(\sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \hat{\theta}_0 \right) \geq 1 \quad (\text{non-support vectors})$$

$$\hat{\alpha}_t \in (0, 1/\lambda) \Rightarrow y^{(t)} \left(\sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \hat{\theta}_0 \right) = 1 \quad (\text{SVs, on the margin})$$

$$\hat{\alpha}_t = 1/\lambda \Rightarrow y^{(t)} \left(\sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \hat{\theta}_0 \right) \leq 1 \quad (\text{SVs, margin violations})$$

Summary

- **Lagrange Multipliers**

- Lagrangian
- Primal-Dual Problems
- Inequality Constraints
- Complementary Slackness

- **Support Vector Machines**

- Maximum Margins
- Dual Problem
- Support Vectors

- **Regularization**

- Slack Variables
- Regularized Hinge Loss
- Bounded Multipliers

Intended Learning Outcomes

Extensions

- Describe the dual problem for the SVM with offset.
- Describe the primal problem for SVM with slack variables. Show that the primal is equivalent to regularized hinge loss. Explain how the regularizing parameter λ affects the margins. Describe the dual problem in terms of box constraints.