# Support Vector Machine
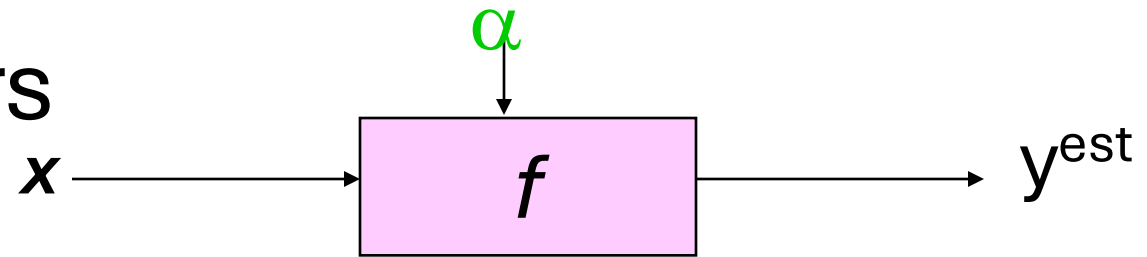
Thushari Silva, PhD

Professor in AI

Department of Computational Mathematics
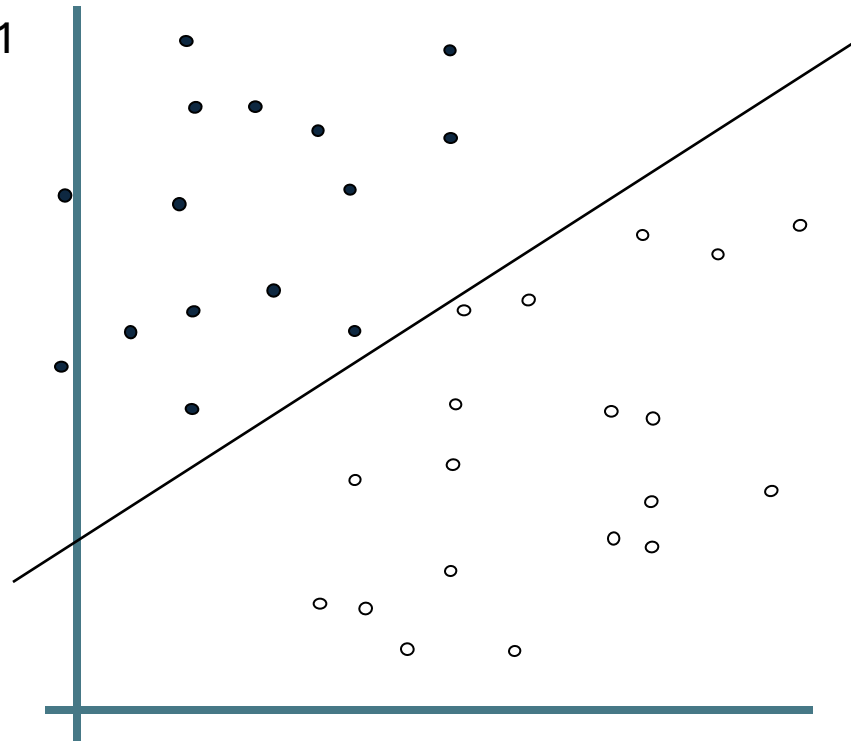
University of Moratuwa

# Linear Classifiers
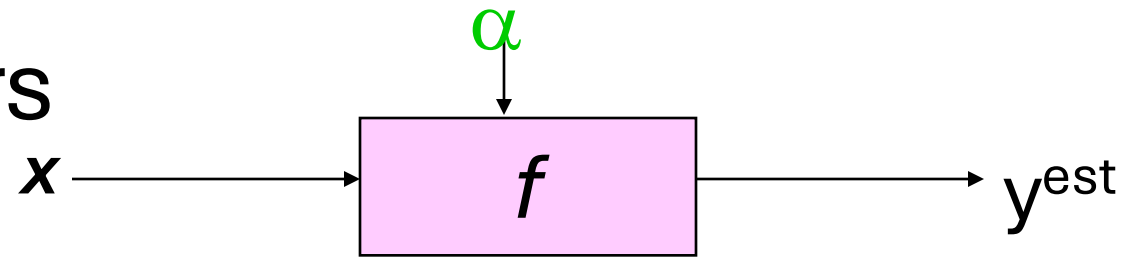
$$\alpha$$

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

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

• denotes +1

∘ denotes -1

How would you classify this data?

# Linear Classifiers

$$\alpha$$

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

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

- denotes +1
- denotes -1

How would you classify this data?

# Linear Classifiers

$\alpha$

$x \longrightarrow$ | $f$ | $\longrightarrow y^{est}$
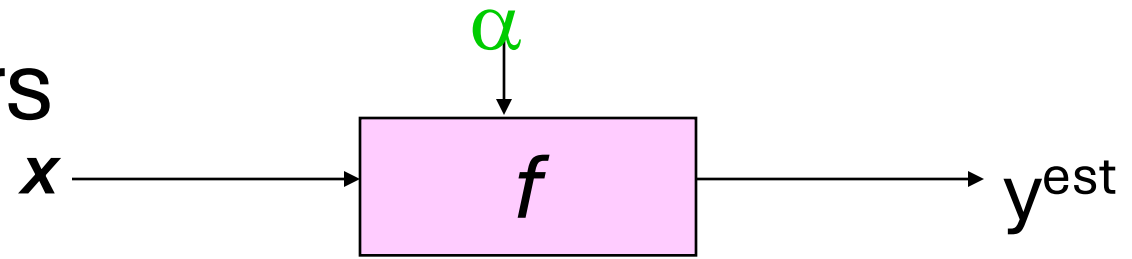
$f(x, w, b) = sign(w \cdot x - b)$

- denotes +1
- denotes -1

Any of these would be fine..

..but which is best?

# Classifier Margin



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

- • denotes +1
- ○ denotes -1

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

# Maximum Margin

$\alpha$

$x \longrightarrow$ [ $f$ ] $\longrightarrow y^{est}$

$f(x,w,b) = sign(w. x - b)$

- denotes +1
○ denotes -1

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

# Maximum Margin



$\alpha$

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

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

- denotes +1
- denotes -1

**Support Vectors** are those datapoints that the margin pushes up against

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

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. LOOCV is easy since the model is immune to removal of any non-support-vector datapoints.

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

5. Empirically it works very very well.

# Specifying a line and margin



- How do we represent this mathematically?
- ...in $m$ input dimensions?

# Specifying a line and margin

"Predict Class = +1" zone

"Predict Class = -1" zone

Plus-Plane

Classifier Boundary

Minus-Plane

$wx+b=1$

$wx+b=0$

$wx+b=-1$

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

Classify as..

| | | |
|---|---|---|
| +1 | if | $w \cdot x + b >= 1$ |
| -1 | if | $w \cdot x + b <= -1$ |
| Universe explodes | 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 $\boldsymbol{w}$ and $b$?

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

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

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

# 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 $\boldsymbol{w}$ and $b$?

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

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

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

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

# Computing the margin width



$M$ = Margin Width

"Predict Class = +1" zone

"Predict Class = -1" zone

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

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

- Plus-plane $= \{ x : w . x + b = +1 \}$
- Minus-plane $= \{ x : w . x + b = -1 \}$
- The vector **w** is perpendicular to the Plus Plane
- Let $x^-$ be any point on the minus plane
- Let $x^+$ be the closest plus-plane-point to $x$.

Any location in $R^m$: not necessarily a datapoint

# Computing the margin width



$M$ = Margin Width

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

"Predict Class = +1" zone

"Predict Class = -1" zone
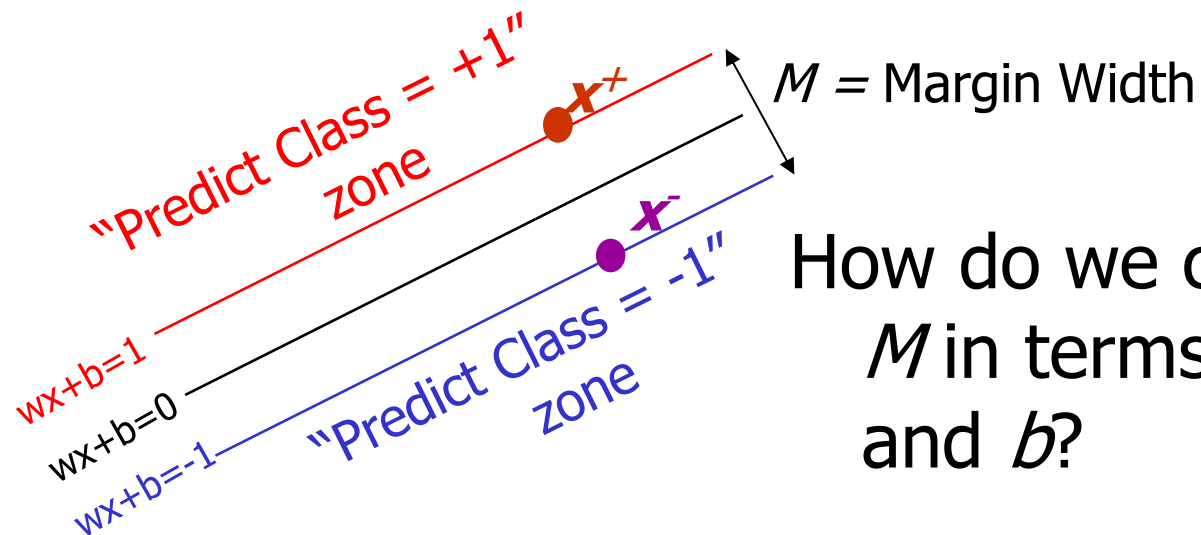
$\mathbf{x}^+$

$\mathbf{x}^-$

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

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

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

- The vector **w** is perpendicular to the Plus Plane

- Let $\mathbf{x}^-$ be any point on the minus plane

- Let $\mathbf{x}^+$ be the closest plus-plane-point to $\mathbf{x}^-$.

Any location in $R^m$: not necessarily a datapoint

# Computing the margin width



$M$ = Margin Width

How do we compute $M$ in terms of $\boldsymbol{w}$ and $b$?

- Plus-plane = $\{\boldsymbol{x} : \boldsymbol{w} \cdot \boldsymbol{x} + b = +1\}$
- Minus-plane = $\{\boldsymbol{x} : \boldsymbol{w} \cdot \boldsymbol{x} + b = -1\}$
- The vector **w** is perpendicular to the Plus Plane
- Let $\boldsymbol{x}^-$ be any point on the minus plane
- Let $\boldsymbol{x}^+$ be the closest plus-plane-point to $\boldsymbol{x}^-$.
- Claim: $\boldsymbol{x}^+ = \boldsymbol{x}^- + \lambda \boldsymbol{w}$ for some value of $\lambda$. Why?

# Computing the margin width



$M$ = Margin Width

The line from $x^-$ to $x^+$ is perpendicular to the planes.

So to get from $x^-$ to $x^+$ travel some distance in direction $w$.

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

# Computing the margin width

"Predict Class = +1" zone

$\boldsymbol{x^+}$

M = Margin Width

$wx+b=1$

$wx+b=0$

$wx+b=-1$

$\boldsymbol{x^-}$

"Predict Class = -1" zone

What we know:

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

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

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

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

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

# Computing the margin width



"Predict Class = +1" zone

$M$ = Margin Width

$x^+$

wx+b=1

wx+b=0

wx+b=-1

$x^-$

"Predict Class = -1" zone

What we know:

- $w \cdot x^+ + b = +1$

- $w \cdot x^- + b = -1$

- $x^+ = x^- + \lambda w$

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

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

$w \cdot (x^- + \lambda w) + b = 1$

=>

$w \cdot x^- + b + \lambda w \cdot w = 1$

=>

$-1 + \lambda w \cdot w = 1$

=>

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

# Computing the margin width

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

"Predict Class = +1" zone
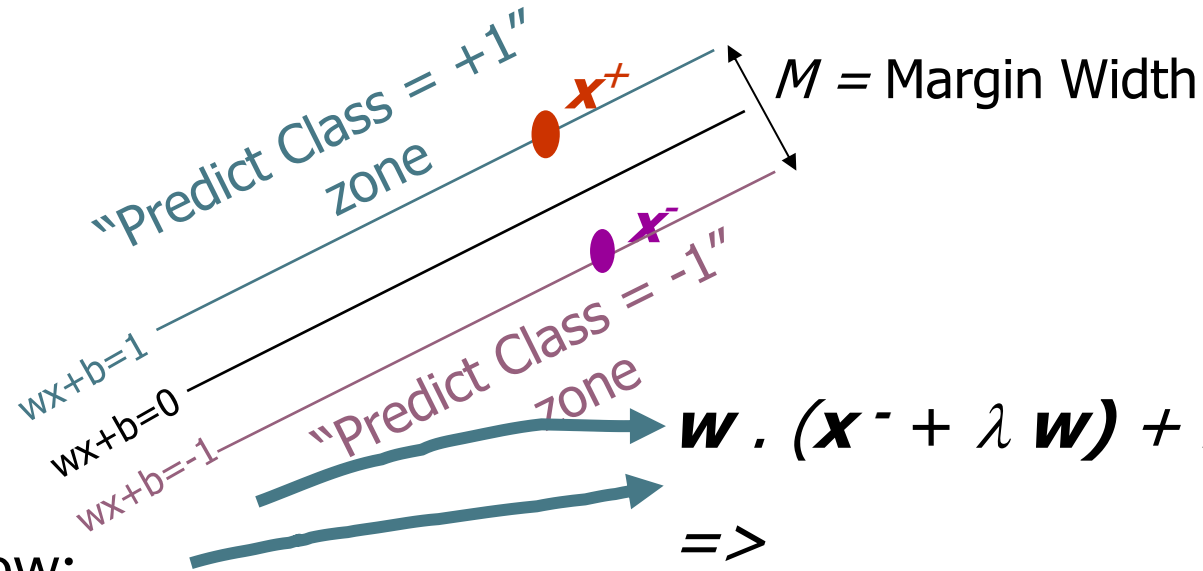
"Predict Class = -1" zone

wx+b=1

wx+b=0

wx+b=-1

**What we know:**

- $\mathbf{w} . \mathbf{x}^+ + b = +1$
- $\mathbf{w} . \mathbf{x}^- + b = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M$
- $\lambda = \dfrac{2}{\mathbf{w}.\mathbf{w}}$

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

$= \lambda |\mathbf{w}| = \lambda \sqrt{\mathbf{w}.\mathbf{w}}$

$= \dfrac{2\sqrt{\mathbf{w}.\mathbf{w}}}{\mathbf{w}.\mathbf{w}} = \dfrac{2}{\sqrt{\mathbf{w}.\mathbf{w}}}$

# Learning the Maximum Margin Classifier

"Predict Class = +1" zone

$\bullet$ $\boldsymbol{x}^+$

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

$wx+b=1$

$wx+b=0$

$wx+b=-1$

$\bullet$ $\boldsymbol{x}^-$

"Predict Class = -1" zone

Given a guess of $\boldsymbol{w}$ and $b$ we can

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

- Compute the width of the margin

So now we just need to write a program to search the space of $\mathbf{w}$'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?

# Learning via Quadratic Programming

- QP is a well-studied class of optimization algorithms to maximize a quadratic function of some real-valued variables subject to linear constraints.

# Quadratic Programming

Find $\quad \underset{\mathbf{u}}{\arg\max} \quad c + \mathbf{d}^T\mathbf{u} + \dfrac{\mathbf{u}^T R\mathbf{u}}{2}$ ← Quadratic criterion

Subject to

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

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

$$\vdots$$

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

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

# Learning the Maximum Margin Classifier

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

"Predict Class = +1" zone

"Predict Class = -1" zone

wx+b=1

wx+b=0

wx+b=-1

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

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

- Compute the margin width

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

What should our quadratic optimization criterion be?

How many constraints will we have?

What should they be?

# Learning the Maximum Margin Classifier

"Predict Class = +1" zone

wx+b=1

wx+b=0

wx+b=-1

"Predict Class = -1" zone

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

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

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

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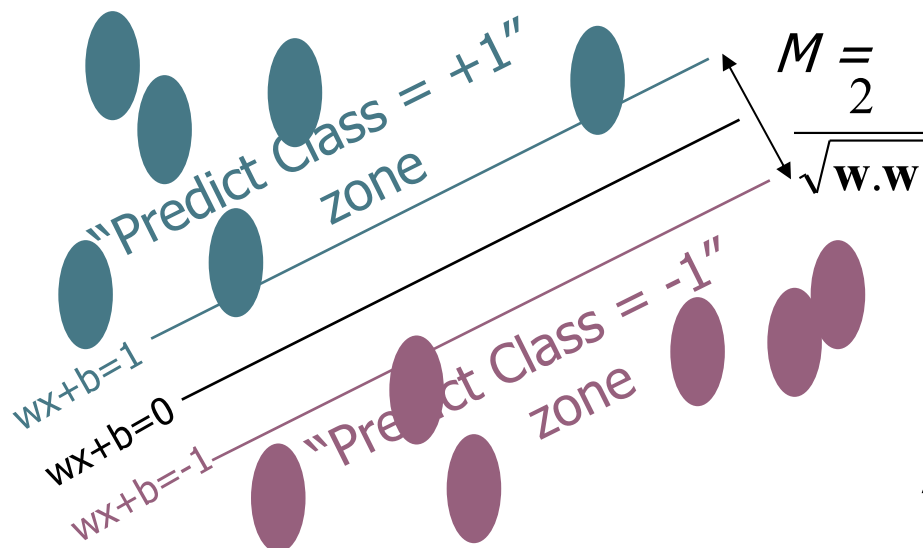
Assume $R$ datapoints, each $(\mathbf{x}_k, y_k)$ where $y_k = +/- 1$

What should our quadratic optimization criterion be?

Minimize **w.w**

How many constraints will we have? $R$

What should they be?

$\mathbf{w} . \mathbf{x}_k + b \geq 1$ if $y_k = 1$

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

# Uh-oh!

This is going to be a problem!
What should we do?

# Uh-oh!



• denotes +1
○ denotes -1

This is going to be a problem!
What should we do?

Idea 1:

Find minimum **w.w**, while minimizing number of training set errors.

Problemette: Two things to minimize makes for an ill-defined optimization

# Uh-oh!

This is going to be a problem!
What should we do?

Idea 1.1:

Minimize

$$w.w + C \text{ (\#train errors)}$$

Tradeoff parameter

There's a serious practical problem that's about to make us reject this approach. Can you guess what it is?

- denotes +1
- denotes -1

# Uh-oh!

This is going to be a problem!
What should we do?

Idea 1.1:
Minimize
$w.w$ + C (#train errors)

- denotes +1

Can't be expressed as a Quadratic Programming problem.

Solving it may be too slow.

(Also, doesn't distinguish between disastrous errors and near misses)

Tradeoff parameter

...ere's a serious practical ...blem that's about to make us reject this approach. Can you guess what it is?

# Uh-oh!



- denotes +1
- denotes -1

This is going to be a problem!
What should we do?

Idea 2.0:

Minimize

$w.w$ + C (distance of error points to their correct place)

# Learning Maximum Margin with Noise



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

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

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

- Compute sum of distances of points to their correct zones

- Compute the margin width

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

What should our quadratic optimization criterion be?

How many constraints will we have?

What should they be?

# Learning Maximum Margin with Noise

$\varepsilon_{11}$

$\varepsilon_2$

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

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

$\varepsilon_7$

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

- Compute sum of distances of points to their correct zones

- Compute the margin width

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

What should our quadratic optimization criterion be?

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

How many constraints will we have? $R$

What should they be?

$\boldsymbol{w} . \boldsymbol{x}_k + b >= 1-\varepsilon_k$ if $y_k = 1$

$\boldsymbol{w} . \boldsymbol{x}_k + b <= -1+\varepsilon_k$ if $y_k = -1$

# Learning Maximum Margin with Noise



$m$ = # input dimensions

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

Given gu... e can

- Compute sum of ...stances of ...to their correct...

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

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

$R$= # records

What should our quadratic optimization criterion be?

Minimize

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

How many constraints ...l... have? $R$

What should they be?

$\mathbf{w} \cdot \mathbf{x}_k + b >= 1 - \varepsilon_k \text{ if } y_k = 1$

$\mathbf{w} \cdot \mathbf{x}_k + b <= -1 + \varepsilon_k \text{ if } y_k = -1$

# Learning Maximum Margin with Noise



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

Given guess of **w**, $b$ we can

- Compute sum of distances of points to their correct zones
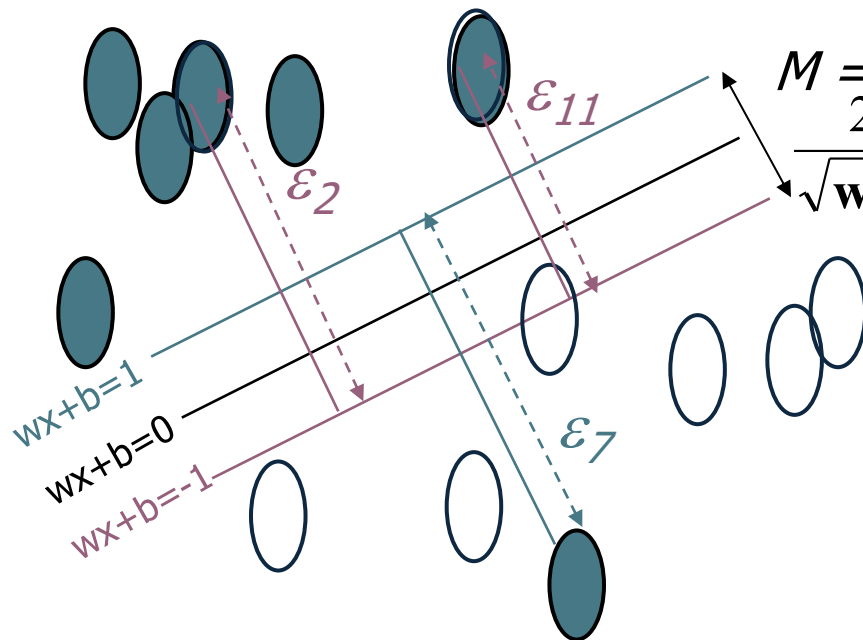
- Compute the margin width

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

What should our quadratic optimization criterion be?

Minimize

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

How many constraints will we have? *R*

What should they be?

$\boldsymbol{w}.\boldsymbol{x}_k + b >= 1-\varepsilon_k$ if $y_k = 1$

$\boldsymbol{w}.\boldsymbol{x}_k + b <= -1+\varepsilon_k$ if $y_k = -1$

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

# Learning Maximum Margin with Noise



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

- Compute sum of distances of points to their correct zones

- Compute the margin width

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

What should our quadratic optimization criterion be?

Minimize

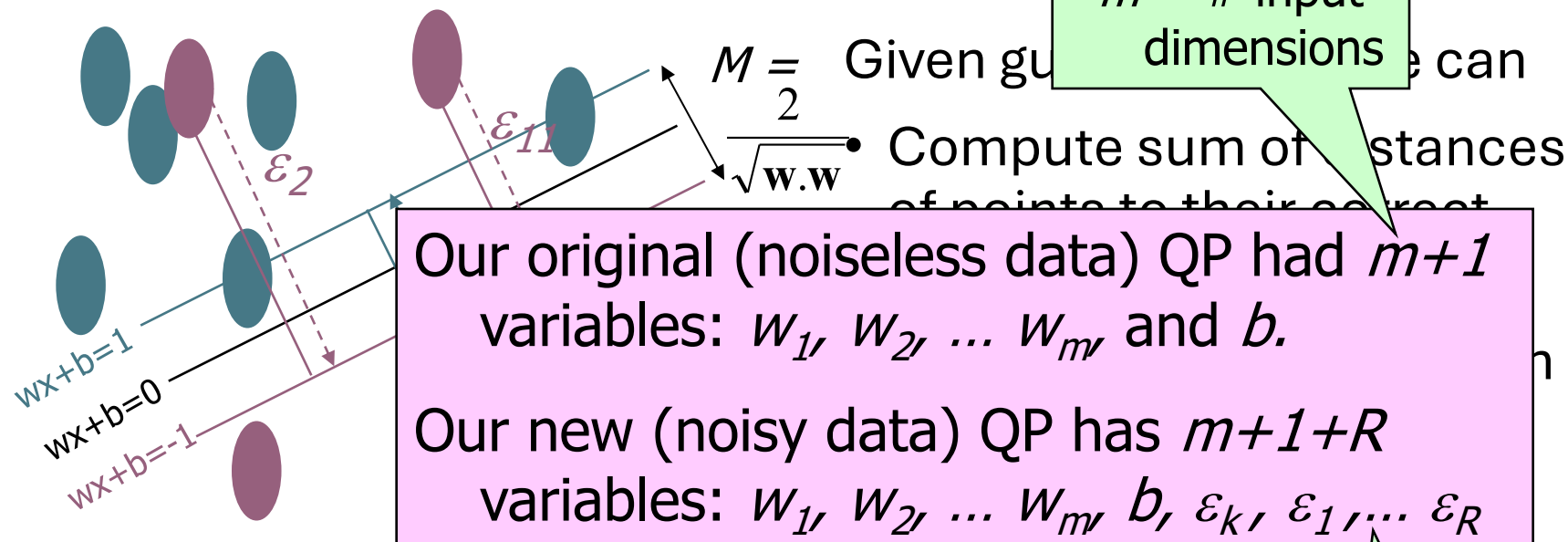$$\frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k$$
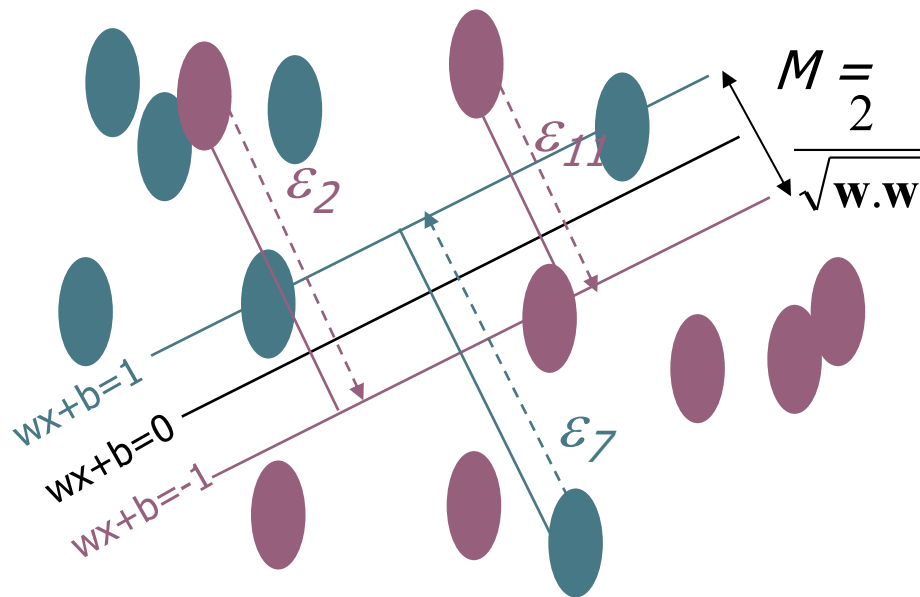
How many constraints will we have? *2R*

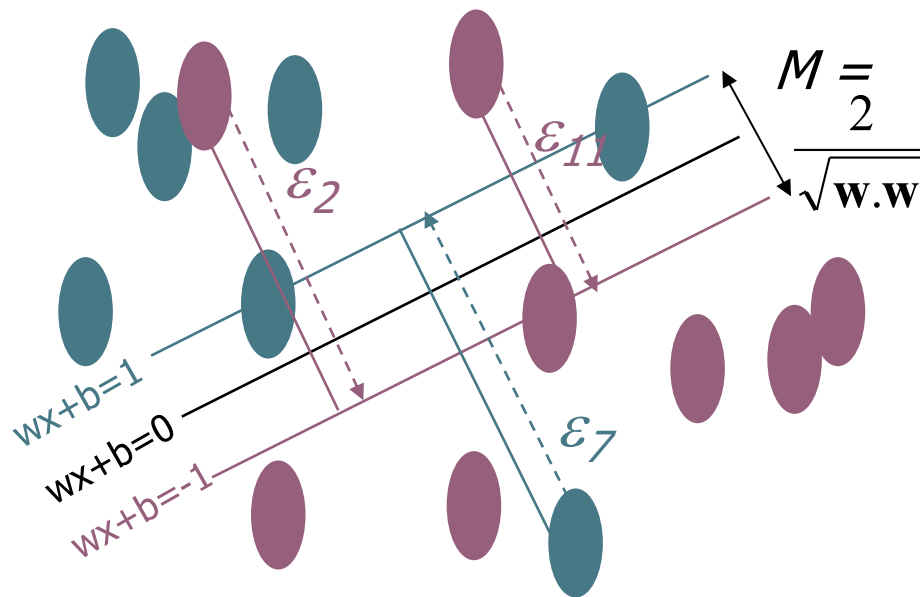What should they be?

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

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

$\varepsilon_k >= 0$ for all $k$

# QP Problems Nature

$$\max_{\alpha} \quad \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} y_n y_m \, \alpha_n \alpha_m \, \mathbf{x}_n^\mathsf{T} \mathbf{x}_m$$

$$\min_{\alpha} \quad \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} y_n y_m \, \alpha_n \alpha_m \, \mathbf{x}_n^\mathsf{T} \mathbf{x}_m - \sum_{n=1}^{N} \alpha_n$$

$$\min_{\alpha} \quad \frac{1}{2} \alpha^\mathsf{T} \underbrace{\begin{bmatrix} y_1 y_1 \, \mathbf{x}_1^\mathsf{T} \mathbf{x}_1 & y_1 y_2 \, \mathbf{x}_1^\mathsf{T} \mathbf{x}_2 & \cdots & y_1 y_N \, \mathbf{x}_1^\mathsf{T} \mathbf{x}_N \\ y_2 y_1 \, \mathbf{x}_2^\mathsf{T} \mathbf{x}_1 & y_2 y_2 \, \mathbf{x}_2^\mathsf{T} \mathbf{x}_2 & \cdots & y_2 y_N \, \mathbf{x}_2^\mathsf{T} \mathbf{x}_N \\ \cdots & \cdots & \cdots & \cdots \\ y_N y_1 \, \mathbf{x}_N^\mathsf{T} \mathbf{x}_1 & y_N y_2 \, \mathbf{x}_N^\mathsf{T} \mathbf{x}_2 & \cdots & y_N y_N \, \mathbf{x}_N^\mathsf{T} \mathbf{x}_N \end{bmatrix}}_{\text{quadratic coefficients}} \alpha + \underbrace{(-\mathbf{1}^\mathsf{T})}_{\text{linear}} \alpha$$

subject to $\quad \underbrace{\mathbf{y}^\mathsf{T} \alpha = 0}_{\text{linear constraint}}$

$$\underbrace{\mathbf{0}}_{\text{lower bounds}} \le \alpha \le \underbrace{\infty}_{\text{upper bounds}}$$

# Solving the Optimization Problem

Find $\mathbf{w}$ and b such that
$\Phi(\mathbf{w}) = \mathbf{w}^T\mathbf{w}$ is minimized
and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ : $\quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1$

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

- Need to optimize a *quadratic* function subject to *linear* constraints.

- Quadratic optimization problems are a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist.

- The solution involves constructing a *dual problem* where a *Lagrange multiplier* $\alpha_i$ is associated with every inequality constraint in the primal (original) problem:

Find $\alpha_1...\alpha_n$ such that
$\mathbf{Q}(\boldsymbol{\alpha}) = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j$ is maximized and
(1) $\Sigma\alpha_i y_i = 0$
(2) $\alpha_i \geq 0$ for all $\alpha_i$

# The Optimization Problem Solution

- Given a solution $\alpha_1 \ldots \alpha_n$ to the dual problem, solution to the primal is:

$$\mathbf{w} = \Sigma \alpha_i y_i \mathbf{x}_i \qquad b = y_k - \Sigma \alpha_i y_i \mathbf{x}_i{}^{\mathbf{T}} \mathbf{x}_k \quad \text{for any } \alpha_k > 0$$

- Each non-zero $\alpha_i$ indicates that corresponding $\mathbf{x_i}$ is a support vector.
- Then the classifying function is (note that we don't need $\mathbf{w}$ explicitly):

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \mathbf{x}_i{}^{\mathbf{T}} \mathbf{x} + b$$

- Notice that it relies on an *inner product* between the test point $\mathbf{x}$ and the support vectors $\mathbf{x}_i$ – we will return to this later.
- Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x}_i{}^{\mathbf{T}} \mathbf{x}_j$ between all training points.

## Testing

- For testing with a new data z:

- Compute $WZ + b = \sum_{i=1}^{N_1} \alpha_i y_i (X_i.Z) + b$ and classify $Z$ as $y_i = +1$ if the sum is positive, $y_i = -1$ otherwise.

- Note that we do not need to form $W$ explicitly.

- Suppose we are given the following positively labeled data points in $\Re^2$:

$$\begin{pmatrix} 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 1 \end{pmatrix}, \begin{pmatrix} 6 \\ -1 \end{pmatrix}$$

and

- the following negatively labeled data points in $\Re^2$:

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}$$
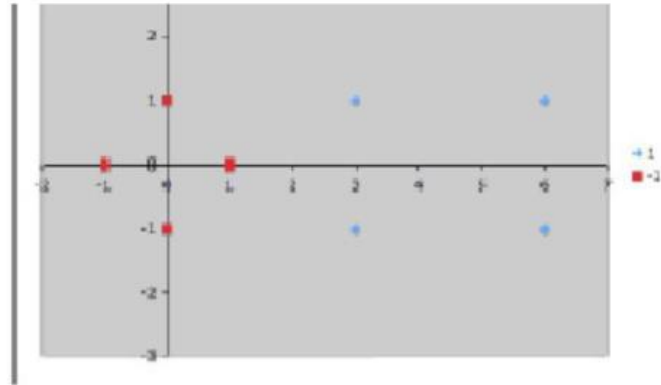


Figure 1: Sample data points in $\Re^2$. Blue diamonds are positive examples red squares are negative examples.

Lets define simple SVM that accurately discriminates the two classes. Since the data is linearly separable, we can use a linear SVM. it should be obvious that there are three support vectors:

$$s_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, s_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, s_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix}$$

- We will use vectors augmented with a 1 as a bias input, and for clarity we will differentiate these with an over-tilde. So, if s1 = (10), then $\tilde{s_1}$ = (101).

$$\tilde{s}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \tilde{s}_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix}, \tilde{s}_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

- Our task is to find values for the $\alpha_i$ such that,

$$\alpha_i \tilde{s}_1 . \tilde{s}_1 + \alpha_2 \tilde{s}_2 . \tilde{s}_1 + \alpha_3 \tilde{s}_3 . \tilde{s}_1 = -1$$
$$\alpha_i \tilde{s}_1 . \tilde{s}_2 + \alpha_2 \tilde{s}_2 . \tilde{s}_2 + \alpha_3 \tilde{s}_3 . \tilde{s}_2 = +1$$
$$\alpha_i \tilde{s}_1 . \tilde{s}_3 + \alpha_2 \tilde{s}_2 . \tilde{s}_3 + \alpha_3 \tilde{s}_3 . \tilde{s}_3 = +1$$

**Example**

- Our task is to find values for the $\alpha_i$ such that,

$$\alpha_i \tilde{s}_1.\tilde{s}_1 + \alpha_2 \tilde{s}_2.\tilde{s}_1 + \alpha_3 \tilde{s}_3.\tilde{s}_1 = -1$$
$$\alpha_i \tilde{s}_1.\tilde{s}_2 + \alpha_2 \tilde{s}_2.\tilde{s}_2 + \alpha_3 \tilde{s}_3.\tilde{s}_2 = +1$$
$$\alpha_i \tilde{s}_1.\tilde{s}_3 + \alpha_2 \tilde{s}_2.\tilde{s}_3 + \alpha_3 \tilde{s}_3.\tilde{s}_3 = +1$$

- Computing the dot product

For example, $\tilde{s}_1.\tilde{s}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} . \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = 1 \times 1 + 0 \times 0 + 1 \times 1 = 2$

$$2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$
$$4\alpha_2 + 11\alpha_2 + 9\alpha_3 = +1$$
$$4\alpha_1 + 9\alpha_2 + 11\alpha_3 = +1$$

$$\alpha_1 = -3.5 \text{ and } \alpha_2 = 0.75 \text{ and } \alpha_3 = 0.75$$

**Example**

- How to find the hyper-plane that discriminates the positive values?

$$\tilde{w} = \sum_{i=1}^{3} \alpha_i \tilde{s}_i$$

$$= -3.5 \times \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 0.75 \times \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + 0.75 \times \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix}$$
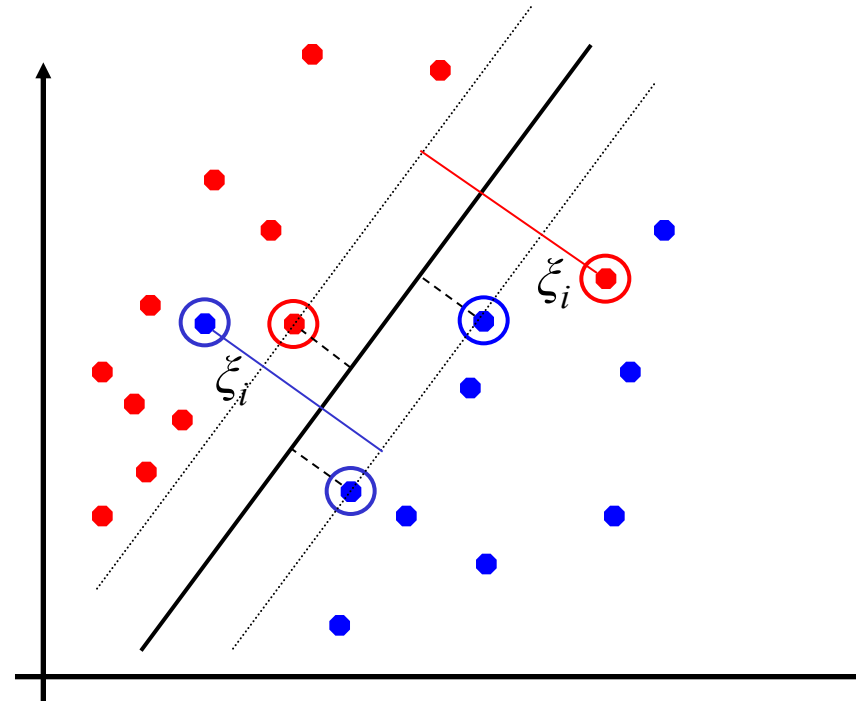
- The bias $b$ and $w$ are:

$$w = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \text{ and } b = -2$$

# Soft Margin Classification

- What if the training set is not linearly separable?
- *Slack variables* $\xi_i$ can be added to allow misclassification of difficult or noisy examples, resulting margin called *soft*.

# Soft Margin Classification Mathematically

- The old formulation:

> Find **w** and b such that
> $\Phi(\mathbf{w}) = \mathbf{w^T w}$  is minimized
> and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ :     $y_i(\mathbf{w^T x}_i + b) \geq 1$

- Modified formulation incorporates slack variables:

> Find **w** and b such that
> $\Phi(\mathbf{w}) = \mathbf{w^T w} + C\Sigma\xi_i$   is minimized
> and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ :     $y_i(\mathbf{w^T x}_i + b) \geq 1 - \xi_{i,}$ ,    $\xi_i \geq 0$

- Parameter $C$ can be viewed as a way to control overfitting:  it "trades off" the relative importance of maximizing the margin and fitting the training data.

# Soft Margin Classification – Solution

- Dual problem is identical to separable case (would *not* be identical if the 2-norm penalty for slack variables $C\Sigma\xi_i^2$ was used in primal objective, we would need additional Lagrange multipliers for slack variables):

> Find $\alpha_1...\alpha_N$ such that
>
> $\mathbf{Q(\alpha)} = \Sigma\alpha_i - \tfrac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j$ is maximized and
>
> (1) $\Sigma\alpha_i y_i = 0$
>
> (2) $0 \leq \alpha_i \leq C$ for all $\alpha_i$

- Again, $\mathbf{x}_i$ with non-zero $\alpha_i$ will be support vectors.
- Solution to the dual problem is:

> $\mathbf{w} = \Sigma\alpha_i y_i \mathbf{x}_i$
>
> $b = y_k(1 - \xi_k) - \Sigma\alpha_i y_i \mathbf{x}_i^T\mathbf{x}_k$    for any $k$ s.t. $\alpha_k > 0$

Again, we don't need to compute $\mathbf{w}$ explicitly for classification:

> $f(\mathbf{x}) = \Sigma\alpha_i y_i \mathbf{x}_i^T\mathbf{x} + b$