

# Support Vector Machines

#### Instructor: Jessica Wu -- Harvey Mudd College

The instructor gratefully acknowledges Andrew Ng (Stanford), Eric Eaton (UPenn), David Sontag (NYU), Piyush Rai (Utah) and the many others who made their course materials freely available online.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

## **SVM Basics**

Learning Goals

- Describe the goal of SVMs
  - Start from perceptron and build graphical intuition
- Setup formal SVM optimization problem

## Recall the Perceptron

Assume there is a linear classifier Start with simple learner and analyze with it does

Let 
$$y \in \{-1, +1\}$$
 use explicit bias term rather than folding into features  $x_{theta}$   
 $h_{\theta}(x) = g(\theta^T x + b)$ 

where

$$g(z) = \operatorname{sgn}(z) \begin{cases} +1, \ z \ge 0\\ -1, \ z < 0 \end{cases}$$



So if  $\theta^T x + b \ge 0$ , then y = +1if  $\theta^T x + b < 0$ , then y = -1





# Summary

- Perceptron finds one of many possible hyperplanes separating data (if one exists)
- Of possible choices, find one with maximum margin ⇒ Support Vector Machines (SVMs)
- Why?
  - Good according to intuition, theory, and practice
- Some history
  - SVM became famous when, in handwriting recognition task using images as input, it gave accuracy comparable to neural network with hand-designed features

## **Review: Margins**

For training set  $\left\{\left(m{x}^{(i)}, y^{(i)}
ight)
ight\}_{i=1}^n$  and parameter vector  $m{ heta}$ :

• functional margin of *i*<sup>th</sup> example

$$\hat{\gamma}^{(i)} = y^{(i)} \left( oldsymbol{ heta}^T oldsymbol{x} + b 
ight)$$

- positive if  $\hat{ heta}$  classifies  $x^{(i)}$  correctly
- absolute value = "confidence" in predicted label (or "mis-confidence")

$$\gamma^{(i)} = rac{\hat{\gamma}^{(i)}}{\|oldsymbol{ heta}\|} = rac{y^{(i)} \left(oldsymbol{ heta}^T oldsymbol{x} + b
ight)}{\|oldsymbol{ heta}\|}$$

— signed distance of  $m{x}^{(i)}$  to hyperplane (positive if classified correctly)

#### • margin over training set

 $\begin{array}{ll} - \mbox{ minimum functional margin } & \hat{\gamma} = \min_{i} \hat{\gamma}^{(i)} \\ - \mbox{ minimum geometric margin } & \gamma = \min_{i} \gamma^{(i)} \end{array}$ 



# **SVM Exercise**

Consider the following training data:

$x_1$	$x_2$	label
1	1	+
2	2	+
2	0	+
0	0	_
1	0	—
0	1	_

- Plot these points. Are the classes  $\{+, -\}$  linearly separable?
- Plot the maximum-margin hyperplane by inspection. What is the associated weight vector  $\theta$  of this hyperplane? Identify the support vectors.
- If you remove one of the support vectors, does the size of the optimal margin decrease, stay the same, or increase?
- *Extra*: Is your answer above true for any data set? Provide a counterexample or give a short proof.

Example adapted from Russell Norvig



(This slide intentionally left blank.)



# **Optimization Problem**

 $\min_{\boldsymbol{\theta}, b} \frac{1}{2} \|\boldsymbol{\theta}\|^2 \qquad \qquad \max \operatorname{inmum margin} = \max \frac{1}{||\boldsymbol{\theta}||} \\ \Rightarrow \min_{\boldsymbol{\theta}, b} \frac{||\boldsymbol{\theta}||^2}{\Rightarrow \min_{\boldsymbol{\theta}, b} \frac{1}{2} ||\boldsymbol{\theta}||^2} \\ \text{s.t. } y^{(i)} \left(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + b\right) \ge 1 \text{ for } i = 1, \dots, n$ 

Note

• convex quadratic objective with linear constraints (*n* of them)

⇒ Quadratic program (QP) Polynomial-time algorithms exist for solving QPs

#### **SVM**s

Learning Goals

- Solving the SVM optimization problem
  - using Lagrange multipliers (leads to dual problem)
- Allowing misclassified examples
  - using slack variables (leads to soft-margin SVM)
- Describe the SVM loss function
- Allowing non-linear decision boundaries
  - using kernels



# Math Details

For math to work,

Karush-Kuhn-Tucker (KKT) conditions must hold:

for optimal  $\theta$ , b,  $\alpha$ ,  $\alpha_i \left[ 1 - y^{(i)} \left( \theta^T x^{(i)} + b \right) \right] = 0$   $\Rightarrow$  if  $\alpha_i > 0$ , then  $y^{(i)}(\theta^T x^{(i)} + b) = 1$ i.e. if weight for example i > 0, then  $x^{(i)}$  lies on a margin boundary ( $x^{(i)}$  is a SV)  $\Rightarrow$  if  $y^{(i)}(\theta^T x^{(i)} + b) > 1$ , then  $\alpha_i = 0$ i.e. if  $x^{(i)}$  does not lie a margin boundary ( $x^{(i)}$  is not a SV) then weight for example i = 0

# Solving SVM Problem

Solve dual in lieu of primal

- Solve for  $\alpha$  directly
- How?

- coordinate descent, sequential minimal optimization (SMO)

Compute primal solution from dual

• to obtain  $\theta$ ,

$$oldsymbol{ heta} = \sum_{i=1}^n lpha_i y^{(i)} oldsymbol{x}^{(i)}$$

- can sum over SVs since only these examples have  $\alpha_i > 0$
- remember, we have not added  $x_0 = 1$
- to obtain b, find a SV  $\boldsymbol{x}^{(i)}$ , solve

$$y^{(i)}\left(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + b\right) = 1$$

$$y^{(i)}y^{(i)}\left(\boldsymbol{\theta}^{T}\boldsymbol{x}^{(i)}+b\right) = y^{(i)}$$
$$b = y^{(i)} - \boldsymbol{\theta}^{T}\boldsymbol{x}^{(i)}$$

# Sparsity and Generalization

 $lpha_i > 0$  only for SVs # SVs <<< # training examples

 $\Rightarrow$  sparseness leads to better generalization Why?

What is maximum LOOCV error?

(This slide intentionally left blank.)

(This slide intentionally left blank.)

## **SVMs**

#### Learning Goals

- ✓ Solving the SVM optimization problem
  - using Lagrange multipliers (leads to dual problem)
- Allowing misclassified examples
  - using slack variables (leads to soft-margin SVM)
- Describe the SVM loss function
- Allowing non-linear decision boundaries
  - using kernels

Ideas?









- margin boundaries still at  ${\pmb heta}^T {\pmb x} + b = \pm 1$  large  $C[\min \Sigma_i \, \xi_i]$
- training example ... within margin region when  $0 < \xi_i < 1$ misclassified when  $\xi_i > 1$

- examples can have (functional) margin < 1
- if  $\xi_i > 0$ , pay penalty  $C\xi_i$
- C is slack penalty [C > 0]trade-off between...
  - minimizing  $||\boldsymbol{\theta}||^2$  (maximizing margin)
  - ensuring most examples have functional margin of at least 1
- small  $C[\min \frac{1}{2}||\boldsymbol{\theta}||^2]$



# $\begin{aligned} & \max_{\alpha} \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)} \rangle \\ & \text{s.t. } 0 \leq \alpha_i \leq C \quad \text{for } i = 1, \dots, n \\ & \max_{n \text{ ew constraint: add upper bound } \alpha_i \leq C \\ & \sum_{i=1}^{n} \alpha_i y^{(i)} = 0 \end{aligned}$

#### intuition

• w/o slack

 $\alpha_i$  can increase without bound to prevent misclassification

• w/ slack

upper bound of *C* limits  $\alpha_i$  to allow misclassifications





Deriving the Loss Function  
Recall soft-margin SVM optimization problem  

$$\begin{split} & \underset{\theta,b,\xi}{\min} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ & \underset{\theta,b,\xi}{\inf} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ & \underset{\xi_i \ge 0}{\inf} \xi_i \\ & \underset{\xi_i \ge 0}{\inf} \frac{1}{2} - \xi_i \quad \text{for } i = 1, \dots, n \\ & \underset{\xi_i \ge 0}{\inf} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_i \\ & \underset{\xi_i \ge 0}{ \underbrace} \frac{1}{2} - \xi_$$



## Soft-Margin SVM as Regularization

Let note:  $L\left(y,\hat{y}\right) = \max\left(0,1-y\hat{y}\right)$  $\hat{y} = \boldsymbol{\theta}^T \boldsymbol{x} + b$  is "raw" output, not predicted class So  $\min_{\boldsymbol{\theta},b} \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{i=1}^n L\left(y^{(i)}, \hat{y}^{(i)}\right)$ empirical risk regularization using hinge loss Standard form of regularization So  $C \propto \frac{1}{\lambda}$  $\min_{\boldsymbol{\theta},b} \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 + \frac{1}{n} \sum_{i=1}^{n} L\left(y^{(i)}, \hat{y}^{(i)}\right)$ low regularization small  $\lambda$ , large C high regularization large  $\lambda$ , small C