

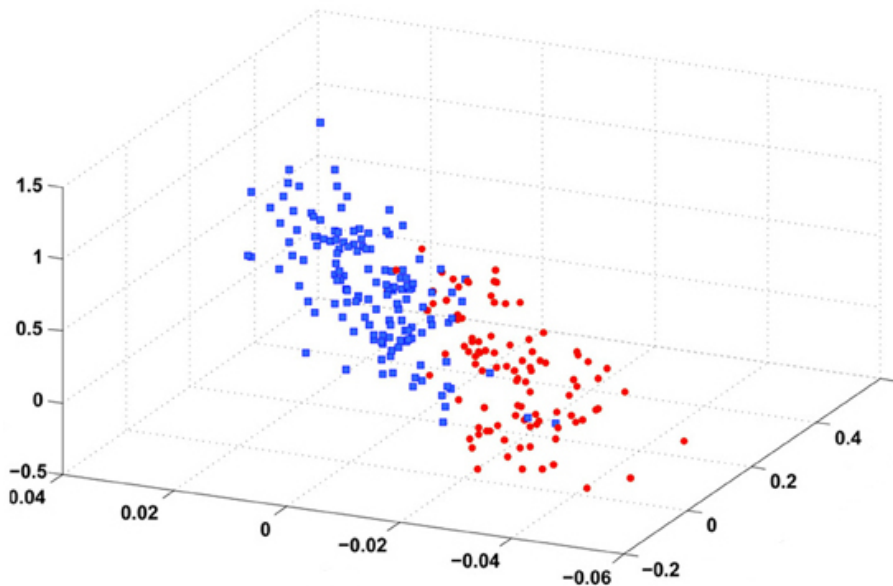
CS189/CS289A
Introduction to Machine Learning
Lecture 2: Linear classifiers

Peter Bartlett

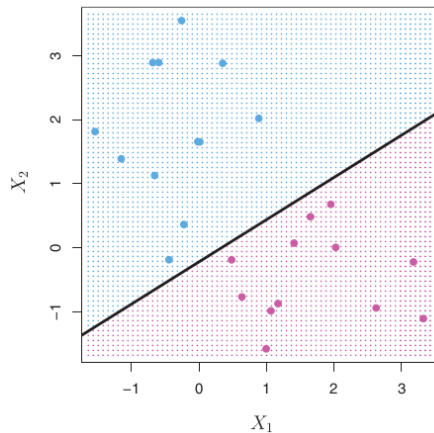
January 22, 2015

Linear Classifiers:

$$x \in \mathbb{R}^d, y \in \{-1, 1\}$$

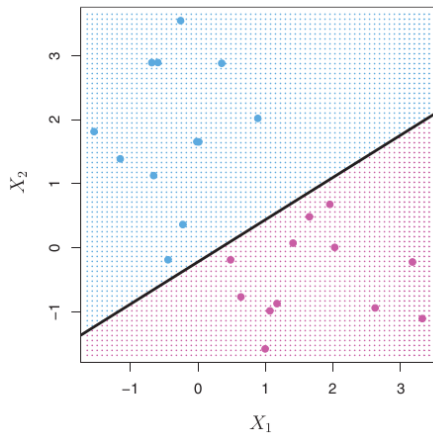


Linear Classifiers



1 Training

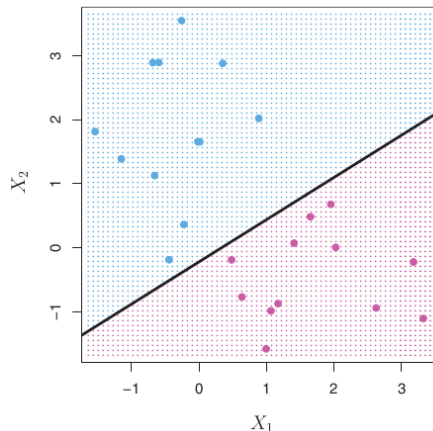
Linear Classifiers



1 Training

- Collect labeled data.

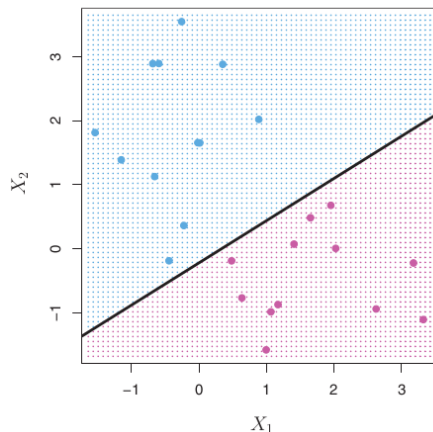
Linear Classifiers



1 Training

- Collect labeled data.
- Find a good separating hyperplane.

Linear Classifiers

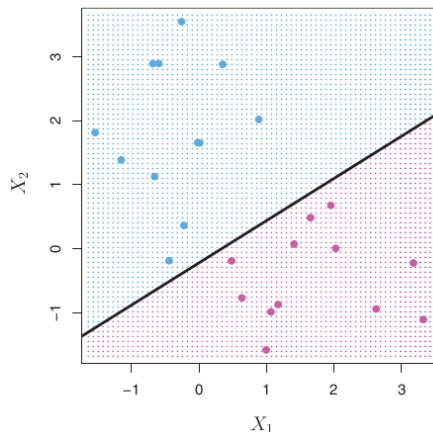


1 Training

- Collect labeled data.
- Find a good separating hyperplane.

2 Test

Linear Classifiers



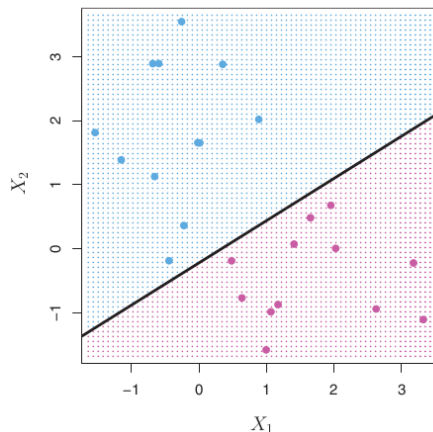
1 Training

- Collect labeled data.
- Find a good separating hyperplane.

2 Test

- Given an unlabeled point, predict according to which side of the hyperplane it's on.

Linear Classifiers



1 Training

- Collect labeled data.
- Find a good separating hyperplane.

2 Test

- Given an unlabeled point, predict according to which side of the hyperplane it's on.

Linear Classifiers: Outline for today

Linear Classifiers: Outline for today

- Properties of linear classifiers.

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.
 - The perceptron algorithm as a stochastic gradient method.

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.
 - The perceptron algorithm as a stochastic gradient method.
 - Support vector machines (Vapnik, 1990s).

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.
 - The perceptron algorithm as a stochastic gradient method.
 - Support vector machines (Vapnik, 1990s).
 - Hard margin SVM.

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.
 - The perceptron algorithm as a stochastic gradient method.
 - Support vector machines (Vapnik, 1990s).
 - Hard margin SVM.
 - Soft margin SVM.

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.
 - The perceptron algorithm as a stochastic gradient method.
 - Support vector machines (Vapnik, 1990s).
 - Hard margin SVM.
 - Soft margin SVM.
 - Both are quadratic programs.

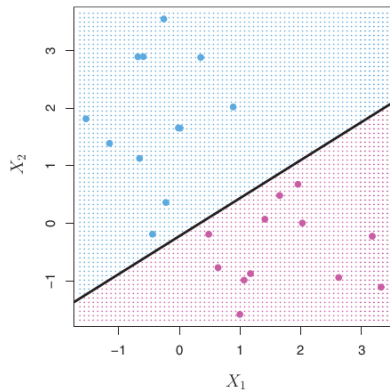
Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.
 - The perceptron algorithm as a stochastic gradient method.
 - Support vector machines (Vapnik, 1990s).
 - Hard margin SVM.
 - Soft margin SVM.
 - Both are quadratic programs.
 - Only use inner products.

Linear threshold functions

Linear classifier:

For a pattern $x \in \mathbb{R}^d$



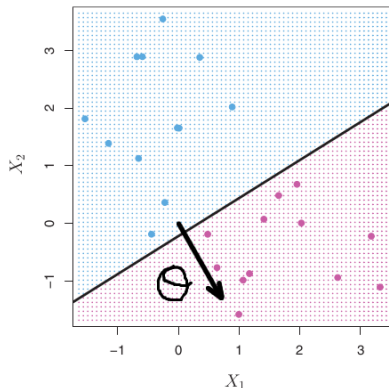
Linear threshold functions

Linear classifier:

For a pattern $x \in \mathbb{R}^d$ and parameters $\theta \in \mathbb{R}^d$, $\theta_0 \in \mathbb{R}$, define

$$f(x) = \theta \cdot x + \theta_0 = \sum_{i=1}^d \theta_i x_i + \theta_0,$$

$$\hat{y} = \begin{cases} 1 & \text{if } f(x) \geq 0, \\ -1 & \text{if } f(x) < 0. \end{cases}$$



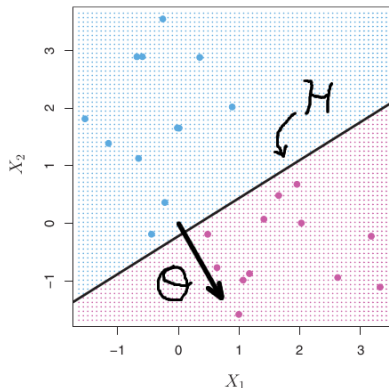
Linear threshold functions

Linear classifier:

For a pattern $x \in \mathbb{R}^d$ and parameters $\theta \in \mathbb{R}^d$, $\theta_0 \in \mathbb{R}$, define

$$f(x) = \theta \cdot x + \theta_0 = \sum_{i=1}^d \theta_i x_i + \theta_0,$$

$$\hat{y} = \begin{cases} 1 & \text{if } f(x) \geq 0, \\ -1 & \text{if } f(x) < 0. \end{cases}$$



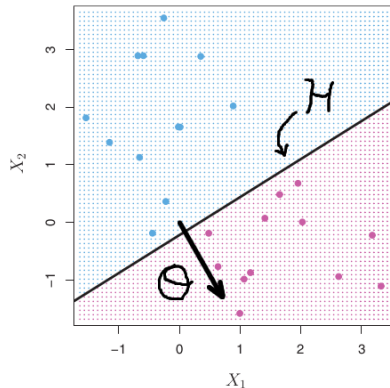
Decision boundary:

$$H = \{x \in \mathbb{R}^d : f(x) = 0\} = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

Linear threshold functions

Decision boundary:

$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

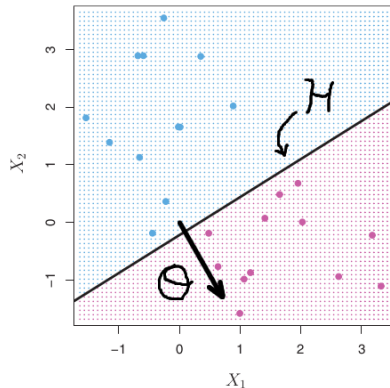


Linear threshold functions

Decision boundary:

$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

In \mathbb{R}^2 , this is a line.

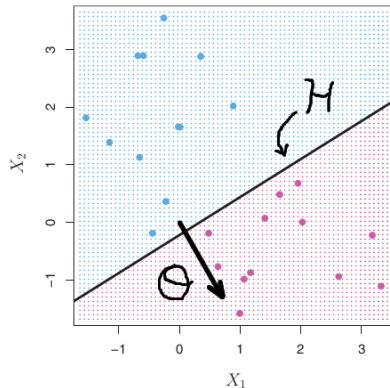


Linear threshold functions

Decision boundary:

$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

In \mathbb{R}^d , it is a hyperplane.



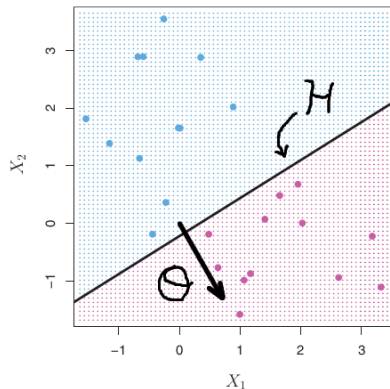
Linear threshold functions

Decision boundary:

$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

In \mathbb{R}^d , it is a hyperplane.

- It is normal to θ :



Linear threshold functions

Decision boundary:

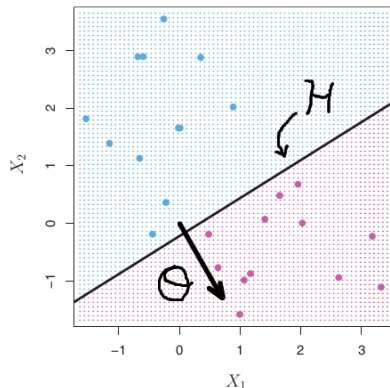
$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

In \mathbb{R}^d , it is a hyperplane.

- It is normal to θ :

For any $x, \tilde{x} \in H$, we have

$$\begin{aligned} 0 &= (\theta \cdot x + \theta_0) - (\theta \cdot \tilde{x} + \theta_0) \\ &= \theta \cdot (x - \tilde{x}). \end{aligned}$$



Linear threshold functions

Decision boundary:

$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

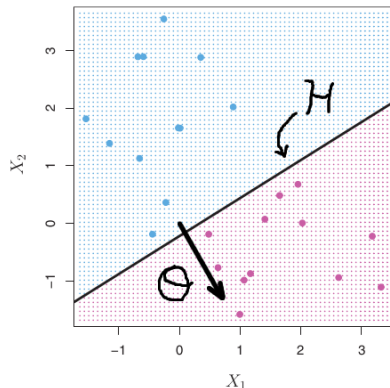
In \mathbb{R}^d , it is a hyperplane.

- It is normal to θ :

For any $x, \tilde{x} \in H$, we have

$$\begin{aligned} 0 &= (\theta \cdot x + \theta_0) - (\theta \cdot \tilde{x} + \theta_0) \\ &= \theta \cdot (x - \tilde{x}). \end{aligned}$$

- θ_0 determines the offset:



Linear threshold functions

Decision boundary:

$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

In \mathbb{R}^d , it is a hyperplane.

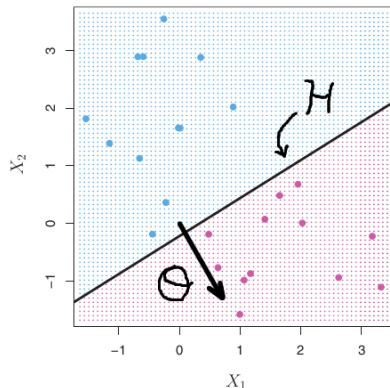
- It is normal to θ :

For any $x, \tilde{x} \in H$, we have

$$\begin{aligned} 0 &= (\theta \cdot x + \theta_0) - (\theta \cdot \tilde{x} + \theta_0) \\ &= \theta \cdot (x - \tilde{x}). \end{aligned}$$

- θ_0 determines the offset:

The signed distance (in the direction of θ) from the origin to H is $-\theta_0/\|\theta\|$.



Linear threshold functions

Decision boundary:

$$H = \{x \in \mathbb{R}^d : \theta \cdot x + \theta_0 = 0\}.$$

In \mathbb{R}^d , it is a hyperplane.

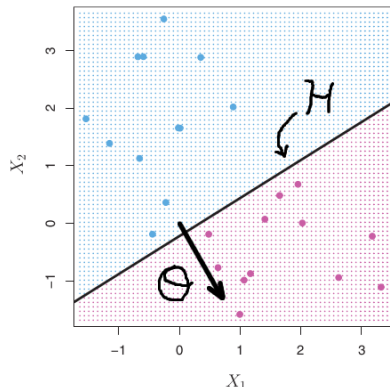
- It is normal to θ :

For any $x, \tilde{x} \in H$, we have

$$\begin{aligned} 0 &= (\theta \cdot x + \theta_0) - (\theta \cdot \tilde{x} + \theta_0) \\ &= \theta \cdot (x - \tilde{x}). \end{aligned}$$

- θ_0 determines the offset:

The signed distance (in the direction of θ) from the origin to H is $-\theta_0/\|\theta\|$. (Why?)



Aside: detail

How far away from the origin is the hyperplane? Since θ is normal to H , we can look at the ray from the origin towards H in the direction of θ . That ray hits H at $c\theta$, where

$$0 = f(c\theta) = c\theta \cdot \theta + \theta_0 = c\|\theta\|^2 + \theta_0.$$

So $c = -\theta_0/\|\theta\|^2$, and the signed distance to H is

$$c\|\theta\| = -\theta_0 \frac{\|\theta\|}{\|\theta\|^2} = -\frac{\theta_0}{\|\theta\|}.$$

Linear threshold functions

We can simplify our notation by augmenting the pattern $x \in \mathbb{R}^d$ with a constant component. This allows us to dispense with the offset θ_0 :

For a pattern $x \in \mathbb{R}^d$, and parameters $\theta \in \mathbb{R}^d$, $\theta_0 \in \mathbb{R}$, define

$$\tilde{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}, \quad \tilde{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{pmatrix}.$$

Then define

$$\tilde{f}(\tilde{x}) := \tilde{\theta} \cdot \tilde{x} = \sum_{i=1}^d \theta_i x_i + \theta_0 = \theta \cdot x + \theta_0 = f(x),$$
$$\hat{y} = \begin{cases} 1 & \text{if } f(x) \geq 0, \\ -1 & \text{if } f(x) < 0. \end{cases} = \begin{cases} 1 & \text{if } \tilde{f}(\tilde{x}) \geq 0, \\ -1 & \text{if } \tilde{f}(\tilde{x}) < 0. \end{cases}$$

Linear threshold functions

So we can always consider linear classifiers of this simpler form:

$$f(x) = \theta \cdot x,$$
$$\hat{y} = \begin{cases} 1 & \text{if } f(x) \geq 0, \\ -1 & \text{if } f(x) < 0. \end{cases}$$

The decision boundary is the hyperplane

$$H = \{x : \theta \cdot x = 0\},$$

passing through the origin, normal to the vector θ .

Linear classifiers

Now that we understand the functions computed by linear classifiers, let's consider how we might use labeled data $(x_1, y_1), \dots, (x_n, y_n)$ to choose a good classifier.

Linear classifiers

Now that we understand the functions computed by linear classifiers, let's consider how we might use labeled data $(x_1, y_1), \dots, (x_n, y_n)$ to choose a good classifier.

A good classifier is one that classifies subsequent (x, y) pairs accurately.

Linear classifiers

Now that we understand the functions computed by linear classifiers, let's consider how we might use labeled data $(x_1, y_1), \dots, (x_n, y_n)$ to choose a good classifier.

A good classifier is one that classifies subsequent (x, y) pairs accurately.

A reasonable approach:

Find a linear classifier that fits the training data well.

Linear classifiers

Now that we understand the functions computed by linear classifiers, let's consider how we might use labeled data $(x_1, y_1), \dots, (x_n, y_n)$ to choose a good classifier.

A good classifier is one that classifies subsequent (x, y) pairs accurately.

A reasonable approach:

Find a linear classifier that fits the training data well.
For instance, we might aim to *minimize the number of misclassifications* on the training data.

Linear classifiers

Now that we understand the functions computed by linear classifiers, let's consider how we might use labeled data $(x_1, y_1), \dots, (x_n, y_n)$ to choose a good classifier.

A good classifier is one that classifies subsequent (x, y) pairs accurately.

A reasonable approach:

Find a linear classifier that fits the training data well.

For instance, we might aim to *minimize the number of misclassifications* on the training data.

This is known as *empirical risk minimization*.

Perceptron algorithm

Perceptron algorithm:

Input: $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{\pm 1\}$

while some $y^i \neq \text{sign}(\theta \cdot x^i)$

 improve θ

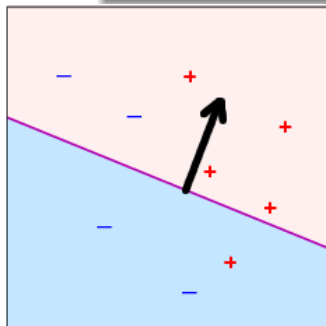
Return θ .

Perceptron algorithm

Perceptron algorithm:

Input: $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{\pm 1\}$
while some $y^i \neq \text{sign}(\theta \cdot x^i)$
 improve θ

Return θ .



Perceptron algorithm

Perceptron algorithm:

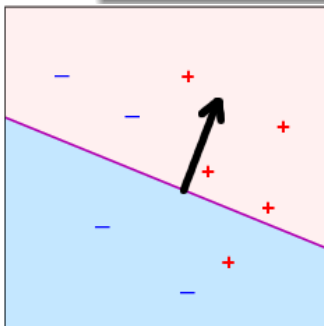
Input: $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{\pm 1\}$

while some $y^i \neq \text{sign}(\theta \cdot x^i)$

 pick some misclassified (x^i, y^i)

 update θ

Return θ .



Perceptron algorithm

Perceptron algorithm:

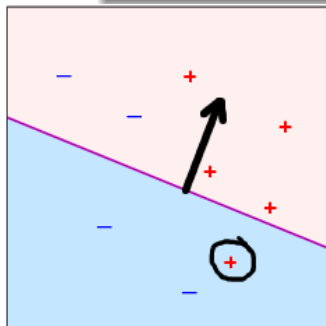
Input: $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{\pm 1\}$

while some $y^i \neq \text{sign}(\theta \cdot x^i)$

 pick some misclassified (x^i, y^i)

 update θ

Return θ .



Perceptron algorithm

Perceptron algorithm:

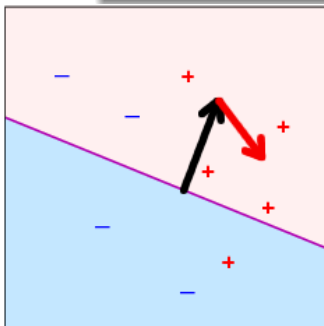
Input: $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{\pm 1\}$

while some $y^i \neq \text{sign}(\theta \cdot x^i)$

 pick some misclassified (x^i, y^i)

$\theta \leftarrow \theta + x^i$

Return θ .



Perceptron algorithm

Perceptron algorithm:

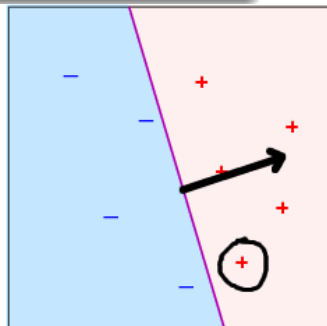
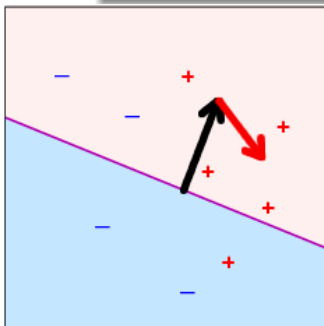
Input: $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{\pm 1\}$

while some $y^i \neq \text{sign}(\theta \cdot x^i)$

 pick some misclassified (x^i, y^i)

$\theta \leftarrow \theta + x^i$

Return θ .



Perceptron algorithm

Perceptron algorithm:

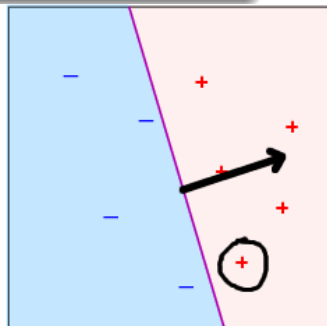
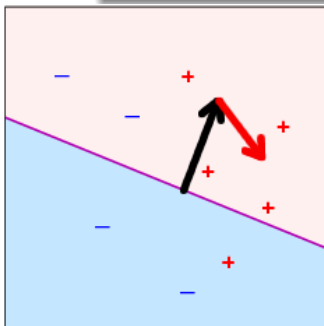
Input: $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \{\pm 1\}$

while some $y^i \neq \text{sign}(\theta \cdot x^i)$

 pick some misclassified (x^i, y^i)

$\theta \leftarrow \theta + y^i x^i$

Return θ .



Linear threshold functions

Perceptron convergence theorem

Given *linearly separable data* (i.e., there is a $\theta \in \mathbb{R}^d$ such that for all i , $y^i \theta \cdot x^i > 0$), for any choices of updates, the perceptron algorithm terminates with all data correctly classified.

Linear threshold functions

Perceptron convergence theorem

Given *linearly separable data* (i.e., there is a $\theta \in \mathbb{R}^d$ such that for all i , $y^i \theta \cdot x^i > 0$), for any choices of updates, the perceptron algorithm terminates with all data correctly classified.

Furthermore, it makes no more than $\frac{R^2}{\gamma^2}$ updates, where

$$R = \max_i \|x^i\|, \quad (\text{radius of data})$$

$$\gamma = \min_i \frac{y^i(\theta \cdot x^i)}{\|\theta\|}. \quad (\text{margin})$$

NB: $\frac{1}{\|\theta\|} \theta \cdot x$ is the signed distance from H to x (in the direction θ).

Linear threshold functions

Perceptron convergence theorem

Given *linearly separable data* (i.e., there is a $\theta \in \mathbb{R}^d$ such that for all i , $y^i \theta \cdot x^i > 0$), for any choices of updates, the perceptron algorithm terminates with all data correctly classified.

Furthermore, it makes no more than $\frac{R^2}{\gamma^2}$ updates, where

$$R = \max_i \|x^i\|, \quad (\text{radius of data})$$

$$\gamma = \min_i \frac{y^i(\theta \cdot x^i)}{\|\theta\|}. \quad (\text{margin})$$

Proof idea: Fix a separating θ^* . Each update increases $\theta \cdot \theta^*$ a lot, but only increases $\|\theta\|$ a little. So there can't be too many updates.

Perceptron algorithm: only uses inner products

Properties:

- $\theta = \sum_i \alpha^i y^i x^i$.

Perceptron algorithm: only uses inner products

Properties:

- $\theta = \sum_i \alpha^i y^i x^i$.
- The only properties of the data that we use are inner products:

$$\theta \cdot x = \sum_i \alpha^i y^i x^i \cdot x.$$

Perceptron algorithm: only uses inner products

Properties:

- $\theta = \sum_i \alpha^i y^i x^i$.
- The only properties of the data that we use are inner products:

$$\theta \cdot x = \sum_i \alpha^i y^i x^i \cdot x.$$

- So we can work with data in any inner product space (not just finite-dimensional vectors).

Perceptron algorithm

We can also view the perceptron algorithm as a stochastic gradient method to minimize a cost function.

Perceptron algorithm

We can also view the perceptron algorithm as a stochastic gradient method to minimize a cost function.

Margin cost function

- We want the sign of $\theta \cdot x^i$ to be the same as y^i .

Perceptron algorithm

We can also view the perceptron algorithm as a stochastic gradient method to minimize a cost function.

Margin cost function

- We want the sign of $\theta \cdot x^i$ to be the same as y^i .
- We want $y^i(\theta \cdot x^i) > 0$.

Perceptron algorithm

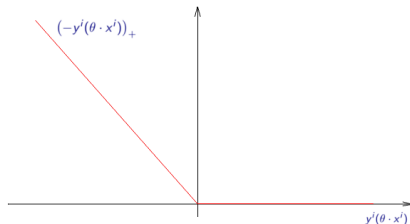
We can also view the perceptron algorithm as a stochastic gradient method to minimize a cost function.

Margin cost function

- We want the sign of $\theta \cdot x^i$ to be the same as y^i .
- We want $y^i(\theta \cdot x^i) > 0$.
- Define

$$J(\theta) = \sum_i (-y^i(\theta \cdot x^i))_+$$

$$(\alpha)_+ = \begin{cases} \alpha & \text{if } \alpha > 0, \\ 0 & \text{otherwise.} \end{cases}$$



Gradients

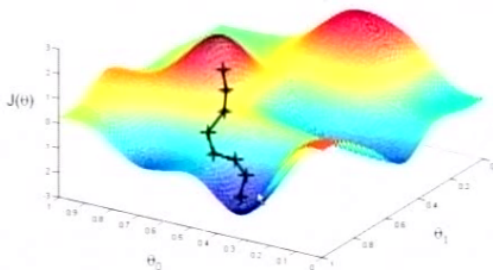
$$\nabla J(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} J(\theta) \\ \frac{\partial}{\partial \theta_2} J(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_d} J(\theta) \end{pmatrix}.$$

Example:

$$\nabla (\theta \cdot x) = \nabla \left(\sum_{i=1}^d \theta_i x_i \right) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} \theta \cdot x \\ \frac{\partial}{\partial \theta_2} \theta \cdot x \\ \vdots \\ \frac{\partial}{\partial \theta_d} \theta \cdot x \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} = x.$$

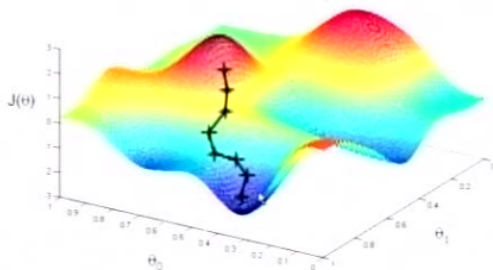
Perceptron algorithm

Gradient descent



Perceptron algorithm

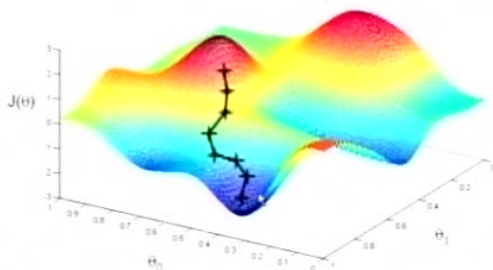
Gradient descent



$$\theta \leftarrow \theta - \eta \nabla J(\theta)$$

Perceptron algorithm

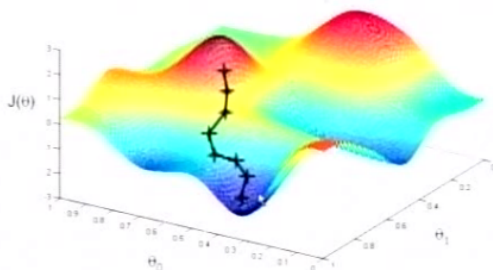
Gradient descent



$$\theta \leftarrow \theta - \underbrace{\eta \nabla J(\theta)}_{\text{uphill}}$$

Perceptron algorithm

Gradient descent



$$\theta \leftarrow \theta - \underbrace{\eta}_{\text{stepsize}} \underbrace{\nabla J(\theta)}_{\text{uphill}}$$

Perceptron algorithm

$$\begin{aligned}J(\theta) &= \sum_i J_i(\theta) = \sum_i (-y^i(\theta \cdot x^i))_+ \\ \nabla J(\theta) &= \sum_i \nabla J_i(\theta) \\ \nabla J_i(\theta) &= \begin{cases} -y^i x^i & \text{if } y^i(\theta \cdot x^i) < 0 \\ 0 & \text{otherwise.} \end{cases}\end{aligned}$$

Stochastic gradient

- Pick a component J_i of the cost function J (i.e., an (x^i, y^i) pair)
- Move downhill wrt that component:
If $y^i(\theta \cdot x^i) > 0$, don't change θ (because the gradient is zero.)
Otherwise: $\theta \leftarrow \theta - \nabla J_i(\theta) = \theta + y^i x^i$.

Perceptron algorithm

$$J(\theta) = \sum_i J_i(\theta) = \sum_i (-y^i(\theta \cdot x^i))_+$$
$$\nabla J(\theta) = \sum_i \nabla J_i(\theta)$$
$$\nabla J_i(\theta) = \begin{cases} -y^i x^i & \text{if } y^i(\theta \cdot x^i) < 0 \\ 0 & \text{otherwise.} \end{cases}$$

Stochastic gradient

- Pick a component J_i of the cost function J (i.e., an (x^i, y^i) pair)
- Move downhill wrt that component:

If $y^i(\theta \cdot x^i) > 0$, don't change θ (because the gradient is zero.)

Otherwise: $\theta \leftarrow \theta - \nabla J_i(\theta) = \theta + y^i x^i$.

Perceptron algorithm

Perceptron algorithm

Issues:

- Converges only if the data are *separable*.

Perceptron algorithm

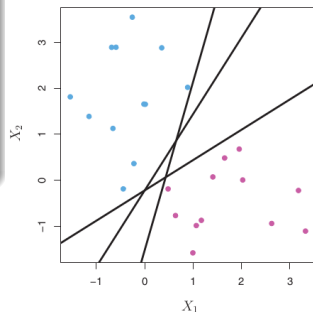
Issues:

- Converges only if the data are *separable*.
- Time to convergence depends on *margin*.

Perceptron algorithm

Issues:

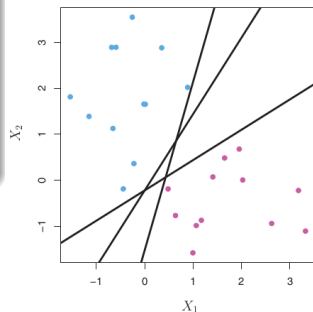
- Converges only if the data are *separable*.
- Time to convergence depends on *margin*.
- The solution depends on starting point.



Perceptron algorithm

Issues:

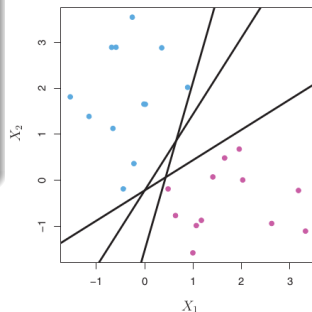
- Converges only if the data are *separable*.
- Time to convergence depends on *margin*.
- The solution depends on starting point.
- Will not converge if data are not separable.



Perceptron algorithm

Issues:

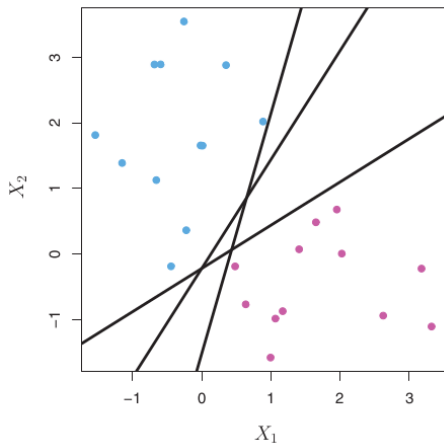
- Converges only if the data are *separable*.
- Time to convergence depends on *margin*.
- The solution depends on starting point.
- Will not converge if data are not separable.



(Looks like a toy example. But stochastic gradient methods turn out to be very useful for large scale problems.)

Support vector machines

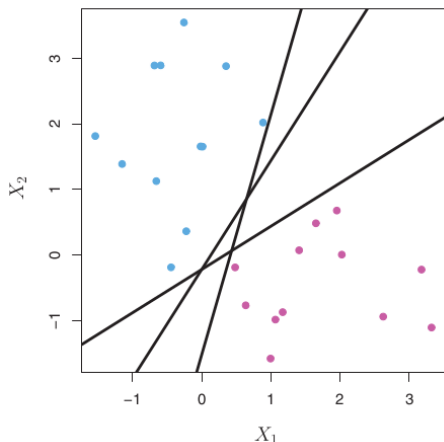
- There are always many linear classifiers that give identical classifications of the training data. Which is better?



Support vector machines

- There are always many linear classifiers that give identical classifications of the training data. Which is better?
- One option is to choose the classifier that maximizes the *margin* on the training data, where the margin is the minimum over (x^i, y^i) pairs of the signed distance to the decision boundary,

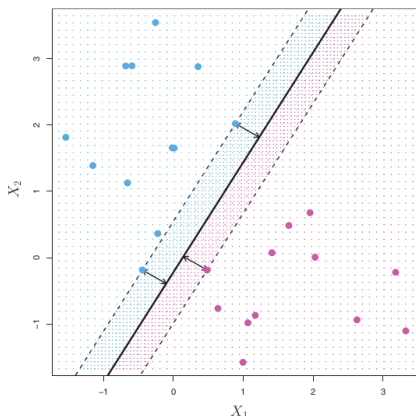
$$\text{margin of } (x^i, y^i) := y^i \frac{\theta \cdot x^i}{\|\theta\|}.$$



Support vector machines

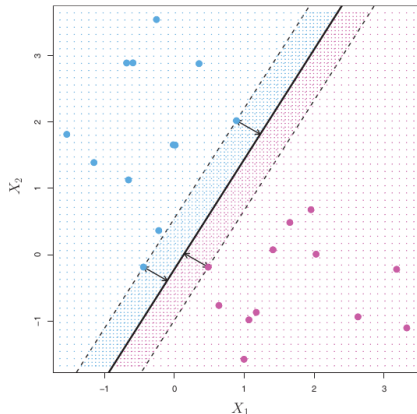
- There are always many linear classifiers that give identical classifications of the training data. Which is better?
- One option is to choose the classifier that maximizes the *margin* on the training data, where the margin is the minimum over (x^i, y^i) pairs of the signed distance to the decision boundary,

$$\text{margin of } (x^i, y^i) := y^i \frac{\theta \cdot x^i}{\|\theta\|}.$$



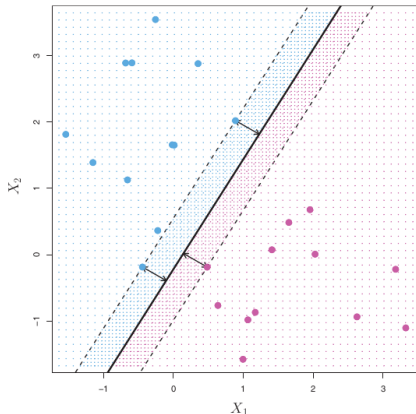
Support vector machines

- Maximizing the margin means choosing the separating hyperplane so that the sandwich around it containing no data is as thick as possible.



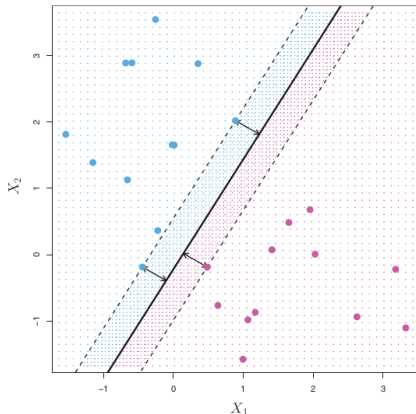
Support vector machines

- Maximizing the margin means choosing the separating hyperplane so that the sandwich around it containing no data is as thick as possible.
- (The texts emphasize that this provides a unique solution. But why should we care?)



Support vector machines

- Maximizing the margin means choosing the separating hyperplane so that the sandwich around it containing no data is as thick as possible.
- (The texts emphasize that this provides a unique solution. But why should we care?)
- Much more importantly, we can expect that a large margin between training examples and the decision boundary will lead to good separation on the test data. (And there are theorems that show this.)



Support vector machines

For linearly separable data, we can find a θ such that for all i ,

$$y^i \theta \cdot x^i > 0.$$

Support vector machines

For linearly separable data, we can find a θ such that for all i ,

$$y^i \theta \cdot x^i > 0.$$

Let δ be the smallest of these values, so that for all i ,

$$y^i \theta \cdot x^i \geq \delta.$$

Support vector machines

For linearly separable data, we can find a θ such that for all i ,

$$y^i \theta \cdot x^i > 0.$$

Let δ be the smallest of these values, so that for all i ,

$$y^i \theta \cdot x^i \geq \delta.$$

Then we can choose a new value of θ (scaled by $1/\delta$) so that, for all i ,

$$y^i \theta \cdot x^i \geq 1.$$

Support vector machines

For linearly separable data, we can find a θ such that for all i ,

$$y^i \theta \cdot x^i > 0.$$

Let δ be the smallest of these values, so that for all i ,

$$y^i \theta \cdot x^i \geq \delta.$$

Then we can choose a new value of θ (scaled by $1/\delta$) so that, for all i ,

$$y^i \theta \cdot x^i \geq 1.$$

With this scaling of θ , the margin is

$$\min_i y^i \frac{\theta \cdot x^i}{\|\theta\|} = \frac{1}{\|\theta\|}.$$

Support vector machines

Maximizing the margin:

$$\begin{array}{ll}\max_{\theta} & \frac{1}{\|\theta\|} \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)\end{array}$$

Maximizing the margin:

$$\min_{\theta} \quad \|\theta\|$$

$$\text{s.t.} \quad y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)$$

Maximizing the margin:

Quadratic Program

$$\min_{\theta} \quad \|\theta\|^2$$

$$\text{s.t.} \quad y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)$$

Maximizing the margin: Quadratic Program

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)\end{array}$$

- A quadratic program (linear constraints and a quadratic criterion) is a convex optimization problem.

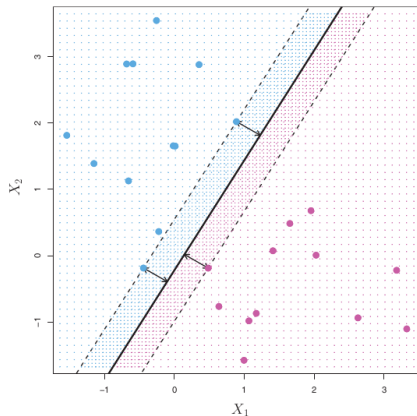
Maximizing the margin: Quadratic Program

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)\end{array}$$

- A quadratic program (linear constraints and a quadratic criterion) is a convex optimization problem.
- There are efficient algorithms for solving QPs.

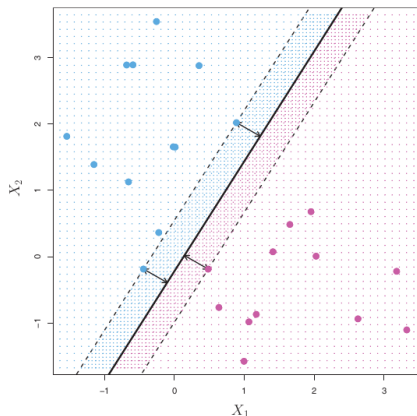
Support vector machines

- The points that satisfy these constraints with equality (*active constraints*) are called *support vectors*.



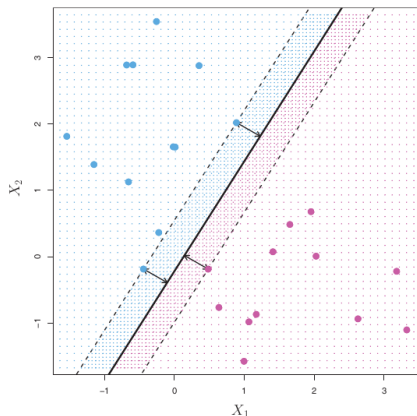
Support vector machines

- The points that satisfy these constraints with equality (*active constraints*) are called *support vectors*.
- (Sufficiently small) changes to any other data points will not affect the maximal margin hyperplane.



Support vector machines

- The points that satisfy these constraints with equality (*active constraints*) are called *support vectors*.
- (Sufficiently small) changes to any other data points will not affect the maximal margin hyperplane.
- We can think of the set of support vectors as a *compressed* version of the training data.



Support vector machines: only uses inner products

Properties:

- $\theta = \sum_i \alpha^i y^i x^i$.
($\alpha^i \neq 0$ only for support vectors.)

Support vector machines: only uses inner products

Properties:

- $\theta = \sum_i \alpha^i y^i x^i$.
($\alpha^i \neq 0$ only for support vectors.)
- The only properties of the data that we use are inner products:

$$\theta \cdot x = \sum_i \alpha^i y^i x^i \cdot x,$$

$$\|\theta\|^2 = \sum_{i,j} \alpha^i \alpha^j y^i y^j x^i \cdot x^j.$$

Support vector machines: only uses inner products

Properties:

- $\theta = \sum_i \alpha^i y^i x^i$.
($\alpha^i \neq 0$ only for support vectors.)
- The only properties of the data that we use are inner products:

$$\theta \cdot x = \sum_i \alpha^i y^i x^i \cdot x,$$

$$\|\theta\|^2 = \sum_{i,j} \alpha^i \alpha^j y^i y^j x^i \cdot x^j.$$

- So we can work with data in any inner product space (not just finite-dimensional vectors).

Hard margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)\end{array}$$

- What if the data are not linearly separable?
There will not be any feasible value of θ .
(Cannot satisfy all of the inequalities $y^i \theta \cdot x^i \geq 1$.)

Hard margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)\end{array}$$

- What if the data are not linearly separable?
There will not be any feasible value of θ .
(Cannot satisfy all of the inequalities $y^i \theta \cdot x^i \geq 1$.)
- We can relax these inequalities $y^i \theta \cdot x^i \geq 1$:

Hard margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)\end{array}$$

- What if the data are not linearly separable?
There will not be any feasible value of θ .
(Cannot satisfy all of the inequalities $y^i \theta \cdot x^i \geq 1$.)
- We can relax these inequalities $y^i \theta \cdot x^i \geq 1$:
 - Introduce slack variables: $\xi_i \geq 0$,

Hard margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n)\end{array}$$

- What if the data are not linearly separable?
There will not be any feasible value of θ .
(Cannot satisfy all of the inequalities $y^i \theta \cdot x^i \geq 1$.)
- We can relax these inequalities $y^i \theta \cdot x^i \geq 1$:
 - Introduce slack variables: $\xi_i \geq 0$,
 - Enforce the milder constraint: $y^i \theta \cdot x^i \geq 1 - \xi_i$,

Hard margin SVM

$$\begin{array}{ll} \min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 \quad (i = 1, \dots, n) \end{array}$$

- What if the data are not linearly separable?
There will not be any feasible value of θ .
(Cannot satisfy all of the inequalities $y^i \theta \cdot x^i \geq 1$.)
- We can relax these inequalities $y^i \theta \cdot x^i \geq 1$:
 - Introduce slack variables: $\xi_i \geq 0$,
 - Enforce the milder constraint: $y^i \theta \cdot x^i \geq 1 - \xi_i$,
 - Add the ξ_i to the criterion.

Hard margin SVM

Quadratic Program

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1\end{array}$$

$$(i = 1, \dots, n)$$

Hard→Soft margin SVM

$$\min_{\theta} \quad \|\theta\|^2$$

$$\text{s.t.} \quad y^i \theta \cdot x^i \geq 1$$

$$(i = 1, \dots, n)$$

Hard→Soft margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 - \xi_i \\ & (i = 1, \dots, n)\end{array}$$

Hard→Soft margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad (i = 1, \dots, n)\end{array}$$

Hard→Soft margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 + \sum_{i=1}^n \xi_i \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad (i = 1, \dots, n)\end{array}$$

Hard→Soft margin SVM

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad (i = 1, \dots, n)\end{array}$$

Soft margin SVM

Quadratic Program

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad (i = 1, \dots, n)\end{array}$$

Soft margin SVM

Quadratic Program

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad (i = 1, \dots, n)\end{array}$$

- The parameter C adjusts the trade-off: $\|\theta\|^2$ versus fit to the data.

Soft margin SVM

Quadratic Program

$$\begin{array}{ll}\min_{\theta} & \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & y^i \theta \cdot x^i \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad (i = 1, \dots, n)\end{array}$$

- The parameter C adjusts the trade-off: $\|\theta\|^2$ versus fit to the data.
- The inequalities imply $\xi_i = (1 - y^i \theta \cdot x^i)_+$.

Support vector machines

Soft margin SVM

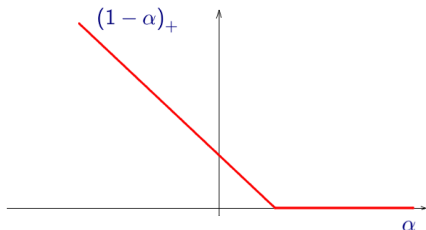
$$\min_{\theta} \quad \|\theta\|^2 + C \sum_{i=1}^n (1 - y^i \theta \cdot x^i)_+.$$

Support vector machines

Soft margin SVM

$$\min_{\theta} \quad \|\theta\|^2 + C \sum_{i=1}^n (1 - y^i \theta \cdot x^i)_+.$$

SVM margin cost function:

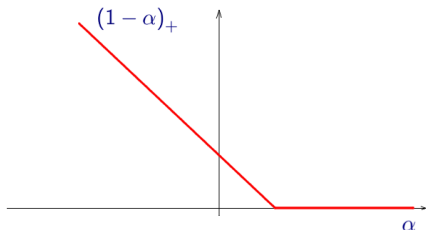


Support vector machines

Soft margin SVM

$$\min_{\theta} \quad \|\theta\|^2 + C \sum_{i=1}^n (1 - y^i \theta \cdot x^i)_+.$$

SVM margin cost function:



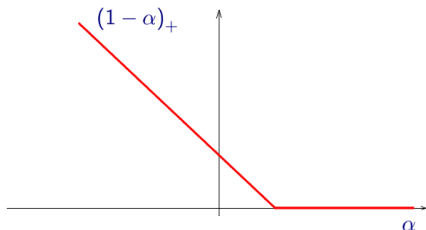
- Convex relaxation (upper bound) for misclassifications:

Support vector machines

Soft margin SVM

$$\min_{\theta} \quad \|\theta\|^2 + C \sum_{i=1}^n (1 - y^i \theta \cdot x^i)_+.$$

SVM margin cost function:



- Convex relaxation (upper bound) for misclassifications:

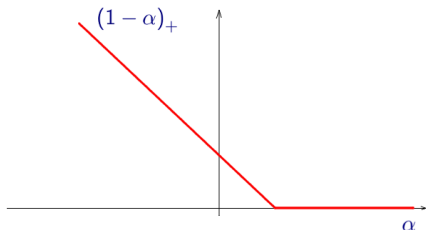
$$1[\alpha \leq 0] = \begin{cases} 1 & \text{if } \alpha \leq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Support vector machines

Soft margin SVM

$$\min_{\theta} \quad \|\theta\|^2 + C \sum_{i=1}^n (1 - y^i \theta \cdot x^i)_+.$$

SVM margin cost function:

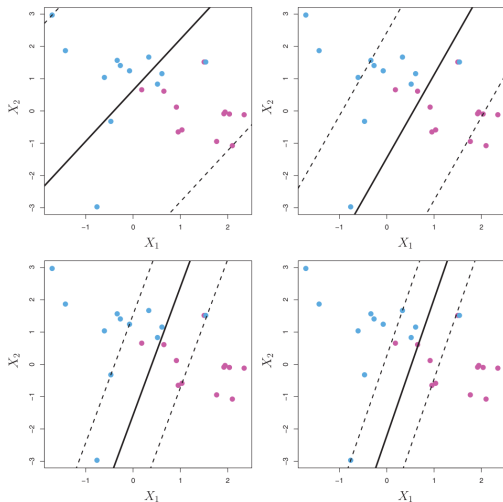


- Convex relaxation (upper bound) for misclassifications:

$$1[\alpha \leq 0] = \begin{cases} 1 & \text{if } \alpha \leq 0, \\ 0 & \text{otherwise.} \end{cases} \leq (1 - \alpha)_+.$$

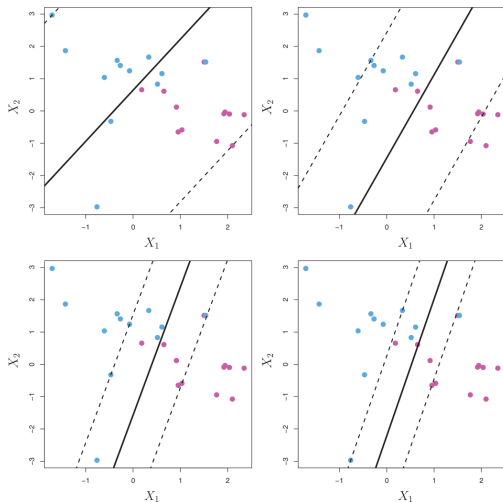
Support vector machines

Varying C :



Support vector machines

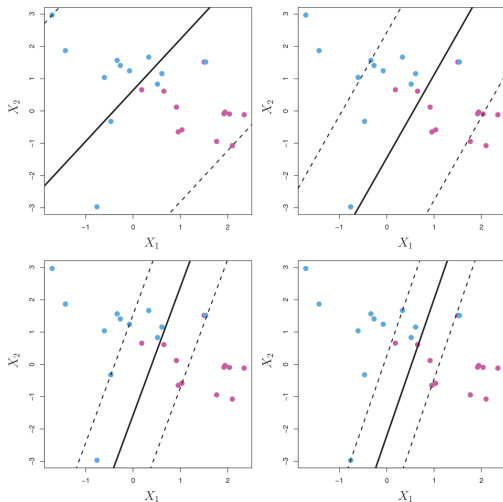
Varying C :



Choosing C :

Support vector machines

Varying C :

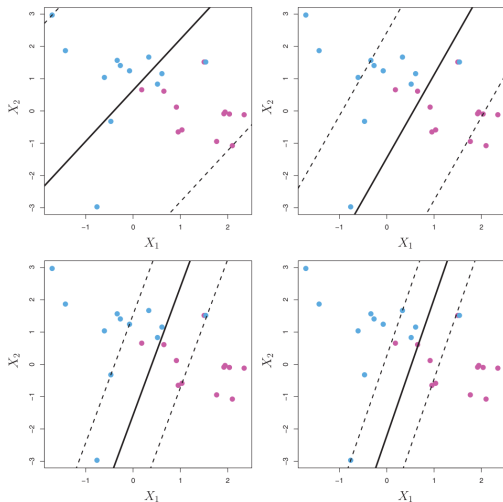


Choosing C :

- *Regularization* parameter.

Support vector machines

Varying C :

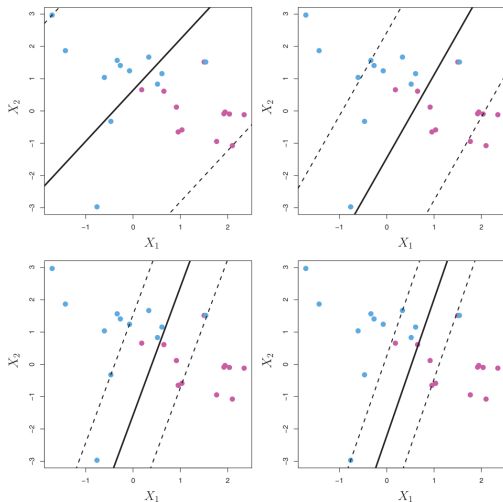


Choosing C :

- *Regularization* parameter.
- Small $C \rightarrow$ small $\|\theta\|^2$.

Support vector machines

Varying C :

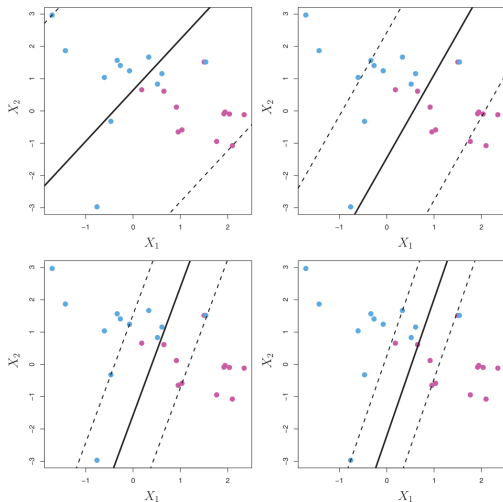


Choosing C :

- *Regularization* parameter.
- Small $C \rightarrow$ small $\|\theta\|^2$.
- Large $C \rightarrow$ small misclassification term.

Support vector machines

Varying C :

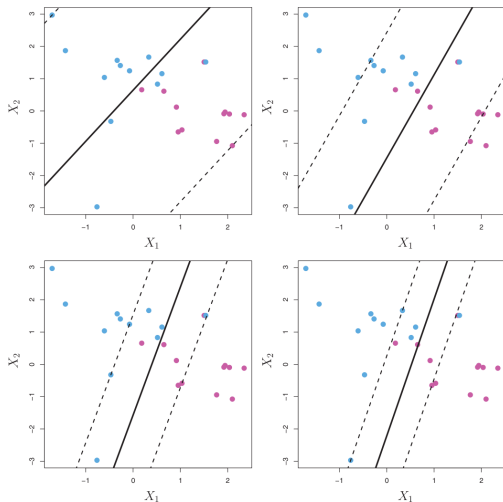


Choosing C :

- *Regularization* parameter.
- Small $C \rightarrow$ small $\|\theta\|^2$.
- Large $C \rightarrow$ small misclassification term.
- Choose C using cross-validation.

Support vector machines

Varying C :



Choosing C :

- *Regularization* parameter.
- Small $C \rightarrow$ small $\|\theta\|^2$.
- Large $C \rightarrow$ small misclassification term.
- Choose C using cross-validation.
- Alternative formulations allow a more intuitive parameterization: ν -SVM. Parameter ν is (roughly) proportion of support vectors.

Soft margin SVM: only uses inner products

Properties:

- We'll see that $\theta = \sum_i \alpha^i y^i x^i$.

Soft margin SVM: only uses inner products

Properties:

- We'll see that $\theta = \sum_i \alpha^i y^i x^i$.
- $\alpha^i \neq 0$ only for *support vectors*, which satisfy $y^i \theta \cdot x^i \leq 1$.

Soft margin SVM: only uses inner products

Properties:

- We'll see that $\theta = \sum_i \alpha^i y^i x^i$.
- $\alpha^i \neq 0$ only for *support vectors*, which satisfy $y^i \theta \cdot x^i \leq 1$.
- The only properties of the data that we use are inner products:

$$\theta \cdot x = \sum_i \alpha^i y^i x^i \cdot x,$$

$$\|\theta\|^2 = \sum_{i,j} \alpha^i \alpha^j y^i y^j x^i \cdot x^j.$$

Soft margin SVM: only uses inner products

Properties:

- We'll see that $\theta = \sum_i \alpha^i y^i x^i$.
- $\alpha^i \neq 0$ only for *support vectors*, which satisfy $y^i \theta \cdot x^i \leq 1$.
- The only properties of the data that we use are inner products:

$$\theta \cdot x = \sum_i \alpha^i y^i x^i \cdot x,$$

$$\|\theta\|^2 = \sum_{i,j} \alpha^i \alpha^j y^i y^j x^i \cdot x^j.$$

- So we can work with data in any inner product space (not just finite-dimensional vectors).

Linear Classifiers: Outline for today

- Properties of linear classifiers.
- Finding a good separating hyperplane.
 - The perceptron algorithm (Rosenblatt, 1950s).
 - The perceptron convergence theorem.
 - The perceptron algorithm as a stochastic gradient method.
 - Support vector machines (Vapnik, 1990s).
 - Hard margin SVM.
 - Soft margin SVM.
 - Both are quadratic programs.
 - Only use inner products.