# COMS 4771 Lecture 8

1. Support vector machines
2. Support vectors and the SVM dual
3. Kernelizing SVMs (maybe next time)
4. Soft-margin SVMs and surrogate losses (maybe next time)

# SUPPORT VECTOR MACHINES

### Setting: two linearly separable classes

Assume there is a linear classifier with **zero training error** on $S$: for some $\boldsymbol{w}_\star \in \mathbb{R}^d$ and $\theta_\star \in \mathbb{R}$,

$$y(\langle \boldsymbol{w}_\star, \boldsymbol{x} \rangle - \theta_\star) > 0, \quad \text{for all } (x, y) \in S.$$

# LINEAR CLASSIFIERS

## Setting: two linearly separable classes

Assume there is a linear classifier with **zero training error** on $S$: for some $\boldsymbol{w}_\star \in \mathbb{R}^d$ and $\theta_\star \in \mathbb{R}$,

$$y(\langle \boldsymbol{w}_\star, \boldsymbol{x} \rangle - \theta_\star) > 0, \quad \text{for all } (x, y) \in S.$$

## Linear programming

Solve linear feasibility problem: find $\boldsymbol{w} \in \mathbb{R}^d$ and $\theta \in \mathbb{R}$ such that

$$y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0, \quad \text{for all } (x, y) \in S.$$

[$d + 1$ variables, $|S|$ constraints.] Can find *some* solution in polynomial time.

# LINEAR CLASSIFIERS

## Setting: two linearly separable classes

Assume there is a linear classifier with **zero training error** on $S$: for some $\boldsymbol{w}_\star \in \mathbb{R}^d$ and $\theta_\star \in \mathbb{R}$,

$$y(\langle \boldsymbol{w}_\star, \boldsymbol{x} \rangle - \theta_\star) > 0, \quad \text{for all } (x, y) \in S.$$

## Linear programming

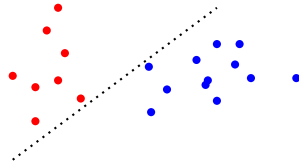Solve linear feasibility problem: find $\boldsymbol{w} \in \mathbb{R}^d$ and $\theta \in \mathbb{R}$ such that

$$y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0, \quad \text{for all } (x, y) \in S.$$

[$d + 1$ variables, $|S|$ constraints.] Can find *some* solution in polynomial time.

## Perceptron algorithm     as if there is a margin

Assuming $\gamma(S) > 0$ finds *some* linear classifier that separates the classes in $S$ *very quickly*.

## Motivation

► Ambiguity in what LP and Perceptron returns, even when $\gamma(S) > 0$.

► Returning an arbitrary linear classifier that separates the classes in $S$ is not very *stable*.

► Unclear what to do (or if algorithms work) when $S$ is not linearly separable, but you still hope to find a good linear separator with small classification error.
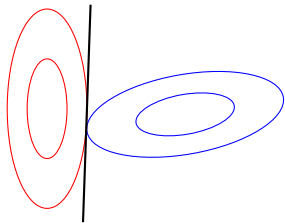
# Support vector machines (SVMs)

## Motivation

also include different data

- ► Ambiguity in what LP and Perceptron returns, even when $\gamma(S) > 0$.
- ► Returning an arbitrary linear classifier that separates the classes in $S$ is not very *stable*.   depends on input sequence, may not the maximum margin solution
- ► Unclear what to do (or if algorithms work) when $S$ is not linearly separable, but you still hope to find a good linear separator with small classification error.

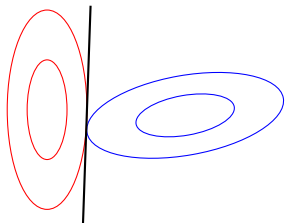## Support vector machines (Vapnik and Chervonenkis, 1963)

- ► Characterize a *stable* solution for linearly separable problems—the **maximum margin solution**.
- ► SVM specified as solution to a **convex optimization problem** that can be solved in polynomial time.
- ► Kernelizable via **convex duality**. (SVM gets its name from its dual form.)
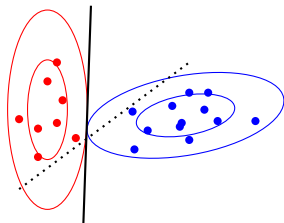- ► Slight alteration to optimization problem gives natural way to handle non-separable cases via **convex surrogate losses**.

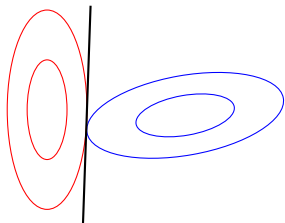Best linear classifier on
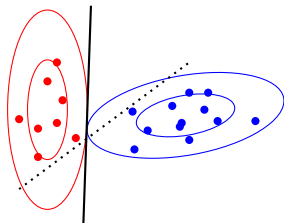population

Best linear classifier on
population

Possible Perceptron or LP
solution on training data $S$

# MAXIMUM MARGIN SOLUTION



Best linear classifier on population

Possible Perceptron or LP solution on training data $S$

"Maximum margin" solution on training data $S$
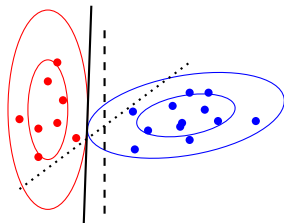
# Maximum margin solution



Best linear classifier on population

Possible Perceptron or LP solution on training data $S$

"Maximum margin" solution on training data $S$

**Why use the "maximum margin" solution?**

(i) Uniquely determined by $S$ (except in degenerate cases), unlike LP's/Perceptron's.

(ii) It is a particular "learning bias"—i.e., an assumption about the problem—that seems to be commonly useful.

Best linear classifier on
population

Possible Perceptron or LP
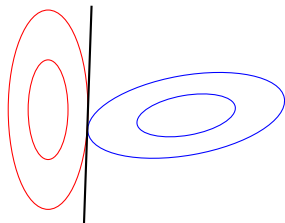solution on training data $S$

"Maximum margin" solution on
training data $S$
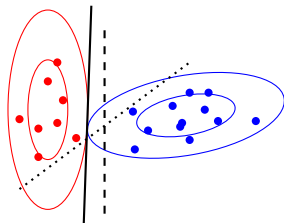
Maximize minimum
distance

**Why use the "maximum margin" solution**?

(i) Uniquely determined by $S$ (except in degenerate cases), unlike LP's/Perceptron's.

(ii) It is a particular "learning bias"—i.e., an assumption about the problem—that seems to be commonly useful.

**Next**: mathematically characterize the "maximum margin" solution.

### Definition
The **distance between a point $x$ and a set $A$** is the Euclidean distance
between $x$ and the closest point in $A$:

$$\operatorname{dist}(\boldsymbol{x}, A) := \min_{\boldsymbol{x}' \in A} \|\boldsymbol{x} - \boldsymbol{x}'\|$$

# DISTANCES TO SETS

### Definition

The **distance between a point $x$ and a set $A$** is the Euclidean distance between $x$ and the closest point in $A$:

$$\text{dist}(\boldsymbol{x}, A) := \min_{\boldsymbol{x}' \in A} \|\boldsymbol{x} - \boldsymbol{x}'\|$$

# MARGIN

If $f_{\boldsymbol{w},\theta}$ linearly separates $S$, then we say its **margin on $S$** is

$$\min_{\boldsymbol{x} \in S_\oplus \cup S_\ominus} \operatorname{dist}(\boldsymbol{x}, H)$$

where

$$H := \text{affine hyperplane corresponding to } f_{\boldsymbol{w},\theta},$$
$$S_\oplus := \{\boldsymbol{x} : (\boldsymbol{x}, +1) \in S\} \quad \text{(positive points)},$$
$$S_\ominus := \{\boldsymbol{x} : (\boldsymbol{x}, -1) \in S\} \quad \text{(negative points)}.$$

Consider linear classifier $f_{\boldsymbol{w},\theta}$ (where $\boldsymbol{w} \in \mathbb{R}^d \setminus \{\boldsymbol{0}\}$ and $\theta \in \mathbb{R}$).

Consider linear classifier $f_{\boldsymbol{w},\theta}$ (where $\boldsymbol{w} \in \mathbb{R}^d \setminus \{\boldsymbol{0}\}$ and $\theta \in \mathbb{R}$).



▶ Correct classification on $(\boldsymbol{x}, y)$:

$$f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y \quad \text{iff} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0.$$

(Ignore boundary case.)

# DISTANCE TO THE BOUNDARY

Consider linear classifier $f_{\boldsymbol{w},\theta}$ (where $\boldsymbol{w} \in \mathbb{R}^d \setminus \{\boldsymbol{0}\}$ and $\theta \in \mathbb{R}$).



▶ Correct classification on $(\boldsymbol{x}, y)$:

$$f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y \quad \text{iff} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0.$$

(Ignore boundary case.)

▶ Projection of $\boldsymbol{x}$ onto $\mathrm{span}(\boldsymbol{w})$ is $\dfrac{\langle \boldsymbol{w}, \boldsymbol{x} \rangle}{\|\boldsymbol{w}\|_2^2} \boldsymbol{w}$.

Consider linear classifier $f_{\boldsymbol{w},\theta}$ (where $\boldsymbol{w} \in \mathbb{R}^d \setminus \{\boldsymbol{0}\}$ and $\theta \in \mathbb{R}$).



- Correct classification on $(\boldsymbol{x}, y)$:

  $$f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y \quad \text{iff} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0.$$

  (Ignore boundary case.)

- Projection of $\boldsymbol{x}$ onto $\mathrm{span}(\boldsymbol{w})$ is $\dfrac{\langle \boldsymbol{w}, \boldsymbol{x} \rangle}{\|\boldsymbol{w}\|_2^2} \boldsymbol{w}$.

- Distance to affine hyperplane $H$ is

  $$\mathrm{dist}(\boldsymbol{x}, H) = \frac{|\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta|}{\|\boldsymbol{w}\|_2}.$$

# DISTANCE TO THE BOUNDARY

Consider linear classifier $f_{\boldsymbol{w},\theta}$ (where $\boldsymbol{w} \in \mathbb{R}^d \setminus \{\boldsymbol{0}\}$ and $\theta \in \mathbb{R}$).



▶ Correct classification on $(\boldsymbol{x}, y)$:

$$f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y \quad \text{iff} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0.$$

(Ignore boundary case.)

▶ Projection of $\boldsymbol{x}$ onto $\mathrm{span}(\boldsymbol{w})$ is $\dfrac{\langle \boldsymbol{w}, \boldsymbol{x} \rangle}{\|\boldsymbol{w}\|_2^2} \boldsymbol{w}$.

▶ Distance to affine hyperplane $H$ is

$$\mathrm{dist}(\boldsymbol{x}, H) = \frac{|\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta|}{\|\boldsymbol{w}\|_2}.$$

▶ If $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y$, then

$$\mathrm{dist}(\boldsymbol{x}, H) = \frac{y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)}{\|\boldsymbol{w}\|_2}.$$

### Distance of a correctly classified point to the boundary

If $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y$  (i.e., $y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0$), then

$$\text{dist}(\boldsymbol{x}, H) = \frac{y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)}{\|\boldsymbol{w}\|_2}.$$

### Distance of a correctly classified point to the boundary

If $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y$ (i.e., $y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0$), then

$$\mathrm{dist}(\boldsymbol{x}, H) = \frac{y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)}{\|\boldsymbol{w}\|_2}.$$

### Maximizing the distances to the boundary

Therefore, find $(\boldsymbol{w}, \theta)$ such that

$$y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0 \quad \text{for all } (\boldsymbol{x}, y) \in S$$

with $\|\boldsymbol{w}\|_2$ as small as possible.

## Distance of a correctly classified point to the boundary

If $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = y$ (i.e., $y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0$), then

$$\mathrm{dist}(\boldsymbol{x}, H) = \frac{y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)}{\|\boldsymbol{w}\|_2}.$$

## Maximizing the distances to the boundary

Therefore, find $(\boldsymbol{w}, \theta)$ such that

$$y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0 \quad \text{for all } (\boldsymbol{x}, y) \in S$$

with $\|\boldsymbol{w}\|_2$ as small as possible.

**Note**: If we have *any* $(\boldsymbol{w}, \theta)$ such that

$$y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) > 0 \quad \text{for all } (\boldsymbol{x}, y) \in S,$$

then can rescale $(\boldsymbol{w}, \theta)$ so that

$$y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \geq 1 \quad \text{for all } (\boldsymbol{x}, y) \in S. \quad \text{how to scale?}$$

# Maximum margin linear classifier

The solution $(\hat{\boldsymbol{w}}, \hat{\theta})$ to the following optimization problem:

$$
\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2
$$
$$
\text{s.t.} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \geq 1 \qquad \text{for all } (\boldsymbol{x}, y) \in S
$$

gives the **linear classifier with the maximum margin on $S$**.
(The $1/2$ and squaring is for mathematical convenience and now the convention.)

# MAXIMUM MARGIN LINEAR CLASSIFIER

The solution $(\hat{\boldsymbol{w}}, \hat{\theta})$ to the following optimization problem:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \geq 1 \qquad \text{for all } (\boldsymbol{x}, y) \in S$$

gives the **linear classifier with the maximum margin on $S$**.
(The $1/2$ and squaring is for mathematical convenience and now the convention.)

The linear classifier obtained by solving this optimization problem is called a **support vector machine**.

# Maximum margin linear classifier

The solution $(\hat{\boldsymbol{w}}, \hat{\theta})$ to the following optimization problem:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{2}\|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \geq 1 \quad \text{for all } (\boldsymbol{x}, y) \in S$$

gives the **linear classifier with the maximum margin on $S$**.
(The $1/2$ and squaring is for mathematical convenience and now the convention.)

The linear classifier obtained by solving this optimization problem is called a **support vector machine**.

The optimization problem is a **convex optimization problem** that can be solved in polynomial time. (Actual algorithm to come later.)

# Maximum margin linear classifier

The solution $(\hat{\boldsymbol{w}}, \hat{\theta})$ to the following optimization problem:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \geq 1 \qquad \text{for all } (\boldsymbol{x}, y) \in S$$

gives the **linear classifier with the maximum margin on $S$**.
(The $1/2$ and squaring is for mathematical convenience and now the convention.)

The linear classifier obtained by solving this optimization problem is called a **support vector machine**.

The optimization problem is a **convex optimization problem** that can be solved in polynomial time. (Actual algorithm to come later.)

If there is a solution (i.e., the problem is separable), then the solution is *unique*. (Compare to LP's and Perceptron's lack of determinism from $S$.)

# SUPPORT VECTORS AND THE SVM DUAL

- Explain how the SVM solution is *entirely determined by certain data points called the "support vectors"*

  In other words: can throw away all data except the support vectors, re-solve SVM optimization problem, and get the same solution.

- Illustrate power of *convex duality*—a useful analytic (& algorithmic) tool.

  (Linear programming duality is a special case.)

## Observation

Where a separating hyperplane may be placed depends on the "outer" points on the sets. Points in the interior do not matter.

# Substituting convex sets

## Observation

Where a separating hyperplane may be placed depends on the "outer" points on the sets. Points in the interior do not matter.



## In geometric terms

For each class, replace all points from the class with *the smallest* **convex set** *which contains all these points*:

# Substituting convex sets

## Observation

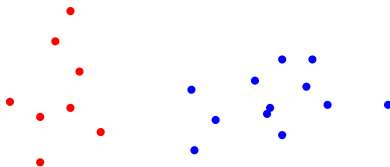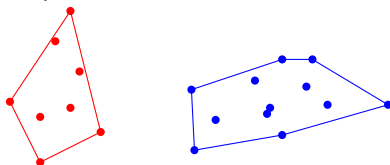Where a separating hyperplane may be placed depends on the "outer" points on the sets. Points in the interior do not matter.



## In geometric terms

For each class, replace all points from the class with *the smallest* **convex set** *which contains all these points*:



a line to separate boudary is ok

A set $C$ is **convex** if it contains the line segments between all pairs of points in $C$.

## Definition

If $A \subset \mathbb{R}^d$ is a set of points, the smallest convex set containing all points in $A$ is called the **convex hull** of $A$, denoted $\mathrm{conv}(A)$.



"Corner points" of the convex set are called **extreme points**.

# Substituting convex sets

## Definition

If $A \subset \mathbb{R}^d$ is a set of points, the smallest convex set containing all points in $A$ is called the **convex hull** of $A$, denoted $\mathrm{conv}(A)$.



"Corner points" of the convex set are called **extreme points**.

## Barycentric coordinates

If convex set $C$ has extreme points $\boldsymbol{v}_1, \ldots, \boldsymbol{v}_m$, then can write any $\boldsymbol{x} \in C$ as

$$\boldsymbol{x} = \sum_{i=1}^{m} \alpha_i \boldsymbol{v}_i \quad \text{where } \alpha_i \geq 0, \ \sum_{i=1}^{m} \alpha_i = 1.$$

Called the **barycentric coordinates** of $\boldsymbol{x}$.

## Key idea

A hyperplane separates two classes if and only if it separates their convex hulls.

If $f_{\boldsymbol{w},\theta}$ linearly separates $S$, then we say its **margin on** $S$ is

$$\min\{\mathrm{dist}(H, \mathrm{conv}(S_\oplus)),\ \mathrm{dist}(H, \mathrm{conv}(S_\ominus))\}.$$



Distance between two sets $A$ and $B$: $\mathrm{dist}(A, B) := \min_{\boldsymbol{x} \in A} \mathrm{dist}(\boldsymbol{x}, B).$

# SUPPORT VECTORS

Extreme points of the convex hulls
closest to the hyperplane are the
**support vectors**.



## Implications

- ▶ SVM optimization problem focuses attention to the area closest to the
  decision surface.
- ▶ SVM classifier entirely determined by support vectors.
  (Much like mistake examples in Perceptron—hope for kernelization?)

### Questions

- Can we determine which points are the support vectors?
- Can we kernelize SVM like we did Perceptron?

# THE DUAL OF THE SVM OPTIMIZATION PROBLEM

### Questions

▶ Can we determine which points are the support vectors?
▶ Can we kernelize SVM like we did Perceptron?

These questions are answered by looking at a corresponding *dual optimization problem*, which we will derive from geometric principles.

## Questions

- Can we determine which points are the support vectors?
- Can we kernelize SVM like we did Perceptron?

These questions are answered by looking at a corresponding *dual optimization problem*, which we will derive from geometric principles.

(Can also derive the dual problem using Lagrangian duality.)

**Fact**: The closest distance between a point $x$ and a convex set $A$ is the maximum over the distances between $x$ and all (affine) hyperplanes that separate $x$ and $A$.

$$\mathrm{dist}(\boldsymbol{x}, A) \;=\; \max_{\substack{H \text{ separating} \\ \boldsymbol{x} \text{ from } A}} \mathrm{dist}(\boldsymbol{x}, H)$$

Convert Distance between convex to Margin!!!

$A$

$x$

The maximum between between x and all H.

$H$

$H'$

## Main idea

Equivalent definition of max-margin affine hyperplane:

1. Find shortest line segment connecting the convex hulls.

2. Place affine hyperplane orthogonal to line segment at midpoint.

## The dual problem

$$\min_{\substack{\boldsymbol{u} \in \mathrm{conv}(S_\oplus), \\ \boldsymbol{v} \in \mathrm{conv}(S_\ominus)}} \frac{1}{2} \|\boldsymbol{u} - \boldsymbol{v}\|_2^2$$

The closest distance between x with set A. is equal to the maximum distance(Margin) over x to all hyperplanes

## Dual problem

$$\min_{\substack{\boldsymbol{u}\in\mathrm{conv}(S_\oplus),\\ \boldsymbol{v}\in\mathrm{conv}(S_\ominus)}} \frac{1}{2}\|\boldsymbol{u}-\boldsymbol{v}\|_2^2.$$

### Dual problem

$$\min_{\substack{\boldsymbol{u}\in\mathrm{conv}(S_\oplus),\\ \boldsymbol{v}\in\mathrm{conv}(S_\ominus)}} \frac{1}{2}\|\boldsymbol{u}-\boldsymbol{v}\|_2^2.$$

### New parameterization

**Express $\boldsymbol{u}$ and $\boldsymbol{v}$ in barycentric coordinates**.
Let $S_\oplus = \{\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(m)}\}$, $S_\ominus = \{\boldsymbol{x}^{(m+1)}, \boldsymbol{x}^{(m+2)}, \ldots, \boldsymbol{x}^{(n)}\}$,

$$\boldsymbol{u} = \sum_{i=1}^{m} \alpha_i \boldsymbol{x}^{(i)}, \qquad \boldsymbol{v} = \sum_{i=m+1}^{n} \alpha_i \boldsymbol{x}^{(i)}.$$

### Dual problem

$$\min_{\substack{\boldsymbol{u}\in\text{conv}(S_\oplus),\\ \boldsymbol{v}\in\text{conv}(S_\ominus)}} \frac{1}{2}\|\boldsymbol{u}-\boldsymbol{v}\|_2^2.$$

### New parameterization

**Express $\boldsymbol{u}$ and $\boldsymbol{v}$ in barycentric coordinates**.
Let $S_\oplus = \{\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(m)}\}$, $S_\ominus = \{\boldsymbol{x}^{(m+1)}, \boldsymbol{x}^{(m+2)}, \ldots, \boldsymbol{x}^{(n)}\}$,

some ai may equal to 0.

$$\boldsymbol{u} = \sum_{i=1}^{m} \alpha_i \boldsymbol{x}^{(i)}, \qquad \boldsymbol{v} = \sum_{i=m+1}^{n} \alpha_i \boldsymbol{x}^{(i)}.$$

**Dual problem in terms of $\alpha_1, \ldots, \alpha_n$** (called "dual variables"):

$$\min_{\alpha_1,\ldots,\alpha_n} \quad \frac{1}{2}\left\|\sum_{i=1}^{m} \alpha_i \boldsymbol{x}^{(i)} - \sum_{i=m+1}^{n} \alpha_i \boldsymbol{x}^{(i)}\right\|_2^2$$

$$\text{s.t.} \quad \sum_{i=1}^{m} \alpha_i = \sum_{i=m+1}^{n} \alpha_i = 1, \quad \alpha_i \geq 0.$$

### Dual problem

$$\min_{\substack{\boldsymbol{u}\in\text{conv}(S_\oplus),\\ \boldsymbol{v}\in\text{conv}(S_\ominus)}} \frac{1}{2}\|\boldsymbol{u}-\boldsymbol{v}\|_2^2.$$

### New parameterization

**Express $\boldsymbol{u}$ and $\boldsymbol{v}$ in barycentric coordinates**.
Let $S_\oplus = \{\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(m)}\}$, $S_\ominus = \{\boldsymbol{x}^{(m+1)}, \boldsymbol{x}^{(m+2)}, \ldots, \boldsymbol{x}^{(n)}\}$,

$$\boldsymbol{u} = \sum_{i=1}^{m} \alpha_i \boldsymbol{x}^{(i)}, \qquad \boldsymbol{v} = \sum_{i=m+1}^{n} \alpha_i \boldsymbol{x}^{(i)}.$$

**Dual problem in terms of $\alpha_1, \ldots, \alpha_n$** (called "dual variables"):

$$\min_{\alpha_1,\ldots,\alpha_n} \quad \frac{1}{2}\left\|\sum_{i=1}^{m} \alpha_i y^{(i)} \boldsymbol{x}^{(i)} + \sum_{i=m+1}^{n} \alpha_i y^{(i)} \boldsymbol{x}^{(i)}\right\|_2^2$$

$$\text{s.t.} \quad \sum_{i=1}^{m} \alpha_i = \sum_{i=m+1}^{n} \alpha_i = 1, \quad \alpha_i \geq 0.$$

### Dual problem

$$\min_{\substack{\boldsymbol{u}\in\text{conv}(S_{\oplus}),\\ \boldsymbol{v}\in\text{conv}(S_{\ominus})}} \frac{1}{2}\|\boldsymbol{u}-\boldsymbol{v}\|_2^2.$$

### New parameterization

**Express $u$ and $v$ in barycentric coordinates**.
Let $S_{\oplus} = \{\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(m)}\}$, $S_{\ominus} = \{\boldsymbol{x}^{(m+1)}, \boldsymbol{x}^{(m+2)}, \ldots, \boldsymbol{x}^{(n)}\}$,

$$\boldsymbol{u} = \sum_{i=1}^{m} \alpha_i \boldsymbol{x}^{(i)}, \qquad \boldsymbol{v} = \sum_{i=m+1}^{n} \alpha_i \boldsymbol{x}^{(i)}.$$

**Dual problem in terms of $\alpha_1, \ldots, \alpha_n$** (called "dual variables"):

$$\min_{\alpha_1,\ldots,\alpha_n} \quad \frac{1}{2}\left\|\sum_{i=1}^{n} \alpha_i y^{(i)} \boldsymbol{x}^{(i)}\right\|_2^2 \qquad \begin{array}{l}\text{negative}\\ \text{great!!!}\end{array}$$

$$\text{s.t.} \quad \sum_{i=1}^{m} \alpha_i = \sum_{i=m+1}^{n} \alpha_i = 1, \quad \alpha_i \geq 0.$$

Expanding the square...

$$\frac{1}{2}\left\|\sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}\right\|_2^2 = \frac{1}{2}\left\langle\sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}, \sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}\right\rangle$$

$$= \frac{1}{2}\sum_{i,j=1}^{n}\alpha_i\alpha_j y^{(i)}y^{(j)}\langle\boldsymbol{x}^{(i)},\boldsymbol{x}^{(j)}\rangle.$$

Expanding the square. . .

$$
\begin{aligned}
\frac{1}{2}\left\|\sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}\right\|_2^2 &= \frac{1}{2}\left\langle \sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}, \sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}\right\rangle \\
&= \frac{1}{2}\sum_{i,j=1}^{n}\alpha_i\alpha_j y^{(i)}y^{(j)}\langle \boldsymbol{x}^{(i)},\boldsymbol{x}^{(j)}\rangle.
\end{aligned}
$$

Usual form of the dual problem

$$
\begin{aligned}
\min_{\alpha_1,\ldots,\alpha_n} \quad & \frac{1}{2}\sum_{i,j=1}^{n}\alpha_i\alpha_j y^{(i)}y^{(j)}\langle \boldsymbol{x}^{(i)},\boldsymbol{x}^{(j)}\rangle \\
\text{s.t.} \quad & \sum_{i=1}^{m}\alpha_i = \sum_{i=m+1}^{n}\alpha_i = 1, \quad \alpha_i \geq 0.
\end{aligned}
$$

Expanding the square...

$$\frac{1}{2}\left\|\sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}\right\|_2^2 = \frac{1}{2}\left\langle\sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}, \sum_{i=1}^{n}\alpha_i y^{(i)}\boldsymbol{x}^{(i)}\right\rangle$$

$$= \frac{1}{2}\sum_{i,j=1}^{n}\alpha_i\alpha_j y^{(i)}y^{(j)}\langle\boldsymbol{x}^{(i)},\boldsymbol{x}^{(j)}\rangle.$$

Usual form of the dual problem

$$\max_{\alpha_1,\ldots,\alpha_n} \quad \sum_{i=1}^{n}\alpha_i - \frac{1}{2}\sum_{i,j=1}^{n}\alpha_i\alpha_j y^{(i)}y^{(j)}\langle\boldsymbol{x}^{(i)},\boldsymbol{x}^{(j)}\rangle$$

$$\text{s.t.} \quad \sum_{i=1}^{m}\alpha_i y^{(i)} = 0, \quad \alpha_i \geq 0. \quad \text{(Don't worry about getting to this form.)}$$

# USING THE DUAL SOLUTION

## Getting a linear classifier from the dual solution

Let $\hat{\alpha}_1, \ldots, \hat{\alpha}_n$ be the solution to the dual problem. Then

$$\hat{\boldsymbol{w}} := \hat{\boldsymbol{u}} - \hat{\boldsymbol{v}} = \sum_{i=1}^{n} \hat{\alpha}_i y^{(i)} \boldsymbol{x}^{(i)};$$

$$\hat{\theta} := \frac{\min\limits_{\boldsymbol{x} \in S_\oplus} \langle \hat{\boldsymbol{w}}, \boldsymbol{x} \rangle + \max\limits_{\boldsymbol{x} \in S_\ominus} \langle \hat{\boldsymbol{w}}, \boldsymbol{x} \rangle}{2}.$$

(We want $\hat{\theta} := b\|\hat{\boldsymbol{w}}\|_2^2$ such that $b\hat{\boldsymbol{w}}$ is the midpoint between $\hat{\boldsymbol{u}}$ and $\hat{\boldsymbol{v}}$.)



convex combination

## Getting a linear classifier from the dual solution

Let $\hat\alpha_1, \dots, \hat\alpha_n$ be the solution to the dual problem. Then

$$\hat{\boldsymbol{w}} := \hat{\boldsymbol{u}} - \hat{\boldsymbol{v}} = \sum_{i=1}^{n} \hat\alpha_i y^{(i)} \boldsymbol{x}^{(i)};$$

$$\hat\theta := \frac{\min\limits_{\boldsymbol{x} \in S_\oplus} \langle \hat{\boldsymbol{w}}, \boldsymbol{x}\rangle + \max\limits_{\boldsymbol{x} \in S_\ominus} \langle \hat{\boldsymbol{w}}, \boldsymbol{x}\rangle}{2}.$$

(We want $\hat\theta := b\|\hat{\boldsymbol{w}}\|_2^2$ such that $b\hat{\boldsymbol{w}}$ is the midpoint between $\hat{\boldsymbol{u}}$ and $\hat{\boldsymbol{v}}$.)

$\hat{\boldsymbol{w}}$ depends on $(\boldsymbol{x}^{(i)}, y^{(i)})$ with $\hat\alpha_i > 0$, i.e., the **support vectors**.
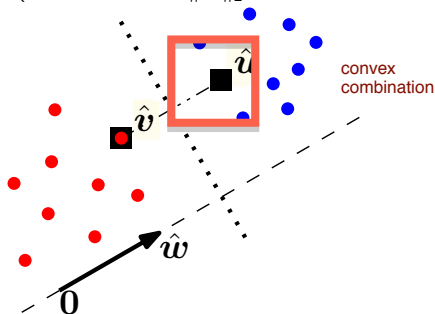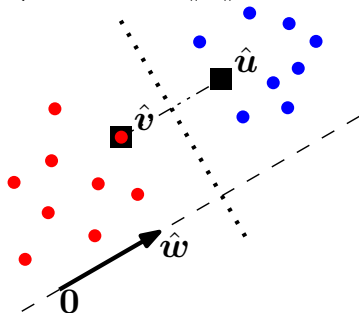
# USING THE DUAL SOLUTION

## Getting a linear classifier from the dual solution

Let $\hat{\alpha}_1, \ldots, \hat{\alpha}_n$ be the solution to the dual problem. Then

$$\hat{w} := \hat{u} - \hat{v} = \sum_{i=1}^n \hat{\alpha}_i y^{(i)} x^{(i)};$$

$$\hat{\theta} := \frac{\min_{x \in S_\oplus} \langle \hat{w}, x \rangle + \max_{x \in S_\ominus} \langle \hat{w}, x \rangle}{2}.$$

(We want $\hat{\theta} := b\|\hat{w}\|_2^2$ such that $b\hat{w}$ is the midpoint between $\hat{u}$ and $\hat{v}$.)



$\hat{w}$ depends on $(x^{(i)}, y^{(i)})$ with $\hat{\alpha}_i > 0$, i.e., the **support vectors**.

These support vector examples have

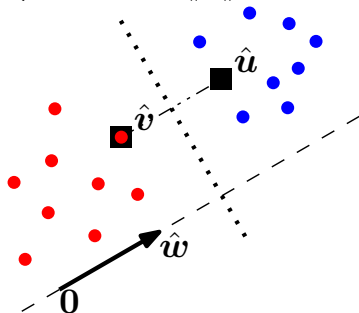$$y^{(i)}\Big(\langle \hat{w}, x^{(i)} \rangle - \hat{\theta}\Big) = \frac{1}{2}\|\hat{w}\|_2^2.$$

Non-support vectors $(x^{(j)}, y^{(j)})$ have

$$y^{(j)}\Big(\langle \hat{w}, x^{(j)} \rangle - \hat{\theta}\Big) > \frac{1}{2}\|\hat{w}\|_2^2.$$

# Solving the primal SVM problem

Technically, to solve the primal SVM problem, need

$$y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq 1 \quad \text{for } i = 1, 2, \ldots, n,$$

but we have

$$y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq \frac{1}{2}\|\boldsymbol{w}\|_2^2 \quad \text{for } i = 1, 2, \ldots, n.$$

# Solving the primal SVM problem

Technically, to solve the primal SVM problem, need

$$y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta\Big) \geq 1 \quad \text{for } i = 1, 2, \ldots, n,$$

but we have

$$y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta\Big) \geq \frac{1}{2}\|\boldsymbol{w}\|_2^2 \quad \text{for } i = 1, 2, \ldots, n.$$

**Solution**: Rescale $\hat{\boldsymbol{w}}$ and $\hat{\theta}$ by dividing by $\frac{1}{2}\|\hat{\boldsymbol{w}}\|_2^2$.

# KERNELIZING SVMs

# KERNELIZING SVMs (BOSER, GUYON, AND VAPNIK, 1992)

SVM dual problem

$$\max_{\alpha_1,\ldots,\alpha_n} \quad \sum_{i=1}^{n}\alpha_i - \frac{1}{2}\sum_{i,j=1}^{n}\alpha_i\alpha_j y^{(i)}y^{(j)}\langle \boldsymbol{x}^{(i)},\boldsymbol{x}^{(j)}\rangle$$

$$\text{s.t.} \quad \sum_{i=1}^{m}\alpha_i y^{(i)} = 0, \quad \alpha_i \geq 0.$$

SVM dual solution

$$\hat{\boldsymbol{w}} := \sum_{i=1}^{n}\hat{\alpha}_i y^{(i)}\boldsymbol{x}^{(i)};$$

$$\hat{\theta} := \frac{\min_{\boldsymbol{x}\in S_\oplus}\langle \hat{\boldsymbol{w}},\boldsymbol{x}\rangle + \max_{\boldsymbol{x}\in S_\ominus}\langle \hat{\boldsymbol{w}},\boldsymbol{x}\rangle}{2}.$$

# KERNELIZING SVMS (BOSER, GUYON, AND VAPNIK, 1992)

SVM dual problem

$$\max_{\alpha_1,\ldots,\alpha_n} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j y^{(i)} y^{(j)} K(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)})$$

$$\text{s.t.} \quad \sum_{i=1}^{m} \alpha_i y^{(i)} = 0, \quad \alpha_i \geq 0.$$

SVM dual solution

$$\langle \hat{\boldsymbol{w}}, \cdot \rangle := \sum_{i=1}^{n} \hat{\alpha}_i y^{(i)} K(\boldsymbol{x}^{(i)}, \cdot);$$

$$\hat{\theta} := \frac{\min_{\boldsymbol{x} \in S_\oplus} K(\hat{\boldsymbol{w}}, \boldsymbol{x}) + \max_{\boldsymbol{x} \in S_\ominus} K(\hat{\boldsymbol{w}}, \boldsymbol{x})}{2}.$$

Just need to keep around the **support vectors** (i.e., examples where $\hat{\alpha}_i > 0$).

# Soft-margin SVMs and surrogate losses

**Non-separable cases**:
No linear classifier has zero training error on $S$.



But if non-separability is only due to a handful of points,
can we still find a good linear classifier?

# SOFT-MARGIN SVMs (CORTES AND VAPNIK, 1995)

When $S = ((\boldsymbol{x}^{(1)}, y^{(1)}), \ldots, (\boldsymbol{x}^{(n)}, y^{(n)}))$ is not linearly separable, the (primal) SVM optimization problem

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad y^{(i)} \Big( \langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta \Big) \geq 1 \qquad \text{for } i = 1, 2, \ldots, n$$

**has no solution**.

# SOFT-MARGIN SVMs (CORTES AND VAPNIK, 1995)

When $S = ((\boldsymbol{x}^{(1)}, y^{(1)}), \ldots, (\boldsymbol{x}^{(n)}, y^{(n)}))$ is not linearly separable, the (primal) SVM optimization problem

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad y^{(i)} \left( \langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta \right) \geq 1 \quad \text{for } i = 1, 2, \ldots, n$$

**has no solution**.

Introduce **slack variables** $\xi_1, \xi_2, \ldots, \xi_n \geq 0$, and a trade-off parameter $C > 0$:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2 + C \sum_{i=1}^{n} \xi_i$$
$$\text{s.t.} \quad y^{(i)} \left( \langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta \right) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \ldots, n$$
$$\xi_i \geq 0 \quad \text{for } i = 1, 2, \ldots, n$$

which is **always feasible**—"soft margin" SVM.

# SOFT-MARGIN SVMs (CORTES AND VAPNIK, 1995)

**Winner of 2008 ACM Paris Kanellakis Award**:

*For "their revolutionary development of a highly effective algorithm known as Support Vector Machines (SVM), a set of related supervised learning methods used for data classification and regression", which is "one of the most frequently used algorithms in machine learning, and is used in medical diagnosis, weather forecasting, and intrusion detection among many other practical applications".*

**Other winners include**: public key cryptography, Lempel-Ziv compression, Splay Trees, interior point method for linear programming, ... and AdaBoost (discussed later in the course).

# SLACK INTERPRETATION

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad \frac{1}{2}\|\boldsymbol{w}\|_2^2 + C\sum_{i=1}^{n} \xi_i$$

$$\text{s.t.} \quad y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, n$$

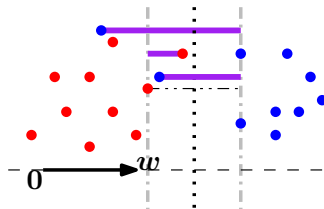$$\xi_i \geq 0 \qquad \text{for } i = 1, 2, \ldots, n$$

# SLACK INTERPRETATION

$$\min_{\boldsymbol{w}\in\mathbb{R}^d,\theta\in\mathbb{R},\boldsymbol{\xi}\in\mathbb{R}^n} \quad \frac{1}{2}\|\boldsymbol{w}\|_2^2 + C\sum_{i=1}^{n}\xi_i$$

$$\text{s.t.} \quad y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, n$$

$$\xi_i \geq 0 \qquad \text{for } i = 1, 2, \ldots, n$$



▶ For a given $(\boldsymbol{w}, \theta)$, $\xi_i/\|\boldsymbol{w}\|_2$ measures distance that $\boldsymbol{x}^{(i)}$ must be moved so that $y^{(i)}\langle \boldsymbol{w}, \boldsymbol{x}^{(i)} - \theta\rangle \geq 1$.

$$\min_{\boldsymbol{w}\in\mathbb{R}^d,\theta\in\mathbb{R},\boldsymbol{\xi}\in\mathbb{R}^n} \quad \frac{1}{2}\|\boldsymbol{w}\|_2^2 + C\sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad y^{(i)}\Big(\langle\boldsymbol{w},\boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, n$$

$$\xi_i \geq 0 \qquad \text{for } i = 1, 2, \ldots, n$$



- For a given $(\boldsymbol{w}, \theta)$, $\xi_i/\|\boldsymbol{w}\|_2$ measures distance that $\boldsymbol{x}^{(i)}$ must be moved so that $y^{(i)}\langle\boldsymbol{w},\boldsymbol{x}^{(i)} - \theta\rangle \geq 1$.
- $C$ controls trade-off between slack penalties and size of margin (which is $1/\|\boldsymbol{w}\|_2$).

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables**: (using $\lambda = 1/(nC)$)

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^{n}\xi_i$$

$$\text{s.t.} \quad y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, n$$

$$\xi_i \geq 0 \qquad \text{for } i = 1, 2, \ldots, n$$

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables**: (using $\lambda = 1/(nC)$)

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^{n} \xi_i$$

$$\text{s.t.} \quad y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta\Big) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, n$$

$$\xi_i \geq 0 \qquad \text{for } i = 1, 2, \ldots, n$$

**Equivalent "hinge loss" form**:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^{n}\Big[1 - y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta\Big)\Big]_+$$

*Notation*: $[a]_+ := \max\{0, a\}$.

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables**: (using $\lambda = 1/(nC)$)

$$\min_{\boldsymbol{w}\in\mathbb{R}^d,\theta\in\mathbb{R},\boldsymbol{\xi}\in\mathbb{R}^n} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, n$$

$$\xi_i \geq 0 \quad \text{for } i = 1, 2, \ldots, n$$

**Equivalent "hinge loss" form**:

$$\min_{\boldsymbol{w}\in\mathbb{R}^d,\theta\in\mathbb{R}} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^n \Big[1 - y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big)\Big]_+$$

*Notation*: $[a]_+ := \max\{0, a\}$.

The **hinge loss** of a linear classifier $f_{\boldsymbol{w},\theta}$ on an example $(\boldsymbol{x}, y)$ is defined to be

$$\mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}, y) := \Big[1 - y(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta)\Big]_+.$$

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables**: (using $\lambda = 1/(nC)$)

$$\min_{\boldsymbol{w}\in\mathbb{R}^d, \theta\in\mathbb{R}, \boldsymbol{\xi}\in\mathbb{R}^n} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^{n}\xi_i$$

$$\text{s.t.} \quad y^{(i)}\Big(\langle \boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\Big) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, n$$

$$\xi_i \geq 0 \quad \text{for } i = 1, 2, \ldots, n$$

**Equivalent "hinge loss" form**:

$$\min_{\boldsymbol{w}\in\mathbb{R}^d, \theta\in\mathbb{R}} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^{n}\mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}^{(i)}, y^{(i)})$$

*Notation*: $[a]_+ := \max\{0, a\}$.

The **hinge loss** of a linear classifier $f_{\boldsymbol{w},\theta}$ on an example $(\boldsymbol{x}, y)$ is defined to be

$$\mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}, y) := \Big[1 - y(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta)\Big]_+.$$

# Zero-one loss vs. hinge loss



$$\mathbb{1}\left\{y\left(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta\right) \leq 0\right\}$$

$$y\left(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta\right)$$

**Zero-one loss**: count if $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) \neq y$.

$$\mathbb{1}\{y(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta) \leq 0\} \ \leq \ \left[1 - y(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta)\right]_{+} \ = \ \mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}, y).$$

# ZERO-ONE LOSS VS. HINGE LOSS



Legend: $\mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}, y)$ ; $\mathbb{1}\{y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \le 0\}$

Axis label: $y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$

**Hinge loss**: an upper-bound on zero-one loss.

$$\mathbb{1}\{y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \le 0\} \ \le \ \left[1 - y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)\right]_+ \ = \ \mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}, y).$$
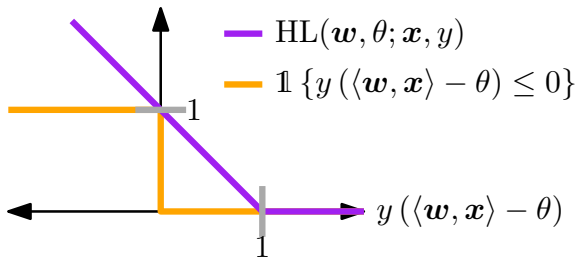
# ZERO-ONE LOSS VS. HINGE LOSS



**Hinge loss**: an upper-bound on zero-one loss.

$$\mathbb{1}\{y(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta) \leq 0\} \ \leq \ \left[1 - y(\langle \boldsymbol{w}, \boldsymbol{x}\rangle - \theta)\right]_{+} \ = \ \text{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}, y).$$

**Soft-margin SVM minimizes an upper-bound on the training error, plus a term that favors large margins.**
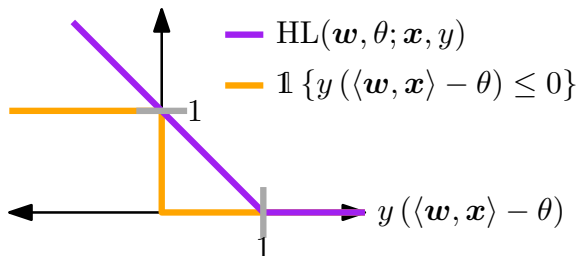
# Zero-one loss vs. hinge loss



Hinge loss: an upper-bound on zero-one loss.

$$\mathbb{1}\{y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta) \leq 0\} \ \leq \ \left[1 - y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)\right]_{+} \ = \ \mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}, y).$$

**Soft-margin SVM minimizes an upper-bound on the training error, plus a term that favors large margins.**

This is **computationally tractable** (unlike minimizing training error) because the hinge loss is a **convex function** of $(\boldsymbol{w}, \theta)$, and so is $\frac{\lambda}{2}\|\boldsymbol{w}\|_2^2$.

# GENERAL FORM

**Empirical risk minimization** (i.e., minimize training error):

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\left\{ y^{(i)}\left( \langle \boldsymbol{w}, \boldsymbol{x}^{(i)} \rangle - \theta \right) \leq 0 \right\}$$

**Soft-margin SVM**:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^{n} \mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}^{(i)}, y^{(i)})$$

# General form

**Empirical risk minimization** (i.e., minimize training error):

$$\min_{\boldsymbol{w}\in\mathbb{R}^d, \theta\in\mathbb{R}} \quad \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}\left\{y^{(i)}\left(\langle\boldsymbol{w}, \boldsymbol{x}^{(i)}\rangle - \theta\right) \leq 0\right\}$$

**Soft-margin SVM**:

$$\min_{\boldsymbol{w}\in\mathbb{R}^d, \theta\in\mathbb{R}} \quad \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2 + \frac{1}{n}\sum_{i=1}^{n}\mathrm{HL}(\boldsymbol{w}, \theta; \boldsymbol{x}^{(i)}, y^{(i)})$$

**Generic learning objective**:

$$\min_{\boldsymbol{w}\in\mathbb{R}^d, \theta\in\mathbb{R}} \quad R(\boldsymbol{w}, \theta) + \frac{1}{n}\sum_{i=1}^{n}\ell(\boldsymbol{w}, \theta; \boldsymbol{x}^{(i)}, y^{(i)})$$

- ▶ **Regularization**: encodes "learning bias" (e.g., preference for large margins), sometimes promotes stability.
- ▶ **Data fitting/empirical loss**: how poorly does the classifier "fit" the data.

# LEARNING VIA OPTIMIZATION

▶ Many different choices for regularization and loss.

  ▶ $R(\boldsymbol{w}) \propto \|\boldsymbol{w}\|_2^2$: encourage large margins
  ▶ $R(\boldsymbol{w}) \propto \|\boldsymbol{w}\|_1$: encourage $\boldsymbol{w}$ to be sparse
  ▶ $R(\boldsymbol{w}) \propto \sum_{i=1}^n w_i \ln w_i$: "maximum entropy" interpretation
  ▶ $\ell(\boldsymbol{w}, \theta; \boldsymbol{x}, y) = [\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta - y]_+^2$
  ▶ $\ell(\boldsymbol{w}, \theta; \boldsymbol{x}, y) = \log(1 + \exp(-y(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)))$
  ▶ . . .
  ▶ Also used beyond classification (e.g., regression, parameter estimation, clustering)

▶ Often want $\ell$ be an upper-bound on zero-one loss—i.e., a **surrogate loss**.

▶ Trade-off parameter $\lambda$: usually determine using hold out error or cross validation error.

▶ Computationally easier when overall objective function is **convex**: possible to efficiently find global minimizer in polynomial time.

▶ **Next**: techniques for analyzing and solving these optimization problems.