Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines

Chapter 6 Linear Models

supplementary slides to Machine Learning Fundamentals <sup>©</sup>Hui Jiang 2020 published by Cambridge University Press

August 2020



supplementary slides to Machine Learning Fundamentals <sup>©</sup>Hui Jiang 2020 published by Cambridge University Press

Perceptron 000	Linear Regression 00	Minimum Classification Error	Logistic Regression 00	Support Vector Machines

### Outline

### 1 Perceptron

- 2 Linear Regression
- 3 Minimum Classification Error
- 4 Logistic Regression
- 5 Support Vector Machines

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines

### Linear Models

- Inear models:  $\mathbf{x} \in \mathbb{R}^d \longmapsto y \in \mathbb{R}$ 
  - $\circ$  linear functions:  $y = \mathbf{w}^\intercal \mathbf{x} \quad (\mathbf{w} \in \mathbb{R}^d)$
  - affine functions:  $y = \mathbf{w}^\intercal \mathbf{x} + b$   $(\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R})$
- binary classification problems: D = {(x<sub>i</sub>, y<sub>i</sub>) | i = 1, 2, ... N} where x<sub>i</sub> ∈ ℝ<sup>d</sup> and binary label y<sub>i</sub> ∈ {+1, -1}
   distinguish linearly separable vs. non-separable cases



Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
•00				

### Perceptron

Use a linear model for 2-class problems (Rosenblatt, 1957):  $y = \operatorname{sign}(\mathbf{w}^{\mathsf{T}}\mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^{\mathsf{T}}\mathbf{x} > 0 \\ -1 & \text{otherwise} \end{cases}$ 

#### Perceptron Algorithm

initialize 
$$\mathbf{w}^{(0)} = 0$$
,  $n = 0$ 

#### loop

randomly choose a sample  $(\mathbf{x}_i, y_i)$  in  $\mathcal{D}$ calculate the actual output  $h_i = \operatorname{sign}(\mathbf{w}^{(n)^{\mathsf{T}}}\mathbf{x}_i)$ if upon a mistake:  $h_i \neq y_i$  then  $\mathbf{w}^{(n+1)} = \mathbf{w}^{(n)} + y_i \mathbf{x}_i$ n = n + 1else if no mistake is found then return  $\mathbf{w}^{(n)}$  and terminate end if end loop

Linear Regression	Logistic Regression	Support vector machines
000 00		

### Perceptron: Convergence

- assume training set *D* is *linearly separable*
- normalizing all input vectors in  $\mathcal{D}$ :  $\|\mathbf{x}_i\| \le 1 \quad \forall i = \{1, 2, \cdots, N\}$
- separation margin (scaling  $\|\hat{\mathbf{w}}\| = 1$ ):

$$\gamma = \min_{\mathbf{x}_i \in \mathcal{D}} \frac{|\hat{\mathbf{w}}^{\mathsf{T}} \mathbf{x}_i|}{\|\hat{\mathbf{w}}\|} = \min_{\mathbf{x}_i \in \mathcal{D}} |\hat{\mathbf{w}}^{\mathsf{T}} \mathbf{x}_i|$$



#### Theorem 1 (Convergence of Perceptron Algorithm)

If the perceptron algorithm is run on a linearly separable training set  $\mathcal{D}$ , the number of mistakes made is at most  $1/\gamma^2$ . In other words, the perceptron algorithm will terminate after at most  $\lceil 1/\gamma^2 \rceil$  updates and return a hyperplane that perfectly separates  $\mathcal{D}$ .

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
000				

### Perceptron: Proof Convergence

#### proof sketch:

$$\begin{array}{l} 1 \quad \text{margin definition: } |\hat{\mathbf{w}}^{\mathsf{T}}\mathbf{x}_{i}| \geq \gamma \implies y_{i}\hat{\mathbf{w}}^{\mathsf{T}}\mathbf{x}_{i} \geq \gamma \quad \left(\forall (\mathbf{x}_{i}, y_{i}) \in \mathcal{D}\right) \\ 2 \quad \text{record all } M \quad \text{mistakes: } \mathcal{M} = \{(\mathbf{x}^{(1)}, y^{(1)}), \cdots, (\mathbf{x}^{(M)}, y^{(M)})\}, \text{ then } \\ \sum_{n \in \mathcal{M}} y^{(n)} \hat{\mathbf{w}}^{\mathsf{T}}\mathbf{x}^{(n)} \geq M \cdot \gamma \\ 3 \quad \text{use Cauchy-Schwarz ineq. and } \|\hat{\mathbf{w}}\| = 1: \\ \sum_{n \in \mathcal{M}} y^{(n)} \hat{\mathbf{w}}^{\mathsf{T}}\mathbf{x}^{(n)} \leq \|\hat{\mathbf{w}}\| \cdot \left\|\sum_{n \in \mathcal{M}} y^{(n)}\mathbf{x}^{(n)}\right\| = \left\|\sum_{n \in \mathcal{M}} y^{(n)}\mathbf{x}^{(n)}\right\| \\ 4 \quad \left\|\sum_{n \in \mathcal{M}} y^{(n)}\mathbf{x}^{(n)}\right\| = \left\|\sum_{n \in \mathcal{M}} (\mathbf{w}^{(n+1)} - \mathbf{w}^{(n)})\right\| = \left\|\mathbf{w}^{(M+1)}\right\| = \\ \sqrt{\left\|\mathbf{w}^{(M+1)}\right\|^{2}} = \sqrt{\sum_{n \in \mathcal{M}} \left(\left\|\mathbf{w}^{(n+1)}\right\|^{2} - \left\|\mathbf{w}^{(n)}\right\|^{2}}\right) \\ 5 \quad \left\|\mathbf{w}^{(n+1)}\right\|^{2} - \left\|\mathbf{w}^{(n)}\right\|^{2} = \left\|\mathbf{w}^{(n)} + y^{(n)}\mathbf{x}^{(n)}\right\|^{2} - \left\|\mathbf{w}^{(n)}\right\|^{2} < \left\|\mathbf{x}^{(n)}\right\|^{2} \leq 1 \\ \left\|\mathbf{w}^{(n)}\right\|^{2} + 2y^{(n)}\mathbf{w}^{(n)^{\mathsf{T}}}\mathbf{x}^{(n)} + (y^{(n)})^{2}\left\|\mathbf{x}^{(n)}\right\|^{2} - \left\|\mathbf{w}^{(n)}\right\|^{2} < \left\|\mathbf{x}^{(n)}\right\|^{2} \leq 1 \\ 6 \quad M \cdot \gamma \leq \left\|\sum_{n \in \mathcal{M}} y^{(n)}\mathbf{x}^{(n)}\right\| < \sqrt{M} \implies M < (1/\gamma)^{2} \\ \end{array}$$

supplementary slides to  $Machine \ Learning \ Fundamentals \ ^{\odot}$  Hui Jiang 2020 published by Cambridge University Press

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
	••			

# Linear Regression (I)

- how about *linearly non-separable* cases?
- find a linear mapping  $y = \mathbf{w}^{\mathsf{T}} \mathbf{x}$  to fit all data in  $\mathcal{D}$ :

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^{\mathsf{T}} \\ \mathbf{x}_2^{\mathsf{T}} \\ \vdots \\ \mathbf{x}_N^{\mathsf{T}} \end{bmatrix}_{N \times d} \qquad \qquad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$

the least square error:

$$E(\mathbf{w}) = \sum_{i=1}^{N} \left( \mathbf{w}^{\mathsf{T}} \mathbf{x}_{i} - y_{i} \right)^{2} = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^{2}$$

linear regression:

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} E(\mathbf{w})$$

supplementary slides to Machine Learning Fundamentals <sup>©</sup> Hui Jiang 2020 published by Cambridge University Press

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
	00			

# Linear Regression (II)

compute the gradient:

$$\frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = 2\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y}$$

derive the closed-form solution by vanishing the gradient:

$$\mathbf{w}^* = \left(\mathbf{X}^\intercal \mathbf{X}
ight)^{-1} \mathbf{X}^\intercal \mathbf{y}$$

with  $\mathbf{X}^{\intercal}\mathbf{X} \in \mathbb{R}^{d imes d}$ 

- may alternatively use gradient descent to derive w\* iteratively to avoid the matrix inversion
- assign a new input x in the test stage:

$$y = \operatorname{sign}(\mathbf{w}^{*\mathsf{T}}\mathbf{x}) = \left\{ \begin{array}{l} +1 & \text{if } \mathbf{w}^{*\mathsf{T}}\mathbf{x} > 0 \\ -1 & \text{otherwise} \end{array} \right.$$

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
		••		

### Minimum Classification Error (I)

- $\blacksquare$  ERM minimizes 0-1 classification errors in  ${\cal D}$
- for each training sample  $(\mathbf{x}_i, y_i)$  in  $\mathcal{D}$

$$-y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i = \begin{cases} > 0 & \Longrightarrow \text{ mis-classification} \\ < 0 & \Longrightarrow \text{ correct classification} \end{cases}$$

count the 0-1 training errors:

$$E_0(\mathbf{w}) = \sum_{i=1}^N H(-y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)$$

• smooth approximation by sigmoid l(x):

$$E_1(\mathbf{w}) = \sum_{i=1}^N l(-y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)$$



	LITON LOgistic Regression	Support vector Machines
00		

### Minimum Classification Error (II)

minimum classification error (MCE): learn w by minimizing the sum of smooth errors E<sub>1</sub>(w)

$$\mathbf{w}_{\mathsf{MCE}} = \arg\min_{\mathbf{w}} E_1(\mathbf{w}) = \arg\min_{\mathbf{w}} \sum_{i=1}^N l(-y_i \mathbf{w}^\mathsf{T} \mathbf{x}_i)$$

- MCE does not have a closed-form solution
- rely on a gradient descent or SGD method:

$$\frac{\partial E_1(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^N y_i \, l(y_i \mathbf{w}^\mathsf{T} \mathbf{x}_i) \, \Big( 1 - l(y_i \mathbf{w}^\mathsf{T} \mathbf{x}_i) \Big) \mathbf{x}_i$$

where  $\frac{dl(x)}{dx} = \frac{e^{-x}}{(1+e^{-x})^2} = l(x)(1-l(x))$ 

supplementary slides to Machine Learning Fundamentals <sup>©</sup> Hui Jiang 2020 published by Cambridge University Press

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
			•0	

# Logistic Regression (I)

• for any  $(\mathbf{x}_i, y_i)$  in  $\mathcal{D}$ , the quantity  $l(-y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i) \in (0, 1)$ :

- MCE interprets it as a soft error to misclassify  $(\mathbf{x}_i, y_i)$
- $\circ\,$  it can be viewed as a probability to misclassify  $(\mathbf{x}_i,y_i)$
- the probability to correctly classify  $(\mathbf{x}_i, y_i)$ :  $1 - l(-u_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i) = l(u_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)$
- $\blacksquare$  logistic regression learns  ${\bf w}$  by maximizing the probability to correctly classify all samples in  ${\cal D}$

$$\mathbf{w}_{LR} = \arg \max_{\mathbf{w}} L(\mathbf{w}) = \arg \max_{\mathbf{w}} \prod_{i=1}^{N} l(y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)$$
$$= \arg \max_{\mathbf{w}} \ln L(\mathbf{w}) = \arg \max_{\mathbf{w}} \sum_{i=1}^{N} \ln l(y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)$$

Perceptron Linear	Regression Minimum C	lassification Error Logis	tic Regression Support	Vector Machines
		00		

# Logistic Regression (II)

rely on gradient descent or SGD:

$$\frac{\partial \ln L(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^{N} y_i \Big( 1 - l(y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i) \Big) \mathbf{x}_i$$

logistic regression vs. MCE

- MCE learning focuses more on the boundary cases
- logistic regression generates significant gradients for all mis-classified samples: prone to outliers but faster convergence
- can be extended to multi-class using the softmax function



Figure: Comparison of the gradient weights of MCE and logistic regression

Logistic Regression

Support Vector Machines

# Support Vector Machines (SVMs)

- linear SVMs
- soft SVMs
- nonlinear SVMs
- solving quadratic programming
- multi-class SVMs

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				0000000000000

# Linear SVMs (1): SVM0

- $\blacksquare$  assume the data  $\mathcal D$  is linearly separable
- how to find the max-margin hyperplane?
- margin:  $\gamma = \min_{\mathbf{x}_i \in \mathcal{D}_N} \frac{y_i(\mathbf{w}^{\mathsf{T}} \mathbf{x}_i + b)}{||\mathbf{w}||}$
- linear SVM  $\implies$  max margin hyperplane:

$$\{\mathbf{w}^*, b^*\} = \arg \max_{\mathbf{w}, b} \gamma$$



#### SVM0

 $\max_{\gamma, \mathbf{w}, b} \quad \gamma$ 

subject to

$$\frac{y_i(\mathbf{w}^{\intercal}\mathbf{x}_i+b)}{||\mathbf{w}||} \ge \gamma \quad \forall i \in \{1, 2, \cdots, N\}$$

Figure: max margin hyperplane



supplementary slides to Machine Learning Fundamentals <sup>©</sup> Hui Jiang 2020 published by Cambridge University Press

ion Minimum Classification Error Logistic Regression Support Vector Machin	es
00 000000000000000000000000000000000000	

# Linear SVMs (2): SVM1

- scale  $\mathbf{w}, b$  to normalize numerators  $\implies$  margin  $\gamma = \frac{1}{\|\mathbf{w}\|}$
- max margin  $\max \gamma \iff \min \|\mathbf{w}\|^2$
- an equivalent formulation of linear SVMs:

#### SVM1

$$\min_{\mathbf{w},b} \quad \frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}$$

subject to

$$y_i(\mathbf{w}^{\mathsf{T}}\mathbf{x}_i+b) \ge 1 \quad \forall i \in \{1, 2, \cdots, N\}$$



Figure: scaling  $\mathbf{w}, b$  does not change the location of the hyperplane  $\mathbf{w}^{\mathsf{T}}\mathbf{x} + b = 0$ 

・ 同 ト ・ ヨ ト ・ ヨ ト

э

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				0000000000000

# Linear SVMs (3): Lagrange Duality

• use the KKT conditions to derive the Lagrange duality 1 introduce Lagrange multipliers  $\alpha_i \ge 0$  to derive the Lagrangian

$$L(\mathbf{w}, b, \{\alpha_i\}) = \frac{1}{2}\mathbf{w}^{\mathsf{T}}\mathbf{w} + \sum_{i=1}^{N} \alpha_i \Big(1 - y_i(\mathbf{w}^{\mathsf{T}}\mathbf{x}_i + b)\Big)$$

2 the Lagrange dual function:

$$\begin{aligned} \mathcal{L}^*(\{\alpha_i\}) &= \inf_{\mathbf{w},b} \quad L(\mathbf{w},b,\{\alpha_i\}) \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^{\mathsf{T}} \mathbf{x}_j \end{aligned}$$

3 subject to the constraint:  $\sum_{i=1}^{N} \alpha_i y_i = 0$ 

supplementary slides to Machine Learning Fundamentals <sup>©</sup> Hui Jiang 2020 published by Cambridge University Press

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				00000000000000

### Linear SVMs (4): Dual Problem

■ linear SVMs are convex optimizatoin  $\implies$  strong duality ■ the unique solution, i.e.  $(\mathbf{w}^*, b^*, \{\alpha_i^*\})$ , is a saddle point

$$\underbrace{\mathsf{SVM1}}_{\mathsf{prime problem}} \iff \underbrace{\max_{\{\alpha_i\}} L^*(\alpha_i) \ \mathsf{s.t.} \sum_{i=1}^N \alpha_i y_i = 0}_{\mathsf{dual problem}}$$

solving the dual problem yields {\$\alpha\_i^\*\$}\$
linear SVMs: \$\mathbf{w}^\* = \sum\_{i=1}^N \alpha\_i^\* y\_i \mathbf{x}\_i\$ and \$b^\* = y\_i - \mathbf{w}^{\*\intercal} \mathbf{x}\_i\$

supplementary slides to Machine Learning Fundamentals <sup>©</sup> Hui Jiang 2020 published by Cambridge University Press

000 00 00 00 00 00 <b>0000000</b>	Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
					000000000000000000

# Linear SVMs (5): SVM2

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_{1} \\ \vdots \\ \alpha_{N} \end{bmatrix}_{N \times 1} \mathbf{y} = \begin{bmatrix} y_{1} \\ \vdots \\ y_{N} \end{bmatrix}_{N \times 1} \mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{N \times 1}$$
$$\mathbf{Q} = \begin{bmatrix} Q_{ij} \end{bmatrix}_{N \times N} = \begin{bmatrix} \mathbf{y} \mathbf{y}^{\mathsf{T}} \end{bmatrix}_{N \times N} \odot \begin{bmatrix} \mathbf{x}_{1}^{\mathsf{T}} \mathbf{x}_{1} & \cdots & \mathbf{x}_{1}^{\mathsf{T}} \mathbf{x}_{N} \\ \vdots & \mathbf{x}_{i}^{\mathsf{T}} \mathbf{x}_{j} & \vdots \\ \mathbf{x}_{N}^{\mathsf{T}} \mathbf{x}_{1} & \cdots & \mathbf{x}_{N}^{\mathsf{T}} \mathbf{x}_{N} \end{bmatrix}_{N \times N}$$

#### SVM2: Quadratic Programming

$$\max_{\boldsymbol{\alpha}} \ \mathbf{1}^{\mathsf{T}} \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^{\mathsf{T}} \mathbf{Q} \boldsymbol{\alpha}$$

subject to

$$\mathbf{y}^{\mathsf{T}} \boldsymbol{\alpha} = 0$$
  
 $\boldsymbol{\alpha} \ge 0$ 

200

supplementary slides to Machine Learning Fundamentals <sup>©</sup> Hui Jiang 2020 published by Cambridge University Press

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				0000000000000

### Soft SVMs

- how about linearly non-separable cases?
- allow non-negative error terms  $\xi$
- soft SVMs: aim to minimize the soft margin

### SVM3

$$\min_{\mathbf{w},b,\xi_i} \quad \frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w} + C \sum_{i=1}^{N} \xi_i$$

subject to

$$y_i(\mathbf{w}^{\mathsf{T}}\mathbf{x}_i + b) \ge 1 - \xi_i \quad \forall i \in \{1, 2, \cdots, N\}$$
$$\xi_i \ge 0 \quad \forall i \in \{1, 2, \cdots, N\}$$



soft margin = margin + linear errors

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				0000000000000

### Soft SVMs: Dual Problem

- use the same technique of Lagrangian
- derive the Lagrange dual function

### SVM4

subject to

$$\max_{\alpha} \mathbf{1}^{\mathsf{T}} \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^{\mathsf{T}} \mathbf{Q} \boldsymbol{\alpha}$$
$$\mathbf{y}^{\mathsf{T}} \boldsymbol{\alpha} = 0$$
$$0 \le \boldsymbol{\alpha} \le C$$

supplementary slides to Machine Learning Fundamentals <sup>©</sup> Hui Jiang 2020 published by Cambridge University Press

< (T) >

-

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				00000000000000

# Nonlinear SVMs

- apply a mapping function to a higher dimensional feature space:  $\mathbf{x} \longmapsto h(\mathbf{x})$
- construct linear SVMs in the feature space
- linear SVMs only depend on  $h^{\intercal}(\mathbf{x}_i)h(\mathbf{x}_j)$
- define the so-called kernel function  $\Phi(\mathbf{x}_i, \mathbf{x}_j) = h^{\mathsf{T}}(\mathbf{x}_i)h(\mathbf{x}_j)$
- choose any kernel function based on the Mercer's condition
  - $\circ~$  polynomial kernel:  $\Phi(\mathbf{x}_i,\mathbf{x}_j) = (\mathbf{x}_i^\mathsf{T}\mathbf{x}_j+1)^p$
  - Gaussian (or RBF) kernel:  $\Phi(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$



$$h(\mathbf{x})$$
 vs.  $\Phi(\mathbf{x}_i, \mathbf{x}_j)$   
 $\circ h(\mathbf{x}): \mathbb{R}^m \to \mathbb{R}^n$ 

• 
$$\Phi(\mathbf{x}_i, \mathbf{x}_j)$$
:  
 $\mathbb{R}^m + \mathbb{R}^m \to \mathbb{R}$ 

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				00000000000000

### Nonlinear SVMs: Kernel Trick

nonlinear SVMs: also solve SVM4 but using

$$\mathbf{Q} = \begin{bmatrix} Q_{ij} \end{bmatrix}_{N \times N} = \begin{bmatrix} \mathbf{y} \mathbf{y}^{\mathsf{T}} \end{bmatrix} \odot \begin{bmatrix} \Phi(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \Phi(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \Phi(\mathbf{x}_i, \mathbf{x}_j) & \vdots \\ \Phi(\mathbf{x}_N, \mathbf{x}_1) & \cdots & \Phi(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

construct the nonlinear SVM:

$$y = \operatorname{sign} \Big( \sum_{i=1}^{N} \alpha_i^* y_i \Phi(\mathbf{x}_i, \mathbf{x}) + b^* \Big)$$

 nonlinear SVMs are very powerful machine learning models

o e.g. a proper RBF kernel is used

kernel tricks in machine learning



Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				000000000000000

### Solving Quadratic Programming in SVMs

Dual problem of all SVMs is a dense quadratic programming:

$$\min_{\boldsymbol{\alpha}} \quad \underbrace{\frac{1}{2} \boldsymbol{\alpha}^{\mathsf{T}} \mathbf{Q} \boldsymbol{\alpha} - \mathbf{1}^{\mathsf{T}} \boldsymbol{\alpha}}_{L(\boldsymbol{\alpha})}$$

subject to  $\mathbf{y}^{\intercal} \boldsymbol{\alpha} = 0$  and  $0 \leq \boldsymbol{\alpha} \leq C$ , where

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_T \end{bmatrix}_{T \times 1} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_T \end{bmatrix}_{T \times 1} \quad \mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{T \times 1}$$
$$\mathbf{Q} = \begin{bmatrix} Q_{ij} \end{bmatrix}_{T \times T} = \begin{bmatrix} \mathbf{y} \mathbf{y}^\mathsf{T} \end{bmatrix}_{T \times T} \odot \begin{bmatrix} \Phi(\mathbf{x}_i, \mathbf{x}_j) \end{bmatrix}_{T \times T}$$

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				000000000000000000000000000000000000000

# Projected Gradient Descent for SVMs

#### Projected Gradient Descent Algorithm

initialize  $\alpha^{(0)} = 0$ , and set n = 0while not converged **do** 

- (1) compute the gradient:  $\nabla L(\boldsymbol{\alpha}^{(n)}) = \mathbf{Q}\boldsymbol{\alpha}^{(n)} \mathbf{1}$
- (2) project the gradient to the hyperplane  $y^{\mathsf{T}} \alpha = 0$ :

$$\tilde{\nabla}L(\boldsymbol{\alpha}^{(n)}) = \nabla L(\boldsymbol{\alpha}^{(n)}) - \frac{\mathbf{y}^{\mathsf{T}} \nabla L(\boldsymbol{\alpha}^{(n)})}{||\mathbf{y}||^2} \mathbf{y}$$

(3) projected gradient descent:  $\alpha^{(n+1)} = \alpha^{(n)} - \eta_n \cdot \tilde{\nabla} L(\alpha^{(n)})$ (4) limit  $\alpha^{(n+1)}$  to [0, C](5) n = n + 1end while

- projected gradient descent: intuitive but inefficient in memory use
- e.g. sequential minimization optimization (SMO) method

Perceptron	Linear Regression	Minimum Classification Error	Logistic Regression	Support Vector Machines
				0000000000000000

### Multi-Class SVM

- SVMs are good for binary classification problems
- how about multi-class problems?
- build multiple binary SVMs:
  - one-vs-one strategy: a binary SVM to separate each pair of classes
  - one-vs-all strategy: a binary SVM to separate each class from all others
  - use majority-voting in the test stage
- directly formulate multi-class SVMs

Perceptron 000	Linear Regression 00	Minimum Classification Error	Logistic Regression 00	Support Vector Machines

# SVM Summary

### SVM Learning Procedure (in a nutshell)

Given a training set as  $\mathcal{D} = \left\{ (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \cdots, (\mathbf{x}_N, y_N) \right\}$ 

- 1. choose a kernel function  $\Phi(\mathbf{x}_i, \mathbf{x}_j)$
- 2. build the matrices  $\mathbf{Q}$ ,  $\mathbf{y}$  and  $\mathbf{1}$  from  $\mathcal{D}$  and  $\Phi(\mathbf{x}_i, \mathbf{x}_j)$
- 3. solve the quadratic programming problem to get:

$$\pmb{lpha}^* = [ \alpha_1^*, \cdots \alpha_N^* ]^{\mathsf{T}}$$
 and  $b^*$ 

4. evaluate the learned model as:

$$y = \operatorname{sign}\left(\sum_{i=1}^{N} \alpha_i^* y_i \Phi(\mathbf{x}_i, \mathbf{x}) + b^*\right)$$