# **Support Vector Machines**

### Machine Learning – Foundations and Algorithms WS 2021/22

Prof. Gerhard Neumann KIT, Institut für Anthrophomatik und Robotik

# Learning Outcomes

#### What will we learn today?

- Understand the concept of Maximum Margin classifiers
- Define the corresponding optimization problem
- How do relax the problem using slack variables
- Connection to the hinge loss
- How do optimize the problem using sub-gradients
- How to do the Kernel Trick with SVMs

# Agenda for Today

#### **Recap: Linear Discriminators**

#### **Support Vector Machines:**

- Maximum Margin
- Optimization Problem
- Soft-Margin
- Hinge-Loss

#### **SVMs with Kernels**

- Constraint Optimization Problem
- Kernel Trick

#### **Basics:**

- Constraint Optimization
- Sub-gradients

### **Binary Classification: Previous Definition**



Given the training data  $(x_i, y_i)$ , i = 1...N, with  $x_i \in \mathbb{R}^d$  and  $y_i \in \{0, 1\}$ , learn a classifier f(x) such that:

$$f(\boldsymbol{x}_i) = \begin{cases} > 0, & \text{if } y_i = 1 \\ < 0, & \text{if } y_i = 0 \end{cases}$$

#### For SVMs, it simplifies notation to use +1 and -1 as class labels

**New definition:** Given the training data  $(x_i, y_i)$ , i = 1...N, with  $x_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$  learn a classifier f(x) such that:

$$f(\boldsymbol{x}_i) = \left\{ \begin{array}{ll} > 0, & \text{if } y_i = 1 \\ < 0, & \text{if } y_i = -1 \end{array} \right.$$

**Or:**  $f(\boldsymbol{x}_i)y_i > 0$  for a correct classification

## **Recap: Linear Classifiers**

A linear classifier is given in the form:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

#### In 2D, the classifier is a line

- w is the normal to the line
- *b* is the bias



# **Basics: Projections of vectors**

The scalar product of 2 vectors can be used to compute the **projection of vector a on vector b**:

- Geometric definition of scalar product  $a^T b = ||a|| ||b|| \cos \theta$
- Angle between 2 vectors

 $\cos \theta = \frac{\boldsymbol{a}^T \boldsymbol{b}}{\|\boldsymbol{a}\| \|\boldsymbol{b}\|}$ 

• Scalar projection of **a** on **b** 

$$a_b = \| oldsymbol{a} \| \cos heta = rac{oldsymbol{a}^T oldsymbol{b}}{\|oldsymbol{b}\|}$$



Observation: If 2 vectors are normal to each other, the projection is 0

# **Recap: Geometrical inspection**

#### **Observations**:

- The decision boundary is normal to w
- Without b, it goes through the origin



# **Recap: Geometrical inspection**



# **Optimal Separation**

Which is the optimal line?



# Maximum Margin

- **Support Vectors:** Data points closes to the decision boundary
  - Other examples can be ignored
- **Margin**  $\rho$  is the distance between the support vectors and the decision boundary
- Margin should be maximized
  - I.e. minimum distance between decision boundary and examples should be maximized



# Maximum Margin

- Maximize distance between hyper-plane and "difficult examples"
  - Examples next to decision boundary
  - Also called Support Vectors
- Intuition:
  - Less examples close to decision boundary
     -> more robust
- Statistical Learning Theory:
  - Maximum Margin Classifier has smaller complexity (VC-dimension)
  - And therefore generalizes better



### **Geometric Inspection**



Gerhard Neumann | Machine Learning | KIT | WS 2021/2022

## **Mathematical Formulation**

#### **Observation:**

- $\mathbf{w}^T \mathbf{x} + b = 0$  and  $c(\mathbf{w}^T \mathbf{x} + b) = 0$ define the same hyper-plane
- Scaling c can be chosen freely

#### Choose scaling such that

- For positive support vectors  $\mathbf{w}^T \mathbf{x}_+ + b = +1$
- For negative support vectors  $\mathbf{w}^T \mathbf{x}_- + b = -1$

#### Margin is then given by

$$\frac{\mathbf{w}^T \mathbf{x}_+ + b}{||\mathbf{w}||} - \frac{\mathbf{w}^T \mathbf{x}_- + b}{||\mathbf{w}||} = \frac{2}{||\mathbf{w}||}$$



# **SVM Optimization**

### **Optimization problem:**

$$\begin{array}{ll} \operatorname{argmax}_{\mathbf{w}} & \frac{2}{||\mathbf{w}||}, \\ \text{s.t.} & \mathbf{w}^{T} \mathbf{x}_{i} + b \left\{ \begin{array}{ll} \geq +1, & \text{if } y_{i} = +1 \\ \leq -1, & \text{if } y_{i} = -1 \end{array} \right. \end{array} \end{array}$$
 Maximize margin Condition for margin

#### **Observations:**

• If the constraints are not satisfied, our definition of the margin would be wrong, i.e,

$$\min_{\mathbf{x}_{+} \in \mathbf{X}_{+}} (\mathbf{w}^{T} \mathbf{x}_{+} + b) = +1 \qquad \max_{\mathbf{x}_{-} \in \mathbf{X}_{-}} (\mathbf{w}^{T} \mathbf{x}_{-} + b) = -1 \qquad \mathbf{X}_{-} : \text{negative examples}$$
Positive support vectors
Negative support vectors
$$\mathbf{X}_{+} : \text{positive examples}$$

- Support vectors have the smallest distance to the decision boundary
- There is at least one positive and one negative data point that satisfy the support vector condition exactly (i.e. equality instead of inequality) from above
  - Why? Because of the argmax! Norm of weight vector could be reduced otherwise

## **SVM Optimization**

#### **Optimization problem:**

$$\begin{array}{ll} \operatorname{argmax}_{\mathbf{w}} & \frac{2}{||\mathbf{w}||}, \\ \text{s.t.} & \mathbf{w}^{T} \mathbf{x}_{i} + b \left\{ \begin{array}{ll} \geq +1, & \text{if } y_{i} = +1 \\ \leq -1, & \text{if } y_{i} = -1 \end{array} \right. \end{array} \end{array}$$
 Maximize margin Condition for margin

Reformulation: Easier to solve, same solution

$$\begin{aligned} \operatorname{argmin}_{\mathbf{w}} & ||\mathbf{w}||^2, \\ \text{s.t.} & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \end{aligned}$$

- Quadratic optimization problem Linear Constraints
- Convex, single optimum We will see later how to solve it...

# Back to linear Separability



### What is the best w?

• Linear separable but: small margin

Large margin but error in classification

We have to choose a trade-off between margin and classification accuracy!

Introduce slack-variables:

 $\xi_i \ge 0$ 

Allows violating the margin conditions

 $y_i(\mathbf{w}^T\mathbf{x}_i+b) \ge 1-\boldsymbol{\xi}_i$ 

- $0 \le \xi_i \le 1$  sample is between margin and decision boundary: margin violation
- $\xi_i > 1$  sample is on the wrong side of the decision boundary: **misclassified**



### Soft Max-Margin

#### **Optimization problem:**

$$\operatorname{argmin}_{\mathbf{w},\boldsymbol{\xi}} \quad ||\mathbf{w}||^2 + C \sum_{i}^{N} \xi_i,$$
  
s.t.  $y_i(\mathbf{w}^T \mathbf{x}_i + b) \ge 1 - \xi_i, \quad \xi_i \ge 0$ 

Punish large slack variables

Condition for soft-margin

#### • C is a (inverse) regularization parameter

- Small C: Constraints have little influence -> large margin -> large regularization
- Large C: Constraints have large influence -> small margin -> small regularization
- C infinite: Constraints are enforced-> hard margin -> no regularization

### Illustration



### Reformulation into an unconstrained problem

#### **Constrained optimization:**

$$\operatorname{argmin}_{\mathbf{w},\boldsymbol{\xi}} ||\mathbf{w}||^2 + C \sum_{i}^{N} \xi_i, \quad \text{s.t.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \ge 1 - \xi_i, \quad \xi_i \ge 0$$

$$Regularization \text{ parameter}$$

SVMs can be reformulated into an unconstrained optimization problem

• Rewrite constraints: 
$$\xi_i \ge 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 - y_i f(\mathbf{x}_i)$$

• Together with  $\xi_i \ge 0$  this results in  $\xi_i = \max(0, 1 - y_i f(x_i))$  (given that  $\xi_i$  should be minimized)

#### **Unconstrained optimization** (over w):

$$\operatorname{argmin}_{\mathbf{w}} \quad \underbrace{||\mathbf{w}||^2}_{\text{regularization}} + \quad C \underbrace{\sum_{i=1}^N \max\left(0, 1 - y_i f(\boldsymbol{x}_i)\right)}_{\text{loss function}}$$

# Hinge Loss



#### Points are in 3 categories:

- $y_i f(\boldsymbol{x}_i) > 1$ : Point outside margin, no contribution to loss
- $y_i f(\boldsymbol{x}_i) = 1$ : Point is on the margin, no contribution to loss as in hard margin
- $y_i f(\boldsymbol{x}_i) \leq 1$ : Point violates the margin, contributes to loss



### Loss function is convex



- There is only one minimum
- We can find it with gradient descent
- However: Hinge loss is not differentiable!

# Comparison to logistic loss function

• SVM-hinge loss

$$\operatorname{argmin}_{\mathbf{w}} \quad \lambda \underbrace{||\mathbf{w}||^{2}}_{\text{regularization}} + \underbrace{\sum_{i=1}^{N} \max\left(0, 1 - y_{i}f(\boldsymbol{x}_{i})\right)}_{\text{data loss}}, \quad \text{with } \lambda = \frac{1}{C}$$

• (Regularized) logistic regression loss (see lecture 2)

$$\operatorname{argmax}_{\mathbf{w}} - \lambda ||\mathbf{w}||^{2} + \sum_{i=1}^{N} c_{i} \log(\sigma(f(\boldsymbol{x}_{i})) + (1 - c_{i}) \log(1 - \sigma(f(\boldsymbol{x}_{i}))), \quad \text{with } c_{i} \in \{0, 1\}$$

$$= \dots$$

$$= \operatorname{argmin}_{\mathbf{w}} \quad \lambda \underbrace{||\mathbf{w}||^{2}}_{\text{regularization}} + \underbrace{\sum_{i=1}^{N} \log(1 + \exp(-y_{i}f(\boldsymbol{x}_{i})))}_{\text{deta lage}}, \quad \text{with } y_{i} \in \{-1, 1\}$$

data loss

#### Both loss functions have similar interpretations

- Keep weights small +  $y_i f(\boldsymbol{x}_i)$  should be large
- Saturates if  $y_i f(\boldsymbol{x}_i)$  gets too large

# Comparison to logistic loss function

#### SVM (hinge) loss:

 $\max\left(0,1-y_if(\boldsymbol{x}_i)\right)$ 

- Outputs -1 or 1 (classlabels)
- Estimates maximum margin solution
- Loss contribution is 0 for correct classification

#### Logistic loss:

 $\log\left(1+\exp(-y_i f(\boldsymbol{x}_i))\right)$ 

- Outputs probabilities
- Contribution never 0
  - Often results in slightly less accurate classification
- Diverges faster than hinge loss
  - More sensitive to outliers



# Comparison to logistic regression (LR)





 SVM is less sensitive to outliers

# Agenda for Today

#### **Recap: Linear Discriminators**

#### **Support Vector Machines:**

- Maximum Margin
- Optimization Problem
- Soft-Margin
- Hinge-Loss

#### **SVMs with Kernels**

- Constraint Optimization Problem
- Kernel Trick

#### **Basics:**

• Sub-gradients

### **Basics: Sub-gradients**

**Remember:** For any convex function  $f : \mathbb{R}^d \to \mathbb{R}$ 

$$f(\boldsymbol{z}) \ge f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^T (\boldsymbol{z} - \boldsymbol{x})$$

• I.e. linear approximation underestimates function



A **subgradient** of a convex function *f* at point *x* is any *g* such that

$$f(\boldsymbol{z}) \ge f(\boldsymbol{x}) + \boldsymbol{g}^T(\boldsymbol{z} - \boldsymbol{x})$$

- Always exists (also if *f* is not differentiable)
- If *f* is differentiable at **x**, then  $\boldsymbol{g} = \nabla f(\boldsymbol{x})$

### Examples

Consider f(x) = |x|

- For  $x \neq 0$  , unique sub-gradient of  $g = \operatorname{sign}(x)$
- For x = 0 , sub-gradient is any element of [-1, 1]



### Examples

Let  $f_1, f_2 : \mathbb{R}^n \to \mathbb{R}$  be convex, differentiable, and consider  $f(x) = \max\{f_1(x), f_2(x)\}$ 

- For  $f_1(x) > f_2(x)$ , unique subgradient  $g = \nabla f_1(x)$
- For  $f_2(x) > f_1(x)$ , unique subgradient  $g = \nabla f_2(x)$
- For  $f_1(x) = f_2(x)$ , subgradient g is any point on the line segment between  $\nabla f_1(x)$  and  $\nabla f_2(x)$



Like gradient descent, but replacing gradients with sub-gradients

#### **Sub-gradient Descent:**

- Given convex *f*, not necessarily differentiable
- Initialize  $m{x}_0$
- Repeat:  $m{x}_{t+1} = m{x}_t + \eta m{g}$  , where  $m{g}$  is any sub-gradient of f at point  $m{x}_t$

## Sub-gradients for hinge loss

$$\mathcal{L}(\mathbf{x}_i, y_i; \mathbf{w}) = \max(0, 1 - y_i f(\mathbf{x}_i))$$
  $f(\mathbf{x}_i) = \mathbf{w}^{\top} \mathbf{x}_i + b$ 



# Sub-gradient descent for SVMs



At each iteration, pick random training sample  $(x_i, y_i)$ 

• If 
$$y_i f(x_i) < 1$$
:  $w_{t+1} = w_t - \eta (2w_t - Cy_i x_i)$ 

• Otherwise: 
$$\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta 2 \boldsymbol{w}_t$$

# **Application: Pedestrian Tracking**

#### **Objective:** Detect (localize) standing humans in images



#### **Detection with a sliding window approach:**

- Reduces object detection to binary classification
- Does an image window contain a person or not?

# **Training Data**

• Positive Data: 1208 examples







• Negative Data: 1218 examples





# Features: Histogram of oriented Gradients (HoG features)



Feature vector dimension =  $16 \times 8$  (for tiling) x 8 (orientations) = 1024

# Example HoG features



# Averaged Positive Example







## **Example detection**



#### Dalal and Triggs, CVPR 2005 Gerhard Neumann | Machine Learning | KIT | WS 2021/2022

### Learned model

Model: 
$$f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} + b$$



# Wrap-up

SVMs have been the "gold standard" in the 90s and 2000s for classification

- SVM have been used successfully in many real-world problems
  - text (and hypertext) categorization
  - image classification
  - bioinformatics (Protein classification, cancer classification)
  - hand-written character recognition
- Can be extended to complex feature spaces using kernels (next part)
- ... and regression problems (support vector regression, not covered)

# In the last 7-10 years, neural networks have outperformed SVMs on most applications

• However, similar insights are still used (e.g, hinge loss is also used for DNNs)

# **SVMs with Kernels**

# Support Vector Machines (continued...)

#### **SVMs with features:**

- Maximum margin principle
- Slack variables allow for margin violation

$$\operatorname{argmin}_{\mathbf{w},\boldsymbol{\xi}} \quad ||\mathbf{w}||^2 + C \sum_{i}^{N} \xi_i,$$
  
s.t.  $y_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + b) \ge 1 - \xi_i, \quad \xi_i \ge 0$ 

Simpler formulation without slack variables

 $\begin{aligned} \operatorname{argmin}_{\mathbf{w}} & ||\mathbf{w}||^2, \\ \text{s.t.} & y_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + b) \geq 1 \end{aligned}$ 



In order to apply the kernel trick, we need to apply Constrained Optimization!

# **Constraint Optimization**

# **Basics: Constrained Optimization**

Simple constrained optimization problem:



How do we solve the constrained optimization problem? Lagrangian Multipliers!

# **Basics: Lagrangian Multipliers**

$$\min_{x} x^2 \quad \text{s.t. } x \ge b$$

#### The Lagrangian:

• L = objective - multiplier \* constraint

$$L(x,\lambda) = \underbrace{x^2}_{-} - \underbrace{\lambda}_{-} \underbrace{(x-b)}_{-}$$

objective multiplier constraint

#### Lagrangian optimization:

$$\min_{x} \max_{\lambda} L(x,\lambda), \quad \text{s.t. } \lambda \ge 0$$

### Why is this equivalent? *Min* fights *max*! • *x* < *b* : $- (x-b) < 0 \to \max_{\lambda} -\lambda(x-b) = \infty$ - *min* won't let that happen • x > b: $- (x-b) > 0, \lambda \ge 0 \to \lambda^* = 0$ L is the same as original objective • x = b: $-\lambda$ can be anything L is the same as original objective *Min* forces *max* to behave such that constraints are satisfied

General Formulation:  $\min_{\boldsymbol{x}} f(\boldsymbol{x})$ ,

s.t. 
$$h_i(\boldsymbol{x}) \ge b_i$$
, for  $i = 1 \dots K$ 

• Several inequality constraints (equality constraints also possible)

**Lagrangian optimization:**  $\min_{\boldsymbol{x}} \max_{\boldsymbol{\lambda}} L(\boldsymbol{x}, \boldsymbol{\lambda}), \quad L(\boldsymbol{x}, \boldsymbol{\lambda}) = f(\boldsymbol{x}) - \sum_{i=1}^{K} \lambda_i (h_i(\boldsymbol{x}) - b_i)$ s.t.  $\lambda_i \ge 0$ , for  $i = 1 \dots K$ 

### **Dual formulation**

Primal optimization problem:

 $\min_{\boldsymbol{x}} f(\boldsymbol{x}),$ 

s.t. 
$$h_i(\boldsymbol{x}) \geq b_i$$
, for  $i = 1 \dots K$ 

Dual optimization problem:
λ\* = arg max g(λ), g(λ) = min L(x, λ) λ
s.t. λ<sub>i</sub> ≥ 0, for i = 1...K
g is also called the dual function of the optimization problem
We essentially swapped min and max in the definition of L

Slaters condition: For a convex objective and convex constraints, solving the dual is equivalent to solving the primal!

Optimal primal parameters can be obtained from optimal dual parameters, i.e.

$$oldsymbol{x}^* = rg\min_{oldsymbol{x}} L(oldsymbol{x},oldsymbol{\lambda}^*)$$

# Lagrangian Optimization

#### Basic "Cookbook":

- 1. Write down Lagrangian  $L(\boldsymbol{x}, \boldsymbol{\lambda}) = f(\boldsymbol{x}) - \sum_{i=1}^{K} \lambda_i (h_i(\boldsymbol{x}) - b_i)$
- 2. Obtain optimal solution for primal parameters
  - Compute derivative, set to zero and solve for x

$$rac{\partial L(oldsymbol{x},oldsymbol{\lambda})}{\partial oldsymbol{x}} = 0 
ightarrow oldsymbol{x}^* = f(oldsymbol{\lambda})$$

3. Set  $x^*$  back into Lagrangian to obtain the dual function

 $g(\boldsymbol{\lambda}) = L(f(\boldsymbol{\lambda}), \boldsymbol{\lambda})$ 

- 4. Obtain optimal solution for the dual function
  - Set derivative to zero or gradient descent

$$\boldsymbol{\lambda}^* = \operatorname{argmax}_{\boldsymbol{\lambda}} g(\boldsymbol{\lambda}), \quad \text{s.t. } \lambda_i \geq 0, \forall i$$

5. Compute optimal primal parameters for given

$$\boldsymbol{x}^* = f(\boldsymbol{\lambda}^*)$$



Example:

# $\min_{x} x^2 \quad \text{s.t. } x \ge 1$

### Example:

$$\min_{x} x^2 \quad \text{s.t. } x \ge 1$$

- 1. Write down Lagrangian
- 2. Find optimal primal parameters  $x^* = f(\lambda)$  $x^* = \operatorname{argmin}_x L(x, \lambda) = f(\lambda)$
- 3. Plug back in Lagrangian to get dual  $g(\lambda) = L(x^*(\lambda), \lambda)$
- 4. Find optimal dual parameters

$$\lambda^* = \operatorname{argmax}_{\lambda} g(\lambda) \text{ s.t. } \lambda \ge 0$$

5. Compute primal solution  $x^* = f(\lambda^*)$ 

$$\begin{split} L(x,\lambda) &= x^2 - \lambda(x-1) \\ \frac{\partial}{\partial x} L(x,\lambda) &= 2x - \lambda = 0 \to x^*(\lambda) = \lambda/2 \\ g(\lambda) &= \lambda^2/4 - \lambda(\lambda/2 - 1) = -\lambda^2/4 + \lambda \\ \frac{\partial}{\partial \lambda} g(\lambda) &= -\lambda/2 + 1 = 0 \to \lambda^* = 2 \\ x^* &= \lambda^*/2 = 1 \end{split}$$

### Dual derivation of the SVM

#### SVM optimization:

• Lagrangian:  $\operatorname{argmin}_{\mathbf{w}} ||\mathbf{w}||^2$ , s.t.  $y_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + b) \ge 1$ 

Compute optimal *w*:  $L(\boldsymbol{w}, \boldsymbol{\lambda}) = \frac{1}{2} \boldsymbol{w}^T \boldsymbol{w} - \sum_i \lambda_i (y_i(\boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_i) + b) - 1)$ 

$$\frac{\partial L}{\partial \boldsymbol{w}} = \boldsymbol{w} - \sum_{i} \lambda_{i} y_{i} \boldsymbol{\phi}(\boldsymbol{x}_{i}) = 0,$$
$$\boldsymbol{w}^{*} = \sum_{i} \lambda_{i} y_{i} \boldsymbol{\phi}(\boldsymbol{x}_{i})$$

- Many of the  $\lambda_i$  will be zero (constraint satisfied)
- If  $\lambda_i$  is not zero,  $\phi(\boldsymbol{x}_i)$  is a support vector
- The optimal weight vector w is a linear combination of the support vectors!

### Dual derivation of the SVM

#### **SVM** optimization:

• Lagrangian:  $\operatorname{argmin}_{\mathbf{w}} ||\mathbf{w}||^2$ , s.t.  $y_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + b) \ge 1$ 

#### Optimality condition for b:

$$L(\boldsymbol{w},\boldsymbol{\lambda}) = \frac{1}{2}\boldsymbol{w}^{T}\boldsymbol{w} - \sum_{i} \lambda_{i} (y_{i}(\boldsymbol{w}^{T}\boldsymbol{\phi}(\boldsymbol{x}_{i}) + b) - 1)$$

- We do not obtain a solution for b
- But an additional condition for the lambdas

$$\frac{\partial L}{\partial b} = -\sum_{i} \lambda_{i} y_{i} \Rightarrow \sum_{i} \lambda_{i} y_{i} = 0$$

#### b can be computed from w:

• If  $\lambda_i > 0$  , then  $oldsymbol{x}_i$  is on the margin, i.e.:

$$y_i(\mathbf{w}^T \boldsymbol{\phi}(\boldsymbol{x}_i) + b) = 1$$
  
 $y_i y_i(\mathbf{w}^T \boldsymbol{\phi}(\boldsymbol{x}_i) + b) = y_i$   
 $b = y_i - \mathbf{w}^T \boldsymbol{\phi}(\boldsymbol{x}_i)$ 

### Kernel Trick in SVMs

Lagrangian: 
$$L(\boldsymbol{w}, \boldsymbol{\lambda}) = \frac{1}{2} \boldsymbol{w}^T \boldsymbol{w} - \sum_i \lambda_i (y_i (\boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_i) + b) - 1), \quad \boldsymbol{w}^* = \sum_i \lambda_i y_i \boldsymbol{\phi}(\boldsymbol{x}_i)$$

**Dualfunction:** 

$$g(\boldsymbol{\lambda}) = L(\boldsymbol{w}^*, \boldsymbol{\lambda})$$

$$= \frac{1}{2} \underbrace{\sum_{i} \sum_{j} \lambda_i \lambda_j y_i y_j \boldsymbol{\phi}(\boldsymbol{x}_i)^T \boldsymbol{\phi}(\boldsymbol{x}_j)}_{\boldsymbol{w}^{*T} \boldsymbol{w}^*} - \underbrace{\sum_{i} \lambda_i y_i (\sum_{j} \lambda_j y_j \boldsymbol{\phi}(\boldsymbol{x}_j))^T \boldsymbol{\phi}(\boldsymbol{x}_i) + \sum_{i} \lambda_i}_{\boldsymbol{w}^*}$$

$$= \underbrace{\sum_{i} \lambda_i - \frac{1}{2} \sum_{i} \sum_{j} \lambda_i \lambda_j y_i y_j \boldsymbol{\phi}(\boldsymbol{x}_i)^T \boldsymbol{\phi}(\boldsymbol{x}_j)}_{j}$$

We just derived the kernel trick for SVMs

$$g(\boldsymbol{\lambda}) = \sum_{i} \lambda_{i} - \frac{1}{2} \sum_{i} \sum_{j} \lambda_{i} \lambda_{j} y_{i} y_{j} \boldsymbol{k}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j})$$

• Scalar products of the feature vectors can be written as kernels

### Kernelized SVM

#### Solve dual optimization problem:

$$\max_{\boldsymbol{\lambda}} \sum_{i} \lambda_{i} - \frac{1}{2} \sum_{i} \sum_{j} \lambda_{i} \lambda_{j} y_{i} y_{j} \boldsymbol{k}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j})$$
  
s.t.  $\lambda_{i} \ge 0, \forall i \in [1 \dots N], \quad \sum_{i} \lambda_{i} y_{i} = 0$ 

#### **Compute primal from dual parameters**

Weight vector (can not be represented): 
$$w^* = \sum_i \lambda_i y_i \phi(x_i)$$
Bias: for any i with  $\lambda_i > 0$ :  $b = y_k - \phi(x_k)^T w^*$   
 $= y_k - \sum_i y_i \lambda_i k(x_i, x_k)$ 
Decision function:  $f(x) = \phi(x)^T w^* + b$ 

$$=\sum_{i} y_i \lambda_i k(\boldsymbol{x}_i, \boldsymbol{x}) + b$$

)

## Relaxed constraints with slack

#### Primal optimization problem:

$$\operatorname{argmin}_{\mathbf{w},\boldsymbol{\xi}} \quad ||\mathbf{w}||^2 + C \sum_{i}^{N} \xi_i,$$
  
s.t.  $y_i(\mathbf{w}^T \mathbf{x}_i + b) \ge 1 - \xi_i, \quad \xi_i \ge 0$ 

#### What changed?

• Added upper bound of C on  $\lambda_i$  !

**Dual optimization problem:** 

$$\max_{\boldsymbol{\lambda}} \sum_{i} \lambda_{i} - \frac{1}{2} \sum_{i} \sum_{j} \lambda_{i} \lambda_{j} y_{i} y_{j} \boldsymbol{k}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j})$$
  
s.t.  $\boldsymbol{C} \ge \lambda_{i} \ge 0, \forall i \in [1 \dots N], \quad \sum_{i} \lambda_{i} y_{i} = 0$ 

• For computing b, we now take an example where  $C > \lambda_i > 0$ 

#### Intuitive explanation:

- Without slack,  $\lambda_i \to \infty$  when constraints are violated (points misclassified)
- Upper bound of C limits the  $\lambda_i$ , so misclassifications are allowed

### Example: SVM with RBF kernel

#### Data is non-linearly separable in original space



### Example: SVM with RBF kernel

#### Data is non-linearly separable in original space



# **Example: Different Cs**

$$\sigma = 1.0$$
  $C = \infty$ 



$$\sigma = 1.0 \quad C = 10$$



$$\sigma = 1.0 \quad C = 100$$



### Example: Different sigma

### $\sigma = 1.0$ $C = \infty$



$$\sigma = 0.25$$
  $C = \infty$ 



 $\sigma = 0.1$   $C = \infty$ 



# Overfitting

#### Huge feature space with kernels: should we worry about overfitting?

- SVM objective seeks a solution with large margin
- Theory says that large margin leads to good generalization
- But everything overfits sometimes!!!

#### Can control overfitting by:

- Setting C (low C -> smaller Complexity)
- Choosing a better Kernel
- Varying parameters of the Kernel (width of Gaussian, etc.)

Model Selection Problems

### Handwritten Digit Clasification

#### US postal service database

• Human performance: 2.5% error

#### Various learning algorithms (pre-deep learning)

- 16.2%: Decision tree (C4.5)
- 5.9%: 2-layer neural network
- 5.1%: LeNet 1 5-layer neural network

#### Various SVM results

- 4.0%: Polynomial kernel (274 support vectors)
- 4.1%: Gaussian kernel

# Handwritten Digit Clasification

#### Very little overfitting due to max-margin

degree of	dimensionality of	support	raw
polynomial	feature space	vectors	error
1	256	282	8.9
2	pprox 33000	227	4.7
3	$pprox 1  imes 10^6$	274	4.0
4	$\approx 1 \times 10^9$	321	4.2
5	$pprox 1  imes 10^{12}$	374	4.3
6	$pprox 1  imes 10^{14}$	377	4.5
7	$pprox 1  imes 10^{16}$	422	4.5

#### **Recent results**

- With more training data, better modeling of invariances, etc.
- Error down to about 0.5% with SVMs and 0.4% with neural networks

# Takeaway messages

#### What have we learned today?

#### Maximum Margin Classifiers:

- A robust formulation for classification
- Margin can be expressed as constrained optimization
- Slack variables allow for constrained violation and regularization
- Can be efficiently optimized using hinge-loss and subgradient descent

#### Sub-gradients:

Use it for non-differentiable convex functions

#### Kernel trick for SVMs:

- Results from the Lagrangian dual formulation
- Optimal solution for w is a linear combination of the support vectors (compare to kernel regression)



# Self-test questions

#### You should understand now:

- Why is it good to use a maximum margin objective for classification?
- How can we define the margin as optimization problem?
- What are slack variables and how can they be used to get a "soft" margin?
- How is the hinge loss defined?
- What is the relation between the slack variables and the hinge loss?
- What are the advantages and disadvantages in comparison to logistic regression?
- What is the difference between gradients and sub-gradients