

# Convex Optimization and Machine Learning

Mengliu Zhao

Machine Learning Reading Group  
School of Computing Science  
Simon Fraser University

March 12, 2014

# Introduction

Formulation of binary SVM problem:  
Given training data set

$$D = \{(x_i, y_i) | x_i \in R^n, y_i \in \{-1, 1\}, i = 1, 2, \dots, m\} \quad (1)$$

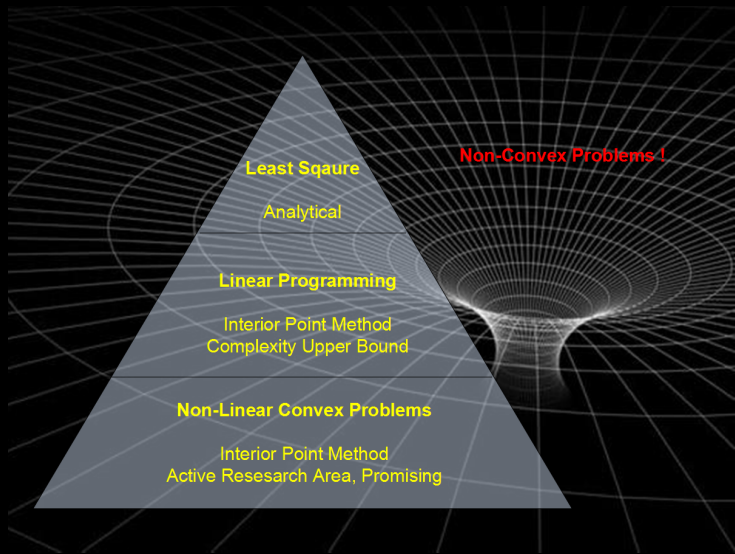
We're trying to find the maximal-margin hyperplane, which can be described by its normal vector  $w$  which satisfies ( $b$  is some offset):

$$\begin{aligned} &\text{minimize} && \|w\|_2 \\ &\text{subject to} && y_i(w x_i - b) \geq 1 \quad i = 1, 2, \dots, m \end{aligned} \quad (2)$$

## Comment

*We encounter a lot of constraint minimization problems in Machine Learning.*

# Why We Want Convex Problems?



# Outline

- ① Lagrange Dual Form
- ② Dual Decomposition, Augmented Lagrangian and ADMM
- ③ SVM and Convex Optimization

# Convex Optimization Problems

General form of convex optimization problem is like following:

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, 2, \dots, m \\ & && h_j(x) = 0, \quad j = 1, 2, \dots, n \end{aligned} \tag{3}$$

where  $f_0, f_i$  are convex functions,  $h_j$  are linear functions.

## Property

*The feasible set of a convex optimization problem is also convex.*

In other words, convex optimization problem is solving a convex function over a convex space.

# General Constraint Problem with Lagrange Duality

However, most constraint problems we optimize are not convex:

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, 2, \dots, m \\ & && h_j(x) = 0, \quad j = 1, 2, \dots, n \end{aligned} \quad (4)$$

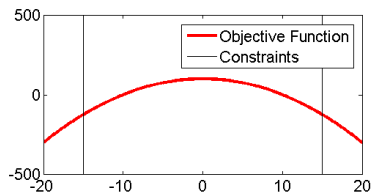
Lagrangian:

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^n \nu_j h_j(x) \quad (5)$$

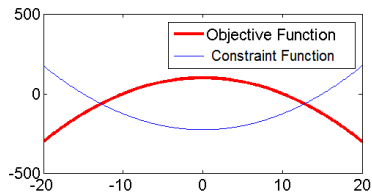
$\lambda_i (> 0)$ ,  $\nu_j$  are called Lagrangian multipliers or dual variables; the Lagrangian dual function is defined as:

$$g(\lambda, \mu) = \inf_x L(x, \lambda, \nu) \quad (6)$$

# Geometric Explanation – Primal Problem

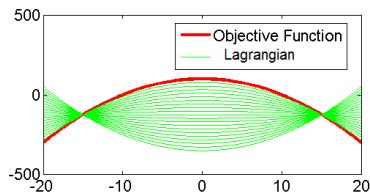


$$\begin{aligned} & \text{minimize} && -x^2 + 15^2 \\ & \text{subject to} && x^2 - 15^2 \leq 0 \end{aligned} \quad (7)$$



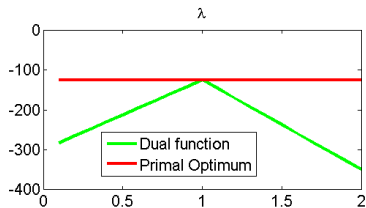
$$(x^2 - 15^2) \quad (8)$$

# Geometric Explanation – Dual Problem



$$-x^2 + 15^2 + \lambda(x^2 - 15^2) \quad (9)$$

$$\lambda = .1 : .1 : 2 \quad (10)$$



$$g(\lambda) = \inf_x \{-x^2 + 15^2 + \lambda(x^2 - 15^2)\} \quad (11)$$
$$\lambda = .1 : .1 : 2$$



## Geometric Explanation – Two Observations

### Observation (I)

*Dual function  $g(\lambda)$  is concave.*

$$\begin{array}{ll} \text{maximize} & g(\lambda) \\ \text{subject to} & \lambda > 0 \end{array} \quad (12)$$

*is a convex optimization problem.*

### Observation (II)

*Let  $p^*$  be the optimal value of the primal problem, then*

$$g(\lambda) \leq p^*, \forall \lambda \quad (13)$$

## Economic Explanation

Company production cost  $f_0$ , with certain limits  $f_i$  below  $a_i$  (rules, resources):

$$\begin{array}{ll} \text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) - a_i \leq 0, \quad i = 1, 2, \dots, m \end{array} \quad (14)$$

However, if , the company can pay a fund rate of  $\lambda_i > 0$  to violate certain rules, which adds back to the total cost:

$$g(\lambda) = \inf_{\lambda} \{f_0(x) + \sum_i \lambda_i (f_i - a_i)\} \quad (15)$$

In this case, the optimal value  $d^*$  for the company is the cost under the least favorable set of prices  $\lambda \rightarrow \max g(\lambda)$ .

## Strong & Weak Duality

How well does the dual problem approximate the original problem?

- 1 Weak Duality: optimal duality gap is always non-negative.

$$p^* - d^* \geq 0 \quad (16)$$

- 2 Strong Duality: duality gap is zero.

$$p^* = d^* \quad (17)$$

Q: When does strong duality hold?

### Theorem (Slater's Theorem)

*D is feasible set. Assume the primal problem is convex:*

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, 2, \dots, m \\ & && h_j(x) = 0, \quad j = 1, 2, \dots, n \end{aligned} \quad (18)$$

*If  $\exists x \in \text{relint } D$ , and  $f_i(x) < 0, i = 0, 1, \dots, m$ , then strong duality holds.*

# KKT Condition

For constrained problem:

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, 2, \dots, m \\ & && h_j(x) = 0, \quad j = 1, 2, \dots, n \end{aligned} \quad (19)$$

If  $x^*$  is the primal minimum, then it satisfies the following necessary condition:

$$\nabla f_0(x^*) + \sum_{i=1}^m \lambda_i \nabla f_i(x^*) + \sum_{j=1}^n \nu_j \nabla h_j(x^*) = 0 \quad (20)$$

# Outline

- ① Lagrange Dual Form
- ② Dual Decomposition, Augmented Lagrangian and ADMM
- ③ SVM and Convex Optimization

## Dual Form, Then What?

Once we get the dual problem, it's easy to solve, e.g., by gradient approach (dual ascent).

### Property

*If  $g(\lambda)$  is a convex (concave) function, then  $\nabla f(\lambda^*) = 0$  iff  $\lambda^*$  is the global minimizer (maximizer).*

### Comment

*A lot of conditions need to be satisfied for a stable gradient method.*

## Dual Ascent for Solving Dual Problem

Let's look at a simplified version of the constrained problem:

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && Ax = b \end{aligned} \tag{21}$$

Its dual form:

$$\max g(\lambda) = \max_{\lambda} \{ \min_x L(x, \lambda) \} \tag{22}$$

$$L(\lambda, x) = f(x) + \lambda(Ax - b) \tag{23}$$

Update  $x, \lambda$  at each iteration:

$$x^{k+1} = \min_x L(x, \lambda^k) \tag{24}$$

$$\lambda^{k+1} = \lambda^k + \alpha^{k+1} \nabla g(x^{k+1}, \lambda^k) \tag{25}$$

### Question

*What if we have a much more complex situation?*

## Dual Decomposition

Suppose the problem is of high dimension,  $\hat{x} = (x, z)$ , and  $f(\hat{x})$  is separable:

$$f(\hat{x}) = f_1(x) + f_2(z) \quad (26)$$

$$A\hat{x} - b = (A_1x - b_1) + (A_2z - b_2) \quad (27)$$

Then we can do dual ascent on each dimension separately:

$$L_1(x) = f_1(x) + \lambda_1(A_1x - b_1) \quad (28)$$

$$L_2(z) = f_2(z) + \lambda_2(A_2z - b_2) \quad (29)$$

$$x^{k+1} = \min_x L_1(x, z^k, \lambda^k) \quad (30)$$

$$z^{k+1} = \min_z L_2(x^{k+1}, z, \lambda^k) \quad (31)$$

$$\lambda^{k+1} = \lambda^k + \alpha^{k+1} \nabla g(x^{k+1}, z^{k+1}, \lambda^k) \quad (32)$$

### Comment

- 1 *Simple dual ascent is usually slow;*



## Alternative to Dual Ascent – Augmented Lagrangian

Primal problem:

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & Ax = b \end{array} \quad (33)$$

Dual problem:

$$L(x, \lambda, \theta) = f(x) + \lambda^T (Ax - b) + \frac{\theta}{2} \|Ax - b\|_2^2 \quad (34)$$

Update by method of multipliers (fixed step):

$$x^{k+1} := \min_x L(x, \lambda^k, \theta) \quad (35)$$

$$\lambda^{k+1} := \lambda^k + \theta(Ax^{k+1} - b) \quad (36)$$

# Method of Multipliers

Comparing to dual ascent:

- ① *Good news*: convergence under more relaxed conditions;
- ② *Bad news*: dual decomposition no longer works (now we have quadratic terms)!

Comment

*ADMM can help!*

# ADMM

Alternating Direction Method of Multipliers

$$\begin{aligned} & \text{minimize} && f(x) + g(z) \\ & \text{subject to} && Ax + Bz = b \end{aligned} \quad (37)$$

Its Lagrangian is:

$$L_{\theta}(x, \lambda, z) = f(x) + g(z) + \lambda^T (Ax + Bz - b) + \frac{\theta}{2} \|Ax + Bz - b\|_2^2 \quad (38)$$

ADMM scheme:

$$\begin{aligned} x^{k+1} & := \min_x L_{\theta}(x, z^k, \lambda^k) \\ z^{k+1} & := \min_z L_{\theta}(x^{k+1}, z, \lambda^k) \\ \lambda^{k+1} & := \lambda^k + \theta(Ax^{k+1} + Bz^{k+1} - b) \end{aligned} \quad (39)$$

## A Closer Look at ADMM

### Comment

*We need more convincing evidence that the scheme will work!*

The thing unnatural here is the new variable  $z$ . We'll check the KKT condition with the constraint problem above:

$$\nabla g(z) + B^T \lambda = 0 \quad (40)$$

We'll check if this could be satisfied by the ADMM scheme. Since  $z^{k+1}$  minimized  $L_\theta(x^{k+1}, z, \lambda^k)$ , then

$$0 = \nabla g(z^{k+1} + B^T \lambda^k + \theta B^T (Ax^{k+1} + Bz^{k+1} - b)) \quad (41)$$

$$= \nabla g(z^{k+1} + B^T \lambda^k) \quad (42)$$

Which means the KKT condition is satisfied.

# Outline

- ① Lagrange Dual Form
- ② Dual Decomposition, Augmented Lagrangian and ADMM
- ③ SVM and Convex Optimization

## Dual Form of SVM

Now let's come back to the constrained version of SVM model:

$$\begin{aligned} & \text{minimize} && \|w\|_2 \\ & \text{subject to} && y_i(wx_i - b) \geq 1 \quad i = 1, 2, \dots, m \end{aligned} \quad (43)$$

It's easy to convert it to Lagrangian dual form as following:

$$\max_{\lambda} \left\{ \min_{w, b} \left\{ \|w\|_2^2 + \sum \lambda_i [1 - y_i(wx_i - b)] \right\} \right\} \quad (44)$$

### Comment

*The formulation is too complex! We can do further to simplify it!*

## Dual Form of SVM

Check KKT condition, taking 1-order derivative of  $w$  and  $b$  on Lagrangian function  $\|w\|_2^2 + \sum \lambda_i [1 - y_i(wx_i - b)]$ :

$$w = \sum_i \lambda_i y_i x_i \quad (45)$$

$$0 = \sum \lambda_i y_i \quad (46)$$

Replace them back in (44), we have:

$$\begin{aligned} \max_{\lambda} g(\lambda) &= \max_{\lambda} \left\{ \sum \lambda_i - \frac{1}{2} \sum y_i y_j \lambda_i \lambda_j (x_i)^T x_j \right\} \\ \text{s.t. } \lambda_i &\geq 0, i = 1, 2, \dots, m \\ \sum \lambda_i y_i &= 0 \end{aligned} \quad (47)$$

# Summary

- ① Lagrangian Duality, KKT condition
- ② Dual Decomposition, Augmented Lagrangian, ADMM
- ③ Example using Lagrangian Duality on SVM



# Thank You!