

# Support Vector Machines

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

# Announcement

- Problem set 3 released.
- Due March 1, 11:59pm. (Note the date!)

# Outline

- 1 Kernel methods (previous lecture)
- 2 Support vector machines – Geometric interpretation
- 3 SVM – Hinge loss

# Kernelized ridge regression

## Setup

- Input:  $\mathbf{x} \in \mathbb{R}^D$  (covariates, predictors, features)
- Transformation:  $\phi(\mathbf{x}) \in \mathbb{R}^M$ .
- Output:  $y \in \mathbb{R}$  (responses, targets, outcomes, outputs)
- Hypotheses/Model:  $h_{\mathbf{w}} : \mathbf{x} \rightarrow y$ , with  $h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$
- Training data:  $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

# Non-linear basis functions

$$\phi(\mathbf{x}) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \vdots \\ \sqrt{2}x_D \\ x_1^2 \\ x_1x_2 \\ \vdots \\ x_1x_D \\ x_2x_1 \\ x_2^2 \\ \vdots \\ x_2x_D \\ \vdots \\ x_Dx_1 \\ \vdots \\ x_D^2 \end{pmatrix}$$

# Two ways to do ridge regression

## Primal

Learn  $\mathbf{w}$  by minimizing the cost function  $J(\mathbf{w})$ :

$$\hat{\mathbf{w}} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y}$$

Predict  $y$  for a new input  $\mathbf{x}$

$$\hat{\mathbf{w}}^T \phi(\mathbf{x})$$

Here:

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M = O(D^2)} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$$

# Two ways to do ridge regression

## Primal

Learn  $\mathbf{w}$  by minimizing the cost function  $J(\mathbf{w})$ :

$$\hat{\mathbf{w}} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

Predict  $y$  for a new input  $\mathbf{x}$

$$\hat{\mathbf{w}}^T \phi(\mathbf{x})$$

- Cost of computing new features :  $O(ND^2)$ .
- Cost of learning (matrix multiplication and matrix inversion):  $O((D^2)^2N + (D^2)^3)$ .
- Cost of prediction  $O(D^2)$ .

# Two ways to do ridge regression

## Dual

Learn  $\alpha$  by minimizing the cost function  $J(\alpha)$ :

$$\hat{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Predict  $y$  for a new input  $x$

$$\mathbf{y}^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{k}_x$$

Here:

$$\mathbf{K} = \Phi \Phi^T$$

$$\mathbf{k}_x = \Phi \phi(x)$$

For the feature vector  $\phi$  that we chose, we can compute them as:

$$K_{m,n} = (1 + \mathbf{x}_m^T \mathbf{x}_n)^2$$

$$k_{x,n} = (1 + \mathbf{x}_n^T \mathbf{x})^2$$



# Two ways to do ridge regression

## Dual

Learn  $\alpha$  by minimizing the cost function  $J(\alpha)$ :

$$\hat{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Predict  $y$  for a new input  $x$

$$\mathbf{y}^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{k}_x$$

- Computing kernel matrix  $\mathbf{K}$ :  $O(N^2D)$ .
- Cost of learning:  $O(N^3)$ .
- Cost of prediction:  $O(N^2)$ .

# Kernel functions

- To use kernelized ridge regression, we need to be able to compute an inner product between a test point and any training point.
- Since we need this for any possible pair of points, we need a function that takes a pair of points and computes an inner product.
- This is the kernel function  $k(\cdot, \cdot)$ .
- Given two inputs,  $k(\cdot, \cdot)$  tells us how “similar” or “close” these inputs are in the space defined by the function  $\phi$ .

# Kernel functions

**Definition:** a kernel function  $k(\cdot, \cdot)$  is a bivariate function that satisfies the following properties. For any  $\mathbf{x}_m$  and  $\mathbf{x}_n$ ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function  $\phi(\cdot)$ .

# Kernel functions

**Mercer theorem** (loosely), a bivariate function  $k(\cdot, \cdot)$  is a kernel function, if and only if, for *any  $N$  and any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \text{and } \mathbf{x}_N$* , the matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is positive semidefinite.

# There are infinite numbers of kernels to use!

## Rules of composing kernels (this is just a partial list)

- if  $k(\mathbf{x}_m, \mathbf{x}_n)$  is a kernel, then  $ck(\mathbf{x}_m, \mathbf{x}_n)$  is also if  $c > 0$ .
- if both  $k_1(\mathbf{x}_m, \mathbf{x}_n)$  and  $k_2(\mathbf{x}_m, \mathbf{x}_n)$  are kernels, then  $\alpha k_1(\mathbf{x}_m, \mathbf{x}_n) + \beta k_2(\mathbf{x}_m, \mathbf{x}_n)$  are also if  $\alpha, \beta \geq 0$
- if both  $k_1(\mathbf{x}_m, \mathbf{x}_n)$  and  $k_2(\mathbf{x}_m, \mathbf{x}_n)$  are kernels, then  $k_1(\mathbf{x}_m, \mathbf{x}_n)k_2(\mathbf{x}_m, \mathbf{x}_n)$  are also.
- if  $k(\mathbf{x}_m, \mathbf{x}_n)$  is a kernel, then  $e^{k(\mathbf{x}_m, \mathbf{x}_n)}$  is also.
- ...

In practice, choosing an appropriate kernel is an “art”

People typically start with polynomial and Gaussian RBF kernels or incorporate domain knowledge.

# Mini-Summary

- Kernels allow us to design algorithms that use rich set of features while being computationally efficient.
- Many machine learning algorithms can be “kernelized”.
- Picking kernels is an art.
- We still need to tune hyperparameters.

# Outline

- 1 Kernel methods (previous lecture)
- 2 Support vector machines – Geometric interpretation
- 3 SVM – Hinge loss

# Support vector machines

- One of the most commonly used classification algorithms
- Good generalization in theory and practice
- Can be kernelized



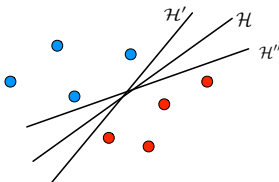
# Perceptron (linear classification)

## Setup

- Input:  $\mathbf{x} \in \mathbb{R}^D$
- Output:  $y \in \{-1, +1\}$
- Hypotheses/Model:  $h_{\mathbf{w},b} : \mathbf{x} \rightarrow y$ , with  $h_{\mathbf{w},b}(\mathbf{x}) = \text{SIGN}(b + \mathbf{w}^T \phi(\mathbf{x}))$   
 $\mathbf{w} = [w_1 \ w_2 \ \cdots \ w_M]^T$ : *weights, parameters, or parameter vector*  
 $b$  is called *bias*.
- Training data:  $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

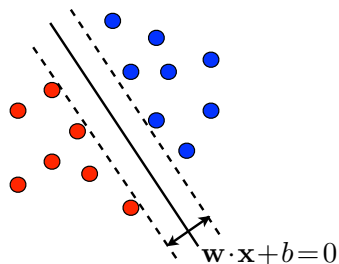
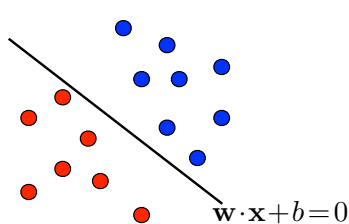
## Intuition: where to put the decision boundary?

Consider the following *separable* training dataset, i.e., we assume there exists a decision boundary that separates the two classes perfectly. There are an *infinite* number of decision boundaries  $\mathcal{H} : \mathbf{w}^T \phi(\mathbf{x}) + b = 0$ !

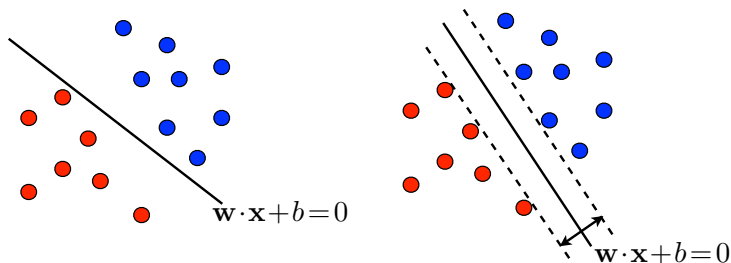


Which one should we pick?

# Intuition: where to put the decision boundary?



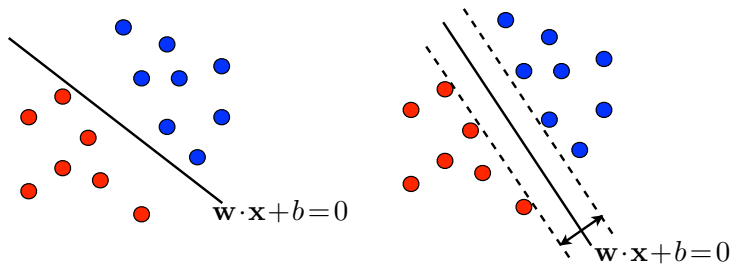
## Intuition: where to put the decision boundary?



Idea: Find a decision boundary in the '*middle*' of the two classes. In other words, we want a decision boundary that:

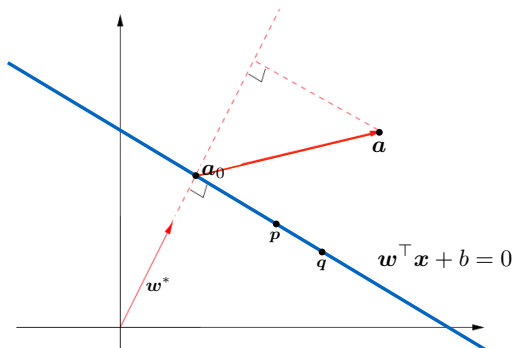
- Perfectly classifies the training data
- Is as far away from every training point as possible

# Review of hyperplane



- General equation is  $w^\top x + b = 0$
- $w \in \mathbb{R}^d$  is a non-zero normal vector,  $b$  is a scalar intercept
- Divides the space in half, i.e.,  $w^\top x + b > 0$  and  $w^\top x + b < 0$
- A hyperplane is a line in 2D and a plane in 3D

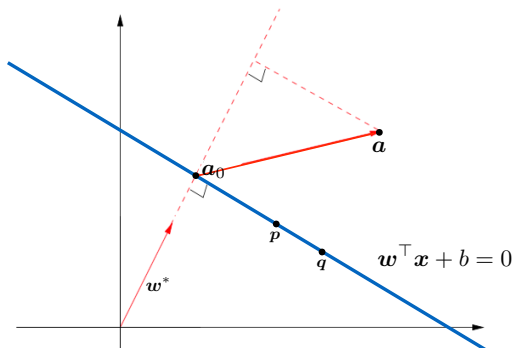
# Properties of hyperplanes



## Why is $w$ normal to this hyperplane?

- If two points,  $p$  and  $q$  are both on the hyperplane, then
$$w^\top(p - q) = w^\top p - w^\top q = b - b = 0$$

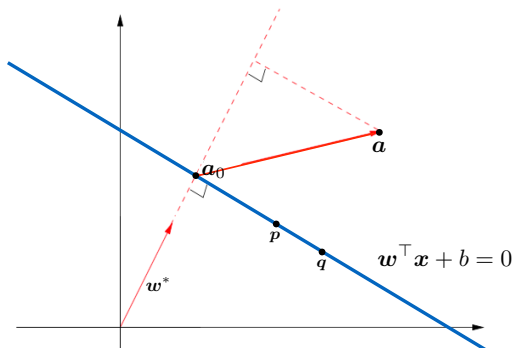
# Properties of hyperplanes



## Why is $w$ normal to this hyperplane?

- If two points,  $p$  and  $q$  are both on the hyperplane, then  $w^\top(p - q) = w^\top p - w^\top q = b - b = 0$
- $p - q$  is an arbitrary vector parallel to the hyperplane, thus  $w$  is orthogonal

# Properties of hyperplanes

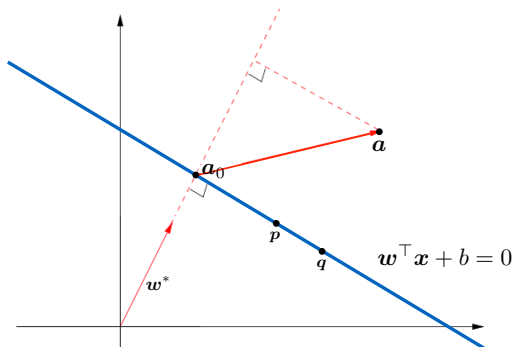


## Why is $w$ normal to this hyperplane?

- If two points,  $p$  and  $q$  are both on the hyperplane, then  $w^\top(p - q) = w^\top p - w^\top q = b - b = 0$
- $p - q$  is an arbitrary vector parallel to the hyperplane, thus  $w$  is orthogonal
- $w^* = \frac{w}{\|w\|}$  is the unit normal vector



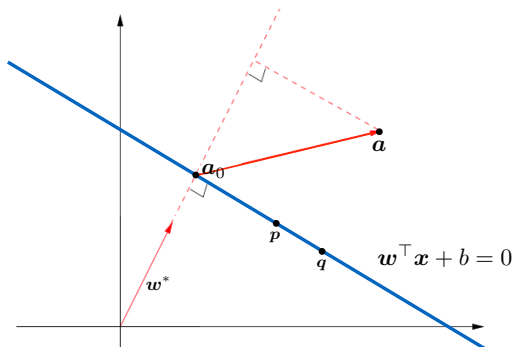
# Properties of hyperplanes



## How to compute signed distance from $a$ to the hyperplane?

- If we define point  $a_0$  on the hyperplane, then this distance corresponds to length of  $a - a_0$  in direction of  $w^*$ , which equals
$$w^{*\top}(a - a_0) = \frac{w}{\|w\|}^\top (a - a_0) = \frac{w^\top a - w^\top a_0}{\|w\|}$$

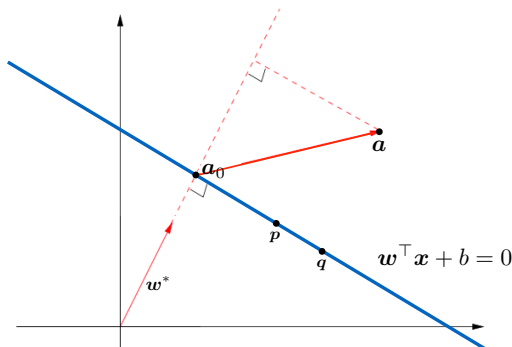
# Properties of hyperplanes



## How to compute signed distance from $a$ to the hyperplane?

- If we define point  $a_0$  on the hyperplane, then this distance corresponds to length of  $a - a_0$  in direction of  $w^*$ , which equals
$$w^{*\top}(a - a_0) = \frac{w}{\|w\|}^\top (a - a_0) = \frac{w^\top a - w^\top a_0}{\|w\|}$$
- Since  $w^\top a_0 = -b$ ,

# Properties of hyperplanes



## How to compute signed distance from $a$ to the hyperplane?

- If we define point  $a_0$  on the hyperplane, then this distance corresponds to length of  $a - a_0$  in direction of  $w^*$ , which equals  $w^{*\top}(a - a_0) = \frac{w}{\|w\|}^\top (a - a_0) = \frac{w^\top a - w^\top a_0}{\|w\|}$
- Since  $w^\top a_0 = -b$ , the distance equals  $\frac{1}{\|w\|}(w^\top a + b)$

# Distance from a point to decision boundary

The *unsigned* distance from a point  $\phi(\mathbf{x})$  to decision boundary (hyperplane)  $\mathcal{H}$  is

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{|\mathbf{w}^T \phi(\mathbf{x}) + b|}{\|\mathbf{w}\|_2}$$

# Distance from a point to decision boundary

The *unsigned* distance from a point  $\phi(\mathbf{x})$  to decision boundary (hyperplane)  $\mathcal{H}$  is

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{|\mathbf{w}^T \phi(\mathbf{x}) + b|}{\|\mathbf{w}\|_2}$$

We can remove the absolute value  $|\cdot|$  by exploiting the fact that the decision boundary classifies every point in the training dataset correctly.

# Distance from a point to decision boundary

The *unsigned* distance from a point  $\phi(\mathbf{x})$  to decision boundary (hyperplane)  $\mathcal{H}$  is

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{|\mathbf{w}^T \phi(\mathbf{x}) + b|}{\|\mathbf{w}\|_2}$$

We can remove the absolute value  $|\cdot|$  by exploiting the fact that the decision boundary classifies every point in the training dataset correctly.

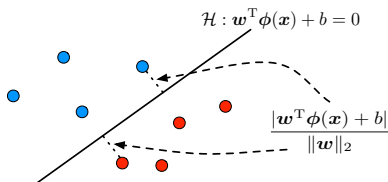
Namely,  $(\mathbf{w}^T \phi(\mathbf{x}) + b)$  and  $\mathbf{x}$ 's label  $y$  must have the same sign, so:

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{y[\mathbf{w}^T \phi(\mathbf{x}) + b]}{\|\mathbf{w}\|_2}$$

# Optimizing the Margin

**Margin** Smallest distance between the hyperplane and all training points

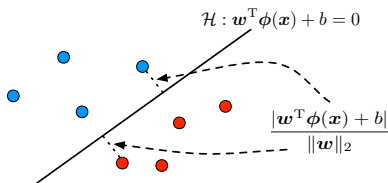
$$\text{MARGIN}(\mathbf{w}, b) = \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2}$$



# Optimizing the Margin

**Margin** Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(\mathbf{w}, b) = \min_n \frac{y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2}$$



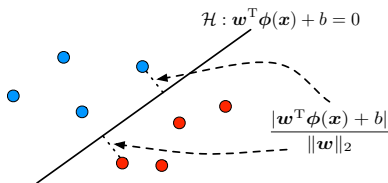
**How should we pick  $(\mathbf{w}, b)$  based on its margin?**



# Optimizing the Margin

**Margin** Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(\mathbf{w}, b) = \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2}$$



## How should we pick $(\mathbf{w}, b)$ based on its margin?

We want a decision boundary that is as far away from all training points as possible, so we want to *maximize* the margin!

$$\max_{\mathbf{w}, b} \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2} = \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \min_n y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]$$

# Scale of $w$

**Margin** Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(w, b) = \min_n \frac{y_n [w^T \phi(x_n) + b]}{\|w\|_2}$$

**Consider three hyperplanes**

- $(w, b)$
- $(2w, 2b)$
- $(.5w, .5b)$

**Which one has the largest margin?**

# Scale of $\mathbf{w}$

**Margin** Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(\mathbf{w}, b) = \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2}$$

## Consider three hyperplanes

- $(\mathbf{w}, b)$
- $(2\mathbf{w}, 2b)$
- $(.5\mathbf{w}, .5b)$

## Which one has the largest margin?

- The MARGIN doesn't change if we scale  $(\mathbf{w}, b)$  by a constant factor  $c$
- $\mathbf{w}^T \phi(\mathbf{x}) + b = 0$  and  $(c\mathbf{w})^T \phi(\mathbf{x}) + (cb) = 0$ : same decision boundary!
- Can we further constrain the problem?

## Rescaled Margin

We can further constrain the problem by scaling  $(\mathbf{w}, b)$  such that

$$\min_n y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1$$

## Rescaled Margin

We can further constrain the problem by scaling  $(\mathbf{w}, b)$  such that

$$\min_n y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1$$

We've fixed the numerator in the  $\text{MARGIN}(\mathbf{w}, b)$  equation, and we have:

$$\text{MARGIN}(\mathbf{w}, b) = \frac{1}{\|\mathbf{w}\|_2}$$

## Rescaled Margin

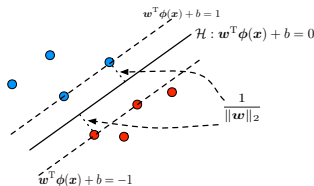
We can further constrain the problem by scaling  $(\mathbf{w}, b)$  such that

$$\min_n y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1$$

We've fixed the numerator in the  $\text{MARGIN}(\mathbf{w}, b)$  equation, and we have:

$$\text{MARGIN}(\mathbf{w}, b) = \frac{1}{\|\mathbf{w}\|_2}$$

Hence the points closest to the decision boundary are at distance  $\frac{1}{\|\mathbf{w}\|_2}$ !



# SVM: max margin formulation for separable data

Assuming separable training data, we thus want to solve:

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{such that} \quad \min_n y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1$$

This is equivalent to

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{such that} \quad y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad n \in \{1, \dots, N\}$$

# SVM: max margin formulation for separable data

Assuming separable training data, we thus want to solve:

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{such that} \quad \min_n y_n [\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b] = 1$$

This is equivalent to

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{such that} \quad y_n [\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b] \geq 1, \quad n \in \{1, \dots, N\}$$

This is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b] \geq 1, \quad n \in \{1, \dots, N\} \end{aligned}$$

Given our geometric intuition, SVM is called a *max margin* (or large margin) classifier. The constraints are called *large margin constraints*.



# SVM for separable data

## Constrained optimization problem

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2}_{\text{Objective function}} \\ \text{s.t.} \quad & \underbrace{y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b]}_{\text{Constraint}} \geq 1, \quad n \in \{1, \dots, N\} \end{aligned}$$

- Variables  $(\mathbf{w}, b)$
- Objective function
- Constraints

# Convex optimization problem (also called convex program)

$$\begin{aligned} \min_{\mathbf{x}} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & \mathbf{a}_i^T \mathbf{x} = b_i, \quad i = 1, \dots, p \end{aligned}$$

where  $\mathbf{x} \in \mathbb{R}^n$ ,  $f_0, \dots, f_m$  are convex functions.

- A point  $\mathbf{x}$  that satisfies all constraints is **feasible**.
- If there exists at least one feasible point, the problem is **feasible**. Otherwise, it is **infeasible**.

# SVM for separable data

## Constrained optimization problem

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & -y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] + 1 \leq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

- Objective function: convex in  $(\mathbf{w}, b)$
- Inequality constraints are linear in  $(\mathbf{w}, b)$ . Hence, convex.

# SVM for non-separable data

## SVM formulation for separable data

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n [\mathbf{w}^\top \phi(\mathbf{x}_n) + b] \geq 1, \quad n \in \{1, \dots, N\} \end{aligned}$$

**Non-separable setting** In practice our training data will not be separable. What issues arise with the optimization problem above when data is not separable?

# SVM for non-separable data

## SVM formulation for separable data

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad n \in \{1, \dots, N\} \end{aligned}$$

**Non-separable setting** In practice our training data will not be separable. What issues arise with the optimization problem above when data is not separable?

- For every  $\mathbf{w}$  there exists a training point  $\mathbf{x}_i$  such that

$$y_i [\mathbf{w}^T \phi(\mathbf{x}_i) + b] \leq 0$$

- There is no feasible  $(\mathbf{w}, b)$  as at least one of our constraints is violated!

# SVM for non-separable data

## Constraints in separable setting

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad n \in \{1, \dots, N\}$$

## Constraints in non-separable setting

Idea: modify our constraints to account for non-separability! Specifically, we introduce **slack variables**  $\xi_n \geq 0$ :

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\}$$

# SVM for non-separable data

## Constraints in separable setting

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad n \in \{1, \dots, N\}$$

## Constraints in non-separable setting

Idea: modify our constraints to account for non-separability! Specifically, we introduce **slack variables**  $\xi_n \geq 0$ :

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\}$$

- For “hard” training points, we can increase  $\xi_n$  until the above inequalities are met
- What does it mean when  $\xi_n$  is very large?

# Soft-margin SVM formulation

We do not want  $\xi_n$  to grow too large, and we can control their size by incorporating them into our optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\} \\ & \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$



# Soft-margin SVM formulation

We do not want  $\xi_n$  to grow too large, and we can control their size by incorporating them into our optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\} \\ & \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

What is the role of  $C$ ?

# Soft-margin SVM formulation

We do not want  $\xi_n$  to grow too large, and we can control their size by incorporating them into our optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\} \\ & \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

What is the role of  $C$ ?

- User-defined hyperparameter
- Trades off between the two terms in our objective
- Same idea as the regularization term in ridge regression, i.e.,  $C = \frac{1}{\lambda}$

# How to solve this problem?

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\} \\ & \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

- This is a *convex problem*: the objective function is quadratic in  $\mathbf{w}$  and linear in  $\xi$  and the constraints are linear (inequality) constraints in  $\mathbf{w}$ ,  $b$  and  $\xi_n$ .

# How to solve this problem?

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\} \\ & \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

- This is a *convex problem*: the objective function is quadratic in  $\mathbf{w}$  and linear in  $\xi$  and the constraints are linear (inequality) constraints in  $\mathbf{w}$ ,  $b$  and  $\xi_n$ .
- Given  $\phi(\cdot)$ , we can solve the optimization problem using general-purpose solvers.

# How to solve this problem?

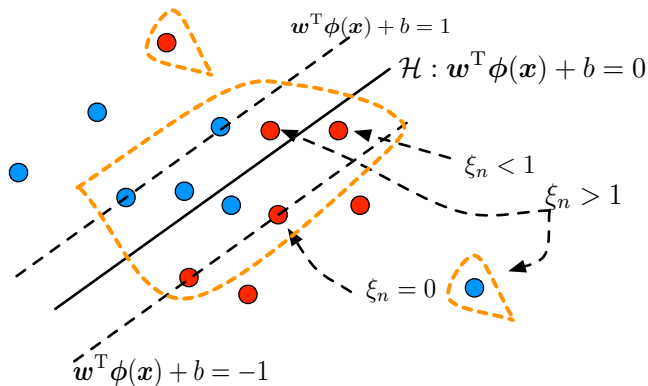
$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad n \in \{1, \dots, N\} \\ & \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

- This is a *convex problem*: the objective function is quadratic in  $\mathbf{w}$  and linear in  $\xi$  and the constraints are linear (inequality) constraints in  $\mathbf{w}$ ,  $b$  and  $\xi_n$ .
- Given  $\phi(\cdot)$ , we can solve the optimization problem using general-purpose solvers.
- There are several specialized methods for solving this problem, taking advantage of the special structure of the objective function and the constraints (we will not discuss them). Most existing SVM implementation/packages leverage these methods.

# Meaning of “support vectors” in SVMs

- The SVM solution is only determined by a subset of the training instances.
- These instances are called *support vectors*
- All other training points do not affect the optimal solution, i.e., if remove the other points and construct another SVM classifier on the reduced dataset, the optimal solution will be the same

# Visualization of how training data points are categorized



*Support vectors* are highlighted by the dotted orange lines

# Outline

- 1 Kernel methods (previous lecture)
- 2 Support vector machines – Geometric interpretation
- 3 SVM – Hinge loss



# A general view of supervised learning

**Definition** Assume  $y \in \{-1, 1\}$  and the decision rule is  $h(\mathbf{x}) = \text{SIGN}(a(\mathbf{x}))$  with  $a(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$

**For classification:** 0/1 loss

$$\ell^{0/1}(y, a(\mathbf{x})) = \begin{cases} 0 & \text{if } ya(\mathbf{x}) \geq 0 \\ 1 & \text{otherwise} \end{cases}$$

# Minimize weighted sum of empirical risk and regularizer

$$\begin{aligned} & \arg \min_{\mathbf{w}, b} R^{\text{EMP}}[h_{\mathbf{w}, b}(x)] + \lambda R(\mathbf{w}, b) \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \ell(y_n, a(\mathbf{x}_n)) + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$

# Minimize weighted sum of empirical risk and regularizer

$$\begin{aligned} & \arg \min_{\mathbf{w}, b} R^{\text{EMP}}[h_{\mathbf{w}, b}(x)] + \lambda R(\mathbf{w}, b) \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \ell(y_n, a(\mathbf{x}_n)) + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$

- Problem with minimizing the 0/1 loss ?

# Hinge loss

**Definition** Assume  $y \in \{-1, 1\}$  and the decision rule is  $h(\mathbf{x}) = \text{SIGN}(a(\mathbf{x}))$  with  $a(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ ,

$$\ell^{\text{HINGE}}(y, a(\mathbf{x})) = \begin{cases} 0 & \text{if } ya(\mathbf{x}) \geq 1 \\ 1 - ya(\mathbf{x}) & \text{otherwise} \end{cases}$$

## Intuition

# Hinge loss

**Definition** Assume  $y \in \{-1, 1\}$  and the decision rule is  $h(\mathbf{x}) = \text{SIGN}(a(\mathbf{x}))$  with  $a(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$ ,

$$\ell^{\text{HINGE}}(y, a(\mathbf{x})) = \begin{cases} 0 & \text{if } ya(\mathbf{x}) \geq 1 \\ 1 - ya(\mathbf{x}) & \text{otherwise} \end{cases}$$

## Intuition

- No penalty if raw output,  $a(\mathbf{x})$ , has same sign and is far enough from decision boundary (i.e., if 'margin' is large enough)
- Otherwise pay a growing penalty, between 0 and 1 if signs match, and greater than one otherwise

# Hinge loss

**Definition** Assume  $y \in \{-1, 1\}$  and the decision rule is  $h(\mathbf{x}) = \text{SIGN}(a(\mathbf{x}))$  with  $a(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$ ,

$$\ell^{\text{HINGE}}(y, a(\mathbf{x})) = \begin{cases} 0 & \text{if } ya(\mathbf{x}) \geq 1 \\ 1 - ya(\mathbf{x}) & \text{otherwise} \end{cases}$$

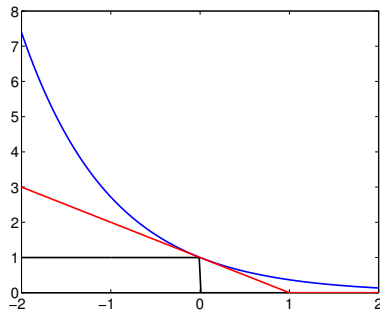
## Intuition

- No penalty if raw output,  $a(\mathbf{x})$ , has same sign and is far enough from decision boundary (i.e., if 'margin' is large enough)
- Otherwise pay a growing penalty, between 0 and 1 if signs match, and greater than one otherwise

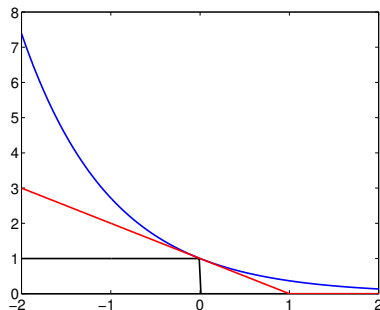
## Convenient shorthand

$$\ell^{\text{HINGE}}(y, a(\mathbf{x})) = \max(0, 1 - ya(\mathbf{x})) = (1 - ya(\mathbf{x}))_+$$

# Visualization and Properties



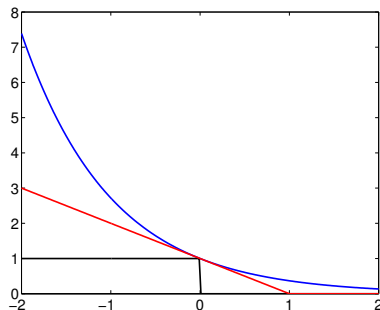
# Visualization and Properties



- Upper-bound for 0/1 loss function (black line)
- We use hinge loss as a *surrogate* to 0/1 loss – Why?

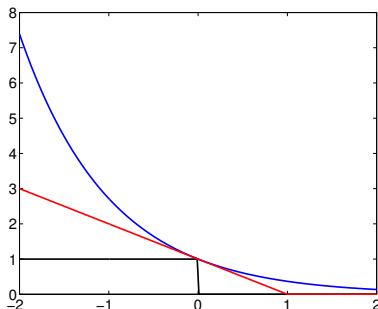


# Visualization and Properties



- Upper-bound for 0/1 loss function (black line)
- We use hinge loss as a *surrogate* to 0/1 loss – Why?
- Hinge loss is convex, and thus easier to work with.

# Visualization and Properties



- Other **surrogate losses** can be used, e.g., exponential loss for Adaboost (in blue), logistic loss (not shown) for logistic regression
- Hinge loss less sensitive to outliers than exponential (or logistic) loss

# Primal formulation of support vector machines (SVM)

**Minimizing the total hinge loss on all the training data with  $l_2$  regularization**

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Analogous to  $l_2$  regularized least squares, as we balance between two terms (the loss and the regularizer).

# Primal formulation of support vector machines (SVM)

## Minimizing the total hinge loss on all the training data with $l_2$ regularization

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Analogous to  $l_2$  regularized least squares, as we balance between two terms (the loss and the regularizer).

Previously, we used geometric arguments to derive:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n \quad \text{and} \quad \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

*Do these yield the same solution?*

# Recovering our previous SVM formulation

## Minimizing the total hinge loss on all the training data

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Define  $C = 1/\lambda$ :

$$\min_{\mathbf{w}, b} C \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

# Recovering our previous SVM formulation

## Minimizing the total hinge loss on all the training data

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Define  $C = 1/\lambda$ :

$$\min_{\mathbf{w}, b} C \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

Define  $\xi_n = \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b])$

$$\min_{\mathbf{w}, b, \xi} C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t.} \quad \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) = \xi_n, \quad n \in \{1, \dots, N\}$$

# Recovering our previous SVM formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) = \xi_n, \quad n \in \{1, \dots, N\} \end{aligned}$$

is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \leq \xi_n, \quad n \in \{1, \dots, N\} \\ & 0 \leq \xi_n, \quad n \in \{1, \dots, N\} \end{aligned}$$

At optimal solution constraints are active so we have equality! Why?



At optimal solution constraints are active so we have equality! Why?

- If  $\xi_n^* > \max(0, 1 - y_n f(\mathbf{x}_n))$ , we could choose  $\bar{\xi}_n < \xi_n^*$  and still satisfy the constraint while reducing our objective function!
- Since  $c \geq \max(a, b) \iff c \geq a, c \geq b$ , we recover previous formulation