Support Vector Machines: Max Margin, ...



R Greiner Department of Computing Science University of Alberta

Much of this is taken from Andrew W. Moore; CMU http://www.cs.cmu.edu/~awm/tutorials

+ Nello Cristianini, Ron Meir, Ron Parr, Barnabas Poczos



Error Function

<u>Jump</u>

Given data $\{[x^{(i)}, y^{(i)}]\}_{i=1..m}$, optimize...

1. Classification error
 Perceptron Training

- $err_{Class}(w) = \frac{1}{m} \sum_{i=1}^{m} I[y^{(i)} \neq sign(w^{T} x^{(i)})]$
- 2. Mean-squared error (LMS)
 Direct (Gradient Descent)

$$err_{MSE}(w) = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} [y^{(i)} - w^T x^{(i)}]^2$$

3. (Log) Conditional Probability (LR)
 MSE Gradient Descent; LCL Gradient Descent

$$LCL(w) = \frac{1}{m} \sum_{i=1}^{m} \log P_{w}(y^{(i)} | x^{(i)})$$

4. (Log) Joint Probability (LDA; FDA)
 Direct Computation

$$LL(w) = \frac{1}{m} \sum_{i=1}^{m} \log P_w(y^{(i)}, x^{(i)})$$

5. Hinge Loss (SVM)

$$HL(w) = \frac{1}{m} \sum_{i} \lambda_i \left(y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} - 1 \right)_+ +$$



A Little History

- Support Vector Machines (SVM)
 - introduced in COLT-92
 - greatly developed since then
- Now, a large and diverse community:
 - machine learning
 - optimization
 - statistics
 - neural networks
 - functional analysis, etc etc etc.
- Successful applications in many fields (bioinformatics, text, handwriting recognition, etc)
- Kernel Machines: large class of learning algs
 - SVM = a particular instance
- http://www.kernel-machines.org

Skip in class

Home

Frequently Asked Questions

Publications

Books

Software

Annual Workshop

JMLR.

Links

News

> Machine Learning Summer School / Course On The Analysis On Patterns

2007-02-12

- > New Kernel-Machines.org server 2007-01-30
- > Call for participation: The 2006 kernel workshop, "10 years of kernel machines" 2006-10-06
- Three Workshops on Kernel Methods at NIPS 2005

2005-09-08

Nips 2004 Workshop on Learning With Structured Outputs

2004-10-22

More news...

Books

Books on SVMs and Other Kernel Machines

last modified 2007-01-31 12:45

Vladimir Vapnik. <u>Estimation of Dependences Based on Empirical Data</u>. Springer Verlag, 2006, 2nd edition.

The second edition of Vapnik's classic on learning theory, including several new chapters on the history of events and on non-inductive inference.

 Grace Wahba. Spline Models for Observational Data. SIAM CBMS-NSF Regional Conference Series in Applied Mathematics vol. 59, Philadelphia, 1990.

Discusses (reproducing) kernel methods in nonparametric regression. Not easy reading for machine learning researchers, but containing fundamental material about precedents of today's kernel machines (169 pages, \$33.5).

Vladimir Vapnik. The Nature of Statistical Learning Theory. Springer, NY, 1995.

An overview of statistical learning theory, containing no proofs, but most of the crucial theorems and milestones of learning theory. With a detailed chapter on SVMs for pattern recognition and regression (1st edition: 188 pages, \$65; 2nd edition: 304 pages, \$70).

Vladimir Vapnik. Statistical Learning Theory. Wiley, NY, 1998.

The comprehensive treatment of statistical learning theory, including a large amount of material on SVMs (768 pages, \$120).

 Bernhard Schölkopf, Chris Burges, and Alex Smola (eds). <u>Advances in Kernel Methods - Support Vector Learning</u> MIT Press, Cambridge, MA, 1999.

A collection of articles written by experts in the field. Includes an introductory tutorial, overviews of the theory of SVMs, contributions on novel algorithms, and three chapters on SVM implementations (392 pages, \$53).

Nello Cristianini and John Shawe-Taylor. <u>An Introduction to Support Vector Machines</u>. Cambridge University Pres. Cambridge, UK, 2000.

An introduction to SVMs which is concise yet comprehensive in its description of the theoretical foundations of large margin algorithms (189 pages, \$45).

 Alex Smola, Peter Bartlett, Bernhard Schölkopf, and Dale Schuurmans (eds). <u>Advances in Large Margin Classifier</u> MIT Press, Cambridge, MA, 2000.

A collection of articles dealing with one of the main ideas of SVMs, large margin regularization. Contains an introduction, articles on new kernels, SVMs, and boosting algorithms (422 pages, \$45)

Home

Frequently Asked Questions

Publications

Books

Software

Annual Workshop

JMLR

Links

News

Machine Learning Summer School / Course On The Analysis On Patterns

2007-02-12

- > New Kernel-Machines.org server 2007-01-30
- > Call for participation: The 2006 kernel workshop, "10 years of kernel machines"

2006-10-06

> Three Workshops on Kernel Methods at NIPS 2005

2005-09-08

Nips 2004 Workshop on Learning With Structured Outputs

2004-10-22

More news...

Software

Kernel-Machines.Org software links

Skip in class

last modified 2007-01-31 12:47

Gaussian Processes

- GP Demo. Demonstration Software for Gaussian Processes by David MacKay (in OCTAVE).
- gpml. Matlab implementations of algorithms from Rasmussen & Williams "Gaussian Processes for Machine Learning the MIT Press 2006.
- <u>LS-SVMlab</u>. Matlab/C toolbox for least squares support vector machines.
- MAP-1. Package for MAP estimation by Carl Rasmussen.
- MC-1. Package for MAP estimation by Carl Rasmussen.
- Flexible Bayesian Modelling. Package by <u>Radford Neal</u>. It includes programs for Neural Networks, Gaussian Processes, and Mixture Models.
- Netlab. Matlab toolbox including Gaussian Process Regression, Mixture models and Neural Networks.
- Sparse Gaussian Processes. Matlab Toolbox for Sparse Inference using Gaussian Processes.
- Tpros and Cpros. Package by Mark Gibbs.

Mathematical Programming

- · CPLEX. Barrier/OP Solver.
- LOQO. Linear and Quadratic Optimization Package by Robert Vanderbei.
- MINOS. Linear and Quadratic Solver.

Support Vectors

- Nearest Point Algorithm. by Sathiya Keerthi (in FORTRAN).
- SVM Java Applet, by Chris Burges et al.
- BSVM. A decomposition method for bound-constrained SVM formulations.
- OP SVM Classification and Regression, Fortran Implementation.
- CLISP/LibSVM. A module for using LibSVM from GNU CLISP (an ANSI Common Lisp implementation).
- <u>Chunking Code</u>. by C. Saunders, M. O. Stitson, J. Weston, L. Bottou, B. Schölkopf, and A. Smola at Royal Holloway, AT&T, and GMD FIRST (Documentation).
- cSVM. SVM for classification tasks with model selection.
- 2D SVM Interactive Demo. runs under Matlab 6 and produces nice pictures useful for courses.
- DTREG. by Phillip H. Sherrod.
- Interior Point Optimizer for SVM Pattern Recognition. by Alex Smola.
- Equbits Foresight, Commercial SVM based Classification and Regression Application Designed for Drug Discovery.
- Gini-SVM. A multi-class Probabilistic regression software for large data sets.
- GiniSVM. Multi-class SVM Probability regression package.
- Gist. Gist contains software tools for support vector machine classification and for kernel principal components analysis. The SVM portion of Gist is available via an interactive web server.
- Parallel GPDT. Parallel and serial training of SVM.
- HeroSvm1.0. A high performance DLL for training SVM on a very large training set efficiently.
- SVM java implementation. This implementation is simple and easy to modify.
- LEARNSC. SVM, NN and FL MATLAB based user-friendly routines.
- LTDCVM An CVM library with a graphic interface

Outline

- Foundations
 - Linear Programming;Primal/Dual
 - Constrained Optimization
 - Lagrange Multipliers
 - KKT
 - Perceptron Factoids
 - Dual Representation
- <u>"Best"</u> Linear Separator: Max Margin!
- Coping with Non-Linearly Separated Data
 - Slack Variables
- Kernel Trick
- Regression



Background: Linear Programming

Linear Programming

- Given $\mathbf{c} \in \mathfrak{R}^n$, $\mathbf{A} \in \mathfrak{R}^{n \times m}$, $\mathbf{b} \in \mathfrak{R}^m$
- find $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} \mathbf{c}^\mathsf{T} \mathbf{x}$
- subject to

```
■ A x \le b -- ie, a_i^T x \le b_i for i = 1 ... m
■ x \ge 0 -- ie, x_j \ge 0 for j = 1 ... n
```

- ∃ fast algorithms for solving linear programs ...including
 - simplex algorithm
 - Karmarkar's algorithm



Duality

• Given $\mathbf{C} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^n$:

Primal

- find $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}}(\mathbf{c}^{\mathsf{T}})\mathbf{x}$
- subject to

Equivalent Dual

- find $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}}(\mathbf{b}^{\mathsf{T}})\mathbf{w}$
- subject to
 - **■**(A) w ≥ (c)
 - w ≥ 0

Skip in class

Strong duality result: If ** is an optimal solution for the primal,

then the dual has optimal solution w* s.t.

If A is $m \times n$ then, $|\mathbf{x}| = m$ $|\mathbf{w}| = n ...$ If m \gg n, then $|\mathbf{x}| \gg |\mathbf{w}|$

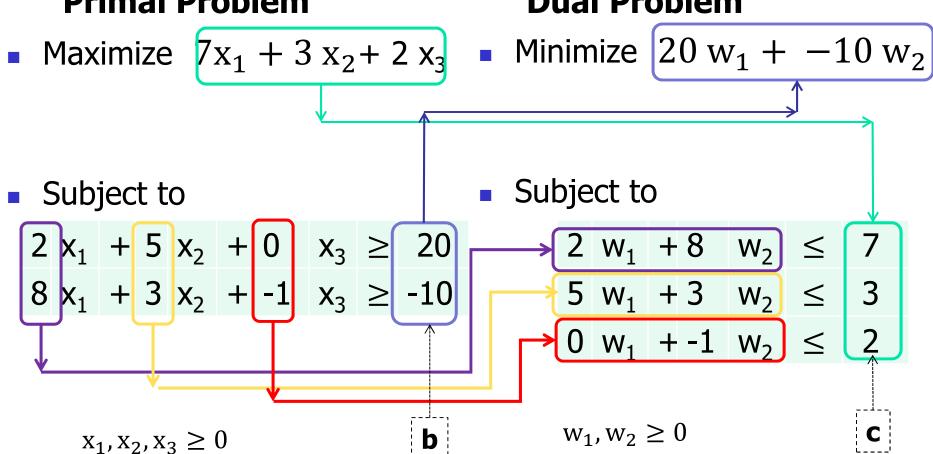
Proof: Lagrange multipliers...



Primal ⇔ Dual

Primal Problem

Dual Problem





Constrained optimization

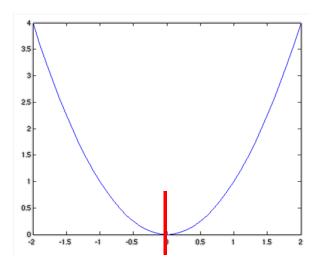
 $min_x x^2$

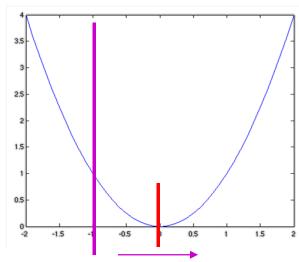
$$min_x x^2$$

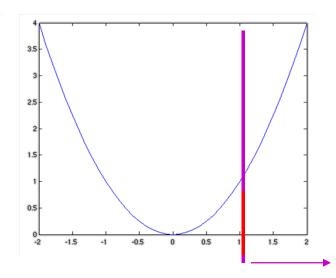
s.t. $x \ge -1$

$$min_x x^2$$

s.t. $x \ge 1$







min @ x=0

Unconstrained

min @ x=0

Constraint **ir**relevant

$$min @ x=+1$$

Constraint relevant

Gradient ... partial derivative

- This vector is TANGENT to $f(\mathbf{x})$
- Eg:
 - $f(x, y) = x^2 + 3xy 5x$

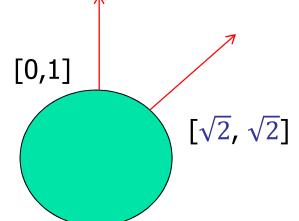
•
$$\nabla f(x,y) = \left[\frac{\partial f(x,y)}{\partial x}, \frac{\partial f(x,y)}{\partial y}\right] = [2x + 3y - 5, 3x]$$

• If evaluated at [x,y] = [1, 3], then this value is $\nabla f([1,3]) = [6, 3]$

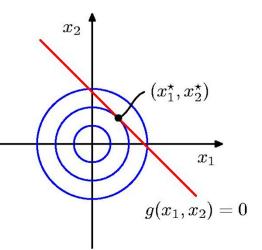
4

Gradient ... partial derivative

- This vector is DIRECTION where $f(\mathbf{x})$ increases fastest
- Eg: $f(x,y) = x^2 + y^2$
 - $\nabla f([x, y]) = [2x, 2y]$
 - ... so $\nabla f([\sqrt{2}, \sqrt{2}]) = [2\sqrt{2}, 2\sqrt{2}]$
 - $\nabla f([0, 1]) = [0, 2]$



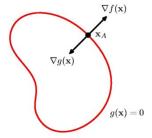
Lagrange Multiplier 101



- Challenge: $\operatorname{argmax}_{\mathbf{x}} f(\mathbf{x}) = -x_1^2 x_2^2$ s.t. $g(\mathbf{x}) = x_1 + x_2 - 1 = 0$
- Consider optimum $\mathbf{x}^* = (\mathbf{x}^*_{1}, \mathbf{x}^*_{2})$
 - On g(x) = 0 line, by def'n!
 - Note $\nabla f(\mathbf{x}^*) \perp g(\mathbf{x}^*)$
 - Otherwise, could walk along g(x)=0 to get larger f(x) value
- In fact, f(x*) is tangent to g(x*)

$$\Rightarrow \nabla f(\mathbf{x}^*)$$
 is $| \mathbf{to} \nabla g(\mathbf{x}^*)$

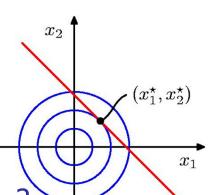
$$\Rightarrow \exists \lambda \text{ s.t. } \nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0$$



- Write $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$
- Note $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = \nabla_{\mathbf{x}} f(\mathbf{x}) + \lambda \nabla_{\mathbf{x}} g(\mathbf{x}) \dots = 0 \otimes \mathbf{x}^*$ $\nabla_{\lambda} L(\mathbf{x}, \lambda) = g(\mathbf{x}) \dots = 0 \Rightarrow g(\mathbf{x}) = 0$ satisfied

4

Lagrange Multiplier 101



- Challenge: $argmax_{\mathbf{x}} f(\mathbf{x}) = -x_1^2 x_2$ s.t. $g(\mathbf{x}) = x_1 + x_2 - 1 = 0$
- $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$

•
$$\nabla_{x_1} L(\mathbf{x}, \lambda) = -2 x_1 + \lambda = 0$$

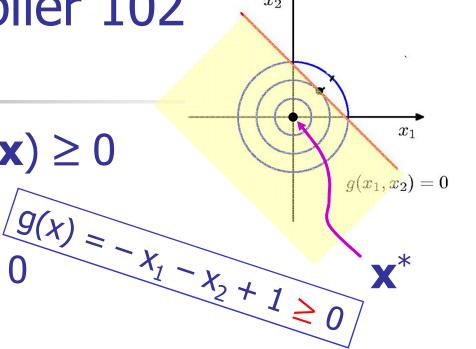
 $\nabla_{x_2} L(\mathbf{x}, \lambda) = -2 x_2 + \lambda = 0$
 $\nabla_{\lambda} L(\mathbf{x}, \lambda) = x_1 + x_2 - 1 = 0$

- Soln: $x_1 = x_2 = \frac{1}{2}$
- Also $\lambda = 1$ (who cares...)

Lagrange Multiplier 102 InEqualities...



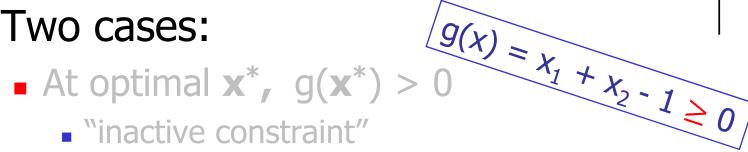
- Two cases:
 - At optimal \mathbf{x}^* , $g(\mathbf{x}^*) > 0$
 - "inactive constraint"
 - Just need $\nabla f(\mathbf{x}) = 0$
 - ... corresponds to $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ when $\lambda = 0$





Lagrange Multiplier 102 InEqualities...

- $argmax_{\mathbf{x}} f(\mathbf{x})$ s.t. $g(\mathbf{x}) \ge 0$
- Two cases:



- "inactive constraint"
- Just need $\nabla f(\mathbf{x}) = 0$
- ... corresponds to $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ when $\lambda = 0$

 $(x_1^{\star}, x_2^{\star})$

 $g(x_1,x_2)=0$

- At optimal \mathbf{x}^* , $g(\mathbf{x}^*) = 0$
 - "active constraint"
 - need $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ with $\lambda \neq 0$
 - Actually... need $\lambda > 0$ as need $\nabla f(\mathbf{x})$ oriented AWAY from $g(\mathbf{x}) > 0$ region... $\nabla f(\mathbf{x}^*) = -\lambda \nabla g(\mathbf{x}^*)$ for some $\lambda > 0$

Lagrange Multiplier 102

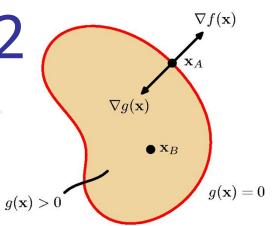
- argmax_x f(x) s.t. $g(x) \ge 0$
- At optimal \mathbf{x}^* , $g(\mathbf{x}^*) > 0$
 - $\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0 \text{ with } \lambda = 0$
- At optimal \mathbf{x}^* , $g(\mathbf{x}^*) = 0$
 - $\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0$ with $\lambda > 0$ Karush-Kuhn-Tucker (KKT):
- Either way... $\lambda g(\mathbf{x}^*) = 0$

$$\lambda g(\mathbf{x}^*) = 0$$

$$\text{Karush-Kuhn-Tucker (KK)}$$

$$\text{If } g(\mathbf{x}) > 0, \text{ then } \lambda = 0$$

- Summary:
 - $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$
 - Solve $\nabla_{x_i} L(x, \lambda) = 0$ $\nabla_{\lambda} L(x, \lambda) = 0$ s.t. $q(\mathbf{x}) \ge 0$ $\lambda \ge 0$ $\lambda q(\mathbf{x}) = 0$



KKT Conditions: Inequality Case

- Karush-Kuhn-Tucker Theorem:
 If
 - function f(x) has a minimum at x* in the feasible set, and
 - $\nabla f(x^*)$ and $\nabla g_i(x^*)$, i=1,2,...,m exist,

then \exists m-dimensional vector λ such that

$$\lambda \ge 0$$

$$\nabla f(x^*) - \sum_{i=1}^{m} \lambda_i \nabla g_i(x^*) = 0$$

$$\lambda_i [g_i(x^*)] = 0, \text{ for } i=1,2,...,m$$

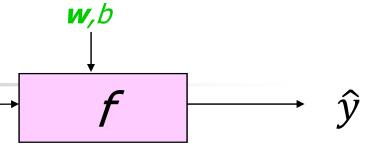
- Each such (x*, λ) is a <u>KKT point;</u>
 λ is the <u>Dual Vector</u> aka the <u>Lagrange Multipliers</u>
- These conditions are sufficient if dealing with a convex programming problem



Outline

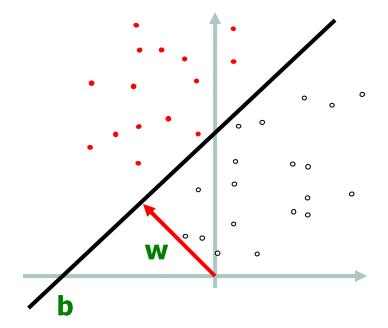
- Foundations
 - Linear Programming;Primal/Dual
 - Constrained Optimization
 - Lagrange Multipliers
 - KKT
 - Perceptron Factoids
 - Dual Representation
- "Best" Linear Separator: Max Margin!
- Coping with Non-Linearly Separated Data
 - Slack Variables
- Kernel Trick
- Regression

Linear Classifiers



Input space	x ∈ X
Output space	$y \in Y = \begin{Bmatrix} +1 \\ -1 \end{Bmatrix}$
Real-valued fn:	f: X → Y
Training Set	$S = \left\{ \begin{bmatrix} \mathbf{x}_1, \mathbf{y}_1 \end{bmatrix} \\ \vdots \\ \begin{bmatrix} \mathbf{x}_m, \mathbf{y}_m \end{bmatrix} \right\}$

$$f(x, w, b) = sign(w^Tx + b)$$



Perceptron Training Rule

```
Initialize \mathbf{w} = 0

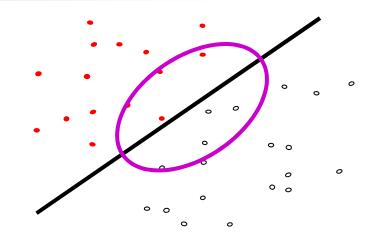
Do until done

Predict "+" iff \mathbf{w}^T \mathbf{x} > 0

else "-"

Mistake on \mathbf{y} = +1: \mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}

Mistake on \mathbf{y} = -1: \mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}
```



$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{y} \mathbf{x}$$

$$\Rightarrow$$
 $\mathbf{w} = \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i}$

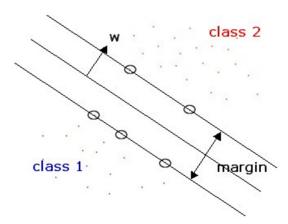
- ... only uses Informative Points (mistake driven)
- Coefficient α_i of point \mathbf{x}_i reflects its 'difficulty'

Mistake Bound Theorem

Theorem: [Rosenblatt 1960]

If data is consistent w/some linear threshold \mathbf{w} , then number of mistakes is $\leq (1/\Delta)^2$, where $\Delta = \min \frac{|\mathbf{w} \cdot \mathbf{x}|}{|\mathbf{w} \cdot \mathbf{x}|}$

- Δ measures "wiggle room" available:
 If |x| = 1, then Δ is max, over all consistent planes, of minimum distance of example to that plane
- w is \perp to separator, as w ' x = 0 at boundary
- So |w x| is projection of x onto plane,
 PERPENDICULAR to boundary line
 ie, is distance from x to that line (once normalized)



≈ margin

4

Dual Representation

$$\mathbf{w} = \sum_{i} \alpha_{i} \mathbf{y}_{i} \mathbf{x}_{i}$$

⇒ can re-write decision function

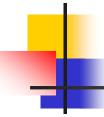
$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b = \sum_{i} \alpha_{i} y_{i} \langle \mathbf{x}_{i}, \mathbf{x} \rangle + b$$

 \Rightarrow can re-write update rule:

If
$$y_j [\sum_i \alpha_i y(\mathbf{x}_i, \mathbf{x}_j) + b] \leq 0$$

Then $\alpha_j \leftarrow \alpha_j + \eta$

In dual representation,
 datapoints x_i appears only inside dot products



Outline

- Foundations
 - Linear Programming; Primal/Dual
 - Constrained Optimization
 - Lagrange Multipliers
 - KKT
 - Perceptron Factoids
 - Dual Representation
- "Best" Linear Separator: Max Margin!
- Coping with Non-Linearly Separated Data
 - Slack Variables
- Kernel Trick
- Regression

Linear Classifiers

 $f \longrightarrow j$

- denotes +1
- ∘ denotes −1



How to classify this data?

Each of these seems fine...

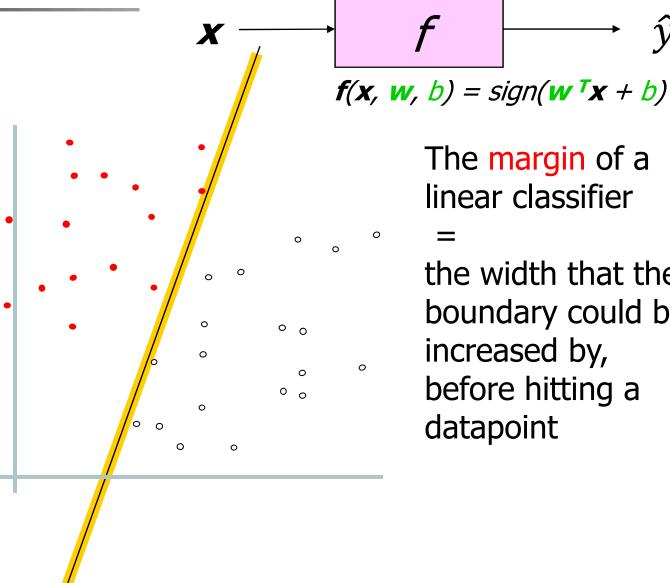
... which is best?

0 0

0



- denotes +1
- denotes -1



The margin of a linear classifier

w,b

the width that the boundary could be increased by, before hitting a datapoint



0 0

 \hat{y} \hat{y}

- denotes +1
- ∘ denotes −1

Support Vectors' are datapoints that "touch" the margin

 $f(x, w, b) = sign(w^{T}x + b)$

Maximum margin linear classifier

=

0 0

linear classifier with maximum margin

Linear SVM

=

simplest kind of SVM

Maximum Margin

- 1. ... this feels safest ...
- 2. If a small error in the location of the boundary (eg, jolt in its perpendicular direction), this gives least chance of causing a misclassification.

w,b

- 3. LOO-CV is easy, since the model is immune to removal of any non-support-vector datapoints.
- ## I some theory (using VC dimension)

 that is related to

 (but not the same as)

 the claim that this is a good thing.
- 5. Empirically it works *very very well*.

- denotes +1
- ∘ denotes −1

Support Vectors' are datapoints that "touch" the margin



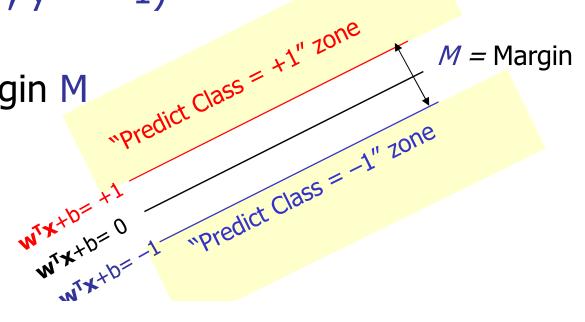
Goal of Max Margin Separator

Want a linear separator

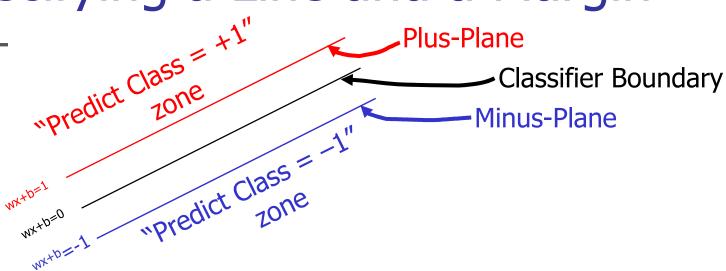
$$\mathbf{w}$$
, \mathbf{b} for $\mathbf{y} = \mathbf{w}^T \mathbf{x} + \mathbf{b}$ s.t.

- For all +points $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)} = + \mathbf{y}^{(i)})$ $\mathbf{w}^{\mathsf{T}} \mathbf{x}^{(i)} + \mathbf{b} \ge (+1)^{\circ}$
- For all -points $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)} = -1)$ $\mathbf{w}^{T} \mathbf{x}^{(i)} + \mathbf{b} \leq -1$
- Maximizes the margin M

Why 1?
Any >0 constant works, as scales.
1 is convenient...



Specifying a Line and a Margin



- Plus-plane = $\{ \mathbf{x} : \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} = -1 \}$

```
Classify as.. +1 if w^Tx + b \ge 1
```

$$-1 if $\mathbf{w}^T \mathbf{x} + b \le -1$ (whatever) if $-1 < \mathbf{w}^T \mathbf{x} \ne b < 1$$$

Never happens

-

Computing the Margin Width

"Predict Class = +1"

"Predict Class = -1"

"Predict Class = -1"

"Predict Class = -1"

How to compute *M* in terms of *w* and *b*?

M = Margin Width

- Plus-plane = $\{ \mathbf{x} : \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} = -1 \}$

Claim: Plus-plane is || to Minus-Plane

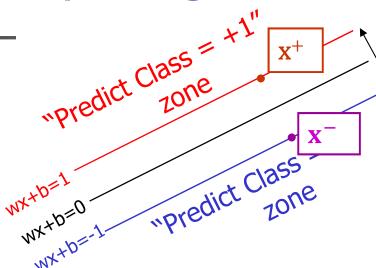
The vector w is perpendicular to the Plus Plane. Why?

- Definitions: "vector" ≡ "point"
- **a** perpendicular to **b** iff $\mathbf{a}^{\mathsf{T}}\mathbf{b} = 0$

Let **u** and **v** be two vectors on the Plus Plane. What is $\mathbf{w}^T(\mathbf{u} - \mathbf{v})$?

4

Computing the Margin Width



How to compute *M* in terms of *w* and *b*?

Any location in \mathfrak{R}^m :

not necessarily

a datapoint

M = Margin Width

- Plus-plane = $\{ \mathbf{x} : \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} = -1 \}$
- The vector w is perpendicular to the Plus Plane
- \mathbf{x}^- = any point on the minus-plane
- x^+ = the point in plus-plane closest to x^-
- Claim: $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$ for some value of $\lambda \in \mathcal{R}^+$. Why?

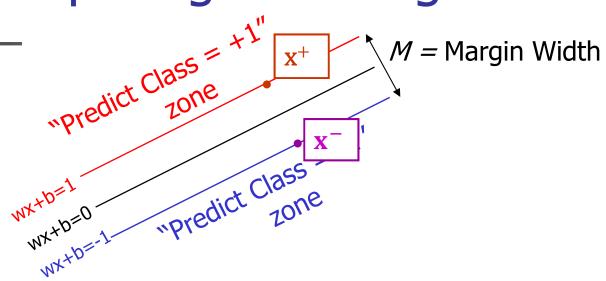
ı

Computing the Margin Width

- "predict Class = +1", *M* = Margin Width Line from \mathbf{x}^- to \mathbf{x}^+ is \perp to the planes. "Predict Class So to get from x^- to x^+ , travel some distance in the direction of w • Plus-plane = $\{x: w^T x + b =$ Minus-plane = $\{x: w^T x + b = -1\}$ So ... $\mathbf{x}^+ - \mathbf{x}^- = \lambda \mathbf{w}$ ular to the Plus Plane átion in Mm. x^- = any point \sqrt{n} e minus plane **Ecessarily Atapoint** • x^+ = the point $\sqrt{}$ plus-plane closest to x^-
- Claim: $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$ for some value of $\lambda \in \mathcal{R}^+$. Why?

-

Computing the Margin Width



Given...

•
$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{\mathsf{+}} + \mathbf{b} = +1$$

•
$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{\mathsf{T}} + \mathbf{b} = -1$$

$$\mathbf{x}^+ - \mathbf{x}^- = \lambda \mathbf{w}$$

$$|x^{+} - x^{-}| = M$$

... easy to get M in terms of **w** and b



Computing the Margin Width

"predict Class" x^{+1} M = Margin Width"predict Class" x^{-1} "predict Class" x^{-1} "predict x^{-1} "predict

Given...

•
$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{\mathsf{+}} + \mathbf{b} = +1$$

•
$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{\mathsf{T}} + \mathbf{b} = -1$$

$$\mathbf{x}^+ - \mathbf{x}^- = \lambda \mathbf{w}$$

$$|x^{+} - x^{-}| = M$$

... easy to get M in terms of w and b

$$(\mathbf{w}^\mathsf{T} \mathbf{x}^- + \mathbf{b}) + \lambda \mathbf{w}^\mathsf{T} \mathbf{w} = 1$$

$$\Rightarrow$$

$$-1 + \lambda \mathbf{w}^\mathsf{T} \mathbf{w} = 1$$

$$\Rightarrow \lambda = \frac{2}{\mathbf{w}^{\mathsf{T}}\mathbf{w}}$$

4

Computing the Margin Width

"Predict Class" $M = Margin Width = \frac{2}{\sqrt{w^T w}}$ "Predict Class" $M = |\mathbf{x}^+ - \mathbf{x}^-| = |\lambda w|$

Given ...

•
$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{\mathsf{+}} + \mathbf{b} = +1$$

•
$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{\mathsf{-}} + \mathbf{b} = -1$$

•
$$x^+ - x^- = \lambda w$$

$$|x^+ - x^-| = M$$

$$\lambda = \frac{2}{\mathbf{w} \cdot \mathbf{w}}$$

$$M = |\mathbf{X}^{+} - \mathbf{X}^{-}| = |\lambda \mathbf{w}|$$

$$= \lambda |\mathbf{w}| = \frac{2}{\mathbf{w}^{T}\mathbf{w}} \sqrt{\mathbf{w}^{T}\mathbf{w}}$$

$$= \frac{2}{\sqrt{\mathbf{w}^{T}\mathbf{w}}}$$

So... to maximize
$$M = \frac{2}{\sqrt{\mathbf{w}^T \mathbf{w}}}$$

just minimize $\mathbf{w}^T \mathbf{w}$

Wait...OMG, I forgot the data!



Goal of Max Margin Separator

Want a linear separator

$$\mathbf{w}$$
, \mathbf{b} for $\mathbf{y} = \mathbf{w}^T \mathbf{x} + \mathbf{b}$ s.t.

• For all +points $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}=+1)$

$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{(\mathsf{i})} + b \geq +1$$

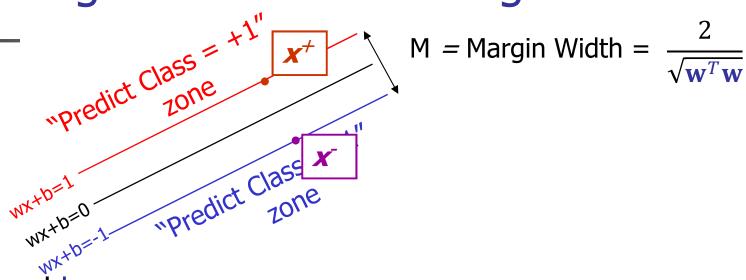
• For all –points $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)} = -1)$

$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{(\mathsf{i})} + b \leq -1$$

Maximizes the margin M

Minimizes w^T w

Learning the Maximum Margin Classifier



Given w and b we can

- Compute whether each data point is in correct half-plane
- Compute the width of the margin

But... need a program to search the space of w's and b's to find the widest margin that matches all the datapoints.

How?

Gradient descent? Simulated Annealing? Matrix Inversion? EM? Newton's Method?



Rewrite Problem

Minimize
$$\frac{1}{2} \mathbf{W}^T \mathbf{W}$$

s.t. $\mathbf{W}^T \mathbf{X}_k + b \ge 1$ if $\mathbf{Y}_k = -1$

$$\mathbf{W}^T \mathbf{X}_k + b \le -1$$
 if $\mathbf{Y}_k = -1$

Lagrange Multiplier

Equivalent optimization...

$$L(\mathbf{w}, b, \lambda) = \frac{1}{2} |\mathbf{w}|^2 - \sum_k \lambda_k [y_k (\mathbf{w}^T \mathbf{x}_k + b) - 1]$$

$$\min_{\mathbf{w},b} L(\mathbf{w}, \mathbf{b}, \lambda) \atop \text{s.t. } \lambda \ge 0 \qquad \begin{array}{c} \mathsf{KKT:} \\ \mathsf{y_k}(\mathbf{w^T} \, \mathbf{x_k} + \mathbf{b}) > 1 \implies \lambda_k = 0 \end{array}$$

Maximum Margin

- denotes +1
- ° denotes -1

Support Vectors'

are datapoints that "touch" the margin

$$\lambda_k = 0$$

$$\lambda_k = 0$$

$$\lambda_k = 0$$

KKT:

$$y_k (\mathbf{w}^T \mathbf{x}_k + \mathbf{b}) > 1 \Rightarrow \lambda_k = 0$$

4

Solving Constrained Optimization

$$\min_{\mathbf{w},b} L(\mathbf{w}, \mathbf{b}, \lambda) = \frac{1}{2} |\mathbf{w}|^2 - \sum_k \lambda_k [y_k (\mathbf{w}^T \mathbf{x}_k + \mathbf{b}) - 1]$$

s.t. $\lambda \ge \mathbf{0}$

Setting derivatives to 0...

$$\mathbf{w} = \sum_{m} \lambda_{m} y_{m} \mathbf{x}_{m}$$

$$0 = \sum_{m} \lambda_{m} y_{m}$$

Substitute back into $L(\cdot,\cdot,\cdot)$:

Find $\lambda \geq 0$ that maximizes

$$L(\lambda) = \sum_{k} \lambda_{k} - \frac{1}{2} \sum_{k} \lambda_{k} y_{k} (\sum_{m} \lambda_{m} y_{m} \mathbf{x}_{m})^{T} \mathbf{x}_{k})$$

$$= \sum_{k} \lambda_{k} - \frac{1}{2} \sum_{k} \sum_{m} \lambda_{k} \lambda_{m} y_{k} y_{m} (\mathbf{x}_{k}^{T} \mathbf{x}_{m})$$



Learning via Quadratic Programming

- QP is a well-studied class of optimization alg's that
 - maximize a quadratic function of some real-valued variables
 - subject to linear constraints
- Popular ML approach:
 - Describe your learning problem as optimization...
 - ...and give it to somebody else to solve!



Quadratic Programming — in general

Find $\underset{\mathbf{w}}{\operatorname{arg\,min}} \quad c + \mathbf{d}^T \mathbf{w} + \frac{\mathbf{w}^T \mathbf{K} \mathbf{w}}{2}$

Quadratic criterion

Subject to

$$\begin{array}{c} a_{11}w_{1} + a_{12}w_{2} + \ldots + a_{1m}w_{m} \leq b_{1} \\ a_{21}w_{1} + a_{22}w_{2} + \ldots + a_{2m}w_{m} \leq b_{2} \\ \vdots \\ a_{n1}w_{1} + a_{n2}w_{2} + \ldots + a_{nm}w_{m} \leq b_{n} \end{array}$$

$$n \text{ additional linear inequality constraints}$$

$$a_{n1}w_{1} + a_{n2}w_{2} + \ldots + a_{nm}w_{m} \leq b_{n}$$

and to

$$a_{(n+1)1}w_1 + a_{(n+1)2}w_2 + \ldots + a_{(n+1)m}w_m = b_{(n+1)}$$

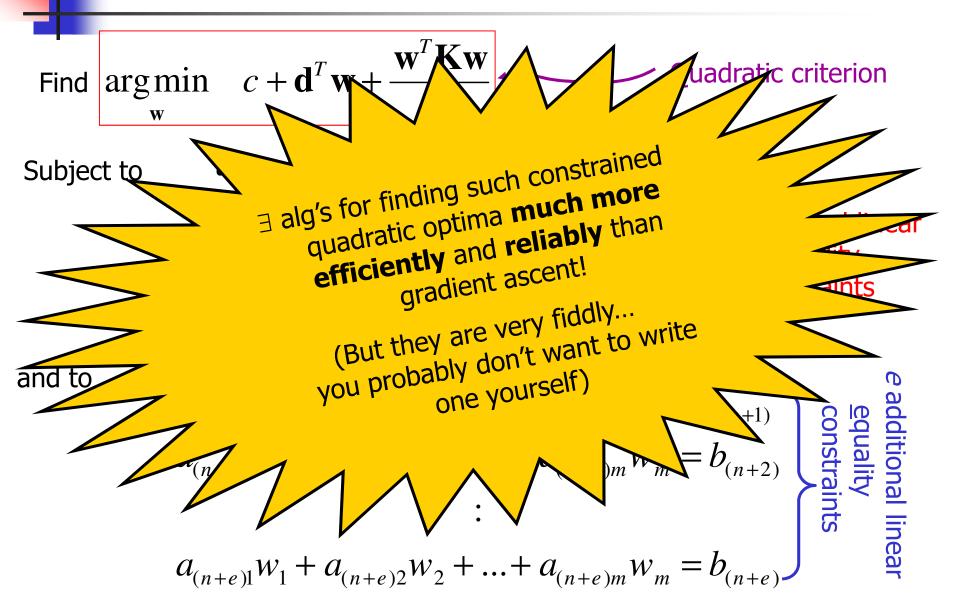
$$a_{(n+2)1}w_1 + a_{(n+2)2}w_2 + \ldots + a_{(n+2)m}w_m = b_{(n+2)}$$

$$\vdots$$

$$a_{(n+e)1}w_1 + a_{(n+e)2}w_2 + \ldots + a_{(n+e)m}w_m = b_{(n+e)}$$

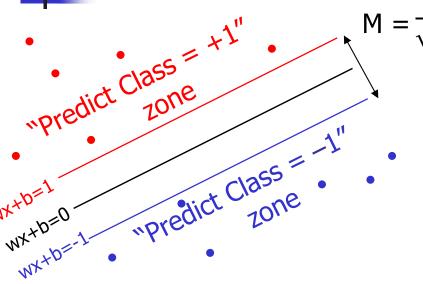
$$a_{(n+e)1}w_1 + a_{(n+e)2}w_2 + \ldots + a_{(n+e)m}w_m = b_{(n+e)}$$

Quadratic Programming – in general



4

Learning the Maximum Margin Classifier



What is quadratic optimization criterion?

Minimize w^Tw

 $M = \overline{\sqrt{\mathbf{w}^T \mathbf{w}}}$ Given guess of \mathbf{w} , \mathbf{b} , can...

- Compute whether each data point is in the correct half-plane
- Compute the margin width
- Leading to R datapoints, $\{[\mathbf{x}_k, y_k]\}$ where $y_k \in \{+1, -1\}$

How many constraints? R

What are they?

$$\mathbf{w}^T \mathbf{x}_k + b \ge 1$$
 if $y_k = 1$
 $\mathbf{w}^T \mathbf{x}_k + b \le -1$ if $y_k = -1$



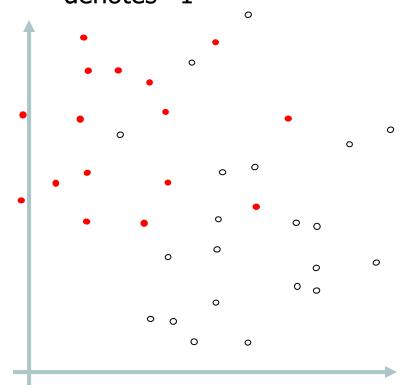
Outline

- Foundations
 - Linear Programming; Primal/Dual
 - Constrained Optimization
 - Lagrange Multipliers
 - KKT
 - Perceptron Factoids
 - Dual Representation
- "Best" Linear Separator: Max Margin!
- Coping with Non-Linearly Separated Data
 - Slack Variables
- Kernel Trick
- Regression



• denotes +1

∘ denotes −1



Earlier analysis does not work!

What should we do?

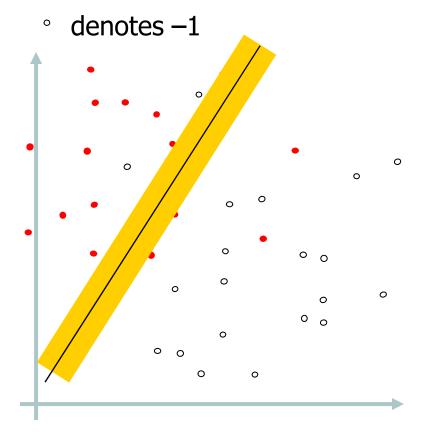
Idea 1: Find w s.t.

- minimum w^Tw
- minimum # of training set errors.

Problem: Minimizing *TWO* objectives is ill-defined!



denotes +1



Earlier analysis does not work!

What should we do?

Idea 1.1: Find w that

Minimizes

 $\mathbf{w}^{\mathsf{T}}\mathbf{w} + C (\# train \ errors)$

Tradeoff parameter

But... a *serious* practical problem dooms this approach

- denotes +1
- denotes -1

Earlier analysis does not work!

What should we do?

Idea 1.1: Find w that

Minimizes

1. Can't be expressed as a Quadratic Programming problem.

- \Rightarrow Solving it is too slow.
- 2. Does not distinguish between disastrous errors and near misses

#train error = 2 + 4

yv^T w + C (#train errors)

Tradeoff parameter



denotes +1

° denotes -1

Earlier analysis does not work!

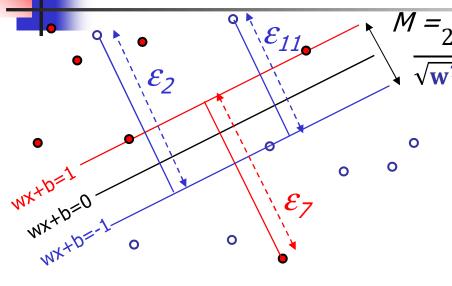
What should we do?

Idea 2: Find w that

Minimizes

w^Tw + C (distance from incorrectly labeled points to their correct place)

Learning Maximum Margin with Noise



Given guess of **w**, b, can...

- $\sqrt{\mathbf{w}^T\mathbf{w}}$ Compute whether each data point is in the correct half-plane
 - Compute the margin width

R datapoints,
$$\{[\mathbf{x}_k, \mathbf{y}_k]\}$$

where $\mathbf{y}_k \in \{+1, -1\}$

What is quadratic optimization criterion?

Minimize
$$\frac{\mathbf{w}^{\mathsf{T}} \mathbf{w}}{2} + C \sum_{k=1}^{R} \varepsilon_k$$

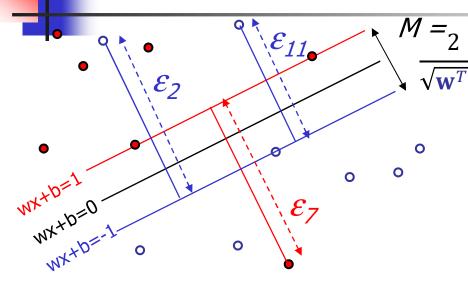
How many constraints?

Minimize
$$\frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{k=1}^{R} \varepsilon_k$$
 What should they be? $\mathbf{w}^T \mathbf{x}_k + b \ge (1 - \varepsilon_k)$ if $y_k = 1$ $\mathbf{w}^T \mathbf{x}_k + b \le (-1 + \varepsilon_k)$ if $y_k = -1$

There's a bug in this QP. Can you spot it?

.

Learning Maximum Margin with Noise



Given guess of **w**, b, can...

- Compute whether each data point is in the correct half-plane
 - Compute the margin width

R datapoints,
$$\{[\mathbf{x}_k, y_k]\}$$

where $y_k \in \{+1, -1\}$

What is quadratic optimization criterion?

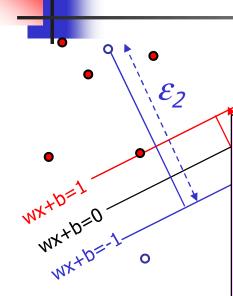
How many constraints? 2R

Minimize
$$\frac{\mathbf{w}^{\mathsf{T}} \mathbf{w}}{2} + C \sum_{k=1}^{R} \varepsilon_k$$

What should they be?



Learning Maximum Margin with Noise



M = 2 Give $\sqrt{\mathbf{w}^T \mathbf{w}} \bullet \text{Compute wheth}$

m = # input dimensions

hidth

Our original (noiseless data) QP had M+1variables: W_1 , W_2 , ... W_m , and b.

Our new (noisy data) QP has m+1+Rvariables: W_1 , W_2 , ... W_m , b, ε_k , ε_1 ,... ε_R

What is quadratic optimization criterion?

Minimize
$$\frac{\boldsymbol{w}^T \boldsymbol{w}}{2} + C \sum_{k=1}^R \varepsilon_k$$

How many constrain R = # instances

What should they be?

$$\mathbf{w}^{\mathsf{T}} \mathbf{x}_{k} + b \ge (1 - \varepsilon_{k})$$
 if $y_{k} = 1$
 $\mathbf{w}^{\mathsf{T}} \mathbf{x}_{k} + b \le (-1 + \varepsilon_{k})$ if $y_{k} = -1$
 $\varepsilon_{k} \ge 0$ for all k

Managing trade-off "C"

• Minimize
$$\frac{\mathbf{w}^{\mathrm{T}} \mathbf{w}}{2} + C \sum_{k=1}^{R} \varepsilon_{k}$$

- Big C ⇒ "Fit the training data as much as possible!" (at the expense of maximizing margin)
- Small C⇒"Maximize the margin as much as possible!" (at the expense of fitting the training data)
- Best value: Internal cross-validation

4

Solving Constrained Optimization

$$\min_{\mathbf{w},b} L(\mathbf{w}, b, \lambda, \varepsilon) = \frac{1}{2} |\mathbf{w}|^2 + C \sum_{k} \varepsilon_{k} - \sum_{k} \lambda_{k} [y_{k}(\mathbf{w}^{\mathsf{T}} \mathbf{x}_{k} + b) - (1 - \varepsilon_{k})]$$

s.t.
$$\lambda$$
, $\varepsilon \geq 0$

Setting derivatives to 0...

•
$$\mathbf{w} = \sum_{m} \lambda_{m} y_{m} \mathbf{x}_{m}$$

$$0 = \sum_{m} \lambda_{m} y_{m}$$

To incorporate slack variables \mathcal{E}_k

Just add constraints:

$$0 \le \lambda_k \le C$$

Substitute back into L(..):

Find $\lambda \ge 0$ that minimizes

Actually:

$$L(\lambda) = ... + \sum_{k} (C - \lambda_{k}) \varepsilon_{k}$$
But $\varepsilon_{k} > 0 \implies \lambda_{k} = C ...$

$$L(\lambda) = \sum_{k} \lambda_{k} - \frac{1}{2} \sum_{k} \sum_{m} \lambda_{k} \lambda_{m} y_{k} y_{m} (\mathbf{x}_{k}^{\mathsf{T}} \mathbf{x}_{m})$$

An Equivalent QP

$$\max_{\lambda} \sum_{k=1}^{R} \lambda_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{r=1}^{R} \lambda_k \lambda_r Q_{kr} \quad \text{where} \quad Q_{kr} = y_k y_r (\mathbf{x}_k^T \mathbf{x}_r)$$

Subject to the constraints:

$$0 \le \lambda_k \le C \quad \forall k \qquad \sum \lambda_k y_k = 0$$

$$\sum_{k=1}^{R} \lambda_k y_k = 0$$

Then define:

$$\mathbf{w} = \sum_{\mathbf{k}} \mathbf{\lambda}_k \mathbf{y}_k \mathbf{x}_k$$

$$k \text{ s.t. } \mathbf{\lambda}_k > 0$$

$$b = y_K (1 - \epsilon_K) - \mathbf{x}_K \cdot \mathbf{w}$$

where
$$K = \arg \max_{k} \lambda_{k}$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, \mathbf{b}) = sign(\mathbf{w}^T \mathbf{x} + \mathbf{b})$$

An Equivalent QP

$$\max_{\lambda} \sum_{k=1}^{R} \lambda_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{r=1}^{R} \lambda_k \lambda_r Q_{kr} \quad \text{where} \quad Q_{kr} = y_k y_r (\mathbf{x}_k^T \mathbf{x}_r)$$

Subject to the constraints:

$$0 \leq \lambda_k \leq C$$

 \mathbf{x}_k only appears in dot product!

k=1

R

Then define:

$$\mathbf{w} = \sum_{k} \mathbf{\lambda}_k \mathbf{y}_k \mathbf{x}_k$$
k s.t. $\mathbf{\lambda}_k > 0$

$$b = y_K (1 - \epsilon_K) - \mathbf{x}_K \cdot \mathbf{w}$$

where $K = \arg \max_{k} \lambda_k$

Datapoints with $\lambda_k > 0$ == support vectors

Then classify with:

 $f(x w b) = sign(w^T x + b)$

...note this sum only needs to be over the support vectors.

(probably << R)

Δ

An Equivalent QP

$$\max_{\lambda} \sum_{k=1}^{R} \lambda_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{r=1}^{R} \lambda_k r Q_{kr} \quad \text{where} \quad Q_{kr} = y_k y_r (\mathbf{x}_k^T \mathbf{x}_r)$$

Subject to

constr

Why use this equivalent QP?

Then as

- QP packages can optimize it more quickly
- $W = \sum_{i=1}^{n} x_i x_i$
- Stay tuned...

k s.t. $\alpha_k > 0$

$$b = y_K (1 - \epsilon_K) - x_K \cdot \mathbf{w}$$

where $K = \arg \max_{k} \alpha_{k}$

needs to be over the support vectors.

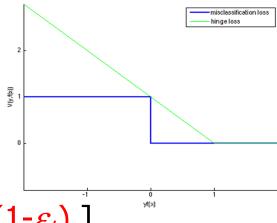
dot product!

4 þ)

(probably << R)



Hinge Loss



$$\min_{\mathbf{w},b} L(\mathbf{w}, b, \lambda, \varepsilon) = \frac{1}{2} |\mathbf{w}|^{2} + C \sum_{k} \varepsilon_{k} - \sum_{k} \lambda_{k} [y_{k}(\mathbf{w}^{\mathsf{T}} \mathbf{x}_{k} + b) - (1 - \varepsilon_{k})]$$
s.t. $\lambda, \varepsilon \geq 0$

■ |w|² is regularizer

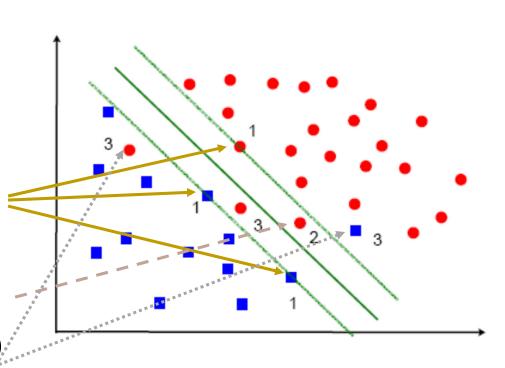
Hidden, as part of HW ...



Types of Support Vectors

Support Vectors (s.v.s)

- 1. Margin s.v. $\xi_i = 0$
 - correct
- 2. Non-Margin s.v. $\xi_i < 0$
 - Correct (in margin)
- 3. Non-Margin s.v. $\xi_i > 0$
 - Error





What do we have?

- Method for learning a maximum-margin linear classifier when the data are ...
 - Linearly separable" –∃ line that gets 0 training error
- But if not linearly separable
 - ie, no such line,

must trade-off between

- maximizing margin vs
- minimizing "stuff-is-on-the-wrong-side-ness"
- OR DO WE?? Kernels!