# EECS 545: Machine Learning

## Lecture 10. Kernel methods: Kernelizing Support Vector Machines

Honglak Lee
02/12/2025

UNIVERSITY OF MICHIGAN

---

## Overview

- Support Vector Machine (SVM)
- Dual optimization
  - General recipe for constrained optimization
  - Hard-margin SVM
  - Soft-margin SVM
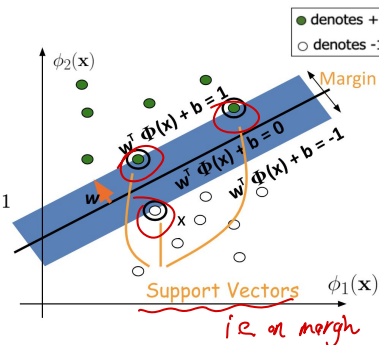
---

## Maximum Margin Classifier

- Optimization problem:
$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|^2$$

subject to

For $y^{(n)} = 1$, $\quad \mathbf{w}^\top \phi\left(\mathbf{x}^{(n)}\right) + b \geq 1$

For $y^{(n)} = -1$, $\quad \mathbf{w}^\top \phi\left(\mathbf{x}^{(n)}\right) + b \leq -1$



- denotes +1
- denotes -1

$\phi_2(\mathbf{x})$

$\mathbf{w}^\top \Phi(\mathbf{x}) + b = 1$
$\mathbf{w}^\top \Phi(\mathbf{x}) + b = 0$
$\mathbf{w}^\top \Phi(\mathbf{x}) + b = -1$

Margin

Support Vectors
i.e. on margin

$\phi_1(\mathbf{x})$

---

## Dual optimization

- So far, we have considered primal optimization which requires a direct access to the feature vectors $\phi\left(\mathbf{x}^{(n)}\right)$
- It is also possible to "kernelize" SVM
  - This formulation is called "Dual" formulation.
  - In this case, you can use any kernel function (such as polynomial, RBF, etc.)

With dual variables $\alpha^{(n)}$, we have the following relations (without proofs)

$$\mathbf{w} = \sum_{n=1}^{N} \alpha^{(n)} y^{(n)} \phi\left(\mathbf{x}^{(n)}\right)$$

$$h(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_{n=1}^{N} \alpha^{(n)} y^{(n)} k\left(\mathbf{x}, \mathbf{x}^{(n)}\right) + b$$

---

## Kernelizing SVM: back to hard-margin case

- Optimization problem:
$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|^2$$
$$\text{subject to} \quad y^{(n)}\left(\mathbf{w}^\top \phi\left(\mathbf{x}^{(n)}\right) + b\right) \geq 1, n = 1, ..., N$$

- This is a constrained optimization problem.
  - We solve this using Lagrange multipliers (convex optimization)
  - Solving dual optimization problem naturally leads to kernalization

---

## Solving Constrained Optimization: General Overview and Recipe

(This section is just a recap, see the supplementary lecture slides for more details)

---

## General (Constrained) Optimization

- General optimization problem:
$$\min_{\mathbf{x}} \quad f(\mathbf{x}) \quad \text{objective (cost) function}$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, i = 1, ..., m \quad \text{inequality constraint functions}$$
$$h_i(\mathbf{x}) = 0, i = 1, ..., p \quad \text{equality constraint functions}$$

- If $\mathbf{x}$ satisfies all the constraints, $\mathbf{x}$ is called underline{feasible} (a feasible solution).
- In general, this is a nontrivial problem to solve, so we use techniques for convex optimization.

---

## Recap: General Recipe

- Given an original optimization
$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, i = 1, ..., p$$

- Solve dual optimization with Lagrangian function:
$$\max_{\lambda,\nu} \min_{\mathbf{x}} \quad \mathcal{L}(\mathbf{x}, \lambda, \nu) = f(\mathbf{x}) + \sum_{i=1}^{m} \lambda_i g_i(\mathbf{x}) + \sum_{i=1}^{p} \nu_i h_i(\mathbf{x})$$
$$\text{subject to} \quad \lambda_i \geq 0, \forall i$$

Add constraint terms with Lagrange multipliers

- Alternatively, solve the dual optimization with underline{Lagrange dual}:
$$\max_{\lambda,\nu} \quad \tilde{\mathcal{L}}(\lambda, \nu) \qquad \text{where} \quad \tilde{\mathcal{L}}(\lambda, \nu) = \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$
$$\text{subject to} \quad \lambda_i \geq 0, \forall i$$

## A Big Picture

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, \, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, \, i = 1, ..., p$$

**Constrained Optimization Problem**

$\mathcal{L}(\mathbf{x}, \lambda, \nu) = f(\mathbf{x}) + \sum_{i=1}^{m} \lambda_i g_i(\mathbf{x}) + \sum_{i=1}^{p} \nu_i h_i(\mathbf{x})$

**Lagrangian**

**Primal Optimization Problem** (min-max)

$\min_{\mathbf{x}} \max_{\nu, \lambda: \lambda_i \geq 0, \forall i} \mathcal{L}(\mathbf{x}, \lambda, \nu)$

e.g. convex optimizations, KKT conditions

**strong duality** (if some conditions are met)

$p^* = d^*$

**weak duality**

$p^* \geq d^*$

**Dual Optimization Problem** (max-min)

$\max_{\nu, \lambda: \lambda_i \geq 0, \forall i} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu)$

9

## Lagrangian Formulation

- The **Lagrangian function** is

$$\mathcal{L}(\mathbf{x}, \lambda, \nu) = f(\mathbf{x}) + \sum_{i=1}^{m} \lambda_i g_i(\mathbf{x}) + \sum_{i=1}^{p} \nu_i h_i(\mathbf{x})$$

$\leq 0$    $= 0$ (when constraints satisfied)

- Here, $\lambda = [\lambda_1, ..., \lambda_m]$ $(\lambda_i \geq 0, \forall i)$ and $\nu = [\nu_1, ..., \nu_p]$ are called Lagrange multipliers (or dual variables)

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, i = 1, ..., p$$

- This leads to **primal optimization problem**

$$\min_{\mathbf{x}} \max_{\nu, \lambda: \lambda_i \geq 0 \, \forall i} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$

otherwise

$\lambda_i g_i(\mathbf{x}) \to \infty$
$\nu_i h_i(\mathbf{x}) \to \pm\infty$
as $\lambda_i, \nu_i \to \pm\infty$

- Difficult to solve directly!

10

## Primal and Feasibility

- Primal optimization problem:

$$p^* = \min_{\mathbf{x}} \max_{\nu, \lambda: \lambda_i \geq 0 \, \forall i} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, \, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, \, i = 1, ..., p$$

where $\mathcal{L}(\mathbf{x}, \lambda, \nu) = f(\mathbf{x}) + \sum_{i=1}^{m} \lambda_i g_i(\mathbf{x}) + \sum_{i=1}^{p} \nu_i h_i(\mathbf{x})$

- Notice that:

$$\mathcal{L}_p(\mathbf{x}) = \max_{\nu, \lambda: \lambda_i \geq 0, \forall i} \mathcal{L}(\mathbf{x}, \lambda, \nu) = \begin{cases} f(\mathbf{x}) & \text{if } \mathbf{x} \text{ is feasible} \\ \infty & \text{otherwise} \end{cases}$$

This eliminates the constraints on **x**, yielding an equivalent optimization problem.

11

## Lagrange Dual

primal vs dual: switching the order of min / max

Note: these are different problems!

- Dual optimization problem:

$$d^* = \max_{\nu, \lambda: \lambda_i \geq 0, \forall i} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$

note: it does not guarantee $\mathcal{L}_d(\mathbf{x}) < \infty$ when $\mathbf{x}$ not feasible!
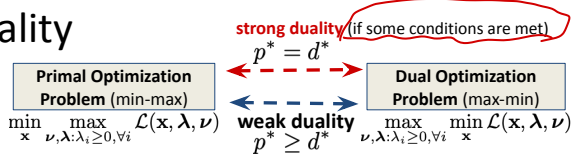
cf) primal optimization problem $p^* = \min_{\mathbf{x}} \max_{\nu, \lambda: \lambda_i \geq 0, \forall i} \mathcal{L}(\mathbf{x}, \lambda, \nu)$

- We can also write as:

(so dual is not equiv to original problem)

$$\max_{\lambda, \nu} \min_{\mathbf{x}} \quad \mathcal{L}(\mathbf{x}, \lambda, \nu)$$
$$\text{subject to} \quad \lambda_i \geq 0, \forall i$$

or

$$\max_{\lambda, \nu} \quad \tilde{\mathcal{L}}(\lambda, \nu)$$
$$\text{subject to} \quad \lambda_i \geq 0, \forall i$$
$$\text{where} \quad \tilde{\mathcal{L}}(\lambda, \nu) = \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$

***Lagrange Dual function***

12

## Weak Duality

**strong duality** (if some conditions are met)

$p^* = d^*$

**Primal Optimization Problem** (min-max)

$\min_{\mathbf{x}} \max_{\nu, \lambda: \lambda_i \geq 0, \forall i} \mathcal{L}(\mathbf{x}, \lambda, \nu)$

**weak duality**

$p^* \geq d^*$

**Dual Optimization Problem** (max-min)

$\max_{\nu, \lambda: \lambda_i \geq 0, \forall i} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu)$

- Claim: 

$$d^* = \max_{\nu, \lambda: \lambda_i \geq 0} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$
$$\leq \min_{\mathbf{x}} \max_{\nu, \lambda: \lambda_i \geq 0} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$
$$= p^*$$

- Difference between $p^*$ and $d^*$ is called the **duality gap**.

- In other words, the dual maximization problem (usually easier) gives a "**lower bound**" for the primal minimization problem (usually more difficult).

13

## Weak Duality

Also see Convex Optimization Review Session

$$d^* = \max_{\nu, \lambda: \lambda_i \geq 0} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu) \leq \min_{\mathbf{x}} \max_{\nu, \lambda: \lambda_i \geq 0} \mathcal{L}(\mathbf{x}, \lambda, \nu) = p^*$$

- Proof: Let $\tilde{\mathbf{x}}$ be feasible. Then for any $\lambda, \nu$ with $\lambda_i \geq 0$,

$$\mathcal{L}(\tilde{\mathbf{x}}, \lambda, \nu) = f(\tilde{\mathbf{x}}) + \sum_{i=1}^{m} \lambda_i g_i(\tilde{\mathbf{x}}) + \sum_{i=1}^{p} \nu_i h_i(\tilde{\mathbf{x}}) \leq f(\tilde{\mathbf{x}})$$

Thus, $\tilde{\mathcal{L}}(\lambda, \nu) = \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu) \leq \mathcal{L}(\tilde{\mathbf{x}}, \lambda, \nu) \leq f(\tilde{\mathbf{x}})$ for any $\lambda, \nu$ with $\lambda_i \geq 0$, any feasible $\tilde{\mathbf{x}}$

Then, maximize LHS (w.r.t. dual variables)

$$d^* = \max_{\nu, \lambda: \lambda_i \geq 0} \tilde{\mathcal{L}}(\lambda, \nu) \leq f(\tilde{\mathbf{x}}) \text{ for any feasible } \tilde{\mathbf{x}}$$

Finally, minimize RHS (w.r.t. primal variable)

$$d^* = \max_{\nu, \lambda: \lambda_i \geq 0} \tilde{\mathcal{L}}(\lambda, \nu) \leq \min_{\tilde{\mathbf{x}}: \text{feasible}} f(\tilde{\mathbf{x}}) = p^*$$

14

## Strong Duality

- If $p^* = d^*$, we say <u>strong duality</u> holds.

- What are the conditions for strong duality?
  - does not hold in general
  - holds for convex problems (under mild conditions)
  - conditions that guarantee strong duality in convex problems are called constraint qualification.

- Two well-known conditions (in convex problems)
  - Slater's constraint qualification (review session)
  - Karush-Kuhn-Tucker (KKT) condition (main focus)

15

## Convex Optimization

Also see Convex Optimization Review Session

- Standard form of **convex problem** has the form:

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, i = 1, ..., p$$

(where *f*, $g_i$ are convex, *and* $h_i$ are affine)

- If **x** satisfies all the constraints, **x** is called <u>feasible</u>.
  - In general, this is a nontrivial problem to solve, so we use techniques for convex optimization.

16

## (Sufficient) Conditions for strong duality: **Slater's constraint qualification**

- Strong duality holds for a **convex** problem

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, i = 1, ..., p$$

(where $f, g_i$ are **convex**, and $h_i$ are **affine**)

*if* the constraint is <u>strictly</u> feasible (by any solution), i.e.,  ✓

$$\exists \mathbf{x}: \quad g_i(\mathbf{x}) < 0, \forall i = 1, ..., m \quad \text{(Not necessarily an optimal solution)}$$
$$h_i(\mathbf{x}) = 0, \forall i = 1, ..., p$$

Slater's condition is a **sufficient** condition for strong duality to hold for a convex problem

17

---

## Karush-Kuhn-Tucker (KKT) condition

(necessary)

Let $\mathbf{x}^*$ be a primal optimal and $\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$ be a dual optimal solution. If the strong duality holds, then we have the following:

$$\nabla_{\mathbf{x}} f(\mathbf{x}^*) + \sum_{i=1}^{m} \lambda_i^* \nabla_{\mathbf{x}} g_i(\mathbf{x}^*) + \sum_{i=1}^{p} \nu_i^* \nabla_{\mathbf{x}} h_i(\mathbf{x}^*) = 0, \quad \boxed{\text{Stationarity (1)}}$$

$$g_i(\mathbf{x}^*) \leq 0, \quad i = 1, \ldots, m, \quad \boxed{\text{Primal feasibility (2)}}$$

$$h_i(\mathbf{x}^*) = 0, \quad i = 1, \ldots, p, \quad \boxed{\text{Primal feasibility (3)}}$$

$$\lambda_i^* \geq 0, \quad i = 1, \ldots, m, \quad \boxed{\text{Dual feasibility (4)}}$$

$$\lambda_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \ldots, m \quad \boxed{\text{Complementary slackness (5)}}$$

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0. \ i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0. \ i = 1, ..., p$$

$$\max_{\lambda, \nu} \min_{\mathbf{x}} \quad \mathcal{L}(\mathbf{x}, \lambda, \nu)$$
$$\text{subject to} \quad \lambda_i \geq 0, \forall i$$
Dual problem

Note: we do **not** assume the optimization problem is necessarily convex for describing KKT condition. However, when the problem is convex (and differentiable), KKT condition ensures strong duality.

18

---

## (Sufficient) Conditions for strong duality: KKT Conditions

- Assume $f, g_i, h_i$ are differentiable

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, i = 1, ..., p$$

- If the original problem is **convex** (where $f, g_i$ are **convex** and $h_i$ are affine), and $\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$ satisfy the KKT conditions, then:
  - $\mathbf{x}^*$ is primal optimal
  - $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$ is dual optimal, and
  - the <u>duality gap is zero</u> (i.e., strong duality holds)

**For convex optimization problems** (+ differentiable objectives/constraints)**, KKT is a sufficient condition for strong duality.**

19

---

## Proof for sufficiency (KKT => Strong duality)

- From (2) and (3), $\mathbf{x}^*$ is primal feasible.
- From (4), $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$ is dual feasible.

**Claim: When KKT (1)-(5) holds, the strong duality holds.**

- $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu})$ is a convex differentiable function. Thus, from (1), $\mathbf{x}^*$ is a minimizer of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$. $\quad \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) = 0$

- Then, 
$$\tilde{\mathcal{L}}(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) = \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$$
$$= \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$$
$$= f(\mathbf{x}^*) + \sum_i \lambda_i^* g_i(\mathbf{x}^*) + \sum_i \nu_i^* h_i(\mathbf{x}^*)$$
$$= f(\mathbf{x}^*)$$

(See also: derivation of complementary slackness)

$= 0$ ∵ (5) complementary slackness

- But, $\tilde{\mathcal{L}}(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) \leq \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} \tilde{\mathcal{L}}(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq \min_{\mathbf{x}: \mathbf{x} \text{ is feasible}} f(\mathbf{x}) \leq f(\mathbf{x}^*) = \tilde{\mathcal{L}}(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$

weak duality

- Then, $\max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} \tilde{\mathcal{L}}(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \min_{\mathbf{x}: \mathbf{x} \text{ is feasible}} f(\mathbf{x})$

which proves that the strong duality holds (i.e., duality gap is zero). 20

---

## KKT conditions: Conclusion

- If a constrained optimization if **differentiable** and has **convex** objective function and constraint sets, then the KKT conditions are **(necessary and) sufficient conditions** for **strong duality** (zero duality gap).

- Thus, the KKT conditions can be used to solve such problems.

21

---

## Applying Constrained Optimization Techniques for solving SVM

22

---

## Kernelizing SVM: back to hard-margin case

- Optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2}||\mathbf{w}||^2 \qquad \text{label is either -1 or +1}$$
$$\text{subject to} \quad y^{(n)} \left( \mathbf{w}^\top \phi\left(\mathbf{x}^{(n)}\right) + b \right) \geq 1, n = 1, ..., N$$

- This is a constrained optimization problem.
  - We solve this using Lagrange multipliers (convex optimization)

23

---

## Back to hard-margin SVM

- Use Lagrange multipliers to enforce constraints while optimizing

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2}||\mathbf{w}||^2 + \sum_{n=1}^{N} \alpha^{(n)} \left\{ 1 - y^{(n)} \left( \mathbf{w}^\top \phi\left(\mathbf{x}^{(n)}\right) + b \right) \right\}$$

- Here, $\alpha^{(n)} \geq 0$ is the Lagrange multiplier (or dual variable) for each constraint (one per data point)

$$y^{(n)} \left( \mathbf{w}^\top \phi\left(\mathbf{x}^{(n)}\right) + b \right) \geq 1 \qquad n = 1, ..., N$$

24

## Lagrangian and Lagrange Dual

- Optimizing the <u>Lagrange</u> dual problem :

$$\max_{\boldsymbol{\alpha}} \min_{\mathbf{w},b} \mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) = \frac{1}{2}\|\mathbf{w}\|^2 + \sum_{n=1}^{N}\alpha^{(n)}\left\{1 - y^{(n)}\left(\mathbf{w}^{\top}\phi\left(\mathbf{x}^{(n)}\right) + b\right)\right\}$$

$$\text{subject to} \quad \alpha^{(n)} \geq 0, \forall n$$

- We first minimize w.r.t. primal variables **w** and b, and get a <u>Lagrange dual problem:</u>

$$\max_{\boldsymbol{\alpha}} \quad \tilde{\mathcal{L}}(\boldsymbol{\alpha})$$

$$\text{subject to} \quad \alpha^{(n)} \geq 0, \forall n$$

$$\text{where} \quad \tilde{\mathcal{L}}(\boldsymbol{\alpha}) = \min_{\mathbf{w},b}\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) \qquad \text{(a.k.a. Lagrange dual function)}$$

(Please see the supplementary material for more explanation about Lagrange Dual)

25

## Maximize the Margin

- Lagrangian function:

$$\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) = \frac{1}{2}\|\mathbf{w}\|^2 + \sum_{n=1}^{N}\alpha^{(n)}\left\{1 - y^{(n)}\left(\mathbf{w}^{\top}\phi\left(\mathbf{x}^{(n)}\right) + b\right)\right\}$$

$$\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) = \mathbf{w} - \sum_{n}\alpha^{(n)}y^{(n)}\phi(x^{(n)})$$

$$:= 0$$

- Set the derivatives of $\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha})$ to zero, to get

$$\mathbf{w} = \sum_{n=1}^{N}\alpha^{(n)}y^{(n)}\phi\left(\mathbf{x}^{(n)}\right) \qquad 0 = \sum_{n=1}^{N}\alpha^{(n)}y^{(n)}$$

c.f. KKT (1) Stationarity
$\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) = 0$
$\nabla_{b}\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) = 0$

$-\sum_{n}\alpha^{(n)}y^{(n)} := 0$

- Substitute in, to eliminate **w** and *b*,

$\mathcal{L}(\mathbf{w}^*, b^*, \alpha)$

$$\max_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}) = \sum_{n=1}^{N}\alpha^{(n)} - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}\alpha^{(n)}\alpha^{(m)}y^{(n)}y^{(m)}\phi\left(\mathbf{x}^{(n)}\right)^{\top}\phi\left(\mathbf{x}^{(m)}\right)$$

$$\text{subject to} \quad \alpha^{(n)} \geq 0, \quad \forall n$$

26

## Dual Representation (with kernel)

- Define a kernel $\quad k\left(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}\right) = \phi\left(\mathbf{x}^{(n)}\right)^{\top}\phi\left(\mathbf{x}^{(m)}\right)$

- Dual optimization is to maximize

$$\max_{\boldsymbol{\alpha}}\tilde{\mathcal{L}}(\boldsymbol{\alpha}) = \sum_{n=1}^{N}\alpha^{(n)} - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}\alpha^{(n)}\alpha^{(m)}y^{(n)}y^{(m)}\underbrace{\phi\left(\mathbf{x}^{(n)}\right)^{\top}\phi\left(\mathbf{x}^{(m)}\right)}_{=k(\mathbf{x}^{(n)},\mathbf{x}^{(m)})}$$

$$\text{subject to} \quad \alpha^{(n)} \geq 0, \quad \forall n$$

- Once we have $\boldsymbol{\alpha}$, we don't need **w**.

- Predict classification for arbitrary input **x** using:

$$h\left(\mathbf{x}\right) = \mathbf{w}^{\top}\phi\left(\mathbf{x}\right) + b = \sum_{n=1}^{N}\alpha^{(n)}y^{(n)}k\left(\mathbf{x},\mathbf{x}^{(n)}\right) + b$$

$$\mathbf{w} = \sum_{n=1}^{N}\alpha^{(n)}y^{(n)}\phi\left(\mathbf{x}^{(n)}\right)$$

27

## Support Vectors

- The KKT conditions are:
$$\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) = 0$$
$$\nabla_{b}\mathcal{L}(\mathbf{w},b,\boldsymbol{\alpha}) = 0$$
$$\alpha^{(n)} \geq 0$$
$$1 - y^{(n)}h\left(\mathbf{x}^{(n)}\right) \leq 0$$
$$\alpha^{(n)}\left\{1 - y^{(n)}h\left(\mathbf{x}^{(n)}\right)\right\} = 0$$

- The last condition (complementary slackness) means:
  - either $\alpha^{(n)} = 0$ or $y^{(n)}h\left(\mathbf{x}^{(n)}\right) = 1$ ·

  support vectors

- That is, only the support vectors matter!
  - To compute $h(\mathbf{x})$ (prediction), sum only over support vectors $\quad h\left(\mathbf{x}\right) = \sum_{m:\text{support vectors}}\alpha^{(m)}y^{(m)}k\left(\mathbf{x},\mathbf{x}^{(m)}\right) + b$

28

## Recovering b

- For any support vector $\mathbf{x}^{(n)} : y^{(n)}h\left(\mathbf{x}^{(n)}\right) = 1$

- Replacing with $\quad h\left(\mathbf{x}\right) = \sum_{m\in S}\alpha^{(m)}y^{(m)}k\left(\mathbf{x},\mathbf{x}^{(m)}\right) + b$

$$y^{(n)}\left(\sum_{m\in S}\alpha^{(m)}y^{(m)}k\left(\mathbf{x}^{(n)},\mathbf{x}^{(m)}\right) + b\right) = 1$$

(index) set of support vectors

- Multiply $y^{(n)}$, and sum over n:

$$b = \frac{1}{N_S}\sum_{n\in S}\left(y^{(n)} - \sum_{m\in S}\alpha^{(m)}y^{(m)}k\left(\mathbf{x}^{(n)},\mathbf{x}^{(m)}\right)\right)$$
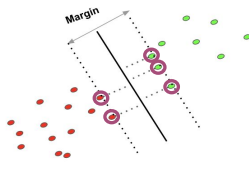
Margin

Image adapted from:
https://www.vebuso.com/2020/02/a-top-machine-learning-algorithm-explained-support-vector-machines-svms/

29

## Formulation of soft-margin SVM

- Maximize the margin, and also penalize for the slack variables

$$C\sum_{n=1}^{N}\xi^{(n)} + \frac{1}{2}\|\mathbf{w}\|^2$$

- The support vectors are now those with

$$y^{(n)}h\left(\mathbf{x}^{(n)}\right) = 1 - \xi^{(n)}$$

30

## Dual formulation of soft-margin SVM

- Lagrangian

$$\mathcal{L}(\mathbf{w},b,\boldsymbol{\xi},\boldsymbol{\alpha},\boldsymbol{\mu}) = \frac{1}{2}\|\mathbf{w}\| + C\sum_{n=1}^{N}\xi^{(n)} + \sum_{n=1}^{N}\alpha^{(n)}\left\{1 - y^{(n)}h(\mathbf{x}^{(n)}) - \xi^{(n)}\right\} + \sum_{n=1}^{N}\mu^{(n)}\left(-\xi^{(n)}\right)$$

$$\text{where} \quad \alpha^{(n)} \geq 0, \quad \mu^{(n)} \geq 0, \quad \xi^{(n)} \geq 0, \quad \forall n$$

- KKT conditions for the constraints

$$\left.\begin{array}{r}1 - y^{(n)}h\left(\mathbf{x}^{(n)}\right) - \xi^{(n)} \leq 0 \\ -\xi^{(n)} \leq 0\end{array}\right\} \text{ Primal variables satisfy the inequality constraints}$$

$$\left.\begin{array}{r}\alpha^{(n)} \geq 0 \\ \mu^{(n)} \geq 0\end{array}\right\} \text{ Dual variables (for above inequalities) are feasible}$$

$$\left.\begin{array}{r}\alpha^{(n)}\left(1 - y^{(n)}h\left(\mathbf{x}^{(n)}\right) - \xi^{(n)}\right) = 0 \\ \mu^{(n)}\xi^{(n)} = 0\end{array}\right\} \text{ Complementary slackness condition}$$

32

## Dual formulation of soft-margin SVM

- Taking derivatives

$$\frac{\partial\mathcal{L}}{\partial\mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w} = \sum_{n=1}^{N}\alpha^{(n)}y^{(n)}\phi\left(\mathbf{x}^{(n)}\right)$$

$$\frac{\partial\mathcal{L}}{\partial b} = 0 \quad \Rightarrow \quad \sum_{n=1}^{N}\alpha^{(n)}y^{(n)} = 0$$

$$\frac{\partial\mathcal{L}}{\partial\xi^{(n)}} = 0 \quad \Rightarrow \quad \alpha^{(n)} = C - \mu^{(n)}$$

33

## Dual formulation of soft-margin SVM

$$\mathbf{w} = \sum_{n=1}^{N} \alpha^{(n)} y^{(n)} \phi\left(\mathbf{x}^{(n)}\right) \qquad \sum_{n=1}^{N} \alpha^{(n)} y^{(n)} = 0 \qquad \alpha^{(n)} = C - \mu^{(n)}$$

- Plug these back into the Lagrangian:

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu}) = \frac{1}{2}\mathbf{w}^\top \mathbf{w} + \sum_{n=1}^{N} \underbrace{(C - \mu^{(n)})}_{\alpha^{(n)}} \xi^{(n)} + \sum_{n=1}^{N} \alpha^{(n)}\{1 - y^{(n)}(\mathbf{w}^\top \phi(\mathbf{x}^{(n)}) + b)) - \xi^{(n)}\}$$

$$= \frac{1}{2}\mathbf{w}^\top \mathbf{w} - \sum_{n=1}^{N} \alpha^{(n)} y^{(n)} \mathbf{w}^\top \phi(\mathbf{x}^{(n)}) - b \underbrace{\sum_{n=1}^{N} \alpha^{(n)} y^{(n)}}_{0} + \sum_{n=1}^{N} \alpha^{(n)}$$

$$= \frac{1}{2}\mathbf{w}^\top \mathbf{w} - \mathbf{w}^\top \underbrace{\left(\sum_{n=1}^{N} \alpha^{(n)} y^{(n)} \phi(\mathbf{x}^{(n)})\right)}_{\mathbf{w}} + \sum_{n=1}^{N} \alpha^{(n)}$$

$$= \sum_{n=1}^{N} \alpha^{(n)} - \frac{1}{2}\mathbf{w}^\top \mathbf{w}$$

$$= \sum_{n=1}^{N} \alpha^{(n)} - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} \alpha^{(n)}\alpha^{(m)} y^{(n)} y^{(m)} \phi(\mathbf{x}^{(n)})^\top \phi(\mathbf{x}^{(m)})$$

## Dual formulation of soft-margin SVM

- Dual optimization (via Lagrange dual)

$$\max_{\boldsymbol{\alpha}} \quad \sum_{n=1}^{N} \alpha^{(n)} - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{M} \alpha^{(n)}\alpha^{(m)} y^{(n)} y^{(m)} k\left(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}\right) \quad \text{Inner product of features replaced with kernel}$$

$$\text{subject to} \quad 0 \leq \alpha^{(n)} \leq C \quad \longleftarrow \mu^{(n)} = C - \alpha^{(n)} \geq 0$$

$$\sum_{n=1}^{N} \alpha^{(n)} y^{(n)} = 0$$

- Solve quadratic problem (convex optimization)

## SVM: practical issues

## Support Vector Machine: Algorithm

1. Choose a kernel function

2. Choose a value for *C*
   (i.e., smaller C → larger regularization)

3. Solve the optimization problem (many software packages available) – primal or dual

4. Construct the discriminant function from the support vectors

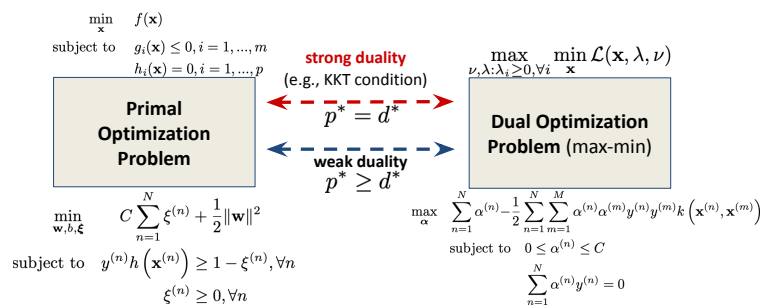## Some Issues

- Linear kernels work fairly well, but can be suboptimal.
- Choice of (nonlinear) kernels
  – Gaussian or polynomial kernel is default
  – If the simple kernels are ineffective, more elaborate kernels are needed
  – Domain experts can give assistance in formulating appropriate similarity measures
- Choice of kernel parameters
  – E.g., Gaussian kernel: $K(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right)$
    - $\sigma$ is the distance between neighboring points whose labels are likely to affect the prediction of the query point.
  – In the absence of reliable criteria, applications rely on the use of a validation set or cross-validation to set such parameters.

## Summary: Support Vector Machine

- Max margin classifier: improved robustness & less over-fitting
- Solved by convex optimization techniques
- Kernel trick can learn complex decision boundaries

$$\min_{\mathbf{x}} \quad f(\mathbf{x})$$
$$\text{subject to} \quad g_i(\mathbf{x}) \leq 0, i = 1, ..., m$$
$$h_i(\mathbf{x}) = 0, i = 1, ..., p$$

**strong duality** (e.g., KKT condition)

$$\max_{\nu, \lambda:\lambda_i \geq 0, \forall i} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \nu)$$

**Primal Optimization Problem**

$$p^* = d^*$$

**Dual Optimization Problem** (max-min)

**weak duality**
$$p^* \geq d^*$$

$$\min_{\mathbf{w}, b, \boldsymbol{\xi}} \quad C\sum_{n=1}^{N} \xi^{(n)} + \frac{1}{2}\|\mathbf{w}\|^2$$
$$\text{subject to} \quad y^{(n)} h\left(\mathbf{x}^{(n)}\right) \geq 1 - \xi^{(n)}, \forall n$$
$$\xi^{(n)} \geq 0, \forall n$$

$$\max_{\boldsymbol{\alpha}} \quad \sum_{n=1}^{N} \alpha^{(n)} - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{M} \alpha^{(n)}\alpha^{(m)} y^{(n)} y^{(m)} k\left(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}\right)$$
$$\text{subject to} \quad 0 \leq \alpha^{(n)} \leq C$$
$$\sum_{n=1}^{N} \alpha^{(n)} y^{(n)} = 0$$

## Additional Resource

- Kernel Methods
  – http://www.kernel-machines.org/

- Convex Optimization
  – http://www.stanford.edu/~boyd/cvxbook/
  – http://www.stanford.edu/class/ee364a/
  – see Chapter 5 (and earlier chapters)

## SVM Implementation

- LIBSVM
  – http://www.csie.ntu.edu.tw/~cjlin/libsvm/
  – One of the most popular generic SVM solver (supports nonlinear kernels)
- Liblinear
  – http://www.csie.ntu.edu.tw/~cjlin/liblinear/
  – One of the fastest <u>linear</u> SVM solver (linear kernel)
- SVMlight
  – http://www.cs.cornell.edu/people/tj/svm_light/
  – Structured outputs, various objective measure (e.g., F1, ROC area), Ranking, etc.
- Scikit-learn
  – https://scikit-learn.org/stable/modules/svm.html

# SVM demo code

- http://www.mathworks.com/matlabcentral/fileexchange/28302-svm-demo

- http://www.alivelearn.net/?p=912