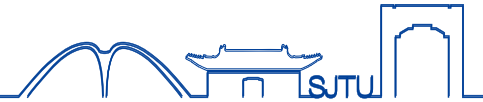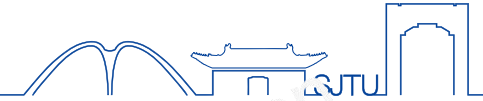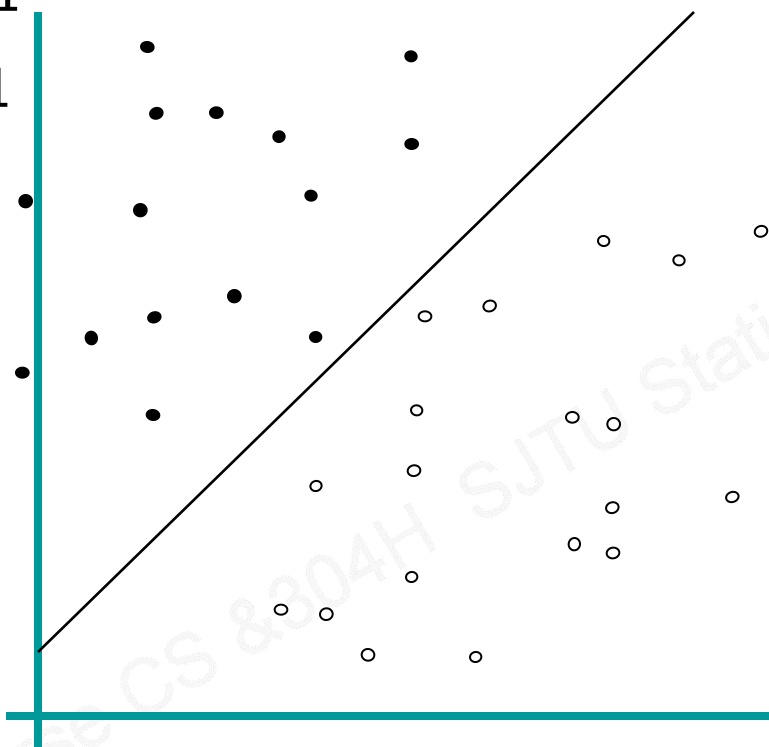# Introduction to SVMs

# SVMs

- Geometric
  - Maximizing Margin
- Kernel Methods
  - Making nonlinear decision boundaries linear
  - Efficiently!
- Capacity
  - Structural Risk Minimization

# Linear Classifiers
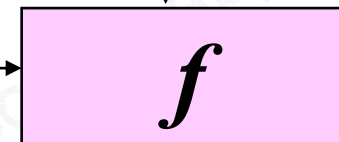
- denotes +1
○ denotes -1

How would you classify this data?

$$\alpha$$

$$\boldsymbol{x} \longrightarrow \boxed{\boldsymbol{f}} \longrightarrow y^{est}$$

$$f(\boldsymbol{x}, \boldsymbol{w}, b) = sign(\boldsymbol{w} \cdot \boldsymbol{x} - b)$$

# Linear Classifiers

- denotes +1
- denotes -1



$$f(x, w, b) = sign(w \cdot x - b)$$

Any of these would be fine..

..but which is best?

# Classifier Margin

- denotes +1
- ○ denotes -1

$\alpha$

$$\mathbf{x} \longrightarrow \boxed{f} \longrightarrow y^{est}$$

$$f(\mathbf{x},\mathbf{w},b) = sign(\mathbf{w}.\mathbf{x} - b)$$

Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin

• denotes +1

○ denotes -1

Support Vectors are those datapoints that the margin pushes up against

$\alpha$

$$f(x, w, b) = sign(w \cdot x - b)$$

$x \longrightarrow \boxed{f} \longrightarrow y^{est}$

The maximum margin linear classifier is the linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

# Why Maximum Margin?



• denotes +1

∘ denotes -1

Support Vectors are those datapoints that the margin pushes up against

$f(\mathbf{x},\mathbf{w},b) = sign(\mathbf{w}.\ \mathbf{x} - b)$

1. Intuitively this feels safest.

2. If we've made a small error in the location of the boundary (it's been jolted in its perpendicular direction) this gives us least chance of causing a misclassification.

3. There's some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing.

4. Empirically it works very very well.

# Maximizing the Margin

# Specifying a line and margin

Classify as..



Plus-Plane

Classifier Boundary

Minus-Plane

"Predict Class = +1" zone

"Predict Class = -1" zone

$wx+b=1$

$wx+b=0$

$wx+b=-1$

- Plus-plane = $\{ x : w \cdot x + b = +1 \}$
- Minus-plane = $\{ x : w \cdot x + b = -1 \}$

$+1$     if    $w \cdot x + b >= 1$

$-1$     if    $w \cdot x + b <= -1$

uncertain    if    $-1 < w \cdot x + b < 1$

# Computing the margin width



"Predict Class = +1" zone

$wx+b=1$
$wx+b=0$
$wx+b=-1$

"Predict Class = -1" zone

$M$ = Margin Width

How do we compute $M$ in terms of $w$ and $b$?

- Plus-plane  = $\{ x : w . x + b = +1 \}$
- Minus-plane = $\{ x : w . x + b = -1 \}$

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

And so of course the vector **w** is also perpendicular to the Minus Plane

Let **u** and **v** be two vectors on the Plus Plane. What is $w . ( u - v )$?

# Computing the margin width



"Predict Class = +1" zone

$wx+b=1$
$wx+b=0$
$wx+b=-1$

$M$ = Margin Width

"Predict Class = -1" zone

$x^+$

$x^-$

Plus-plane  =  $\{ x : w . x + b = +1 \}$
Minus-plane =  $\{ x : w . x + b = -1 \}$

- $x^+ = x^- + \lambda\, w$                    $|x^+ - x^-| = M$

It's now easy to get $M$ in terms of $w$ and $b$

# Computing the margin width



"Predict Class = +1" zone

$M$ = Margin Width

$\boldsymbol{x}^+$

$\boldsymbol{x}^-$

wx+b=1
wx+b=0
wx+b=-1

"Predict Class ..1" zone

**What we know:**

- $\boldsymbol{w^T\,x^+} + b = +1$

- $\boldsymbol{w^T\,x^-} + b = -1$

- $\boldsymbol{x}^+ = \boldsymbol{x}^- + \lambda\,\boldsymbol{w}$

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

It's now easy to get $M$ in terms of $\boldsymbol{w}$ and $b$

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

$$\boldsymbol{w}^{\,T}\boldsymbol{x}^- + b + \overset{\lambda=\frac{2}{\boldsymbol{w}^T\boldsymbol{w}}}{\lambda\,\boldsymbol{w}^{\,T}\boldsymbol{w}} = 1$$

➔ $-1 + \lambda\,\boldsymbol{w^T w} = 1$

# Computing the margin width

$$M = \text{Margin Width} = \frac{2}{\sqrt{\mathbf{w}^T \mathbf{w}}}$$

"Predict Class = +1" zone

$x^+$

wx+b=1

wx+b=0

wx+b=-1

"Predict Class = -1" zone

$x^-$

$$M = |\, \boldsymbol{x}^+ - \boldsymbol{x}^-\, | = |\, \lambda\, \boldsymbol{w}\, | =$$

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

## What we know:

- $\boldsymbol{w} \cdot \boldsymbol{x}^+ + b = +1$
- $\boldsymbol{w} \cdot \boldsymbol{x}^- + b = -1$
- $\boldsymbol{x}^+ = \boldsymbol{x}^- + \lambda \boldsymbol{w}$
- $|\boldsymbol{x}^+ - \boldsymbol{x}^-| = M$

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

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

# Learning the Maximum Margin Classifier

"Predict Class = +1" zone

$\mathbf{x}^{+}$

$M$ = Margin Width = $\dfrac{2}{\sqrt{\mathbf{w.w}}}$

wx+b=1

wx+b=0

wx+b=-1

"Predict Class = -1" zone

$\mathbf{x}^{-}$

Given a guess of **w** and *b* we can

- Compute whether all data points in the correct half-planes

- Compute the width of the margin

- to write 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?

# Quadratic Programming

Find

$$\arg\max_{\mathbf{u}} \quad c + \mathbf{d}^T\mathbf{u} + \frac{\mathbf{u}^T R \mathbf{u}}{2}$$

Quadratic criterion

Subject to

$$a_{11}u_1 + a_{12}u_2 + \ldots + a_{1m}u_m \leq b_1$$

$$a_{21}u_1 + a_{22}u_2 + \ldots + a_{2m}u_m \leq b_2$$

$$\vdots$$

$$a_{n1}u_1 + a_{n2}u_2 + \ldots + a_{nm}u_m \leq b_n$$

$n$ additional linear <u>in</u>equality constraints

And subject to

$$a_{(n+1)1}u_1 + a_{(n+1)2}u_2 + \ldots + a_{(n+1)m}u_m = b_{(n+1)}$$

$$a_{(n+2)1}u_1 + a_{(n+2)2}u_2 + \ldots + a_{(n+2)m}u_m = b_{(n+2)}$$

$$\vdots$$

$$a_{(n+e)1}u_1 + a_{(n+e)2}u_2 + \ldots + a_{(n+e)m}u_m = b_{(n+e)}$$

$e$ additional linear <u>e</u>quality constraints

Find

$$\arg\max_{\mathbf{u}} \quad c + \mathbf{d}^T\mathbf{u} + \frac{\mathbf{u}^T R \mathbf{u}}{2}$$

Quadratic criterion

Subject to

additional linear inequality constraints

And subject to

There exist algorithms for finding such constrained quadratic optima much more efficiently and reliably than gradient ascent.

(But they are very fiddly...you probably don't want to write one yourself)

$$a_{(n+e)1}u_1 + a_{(n+e)2}u_2 + ... + a_{(n+e)m}u_m = b_{(n+e)}$$

*e* additional linear equality constraints

# Learning the Maximum Margin Classifier

$$M = \frac{2}{\sqrt{\mathbf{w.w}}}$$

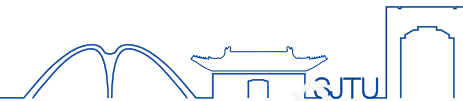"Predict Class = +1" zone

"Predict Class = -1" zone

wx+b=1

wx+b=0

wx+b=-1

Given guess of $\mathbf{w}$, $b$ we can

- Compute whether all data points are in the correct half-planes
- Compute the margin width

Assume $R$ datapoints, each $(\mathbf{x}_k, y_k)$ where $y_k = +/- \mathbf{1}$

What should our quadratic optimization criterion be?

Minimize $\mathbf{w.w}$

How many constraints will we have?
What should they be?     $R$

$\mathbf{w} . \mathbf{x}_k + b >= 1, \quad if \ y_k = 1$

$\mathbf{w} . \mathbf{x}_k + b <= -1, \quad if \ y_k = -1$

# Large-margin Decision Boundary

- The decision boundary should be as far away from the data of both classes as possible
  - We should maximize the margin, *m*
  - Distance between the origin and the line $w^t x = k$ is $k/\|\mathbf{w}\|$

$$m = \frac{2}{\|\mathbf{w}\|}$$

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

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

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

**w**

Class 2

Class 1

*m*

How many constraints will we have?

What should they be?    *R*

$$w \cdot x_k + b >= 1, \quad if \ y_k = 1$$

$$w \cdot x_k + b <= -1, \quad if \ y_k = -1$$

# Finding the Decision Boundary

- Let $\{x_1, ..., x_n\}$ be our data set and let $y_i \in \{1,-1\}$ be the class label of $x_i$
- The decision boundary should classify all points correctly $\Rightarrow$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \qquad \forall i$$

- The decision boundary can be found by solving the following constrained optimization problem

$$\text{Minimize } \frac{1}{2}\|\mathbf{w}\|^2$$

$$\text{subject to } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \qquad \forall i$$

- This is a constrained optimization problem. Solving it requires some new tools
  - Feel free to ignore the following several slides; what is important is the constrained optimization problem above

# Back to the Original Problem

Minimize $\|\mathbf{w}\|^2 / 2$

$\qquad$ subject to $1 - y_i(\mathbf{w}^T\mathbf{x}_i + b) \leq 0,$ $\qquad$ for $i = 1,\ldots,n$

- The Lagrangian is

$$\mathcal{L} = \frac{1}{2}\mathbf{w}^T\mathbf{w} + \sum_{i=1}^{n}\alpha_i\left(1 - y_i(\mathbf{w}^T\mathbf{x}_i + b)\right)$$

  – Setting the gradient of $L$ w.r.t. **w** and b to zero, we have

$$\mathbf{w} + \sum_{i=1}^{n}\alpha_i(-y_i)\mathbf{x}_i = \mathbf{0} \implies \mathbf{w} = \sum_{i=1}^{n}\alpha_i y_i\mathbf{x}_i$$

$$\sum_{i=1}^{n}\alpha_i y_i = 0$$

Illed posed problem?

- The Karush-Kuhn-Tucker conditions,
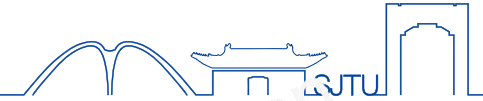
$$\alpha_i \left[ y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \right] = 0, \quad \forall i$$

  - If $\alpha_i > 0$, then $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$, or in other word, $x_i$ is on the boundary of the margin;

  - If $y_i(\mathbf{w}^T \mathbf{x}_i + b) \neq 1$, $x_i$ is not on the boundary of the margin, and $\alpha_i = 0$.

# The Dual Problem

- If we substitute $\mathbf{w} = \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i$ to $\mathcal{L}$, we have

$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i^T \sum_{j=1}^{n} \alpha_j y_j \mathbf{x}_j + \sum_{i=1}^{n} \alpha_i \left( 1 - y_i \left( \sum_{j=1}^{n} \alpha_j y_j \mathbf{x}_j^T \mathbf{x}_i + b \right) \right)$$

$$= \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^{n} \alpha_i - \sum_{i=1}^{n} \alpha_i y_i \sum_{j=1}^{n} \alpha_j y_j \mathbf{x}_j^T \mathbf{x}_i - b \sum_{i=1}^{n} \alpha_i y_i$$

$$= - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^{n} \alpha_i$$

$$\sum_{i=1}^{n} \alpha_i y_i = 0$$

- This is a function of $\alpha_i$ only

- It is known as **the dual problem**: if we know **w**, we know all $\alpha_i$; if we know all $\alpha_i$, we know **w**
- The objective function of the dual problem needs to **be maximized!**

$$\text{max. } W(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^{n} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{subject to } \alpha_i \geq 0, \quad \sum_{i=1}^{n} \alpha_i y_i = 0$$

Properties of $\alpha_i$ when we introduce the Lagrange multipliers

The result when we differentiate the original Lagrangian w.r.t. b

# The Dual Problem

$$\text{max. } W(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^{n} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

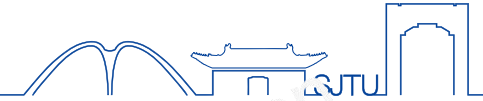$$\text{subject to } \alpha_i \geq 0, \sum_{i=1}^{n} \alpha_i y_i = 0$$

- This is a quadratic programming (QP) problem
  - A global maximum of $\alpha_i$ can always be found

- $\mathbf{w}$ can be recovered by

$$\mathbf{w} = \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i$$

# Characteristics of the Solution

- Many of the $\alpha_i$ are zero
  - **w** is a linear combination of a small number of data points

- $\mathbf{x}_i$ with non-zero $\alpha_i$ are called support vectors (SV)
  - Let $t_j$ $(j=1, ..., s)$ be the indices of the $s$ support vectors. We can write

$$\mathbf{w} = \sum_{j=1}^{S} \alpha_{t_j} y_{t_j} \mathbf{x}_{t_j}$$

- For testing with a new data **z**

  - classify **z** as class 1 if the sum is positive, and class 2 otherwise

$$\mathbf{w}^T \mathbf{z} + b = \sum_{j=1}^{S} \alpha_{t_j} y_{t_j} (\mathbf{x}_{t_j}^T \mathbf{z}) + b$$

# A Geometrical Interpretation



Class 2

$\alpha_{10}=0$

$\alpha_8=0.6$

$\mathbf{W}$

$\alpha_7=0$

$\alpha_2=0$

$\alpha_5=0$

$\alpha_1=0.8$

$\alpha_4=0$

$\alpha_6=1.4$

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

$\alpha_9=0$

$\alpha_3=0$

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

Class 1

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

# Non-linearly Separable Problems

- We allow "error" $\xi_i$ in classification; it is based on the output of the discriminant function $\boldsymbol{w}^T\boldsymbol{x}+b$

- $\xi_i$ approximates the number of misclassified samples

$$\xi_j$$

$$\mathbf{x}_j$$

$$\text{Class 2}$$

$$\mathbf{W}$$

$$\mathbf{x}_i$$

$$\xi_i$$

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

$$\text{Class 1}$$

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

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

# Learning Maximum Margin with Noise



The quadratic optimization criterion:

Minimize

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

How many constraints will we have? *R*

Constraint:

$$\mathbf{w}^T\mathbf{x}_k + b \geq 1 - \varepsilon_k, \quad \text{if } y_k = 1$$

$$\mathbf{w}^T\mathbf{x}_k + b \leq -1 + \varepsilon_k, \text{if } y_k = -1$$

# Learning Maximum Margin with Noise

Our original (noiseless data) QP had $m+1$ variables: $w_1, w_2, \dots w_m,$ and $b.$

Our new (noisy data) QP has $m+1+R$ variables: $w_1, w_2, \dots w_m, b, \varepsilon_k, \varepsilon_1, \dots \varepsilon_R$

The quadratic optimization criterion:

Minimize

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

How many constraints will we have? $2R$

Constraint:

$$y_k\left(\mathbf{w}^T\mathbf{x}_k + b\right) \geq 1 - \varepsilon_k,$$

$$\varepsilon_k \geq 0, \qquad \forall\ k$$

# An Equivalent Dual QP

Minimize

$$\frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k \qquad y_k\left(\mathbf{w}^T\mathbf{x}_k + b\right) \geq 1-\varepsilon_k, \quad \forall\ k$$

$$\varepsilon_k \geq 0, \quad \forall k$$

The Lagrange function:

$$L_p(w,b,\varepsilon) = \frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k - \sum_{k=1}^{R}\mu_k\varepsilon_k -$$

$$- \sum_{k=1}^{R}\alpha_k\left[y_k\left(\mathbf{w}^T\mathbf{x}_k + b\right) - \left(1-\varepsilon_k\right)\right]$$

The Lagrange function:

$$L_p(w,b,\varepsilon) = \frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k - \sum_{k=1}^{R}\mu_k\varepsilon_k$$

$$- \sum_{k=1}^{R}\alpha_k\left[ y_k\left(\mathbf{w}^T\mathbf{x}_k + b\right) - \left(1\text{-}\varepsilon_k\right)\right]$$

- Setting the respective derivatives to zero, we get

$$\mathbf{w} = \sum_{k=1}^{R}\alpha_k y_k \mathbf{x}_k, \qquad \sum_{k=1}^{R}\alpha_k y_k = 0,$$

$$\alpha_i = C - \mu_i, \forall i, \qquad \alpha_i, \mu_i, \varepsilon_i > 0, \forall i$$
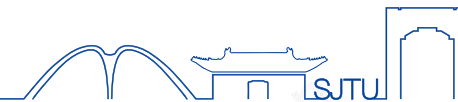
# An Equivalent Dual QP

Minimize

$$\frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k \qquad y_k\left(\mathbf{w}^T\mathbf{x}_k + b\right) \geq 1-\varepsilon_k, \quad \forall\ k$$

$$\varepsilon_k \geq 0, \quad \forall k$$

Dual QP

$$\text{Maximize}\sum_{k=1}^{R}\alpha_k - \frac{1}{2}\sum_{k=1}^{R}\sum_{l=1}^{R}\alpha_k\alpha_l Q_{kl}\text{ where }\quad Q_{kl} = y_k y_l(\mathbf{x}_k.\mathbf{x}_l)$$

Subject to these constraints: $\quad 0 \leq \alpha_k \leq C \quad \forall k \qquad \sum_{k=1}^{R}\alpha_k y_k = 0$

# An Equivalent Dual QP

$$\text{Maximize} \sum_{k=1}^{R} \alpha_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{l=1}^{R} \alpha_k \alpha_l Q_{kl} \text{ where } Q_{kl} = y_k y_l (\mathbf{x}_k . \mathbf{x}_l)$$

Subject to these constraints:
$$0 \le \alpha_k \le C \quad \forall k \qquad \sum_{k=1}^{R} \alpha_k y_k = 0$$

Then define:

$$\mathbf{w} = \sum_{k=1}^{R} \alpha_k y_k \mathbf{x}_k$$

$$b = y_K (1 - \varepsilon_K) - \mathbf{x}_K . \mathbf{w}_K$$

$$\text{where} \quad K = \arg \max_k \alpha_k$$

Then classify with:

$$f(\pmb{x},\pmb{w},b) = sign(\pmb{w}. \pmb{x} - b)$$

Linear separability in threshold perceptrons. Black dots indicate a point in the input space where the value of the function is 1, and white dots indicate a point where the value is 0. The perceptron returns 1 on the region on the non-shaded side of the line. In (c), no such line exists that correctly classifies the inputs.

## Example XOR problem revisited:

Let the nonlinear mapping be :

$$\varphi(\mathbf{x}) = \left(1, \sqrt{2}\,x_1, \sqrt{2}\,x_2, x_1^2, \sqrt{2}\,x_1\,x_2, x_2^2\right)^T ; \mathbf{x} = (x_1, x_2)^T$$

And: $\varphi(\mathbf{x}_i) = \left(1, \sqrt{2}\,x_{i1}, \sqrt{2}\,x_{i2}, x_{i1}^2, \sqrt{2}\,x_{i1}\,x_{i2}, x_{ii2}^2\right)^T$

Therefore the feature space is in 6D with input data in 2D

$\underline{x}_1 = (-1,-1), \quad d_1 = -1$

$\underline{x}_2 = (-1, 1), \quad d_2 = 1$

$\underline{x}_3 = (1,-1), \quad d_3 = 1$

$\underline{x}_4 = (-1,-1), \quad d_4 = -1$

$Q(a) = \Sigma a_i - \frac{1}{2} \Sigma \Sigma a_i a_j d_i d_j \phi(\underline{x}_i)^T \phi(\underline{x}_j)$

$= a_1 + a_2 + a_3 + a_4 - \frac{1}{2}(9 a_1 a_1 - 2a_1 a_2 - 2 a_1 a_3 + 2a_1 a_4$

$+ 9a_2 a_2 + 2a_2 a_3 - 2a_2 a_4 + 9a_3 a_3 - 2a_3 a_4 + 9 a_4 a_4)$

To minimize Q, calculate

$$\frac{\partial Q(a)}{\partial a_i} = 0, i = 1, ..., 4$$

(due to optimality conditions) which gives

$1 = 9 a_1 - a_2 - a_3 + a_4$

$1 = -a_1 + 9 a_2 + a_3 - a_4$

$1 = -a_1 + a_2 + 9 a_3 - a_4$

$1 = a_1 - a_2 - a_3 + 9 a_4$

The solution of which gives the optimal values:

$$a_{0,1} = a_{0,2} = a_{0,3} = a_{0,4} = 1/8$$

$$\underline{w}_0 = \Sigma\, a_{0,i}\, d_i\, \phi(\underline{x}_i)$$
$$= 1/8[\phi(\underline{x}_1) - \phi(\underline{x}_2) - \phi(\underline{x}_3) + \phi(\underline{x}_4)]$$

$$= \frac{1}{8}\left[ -\begin{pmatrix} 1 \\ 1 \\ \sqrt{2} \\ 1 \\ -\sqrt{2} \\ -\sqrt{2} \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ -\sqrt{2} \\ 1 \\ -\sqrt{2} \\ \sqrt{2} \end{pmatrix} - \begin{pmatrix} 1 \\ 1 \\ -\sqrt{2} \\ 1 \\ \sqrt{2} \\ -\sqrt{2} \end{pmatrix} - \begin{pmatrix} 1 \\ 1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \\ \sqrt{2} \end{pmatrix} \right] = \begin{pmatrix} 0 \\ 0 \\ -1\!\!\Big/\!\sqrt{2} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

where the first element of $w_0$ gives the bias b

From earlier we have that the optimal hyperplane is defined by:

$$\underline{w_0}^T \phi(\underline{x}) = 0$$

That is:

$$\underline{w_0}^T \phi(\underline{x}) = \begin{pmatrix} 0 & 0 & -\dfrac{1}{\sqrt{2}} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ x_1^2 \\ \sqrt{2}x_1 x_2 \\ x_2^2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \end{pmatrix} = -x_1 x_2 = 0$$

which is the optimal decision boundary for the XOR problem. Furthermore we note that the solution is unique since the optimal decision boundary is unique

Output for polynomial RBF


Outputs from Linear Rbf Net


Outputs from Gaussian Rbf Net

For a non-linearly separable problem we have to first map data onto feature space so that they are linear separable

$$x_i \rightarrow \varphi(x_i) \qquad \text{sample}: \qquad \{(\varphi(x_i), y_i)\}_{i=1}^{R}$$

Given the training data sample $\{(\underline{x}_i, y_i), i=1, ...,N\}$, find the optimum values of the weight vector w and bias b

$$w = \sum_i a_i y_i \varphi(x_i)$$

where $a_{0,i}$ are the optimal Lagrange multipliers determined by maximizing the following objective function

$$Q(a) = \sum_{i=1}^{N} a_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} a_i a_j d_i d_j \underline{\phi}^T(\underline{x}_i) \underline{\phi}(\underline{x}_j)$$

subject to the constraints

$$\sum_{j=1}^{N} a_i y_i = 0; \quad a_i \geq 0, \forall i$$

**SVM building procedure:**

1. Pick a nonlinear mapping f
2. Solve for the optimal weight vector

However: how do we pick the function $f$?

- In practical applications, if it is not totally impossible to find f, it is very hard
- In the previous example, the function f is quite complex: How would we find it?

Answer: *the Kernel Trick*

Notice that in the dual problem the inner product of input vectors is replaced by an inner-product kernel

$$Q(a) = \sum_{i=1}^{N} a_i - \frac{1}{2} \sum_{i,j=1}^{N} a_i a_j d_i d_j \varphi_j(\underline{x})^T \varphi_j(\underline{x}_i)$$

$$= \sum_{i=1}^{N} a_i - \frac{1}{2} \sum_{i,j=1}^{N} a_i a_j d_i d_j K(\underline{x}_i, \underline{x}_j)$$

How do we relate the output of the SVM to the kernel K?

Look at the equation of the boundary in the feature space and use the optimality conditions derived from the Lagrangian formulations

*Hyperplane is defined by*

$$\sum_{j=1}^{m1} w_j \varphi_j(\underline{x}) + b = 0$$

*or*

$$\sum_{j=0}^{m1} w_j \varphi_j(\underline{x}) = 0; \quad where \; \varphi_0(\underline{x}) = 1$$

*writing* : $\underline{\varphi}(\underline{x}) = [\varphi_0(\underline{x}), \varphi_1(\underline{x}), ..., \varphi_{m1}(\underline{x})]$

*we get* : $\underline{w}^T \underline{\varphi}(\underline{x}) = 0$

*from optimality conditions* : $\underline{w} = \sum_{i=1}^{N} a_i d_i \underline{\varphi}(\underline{x}_i)$

*Thus* the decision boundary:

$$\sum_{i=1}^{N} a_i d_i \underline{\varphi}^T (\underline{x}_i) \underline{\varphi}(\underline{x}) = 0; \quad \text{or} \quad \sum_{i=1}^{N} a_i d_i K(\underline{x}, \underline{x}_i) = 0$$

*Now* the decision output:

$$y = \underline{w}^T \underline{\varphi}(\underline{x}) = \sum_{i=1}^{N} a_i d_i K(\underline{x}, \underline{x}_i)$$

*where* $: K(\underline{x}, \underline{x}_i) = \underline{\varphi}^T (\underline{x}_i) \underline{\varphi}(\underline{x})$

We therefore only need a suitable choice of kernel function cf: Mercer's Theorem:

Let $K(\underline{x},\underline{y})$ be a continuous symmetric kernel that defined in the closed interval [a,b]. The kernel K can be expanded in the form

$$K(\underline{x},\underline{y})=\varphi^T(\underline{y})\varphi(\underline{x})$$

provided it is positive definite. Some of the usual choices for K are:

Polynomial SVM $\quad (\underline{x}^T\underline{x_i}+1)^p \quad$ p specified by user

RBF SVM $\quad \exp(-1/(2\sigma^2)\|\underline{x}-\underline{x_i}\|^2) \quad \sigma$ specified by user

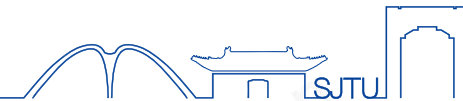MLP SVM $\quad \tanh(s_0\underline{x}^T\underline{x_i}+s_1)$

# Complexity Analysis

- Why the kernel approach provides feasible solution for implementing SVM?

# An Equivalent Dual QP

Maximize $\displaystyle\sum_{k=1}^{R} \alpha_k - \frac{1}{2}\sum_{k=1}^{R}\sum_{l=1}^{R}\alpha_k\alpha_l Q_{kl}$   where   $Q_{kl} = y_k y_l (\mathbf{x}_k . \mathbf{x}_l)$

Constraints:   $0 \le \alpha_k \le C \quad \forall k$   $\displaystyle\sum_{k=1}^{R} \alpha_k y_k = 0$

Then define:

$$\mathbf{w} = \sum_{k=1}^{R} \alpha_k y_k \mathbf{x}_k$$

$$b = y_K (1 - \varepsilon_K) - \mathbf{x}_K . \mathbf{w}_K$$

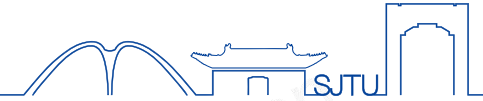$$\text{where} \quad K = \arg\max_k \alpha_k$$

Datapoints with $a_k > 0$ will be the support vectors

Then classify with:

$f(x,\textbf{\textit{w}},b) = sign(\textbf{\textit{w}}. \textbf{\textit{x}} - b)$

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

$$\Phi(\mathbf{x}) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \vdots \\ \sqrt{2}x_m \\ x_1^2 \\ x_2^2 \\ \vdots \\ x_m^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \vdots \\ \sqrt{2}x_1x_m \\ \sqrt{2}x_2x_3 \\ \vdots \\ \sqrt{2}x_1x_m \\ \vdots \\ \sqrt{2}x_{m-1}x_m \end{pmatrix}$$

Constant Term

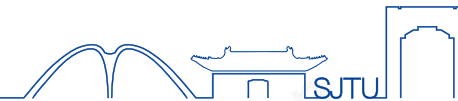Linear Terms

Pure Quadratic Terms

Quadratic Cross-Terms

# Quadratic Basis Functions

Number of terms (assuming m input dimensions) = (m+2)(m+1)/2

= (as near as makes no difference) $m^2/2$

You may be wondering what those $\sqrt{2}$'s are doing.

# QP with basis functions

Maximize $\sum_{k=1}^{R} \alpha_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{l=1}^{R} \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\mathbf{\Phi}(\mathbf{x}_k) . \mathbf{\Phi}(\mathbf{x}_l))$

Constraints: $0 \le \alpha_k \le C \quad \forall k \qquad \sum_{k=1}^{R} \alpha_k y_k = 0$

Then define:

$$\mathbf{w} = \sum_{k \text{ s.t. } \alpha_k > 0} \alpha_k y_k \mathbf{\Phi}(\mathbf{x}_k)$$

Then classify with:

$f(\mathbf{x}, \mathbf{w}, b) = sign(\mathbf{w} . \phi(\mathbf{x}) - b)$

$$b = y_K (1 - \varepsilon_K) - \mathbf{x}_K . \mathbf{w}_K$$

$$\text{where} \quad K = \arg \max_{k} \alpha_k$$

$$\mathbf{\Phi(a)}^T\mathbf{\Phi(b)} = \begin{pmatrix} 1 \\ \sqrt{2}a_1 \\ \sqrt{2}a_2 \\ \vdots \\ \sqrt{2}a_m \\ a_1^2 \\ a_2^2 \\ \vdots \\ a_m^2 \\ \sqrt{2}a_1a_2 \\ \sqrt{2}a_1a_3 \\ \vdots \\ \sqrt{2}a_1a_m \\ \sqrt{2}a_2a_3 \\ \vdots \\ \sqrt{2}a_1a_m \\ \vdots \\ \sqrt{2}a_{m-1}a_m \end{pmatrix} \bullet \begin{pmatrix} 1 \\ \sqrt{2}b_1 \\ \sqrt{2}b_2 \\ \vdots \\ \sqrt{2}b_m \\ b_1^2 \\ b_2^2 \\ \vdots \\ b_m^2 \\ \sqrt{2}b_1b_2 \\ \sqrt{2}b_1b_3 \\ \vdots \\ \sqrt{2}b_1b_m \\ \sqrt{2}b_2b_3 \\ \vdots \\ \sqrt{2}b_1b_m \\ \vdots \\ \sqrt{2}b_{m-1}b_m \end{pmatrix}$$

$$\left.\vphantom{1}\right\} 1$$

$$+$$

$$\left.\vphantom{\sum}\right\} \sum_{i=1}^{m} 2a_i b_i$$

$$+$$

$$\left.\vphantom{\sum}\right\} \sum_{i=1}^{m} a_i^2 b_i^2$$

# Quadratic
$+$ Dot Products

$$\left.\vphantom{\sum}\right\} \sum_{i=1}^{m} \sum_{j=i+1}^{m} 2a_i a_j b_i b_j$$

# Quadratic Dot Products

$$\mathbf{\Phi(a)} \bullet \mathbf{\Phi(b)} =$$

$$1 + 2\sum_{i=1}^{m} a_i b_i + \sum_{i=1}^{m} a_i^2 b_i^2 + \sum_{i=1}^{m}\sum_{j=i+1}^{m} 2a_i a_j b_i b_j$$

Just out of casual, innocent, interest, let's look at another function of $\boldsymbol{a}$ and $\boldsymbol{b}$:

$$(\mathbf{a.b} + 1)^2$$

$$= (\mathbf{a.b})^2 + 2\mathbf{a.b} + 1$$

$$= \left(\sum_{i=1}^{m} a_i b_i\right)^2 + 2\sum_{i=1}^{m} a_i b_i + 1$$

$$= \sum_{i=1}^{m}\sum_{j=1}^{m} a_i b_i a_j b_j + 2\sum_{i=1}^{m} a_i b_i + 1$$

$$= \sum_{i=1}^{m} (a_i b_i)^2 + 2\sum_{i=1}^{m}\sum_{j=i+1}^{m} a_i b_i a_j b_j + 2\sum_{i=1}^{m} a_i b_i + 1$$

# Quadratic Dot Products

Just out of casual, innocent, interest, let's look at another function of **a** and **b**:

$$K(a,b) = (\mathbf{a}.\mathbf{b} + 1)^2$$

$$= (\mathbf{a}.\mathbf{b})^2 + 2\mathbf{a}.\mathbf{b} + 1$$

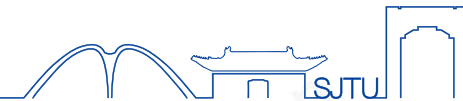$$= \left(\sum_{i=1}^{m} a_i b_i\right)^2 + 2\sum_{i=1}^{m} a_i b_i + 1$$

$$= \sum_{i=1}^{m}\sum_{j=1}^{m} a_i b_i a_j b_j + 2\sum_{i=1}^{m} a_i b_i + 1$$

$$= \sum_{i=1}^{m}(a_i b_i)^2 + 2\sum_{i=1}^{m}\sum_{j=i+1}^{m} a_i b_i a_j b_j + 2\sum_{i=1}^{m} a_i b_i + 1$$

$$\mathbf{\Phi(a)} \bullet \mathbf{\Phi(b)} =$$

$$1 + 2\sum_{i=1}^{m} a_i b_i + \sum_{i=1}^{m} a_i^2 b_i^2 + \sum_{i=1}^{m}\sum_{j=i+1}^{m} 2a_i a_j b_i b_j$$

They're the same!

And this is only O(m) to compute!

# QP with basis functions

Maximize $\sum_{k=1}^{R} \alpha_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{l=1}^{R} \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\mathbf{\Phi}(\mathbf{x}_k).\mathbf{\Phi}(\mathbf{x}_l))$

Constraints: $\qquad 0 \le \alpha_k \le C \quad \forall k \qquad \sum_{k=1}^{R} \alpha_k y_k = 0$

Then define:

$$\mathbf{w} = \sum_{k \text{ s.t. } \alpha_k > 0} \alpha_k y_k \mathbf{\Phi}(\mathbf{x}_k)$$

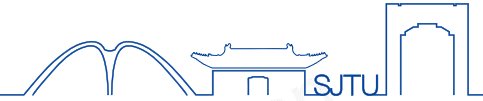$$b = y_K (1 - \varepsilon_K) - \mathbf{x}_K.\mathbf{w}_K$$

$$\text{where} \quad K = \arg \max_{k} \alpha_k$$

We must do $R^2/2$ dot products to get this matrix ready.

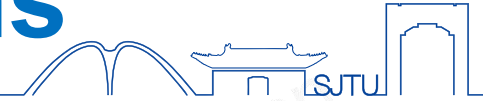Each dot product requires $m^2/2$ additions and multiplications

The whole thing costs $R^2 m^2 /4$.

# Higher Order Polynomials

| Polynomial | $\phi(x)$ | Cost to build $Q_{kl}$ matrix traditionally | Cost $M=100$ | $\phi(a).\phi(b)$ | Cost to build $Q_{kl}$ matrix efficiently | Cost if 100 inputs |
|---|---|---|---|---|---|---|
| Quadratic | All $m^2/2$ terms up to degree 2 | $m^2 R^2 /4$ | $2{,}500\ R^2$ | $(a.b+1)^2$ | $m R^2 / 2$ | $50\ R^2$ |
| Cubic | All $m^3/6$ terms up to degree 3 | $m^3 R^2 /12$ | $83{,}000\ R^2$ | $(a.b+1)^3$ | $m R^2 / 2$ | $50\ R^2$ |
| Quartic | All $m^4/24$ terms up to degree 4 | $m^4 R^2 /48$ | $1{,}960{,}000\ R^2$ | $(a.b+1)^4$ | $m R^2 / 2$ | $50\ R^2$ |

# QP with Quintic basis functions

where $Q_{kl} = y_k y_l \left( \Phi(x_k) \cdot \Phi(x_l) \right)$

We must do $R^2/2$ dot products to get this matrix ready.

In 100-d, each dot product now needs 103 operations instead of 75 million

But there are still worrying things lurking away. What are they?

The use of Maximum Margin magically makes this not a problem
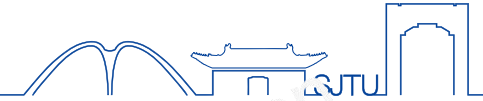
$$\forall k \qquad \sum \alpha_k y_k = 0$$

- The fear of overfitting with this enormous number of terms

- The evaluation phase (doing a set of predictions on a test set) will be very expensive (why?)

Then define:

$$\mathbf{w} = \sum_{k \ \text{s.t.} \ \alpha_k > 0} \alpha_k y_k \Phi(\mathbf{x}_k)$$

Because each **w**. φ(**x**) (see below) needs 75 million operations. *What can be done?*

$$b = y_K (1 - \varepsilon_K) - \mathbf{x}_K^T \mathbf{w}_K$$

$$\text{where } K = \arg \max_k \alpha_k$$

Then classify with:

*f(x,w,b) = sign(w.* **φ***(x) - b)*

# SVM Kernel Functions

- $K(a,b)=(a \cdot b +1)^d$ is an example of an SVM Kernel Function
- Beyond polynomials there are other very high dimensional basis functions that can be made practical by finding the right Kernel Function
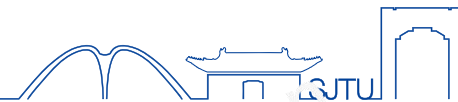
  – Radial-Basis-style Kernel Function:

  $$K(a,b) = exp\left( - \frac{(a-b)^2}{2\sigma^2} \right)$$

  – Neural-net-style Kernel Function:

  $$K(a,b) = tanh(\kappa a.b - \delta)$$

$\sigma$, $\kappa$ and $\delta$ are magic parameters that must be chosen by a model selection method such as CV or VCSRM

*The End of Talk*