# Machine Learning 10-701

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

April 12, 2011

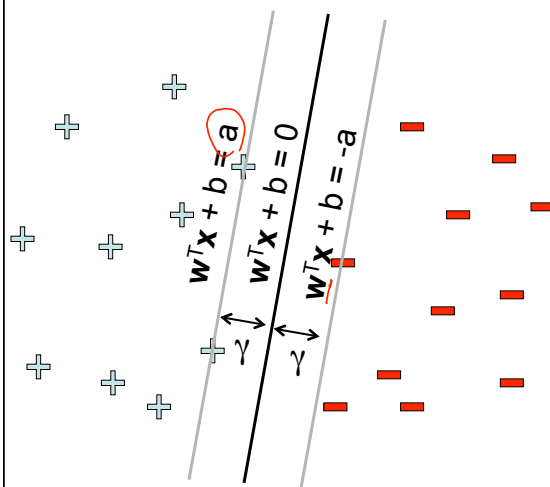| Today: | Readings: |
|---|---|
| • Support Vector Machines<br>• Margin-based learning | Required:<br>SVMs: Bishop Ch. 7, through 7.1.2<br><br>Optional:<br>Remainder of Bishop Ch. 7 |

Thanks to Aarti Singh for several slides
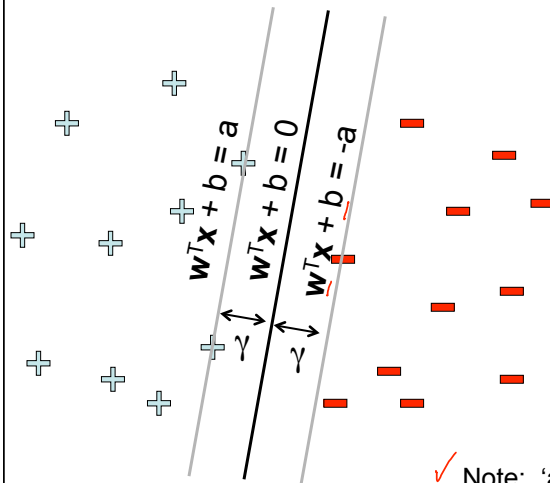
---

# SVM: Maximize the margin

Margin = Distance of
closest examples
from the decision line/
hyperplane



$w^T x + b = a$

$w^T x + b = 0$

$w^T x + b = -a$

$\gamma$   $\gamma$

margin = $\gamma$ = a/‖w‖

# Maximizing the margin

$Y \in \{-1, 1\}$

$\mathbf{w}^T\mathbf{x} + b = a$
$\mathbf{w}^T\mathbf{x} + b = 0$
$\mathbf{w}^T\mathbf{x} + b = -a$

$\gamma$    $\gamma$

Margin = Distance of closest examples from the decision line/ hyperplane

$$\boxed{\text{margin} = \gamma = a/\|\mathbf{w}\|}$$

$a = 1$

$$\max_{\mathbf{w},b} \ \gamma = a/\|\mathbf{w}\|$$

$$\text{s.t.} \ (\mathbf{w}^T\mathbf{x}_j + b) \ y_j \geq a \ \forall j$$

✓ Note: 'a' is arbitrary (can normalize equations by a)

---

# Support Vector Machine (primal form)

$$\max_{\mathbf{w},b} \ \gamma = 1/\|\mathbf{w}\|$$

$$\text{s.t.} \ (\mathbf{w}^T\mathbf{x}_j + b) \ y_j \geq 1 \ \forall j$$

Primal form:

$$\min_{\mathbf{w},b} \ \mathbf{w}^T\mathbf{w}$$

$$\text{s.t.} \ (\mathbf{w}^T\mathbf{x}_j + b) \ y_j \geq 1 \ \forall j$$

$\mathbf{w}^T\mathbf{x} + b = 1$
$\mathbf{w}^T\mathbf{x} + b = 0$
$\mathbf{w}^T\mathbf{x} + b = -1$

$\gamma$    $\gamma$

Solve efficiently by quadratic programming (QP)

– Well-studied solution algorithms

## We can solve either primal or dual forms

Primal form: solve for $\mathbf{w}, b$

$$\min_{\mathbf{w},b} \quad \mathbf{w}^T\mathbf{w}$$
$$\text{s.t.} \quad y_l(\mathbf{w}^T\mathbf{x}_l + b) \geq 1 \quad \forall l \in \text{ training examples}$$

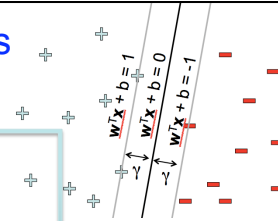Classification test for new $\mathbf{x} : \mathbf{w}^T\mathbf{x} + b > 0$

Dual form: solve for $\alpha_1 ... \alpha_M$

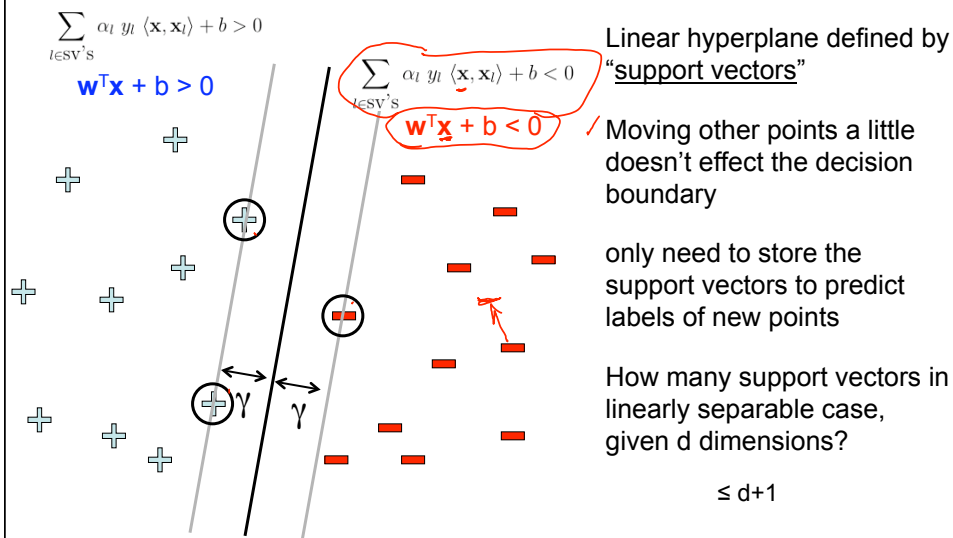$$\max_{\alpha_1...\alpha_M} \quad \sum_{l=1}^{M} \alpha_l - \frac{1}{2} \sum_{j=1}^{M} \sum_{k=1}^{M} \alpha_j \alpha_k y_j y_k \langle \mathbf{x_j}, \mathbf{x_k} \rangle$$
$$\text{s.t.} \quad \alpha_l \geq 0 \quad \forall l \in \text{ training examples}$$
$$\sum_{l=1}^{M} \alpha_l y_l = 0$$

Classification test for new $\mathbf{x} : \sum_{l \in \text{sv's}} \alpha_l y_l \langle \mathbf{x}, \mathbf{x}_l \rangle + b > 0$

$\mathbf{w}^T\mathbf{x} + b = 1$
$\mathbf{w}^T\mathbf{x} + b = 0$
$\mathbf{w}^T\mathbf{x} + b = -1$

both are QP problems with a single local optimum!

---

## Support Vectors

$$\sum_{l \in \text{sv's}} \alpha_l \, y_l \, \langle \mathbf{x}, \mathbf{x}_l \rangle + b > 0$$

$\mathbf{w}^T\mathbf{x} + b > 0$

$$\sum_{l \in \text{sv's}} \alpha_l \, y_l \, \langle \mathbf{x}, \mathbf{x}_l \rangle + b < 0$$

$\mathbf{w}^T\mathbf{x} + b < 0$

Linear hyperplane defined by "support vectors"

Moving other points a little doesn't effect the decision boundary

only need to store the support vectors to predict labels of new points

How many support vectors in linearly separable case, given d dimensions?

≤ d+1

$\gamma$  $\gamma$

3

# Kernel SVM

And because the dual form depends only on inner products, we can apply the kernel trick to work in a (virtual) projected space $\Phi : X \to F$

Primal form: solve for $\mathbf{w}, b$ in the projected higher dim. space

$$\min_{\mathbf{w}, b} \quad \mathbf{w}^T \mathbf{w}$$

$$\text{s.t.} \quad y_l(\mathbf{w}^T \Phi(\mathbf{x}_l) + b) \geq 1 \quad \forall l \in \text{ training examples}$$

Classification test for new $\mathbf{x}$ $\quad \mathbf{w}^T \Phi(\mathbf{x}) + b > 0$

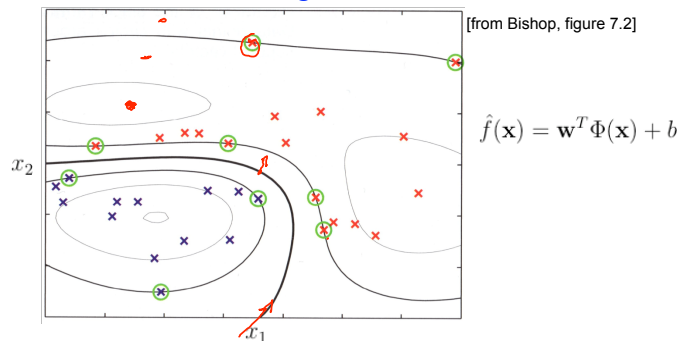Dual form: solve for $\alpha_1 ... \alpha_M$ in the original low dim. space

$$\max_{\alpha_1 ... \alpha_M} \quad \sum_{l=1}^{M} \alpha_l - \frac{1}{2} \sum_{j=1}^{M} \sum_{k=1}^{M} \alpha_j \alpha_k y_j y_k \, \kappa(\mathbf{x_j}, \mathbf{x_k}) \quad = \langle \Phi(x_j), \Phi(x_k) \rangle$$

$$\text{s.t.} \quad \alpha_l \geq 0 \quad \forall l \in \text{ training examples}$$

$$\sum_{l=1}^{M} \alpha_l y_l = 0$$

Classification test for new $\mathbf{x}$ : $\displaystyle\sum_{l \in \text{sv's}} \alpha_l \, y_l \left( \kappa(\mathbf{x}, \mathbf{x}_l) \right) + b > 0$

---

# SVM Decision Surface using Gaussian Kernel



[from Bishop, figure 7.2]

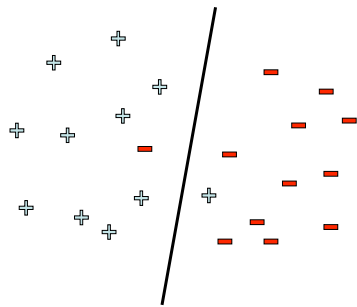$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) + b$$

Circled points are the *support vectors*: training examples with non-zero $\alpha_l$

Points plotted in underline{original} 2-D space.

Contour lines show constant $\hat{f}(\mathbf{x})$

$$\hat{f}(\mathbf{x}) = b + \sum_{l=1}^{M} \alpha_l \, y_l \, \kappa(\mathbf{x}, \mathbf{x}_l) = b + \sum_{l=1}^{M} \alpha_l \, y_l \exp(-\|\mathbf{x} - \mathbf{x}_l\|^2 / 2\sigma^2)$$

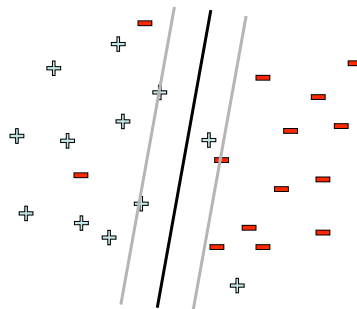## What if data is not linearly separable?

Use features of features
of features of features….

$x_1^2, x_2^2, x_1x_2, ...., exp(x_1)$

But run risk of overfitting!

---

## What if data is still not linearly separable?

Allow "error" in classification

$$\min_{\mathbf{w},b} \; \mathbf{w}^T\mathbf{w} + C \text{ #mistakes}$$
$$\text{s.t. } (\mathbf{w}^T\mathbf{x}_j+b) \, y_j \geq 1 \quad \forall j$$
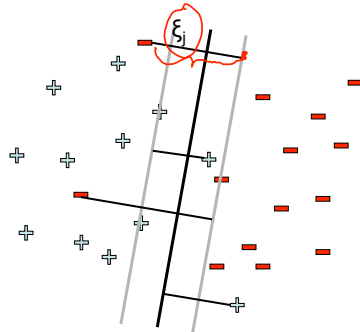
Maximize margin and minimize
# mistakes on training data

C  -  tradeoff parameter

Not QP ☹

0/1 loss (doesn't distinguish between
near miss and bad mistake)
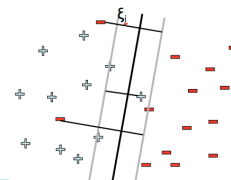
# Support Vector Machine with soft margins

*Max margin* *Mistakes over training examples*

Allow "error" in classification



**Soft margin approach**

$$\min_{\mathbf{w},b} \ \mathbf{w}^T\mathbf{w} + C \ \Sigma_j \xi_j$$

$$\text{s.t.} \ (\mathbf{w}^T\mathbf{x}_j + b)\, y_j \geq 1 - \xi_j \quad \forall j$$

$$\xi_j \geq 0 \qquad \forall j$$

$\xi_j$ - "slack" variables
   = (>1 if $x_j$ misclassifed)
 pay linear penalty if mistake

C - tradeoff parameter (chosen by cross-validation)

Still QP ☺

---

# Primal and Dual Forms for Soft Margin SVM



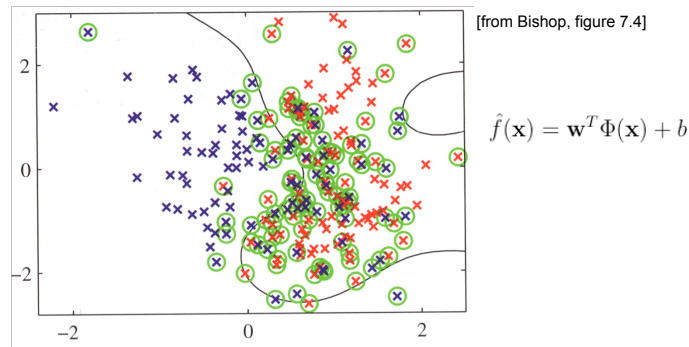Primal form: solve for $\mathbf{w}, b$ in the projected higher dim. space

$$\min_{\mathbf{w},b} \quad \tfrac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{l=1}^{M}\xi_l$$

$$\text{s.t.} \quad y_l(\mathbf{w}^T\Phi(\mathbf{x}_l) + b) \geq 1 - \xi_l \quad \forall l \in \text{ training examples}$$

$$\xi_l \geq 0 \qquad\qquad \forall l \in \text{ training examples}$$

Dual form: solve for $\alpha_1...\alpha_M$ in the original low dim. space

$$\max_{\alpha_1...\alpha_M} \quad \sum_{l=1}^{M}\alpha_l - \tfrac{1}{2}\sum_{j=1}^{M}\sum_{k=1}^{M}\alpha_j\alpha_k y_j y_k \ \kappa(\mathbf{x_j},\mathbf{x_k})$$

$$\text{s.t.} \quad 0 \leq \alpha_l \leq C \qquad \forall l \in \text{ training examples}$$

$$\sum_{l=1}^{M}\alpha_l y_l = 0$$

both are QP problems with a single local optimum ☺

## SVM Soft Margin Decision Surface using Gaussian Kernel

[from Bishop, figure 7.4]

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) + b$$

Circled points are the _support vectors_: training examples with non-zero $\alpha_l$

Points plotted in original 2-D space.

Contour lines show constant $\hat{f}(\mathbf{x})$

$$\hat{f}(\mathbf{x}) = b + \sum_{l=1}^{M} \alpha_l \ y_l \ \kappa(\mathbf{x}, \mathbf{x}_l) = b + \sum_{l=1}^{M} \alpha_l \ y_l \exp(-\|\mathbf{x} - \mathbf{x}_l\|^2 / 2\sigma^2)$$

# SVM Summary

- Objective: maximize margin between decision surface and data
- Primal and dual formulations
    - dual represents classifier decision in terms of _support vectors_
- Kernel SVM's
    - learn linear decision surface in high dimension space, working in original low dimension space
- Handling noisy data: soft margin "slack variables"
    - again primal and dual forms
- SVM algorithm: Quadratic Program optimization
    - single global minimum

# SVM: PAC Results?

---

# VC dimension: examples

What is VC dimension of

- $H_2 = \{ ((w_0 + w_1x_1 + w_2x_2) > 0 \rightarrow y=1) \}$
  - $VC(H_2)=3$
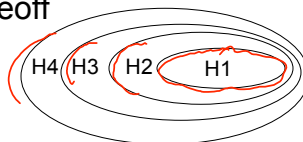- For $H_n$ = linear separating hyperplanes in n dimensions, $VC(H_n)=n+1$

$$m \geq \frac{1}{\epsilon}(4\log_2(2/\delta) + 8VC(H)\log_2(13/\epsilon))$$

# Structural Risk Minimization [Vapnik]

Which hypothesis space should we choose?
* Bias / variance tradeoff



SRM: choose H to minimize bound on expected true error!

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln \frac{2m}{VC(H)} + 1) + \ln \frac{4}{\delta}}{m}}$$

\* unfortunately a somewhat loose bound...

---

# Margin-based PAC Results

[Shawe-Taylor, Langford, McCallester]

Consider a fixed distribution $D$ on pairs $\langle x, y \rangle$ with $x \in R^d$ satisfying $||x|| = 1$ and $y \in \{-1, 1\}$. We are interested in finding a weight vector $w$ with $||w|| = 1$ such that the sign of $w \cdot x$ predicts $y$. For $\gamma > 0$ the error rate of $w$ on distribution $D$ relative to safety margin $\gamma$, denoted $\ell_\gamma(w, D)$ is defined as follows.

$$\ell_\gamma(w, D) \equiv P_{\langle x, y \rangle \sim D} [(w \cdot x)y \leq \gamma]$$

Let $S$ be a sample of $m$ pairs drawn IID from the distribution $D$. The sample $S$ can be viewed as an empirical distribution on pairs. We are interested in bounding $\ell_0(w, D)$ in terms of $\ell_\gamma(w, S)$ and the margin $\gamma$. Bartlett and Shawe-Taylor use fat shattering arguments [2] to show that with probability at least $1 - \delta$ over the choice of the sample $S$ we have the following simultaneously for all weight vectors $w$ with $||w|| = 1$ and margins $\gamma > 0$.

$$\ell_0(w, D) \leq \ell_\gamma(w, S) + 27.18 \sqrt{\frac{\log^2 m + 84}{m\gamma^2}} + O\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (1)$$

recall:

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln \frac{2m}{VC(H)} + 1) + \ln \frac{4}{\delta}}{m}}$$

9

## Maximizing Margin as an Objective Function

- We've talked about many learning algorithms, with different objective functions

- 0-1 loss
- sum sq error
- maximum log data likelihood
- MAP
- maximum margin

How are these all related?
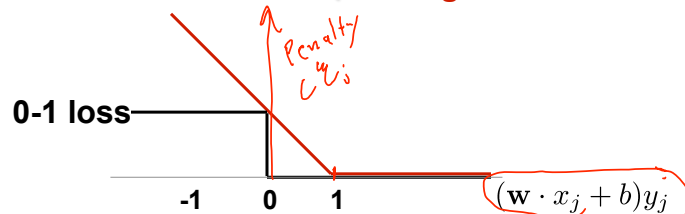
---

## Slack variables – Hinge loss

Complexity penalization

$$\xi_j = \text{loss}(f(x_j), y_j)$$

⬅

$$f(x_j) = \text{sgn}(\mathbf{w} \cdot \mathrm{x_j} + \mathrm{b})$$

$$\min_{\mathbf{w},b} \ \mathbf{w}^\mathsf{T}\mathbf{w} + C \ \Sigma \xi_j$$

$$\text{s.t. } (\mathbf{w}^\mathsf{T}\mathbf{x_j}+b) \ y_j \geq 1\text{-}\xi_j \quad \forall j$$

$$\xi_j \geq 0 \qquad \forall j$$

$$\xi_j = (1 - (\mathbf{w} \cdot x_j + b)y_j))_+$$  ⬅ **Hinge loss**

*Penalty $C\xi_j$*

**0-1 loss**

-1    0    1

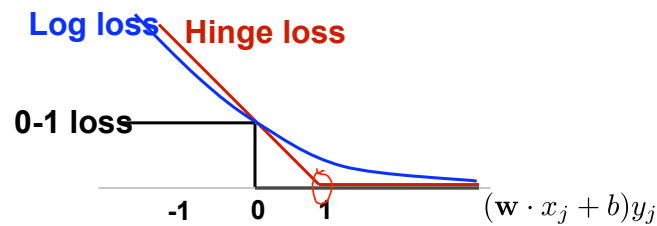$$(\mathbf{w} \cdot x_j + b)y_j$$

10

## SVM vs. Logistic Regression

SVM : **Hinge loss**

$$\text{loss}(f(x_j), y_j) = (1 - (\mathbf{w} \cdot x_j + b)y_j))_+$$

Logistic Regression : **Log loss**  ( -ve log conditional likelihood)

$$\text{loss}(f(x_j), y_j) = -\log P(y_j \mid x_j, \mathbf{w}, b) = \log(1 + e^{-(\mathbf{w} \cdot x_j + b)y_j})$$

**Log loss**   **Hinge loss**

**0-1 loss**

-1    0    1        $(\mathbf{w} \cdot x_j + b)y_j$

---

## What you need to know

Primal and Dual optimization problems

Kernel functions

Support Vector Machines

- Maximizing margin
- Derivation of SVM formulation
- Slack variables and hinge loss
- Relationship between SVMs and logistic regression
  - 0/1 loss
  - Hinge loss
  - Log loss