







Nebiaska Perceptron Training Rule

$$w_j^{t+1} \gets w_j^t + \Delta w_j^t \ , \ \text{where} \ \Delta w_j^t = \eta \left(r^t - y^t\right) x_j^t$$
 and

• r^t is label of training instance t

- y^t is perceptron output on training instance t
- η is small constant (e.g., 0.1) called *learning rate*

I.e., if $(r^t - y^t) > 0$ then increase w_i^t w.r.t. x_i^t , else decrease

Can prove rule will converge if training data is linearly separable and η sufficiently small

Nebiaska Where Does the Training Rule Come From?

• Consider simpler linear unit, where output

$$y^{t} = w_{0}^{t} + w_{1}^{t} x_{1}^{t} + \dots + w_{n}^{t} x_{n}^{t}$$

(i.e., no threshold)

Dutline

- For each example, want to compromise between correctiveness and conservativeness
 - Correctiveness: Tendency to improve on x^t (reduce error)
 - Conservativeness: Tendency to keep w^{t+1} close to w^t (minimize distance)
- Use *cost function* that measures both:

$$U(\mathbf{w}) = dist\left(\mathbf{w}^{t+1}, \mathbf{w}^{t}\right) + \eta \operatorname{error}\left(r^{t}, \overline{\mathbf{w}^{t+1} \cdot \mathbf{x}^{t}}\right)$$







Lincoln	
CSCE 478/878 Lecture 5: Artificial Neural Networks and Support Vector Machines Stephen Scott Introduction	Approximate with $0=2\left(w_{i}^{t+1}-w_{i}^{t} ight)-2\eta\left(r^{t}-\sum_{j=1}^{n}w_{j}^{t}x_{j}^{t} ight)$
Outline The Perceptron	which yields
Training a Perceptron Implementation Approaches Nonlinearly Separable	$w_i^{t+1} = w_i^t + \overbrace{\eta \left(r^t - y^t\right) x_i^t}^{\Delta w_i^t}$

Gradient Descent (cont'd)

・ロト・(部・・ヨ・・ヨ・ のへぐ

 x_i^t











Nebraska

Sigmoid Unit Gradient Descent

Again, use squared error for correctiveness:

$$E(\mathbf{w}^t) = \frac{1}{2} \left(r^t - y^t \right)^2$$

(folding 1/2 of correctiveness into error func)

Thus
$$\frac{\partial E}{\partial w_j^t} = \frac{\partial}{\partial w_j^t} \frac{1}{2} (r^t - y^t)^2$$

= $\frac{1}{2} 2 (r^t - y^t) \frac{\partial}{\partial w_j^t} (r^t - y^t) = (r^t - y^t) \left(-\frac{\partial y^t}{\partial w_j^t}\right)$



Nebiaska
LaterTraining Multilayer Networks
Output UnitsCSCE
478978
Lector 5
Actineat
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
Neural
N

¹This is because all other outputs are constants w.r.t. $w_{ii}^{t} \in \mathbb{R}^{+}$

CSCE 478/678 Lecture 5: Artificial Neural Vector Machines Stephen Scott Introduction Outline The Perceptron Nonlinearly Separable Problems Backgrop Signed Uet Mathewas Backgrop Signed Uet

Nebraska

Training Multilayer Networks

- How can we compute the error term for hidden layers when there is no target output r^t for these layers?
- Instead propagate back error values from output layer toward input layers, scaling with the weights
- Scaling with the weights characterizes how much of the error term each hidden unit is "responsible for"

Nebraska Linden Units (cont'd)

tephen So

Outline

'he 'erceptron

lonlinearly eparable roblems

ckprop

raining Multilay

=

The impact that w_{ji}^t has on E^t is only through net_j^t and units immediately "downstream" of *j*:

$$\frac{\partial E^{t}}{\partial w_{ji}^{t}} = \frac{\partial E^{t}}{\partial net_{j}^{t}} \frac{\partial net_{j}^{t}}{\partial w_{ji}^{t}} = x_{ji}^{t} \sum_{\substack{k \in down(i) \\ k \in down(i)}} \frac{\partial E^{t}}{\partial net_{k}^{t}} \frac{\partial net_{k}^{t}}{\partial net_{j}^{t}}$$

$$= x_{ji}^{t} \sum_{\substack{k \in down(j)}} -\delta_{k}^{t} \frac{\partial net_{k}^{t}}{\partial net_{j}^{t}} = x_{ji}^{t} \sum_{\substack{k \in down(j)}} -\delta_{k}^{t} \frac{\partial net_{k}^{t}}{\partial y_{j}} \frac{\partial y_{j}}{\partial net_{j}^{t}}$$
$$x_{ji}^{t} \sum_{\substack{k \in down(j)}} -\delta_{k}^{t} w_{kj} \frac{\partial y_{j}}{\partial net_{j}^{t}} = x_{ji}^{t} \sum_{\substack{k \in down(j)}} -\delta_{k}^{t} w_{kj} y_{j} (1 - y_{j})$$

Works for arbitrary number of hidden layers

Nebraska **Backpropagation Algorithm**

Stephen Sco Dutline rceptror online



Until termination condition satisfied do

• For each training example (r^t, \mathbf{x}^t) do Input x^t to the network and compute the outputs y^t For each output unit k

 $\delta_k^t \leftarrow y_k^t \left(1 - y_k^t\right) \left(r_k^t - y_k^t\right)$

For each hidden unit h

$$\delta_{h}^{t} \leftarrow y_{h}^{t} \left(1 - y_{h}^{t}\right) \sum_{k \in down(h)} w_{k,h}^{t} \, \delta_{k}^{t}$$

Update each network weight w^t_i

 $w_{j,i}^t \leftarrow w_{j,i}^t + \Delta w_{j,i}^t$

where



Nebraska Remarks

ntroduction Dutline -ceptron

Backpropagation Algorithm

- When to stop training? When weights don't change much, error rate sufficiently low, etc. (be aware of overfitting: use validation set)
- Cannot ensure convergence to global minimum due to myriad local minima, but tends to work well in practice (can re-run with new random weights)
- Generally training very slow (thousands of iterations), use is very fast
- Setting η : Small values slow convergence, large values might overshoot minimum, can adapt it over time

Backpropagation Algorithm Nebraska





tephen Sco

Outline

Backpropagation Algorithm Remarks

• Alternative error function: cross entropy

$$E^{t} = \sum_{k \in outputs} \left(r_{k}^{t} \ln y_{k}^{t} + \left(1 - r_{k}^{t} \right) \ln \left(1 - y_{k}^{t} \right) \right)$$

"blows up" if $r_k^t \approx 1$ and $y_k^t \approx 0$ or vice-versa (vs. squared error, which is always in [0, 1])

• Regularization: penalize large weights to make space more linear and reduce risk of overfitting:

$$E^{t} = \frac{1}{2} \sum_{k \in outputs} \left(r_{k}^{t} - y_{k}^{t} \right)^{2} + \gamma \sum_{i,j} (w_{ji}^{t})^{2}$$

Backpropagation Algorithm Nebraska Remarks (cont'd)

Representational power:

tephen So

Outline

The Perceptron

Ionlinearly eparable

ackprop

- Any boolean function can be represented with 2 layers
- Any bounded, continuous function can be represented with arbitrarily small error with 2 layers
- Any function can be represented with arbitrarily small error with 3 layers

Number of required units may be large

May not be able to find the right weights

Nebraska Hypothesis Space

- Stephen Sco troduction Dutline
- Hyp. space \mathcal{H} is set of all weight vectors (continuous vs. discrete of decision trees)
- Search via Backprop: Possible because error function and output functions are continuous & differentiable
- Inductive bias: (Roughly) smooth interpolation between data points



Support Vector Machines Nebraska Introduction

CSCE 478/878 ecture 5 Artificial Neural tworks a

Stephen Sco

ntroduction

Outline

Similar to ANNs, polynomial classifiers, and RBF networks in that it remaps inputs and then finds a hyperplane

Main difference is how it works

Features of SVMs:

- Maximization of margin
- Duality
 - Use of kernels
 - Use of problem convexity to find classifier (often without local minima)

・ロト・西ト・モー・ 一日・ うへの



• Definition assumes linear separability (more general definitions exist that do not)

Support Vector Machines



Lincoln	The	Perc	eptro	on Alg	orith	n Revi	isited (oartia	al ex	ampl	le)
CSCE 478/878 Lecture 5: Artificial Neural Networks and Support	1 1 2 3	x ₁ 4 5 6	x ₂ 1 3 3	7 +1 +1 +1	w ₁ 0.4 0.4 0.4	w ₂ 0.1 0.1 0.1	b 0.1 0.1 0.1	α 1 0 0	x'_1 4 5 6	x ⁷ /2 1 3 3	
Vector	4	2	2	-1	0.4	0.1	0.1		2	2	
machines	6	3	1	-1	0.4	0.1	0.1	l ö l	3	1	
Stephen Scott Introduction Outline The Perceptron Nonlinearly Separable Problems Backprop	1 2 3 4 5 6 1 2 3 3 4 4 5 6	4 5 6 2 2 3 4 5 6 2 2 3	1 3 1 2 1 3 3 1 2 1 2 1	+1 +1 +1 -1 -1 -1 +1 +1 +1 +1 -1 -1 -1 -1	0.1 0.4 0.4 0.2 0.2 0.2 0.2 0.2 0.2 0.2 0.2 0.2 0.2	0.1 0.1 0.0 0.0 0.0 0.0 0.0 0.0	0.1 0.1 0.1 0.0 0.0 0.0 0.0 0.0	1 0 0 1 0 0 2 0 0 0	4 5 6 2 3 4 5 6 2 3 4 5 6 2 3	1 3 1 2 1 3 1 2 1 3 1 2 1	
Dackprop	1	4	1	+1	0.4	0.0	0.0	2	4	1	
SVMs	3	6	3	+1	0.4	0.0	0.0		6	3	
Margins	4	2	1	-1	0.4	0.0	0.0	2	2	1	
Duality	5	2	2	-1	0.4	0.0	0.0	0	2	2	
Types of Kernels	6	3	1	-1	0.4	0.0	0.0	0	3	1	
svMs 35/50										< 🗆 >	4

Nebraska

t	x_1^l	x2	r	w1	w2	b	α	x_1^{\prime}	x2	r	w1	w2	b	α
1	4	1	+1	0.4	0.1	0.1	1	4	1	+1	0.4	0.0	0.0	2
2	5	3	+1	0.4	0.1	0.1	0	5	3	+1	0.4	0.0	0.0	0
3	6	3	+1	0.4	0.1	0.1	0	6	3	+1	0.4	0.0	0.0	0
4	2	1	-1	0.4	0.1	0.1	0	2	1	-1	0.2	-0.1	-0.1	3
5	2	2	-1	0.4	0.1	0.1	0	2	2	-1	0.2	-0.1	-0.1	0
6	3	1	-1	0.4	0.1	0.1	0	3	1	-1	0.2	-0.1	-0.1	0
1	4	1	+1	0.4	0.1	0.1	1	4	1	+1	0.2	-0.1	-0.1	2
2	5	3	+1	0.4	0.1	0.1	0	5	3	+1	0.2	-0.1	-0.1	0
3	6	3	+1	0.4	0.1	0.1	0	6	3	+1	0.2	-0.1	-0.1	0
4	2	1	-1	0.2	0.0	0.0	1	2	1	-1	0.0	-0.2	-0.2	4
5	2	2	-1	0.2	0.0	0.0	0	2	2	-1	0.0	-0.2	-0.2	0
6	3	1	-1	0.2	0.0	0.0	0	3	1	-1	0.0	-0.2	-0.2	0
	4	1	+1	0.2	0.0	0.0	1	4	1	+1	0.4	-0.1	-0.1	3
2	5	3	+1	0.2	0.0	0.0	0	5	3	+1	0.4	-0.1	-0.1	0
3	6	3	+1	0.2	0.0	0.0	0	6	3	+1	0.4	-0.1	-0.1	0
4	2	1	-1	0.0	-0.1	-0.1	2	2	1	-1	0.4	-0.1	-0.1	4
5	2	2	-1	0.0	-0.1	-0.1	0	2	2	-1	0.4	-0.1	-0.1	0
6	3	1	-1	0.0	-0.1	-0.1	0	3	1	-1	0.4	-0.1	-0.1	0
1	4	1	+1	0.4	0.0	0.0	2	4	1	+1	0.4	-0.1	-0.1	3
2	5	3	+1	0.4	0.0	0.0	0	5	3	+1	0.4	-0.1	-0.1	0
3	6	3	+1	0.4	0.0	0.0	0	6	3	+1	0.4	-0.1	-0.1	0
4	2	1	-1	0.4	0.0	0.0	2	2	1	-1	0.2	-0.2	-0.2	5
5	2	2	-1	0.4	0.0	0.0	0	2	2	-1	0.2	-0.2	-0.2	0
6	3	1	-1	0.4	0.0	0.0	0	3	1	-1	0.2	-0.2	-0.2	0

Lincoln	The Perceptron Algorithm Revisited (partial example)
SCE 8/878 ture 5: ificial eural orks and pport ector chines	At this point, $\mathbf{w} = (0.2, -0.2), b = -0.2, \alpha = (3, 0, 0, 5, 0, 0)$
en Scott	Can compute
uction	$w_1 = \eta(\alpha_1 r^1 x_1^1 + \alpha_4 r^4 x_1^4) = 0.1(3(1)4 + 5(-1)2) = 0.2$
	$w_2 = \eta(\alpha_1 r^1 x_2^1 + \alpha_4 r^4 x_2^4) = 0.1(3(1)1 + 5(-1)1) = -0.2$
early able ms	I.e., $\mathbf{w} = \eta \sum_{t=1}^{N} \alpha_t r^t \mathbf{x}^t$

Support Vector Machines

Nebiaska

47 Lec Ar

Nebraska

Stephen S

Dutline

'he 'erceptror

onlinearl eparable roblems

VMs 41/50

Support Vector Machines

Another way of representing predictor:

$$h(\mathbf{x}) = \operatorname{sgn} \left(\mathbf{w} \cdot \mathbf{x} + b \right) = \operatorname{sgn} \left(\eta \sum_{t=1}^{N} \left(\alpha_t \, r^t \, \mathbf{x}^t \right) \cdot \mathbf{x} + b \right)$$
$$= \operatorname{sgn} \left(\eta \sum_{t=1}^{N} \alpha_t \, r^t \, \left(\mathbf{x}^t \cdot \mathbf{x} \right) + b \right)$$

・ロト・西ト・モー・モー ひゃう

 $(\alpha_t = \# \text{ prediction mistakes on } \mathbf{x}^t)$

Support Vector Machines
Dulity (contid)CSCE
Vector
Neuron
Support
Vector
Neuron
Vector
Neuron
Support
Vector
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron
Neuron<br/





Nebraska	XOR Revisited (cont'd)	Nebraska Lincoln	Kernels
CSCE 478/878 Lecture 5: Artificial Neurai Networks and Support Vector Machines Stephen Scott Introduction Outline The Perceptron Nonlinearly Separable Problems Backprop SVMs Magne Outlar Verse	• Can easily compute the dot product $\phi(\mathbf{x}) \cdot \phi(\mathbf{z})$ (where $\mathbf{x} = [x_1, x_2]$) without first computing ϕ : $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z} + 1)^2 = (x_1z_1 + x_2z_2 + 1)^2$ $= (x_1z_1)^2 + (x_2z_2)^2 + 2x_1z_1x_2z_2 + 2x_1z_1 + 2x_2z_2 + 1$ $= \underbrace{\left[x_1^2, x_2^2, \sqrt{2}x_1 x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1\right]}_{\phi(\mathbf{x})}$ $\cdot \underbrace{\left[z_1^2, z_2^2, \sqrt{2}z_1 z_2, \sqrt{2}z_1, \sqrt{2}z_2, 1\right]}_{\phi(\mathbf{z})}$ • I.e., since we use dot products in new Perceptron algorithm, we can <i>implicitly</i> work in the remapped y space via k	CSCE 478/878 Lecture 5: Artificial Networks and Support Vector Machines Stephen Scott Introduction Outline The Perceptron Nonlinearly Separable Problems Backprop SVMs Magins Dealty Poster of Kemets	 A <i>kernel</i> is a function <i>K</i> such that ∀x, z, <i>K</i>(x, z) = φ(x) ⋅ φ(z) E.g., previous slide (quadratic kernel) In general, for degree-<i>q</i> polynomial kernel, computing (x ⋅ z + 1)^q takes ℓ multiplications + 1 exponentiation for x, z ∈ ℝ^ℓ In contrast, need over (^{ℓ+q-1}/_q) ≥ (^{ℓ+q-1}/_q)^q multiplications if compute φ first

SVMs 42/50

・ロト・日本・モー・モー ひゃう

Nebraska Lincoln Kernels (cont'd)

8

Dutline

ceptro

• Typically start with kernel and take the feature mapping that it yields

• E.g., Let
$$\ell = 1, \mathbf{x} = x, \mathbf{z} = z, K(x, z) = \sin(x - z)$$

$$\sin(x-z) = a_0 + \sum_{n=1}^{\infty} a_n \sin(nx) \sin(nz) + \sum_{n=1}^{\infty} a_n \cos(nx) \cos(nz)$$

for Fourier coeficients a_0, a_1, \ldots

• This is the dot product of two *infinite sequences* of nonlinear functions:

 $\{\phi_i(x)\}_{i=0}^{\infty} = [1, \sin(x), \cos(x), \sin(2x), \cos(2x), \ldots]$

• I.e., there are an infinite number of features in this remapped space!







Nebraska	Types of Kernels Others
CSCE 478/878 Lecture 5: Artificial Neural Networks and Support Vector Machines	Hyperbolic tangent: $K(\mathbf{x}^{t}, \mathbf{x}) = \tanh \left(2\mathbf{x}^{t} \cdot \mathbf{x} + 1 \right)$
Stephen Scott	(not a true kernel)
Introduction Outline The	Also have ones for structured data: e.g., graphs, trees, sequences, and sets of points
Perceptron Nonlinearly Separable Problems	In addition, the sum of two kernels is a kernel, the product of two kernels is a kernel
Backprop SVMs Margins Duality	Finally, note that a kernel is a <i>similarity measure</i> , useful in clustering, nearest neighbor, etc.
Kernels Types of Kernels SVMs 46 / 50	- ロ > - (2)



tephen Scc

Dutline

Support Vector Machines Finding a Hyperplane

Can show that if data linearly separable in remapped space, then get maximum margin classifier by minimizing $\mathbf{w} \cdot \mathbf{w}$ subject to $r^t (\mathbf{w} \cdot \mathbf{x}^t + b) \ge 1$

Can reformulate this in *dual form* as a *convex quadratic program* that can be solved optimally, i.e., *won't encounter local optima*:

$$\begin{array}{ll} \underset{\boldsymbol{\alpha}}{\text{maximize}} & \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \, \alpha_j \, r^i \, r^j \, K(\mathbf{x}^i, \mathbf{x}^j) \\ \text{s.t.} & \alpha_i \geq 0, i = 1, \dots, m \\ & \sum_{i=1}^{N} \alpha_i \, r^i = 0 \end{array}$$

CSCE 4798/878 Lecture 5: Autificial Neural letworks and Support Vector Machines tephen Scott $f(\mathbf{x}) = \text{sgn}\left(\sum_{i=1}^{N} \alpha_i r^t K(\mathbf{x}, \mathbf{x}^t) + b\right)$

Nebraska

troductior

Jonlinearly

eparable

ckprop

/Ms

Outline

Support Vector Machines

Finding a Hyperplane (cont'd)

(Note only need to use \mathbf{x}^t such that $\alpha_t > 0$, i.e., *support* vectors)

Can always find a kernel that will make training set linearly separable, but *beware of choosing a kernel that is too powerful* (overfitting)

· ロ > ・ 酒 > ・ 言 > ・ 言 > ・ 弓 - ・ へ D >

Nebraska

CSCE 478/878 ecture 5 Artificial Neural

Machi

Stephen Sco

Introduction Outline The Perceptron Nonlinearly Separable Problems Nackprop

SVMs Margins Duality

Support Vector Machines Finding a Hyperplane (cont'd)

If kernel doesn't separate, can *soften* the margin with *slack* variables <mark>ξ</mark>i:

$$\begin{array}{ll} \underset{\mathbf{w},b,\boldsymbol{\xi}}{\text{minimize}} & \|\mathbf{w}\|^2 + C \sum_{i=1}^{N} \boldsymbol{\xi}^i \\ \text{s.t.} & r^i((\mathbf{x}^i \cdot \mathbf{w}) + b) \ge 1 - \boldsymbol{\xi}^i, \ i = 1, \dots, N \\ & \boldsymbol{\xi}^i \ge 0, \ i = 1, \dots, N \end{array}$$

The dual is similar to that for hard margin:

$$\begin{array}{ll} \underset{\boldsymbol{\alpha}}{\text{maximize}} & \sum_{i=1}^{N} \alpha_i - \sum_{i,j} \alpha_i \, \alpha_j \, r^i \, r^j \, K(\mathbf{x}^i, \mathbf{x}^j) \\ \text{s.t.} & 0 \le \alpha_i \le C, \ i = 1, \dots, N \\ & \sum_{i=1}^{N} \alpha_i \, r^i = 0 \end{array}$$

a.

Can still solve optimally

Lincoln	Finding a Hyperplane (cont'd)
CSCE 478/878 Lecture 5: Artificial Networks and Support Vector Machines Stephen Scott	If number of training vectors is very large, may opt to approximately solve these problems to save time and space
Introduction Outline	Use e.g., gradient ascent and sequential minimal optimization (SMO)
The Perceptron Nonlinearly	When done, can throw out non-SVs
Separable Problems	
Backprop	
Margins	
Duality Kernels	
Types of Kernels	
SVMs 507/50	< ロ > < 母 > < 臣 > < 臣 > < 臣 > < 臣 > < 臣 > < 臣 > < 臣 > の Q ()

Support Vector Machine