# CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Feb 24, 2021

## Outline

1. Logistics

2. Review of last lecture

3. Support vector machines (primal formulation)

4. Quiz 1 Specifics

## Outline

1. **Logistics**

2. Review of last lecture

3. Support vector machines (primal formulation)

4. Quiz 1 Specifics

## Logistics

- HW 3 was assigned.
- We will discuss quiz specifics at the end of the lecture today.

## Outline

## Kernel functions

**Definition**: a function $k : \mathbb{R}^D \times \mathbb{R}^D \to \mathbb{R}$ is called a *(positive semidefinite) kernel function* if there exists a function $\phi : \mathbb{R}^D \to \mathbb{R}^M$ so that for any $\boldsymbol{x}, \boldsymbol{x}' \in \mathbb{R}^D$,

$$k(\boldsymbol{x}, \boldsymbol{x}') = \phi(\boldsymbol{x})^\mathrm{T} \phi(\boldsymbol{x}')$$

Examples we have seen

$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^\mathrm{T} \boldsymbol{x}')^2$$

$$k(\boldsymbol{x}, \boldsymbol{x}') = \sum_{d=1}^D \frac{\sin(2\pi(x_d - x_d'))}{x_d - x_d'}$$

$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^\mathrm{T} \boldsymbol{x}' + c)^d \qquad \textbf{(polynomial kernel)}$$

$$k(\boldsymbol{x}, \boldsymbol{x}') = e^{-\frac{\|\boldsymbol{x}-\boldsymbol{x}'\|_2^2}{2\sigma^2}} \qquad \textbf{(Gaussian/RBF kernel)}$$

## Kernelizing ML algorithms

Feasible as long as **only inner products are required**:

- regularized linear regression (dual formulation)

$$\phi(\boldsymbol{x})^\mathrm{T} \boldsymbol{w}^* = \phi(\boldsymbol{x})^\mathrm{T} \boldsymbol{\Phi}^\mathrm{T} (\boldsymbol{K} + \lambda \boldsymbol{I})^{-1} \boldsymbol{y} \quad (\boldsymbol{K} = \boldsymbol{\Phi}\boldsymbol{\Phi}^\mathrm{T} \text{ is } \textit{kernel matrix})$$

- nearest neighbor classifier with L2 distance

$$\|\phi(\boldsymbol{x}) - \phi(\boldsymbol{x}')\|_2^2 = k(\boldsymbol{x}, \boldsymbol{x}) + k(\boldsymbol{x}', \boldsymbol{x}') - 2k(\boldsymbol{x}, \boldsymbol{x}')$$

- perceptron, logistic regression, SVM, . . .

## Outline

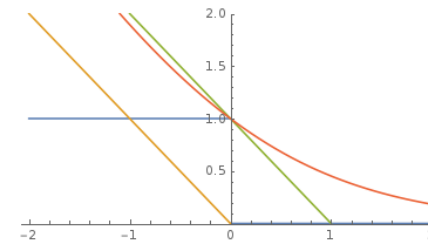## Support vector machines (SVM)

- One of the most commonly used classification algorithms

- Works well with the kernel trick

- Strong theoretical guarantees

We focus on **binary classification** here.

## Primal formulation

In one sentence: linear model with L2 regularized hinge loss. Recall



- perceptron loss $\ell_{\text{perceptron}}(z) = \max\{0, -z\} \rightarrow$ Perceptron
- logistic loss $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z)) \rightarrow$ logistic regression
- hinge loss $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\} \rightarrow$ **SVM**

## Primal formulation

For a linear model $(\boldsymbol{w}, b)$, this means

$$\min_{\boldsymbol{w}, b} \sum_n \max\left\{0, 1 - y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b)\right\} + \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2$$
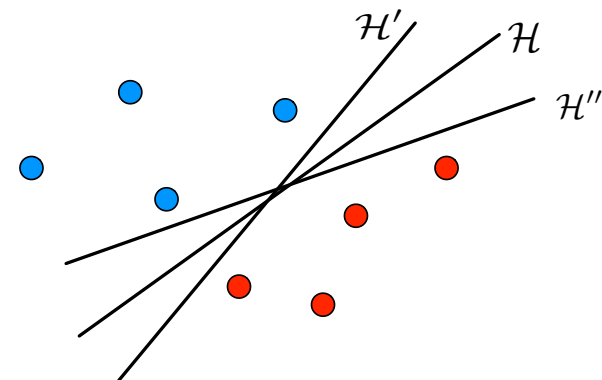
- recall $y_n \in \{-1, +1\}$

- a nonlinear mapping $\phi$ is applied

- the bias/intercept term $b$ is used explicitly (think about why after this lecture)

*So why L2 regularized hinge loss?*

## Geometric motivation: separable case

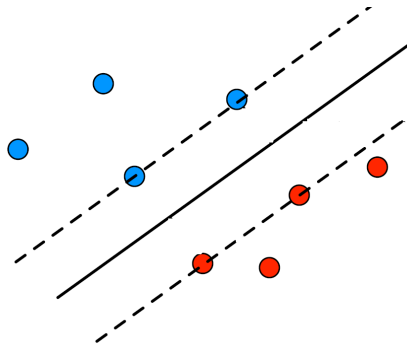When data is **linearly separable**, there are *infinitely many hyperplanes with zero training error*:



So which one should we choose?

## Intuition

The further away from data points the better.



*How to formalize this intuition?*

## Distance to hyperplane

What is the **distance** from a point $x$ to a hyperplane $\{x : w^{\mathrm{T}}x + b = 0\}$?

Assume the **projection** is $x - \ell\frac{w}{\|w\|_2}$, then

$$0 = w^{\mathrm{T}}\left(x - \ell\frac{w}{\|w\|_2}\right) + b = w^{\mathrm{T}}x - \ell\|w\| + b$$

and thus $\ell = \frac{w^{\mathrm{T}}x+b}{\|w\|_2}$.

Therefore the distance is

$$\frac{|w^{\mathrm{T}}x + b|}{\|w\|_2}$$

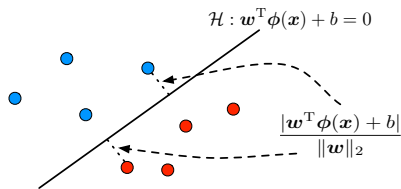For a hyperplane that correctly classifies $(x, y)$, the distance becomes

$$\frac{y(w^{\mathrm{T}}x + b)}{\|w\|_2}$$

## Maximizing margin

**Margin**: the *smallest* distance from all training points to the hyperplane

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



The intuition "**the further away the better**" translates to solving

$$\max_{w,b} \; \min_n \frac{y_n(w^{\mathrm{T}}\phi(x_n) + b)}{\|w\|_2} = \max_{w,b} \frac{1}{\|w\|_2} \min_n y_n(w^{\mathrm{T}}\phi(x_n) + b)$$
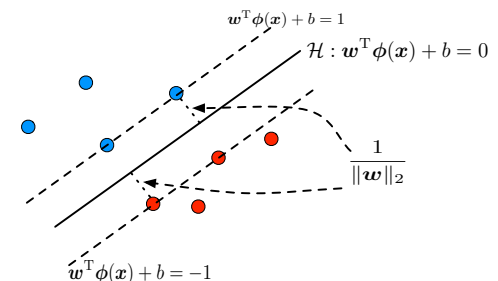
## Rescaling

**Note**: rescaling $(w, b)$ does not change the hyperplane at all.

We can thus always scale $(w, b)$ s.t. $\min_n y_n(w^{\mathrm{T}}\phi(x_n) + b) = 1$

The margin then becomes

$$\text{MARGIN OF } (w, \, b)$$
$$= \frac{1}{\|w\|_2} \min_n y_n(w^{\mathrm{T}}\phi(x_n) + b)$$
$$= \frac{1}{\|w\|_2}$$

## Summary for separable data

For a separable training set, we aim to solve

$$\max_{\boldsymbol{w},b} \frac{1}{\|\boldsymbol{w}\|_2} \quad \text{s.t.} \quad \min_{n} y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b) = 1$$

This is equivalent to

$$\min_{\boldsymbol{w},b} \quad \frac{1}{2}\|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b) \geq 1, \quad \forall\, n$$

SVM is thus also called *max-margin* classifier. The constraints above are called *hard-margin* constraints.

## General non-separable case

If data is not linearly separable, the previous constraint

$$y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b) \geq 1, \quad \forall\, n$$

is obviously *not feasible*.

To deal with this issue, we relax them to **soft-margin** constraints:

$$y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b) \geq 1 - \xi_n, \quad \forall\, n$$

where we introduce **slack variables** $\xi_n \geq 0$.

## SVM Primal formulation

We want $\xi_n$ to be as small as possible too. The objective becomes

$$\min_{\boldsymbol{w},b,\{\xi_n\}} \quad \frac{1}{2}\|\boldsymbol{w}\|_2^2 + C\sum_{n} \xi_n$$
$$\text{s.t.} \quad y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b) \geq 1 - \xi_n, \quad \forall\, n$$
$$\xi_n \geq 0, \quad \forall\, n$$

where $C$ is a hyperparameter to balance the two goals.

## Equivalent form

**Formulation**

$$\min_{\boldsymbol{w},b,\{\xi_n\}} \quad C\sum_{n} \xi_n + \frac{1}{2}\|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad 1 - y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b) \leq \xi_n, \quad \forall\, n$$
$$\xi_n \geq 0, \quad \forall\, n$$

**is equivalent to**

$$\min_{\boldsymbol{w},b,\{\xi_n\}} \quad C\sum_{n} \xi_n + \frac{1}{2}\|\boldsymbol{w}\|_2^2$$
$$\text{s.t.} \quad \max\left\{0, 1 - y_n(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n) + b)\right\} = \xi_n, \quad \forall\, n$$

## Equivalent form

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

$$\text{s.t.} \quad \max\left\{0, 1 - y_n(\boldsymbol{w}^\mathrm{T}\boldsymbol{\phi}(\boldsymbol{x}_n) + b)\right\} = \xi_n, \quad \forall\, n$$

**is equivalent to**

$$\min_{\boldsymbol{w},b} \, C\sum_n \max\left\{0, 1 - y_n(\boldsymbol{w}^\mathrm{T}\boldsymbol{\phi}(\boldsymbol{x}_n) + b)\right\} + \frac{1}{2}\|\boldsymbol{w}\|_2^2$$

and

$$\min_{\boldsymbol{w},b} \, \sum_n \max\left\{0, 1 - y_n(\boldsymbol{w}^\mathrm{T}\boldsymbol{\phi}(\boldsymbol{x}_n) + b)\right\} + \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2$$

with $\lambda = 1/C$. *This is exactly minimizing L2 regularized hinge loss!*

## Optimization

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

$$\text{s.t.} \quad 1 - y_n(\boldsymbol{w}^\mathrm{T}\boldsymbol{\phi}(\boldsymbol{x}_n) + b) \le \xi_n, \quad \forall\, n$$

$$\xi_n \ge 0, \quad \forall\, n$$

- It is a convex (**quadratic** in fact) problem

- thus can apply any convex optimization algorithms, e.g. SGD

- there are **more specialized and efficient** algorithms

- but usually we apply kernel trick, which requires solving the *dual problem*

## Outline

## Logistics

- Quiz 1 is scheduled for March 3, 2021 from 10:00 – 12:00 PM. It is an in-class, open book and notes exam (no other resources are allowed).
- We will be using CrowdMark and WebEx to administer the exam.
- CrowdMark link: `https://app.crowdmark.com/sign-in/usc`
- We'll be releasing some questions (and solutions) for the topics covered in HW3 on Friday using CrowdMark. *Make sure you get familiar with the platform*.
- **Topics:** All topics covered till the next lecture.

## On Quiz day

- Join ∼15 min prior to the class time.
- We'll assign the exam 5 minutes before 10:00 AM on CrowdMark.
- You will have 10:00 – 11:45 AM for the exam, and the last 15 minutes are for you to upload your solutions.
- You will upload the pictures for each question separately.
- Join via the WebEx link on DEN@USC, *required to have video ON*.
- We'll be recording the video via WebEx.
- You may ask your questions privately to the teaching staff using WebEx chat, cannot communicate with fellow students in any way.