# 不確実性を考慮した最適化手法

統計数理研究所 / 理化学研究所AIP センター
武田朗子

# 講義の構成

<u>不確実な最適化問題に対する定式化と解法</u>

- ●第１部： ロバスト最適化 (10:30 – 11:20)
- ●第２部： 確率計画法 (11:40 – 12:30)

- ●第３部： ロバスト最適化や確率計画法の機械学習
    問題への適用 (14:00 – 15:00)
- ●第４部： 演習
- ●第５部： 総括

# Mathematical Optimization

It helps to select a best element (with regard to some criteria) from some set of available alternatives.

**Mathematical Optimization Problem:**

$$\min_{x} \quad f(x)$$
$$\text{s.t.} \quad g_i(x) \geq 0, \quad i = 1, \ldots, m$$

- $f(x), g_1(x), \ldots, g_m(x) : \mathbb{R}^n \rightarrow \mathbb{R}$

- If $f(x), g_1(x), \ldots, g_m(x)$ are linear in $x$, the problem is called a linear programming problem.

**math.**

**power engineering**

**Solving a system of polynomial equations**

$$xy = 1$$
$$2i\, xy^2 + y^2 + x = 1$$

**(Linear Programming)**

**Scheduling of generators**



**Optimal size of solar panel**



**(Robust Optimization)**

**Optimization Problem**

**Min:** $f(\boldsymbol{x})$

**subj.to**: $g_1(\boldsymbol{x}) \geq 0$
$g_2(\boldsymbol{x}) \geq 0$
.......

**machine learning**
**(Global Optimization)**



**Support Vector Machine**

**finance**

**portfolio allocation**



**mathematical optimization**

**Solution Method**
• **Nonconvex Opt.**
• **Robust Optimization**

4

# Various Optimization Problems

## Continuous Optimization

Linear Program.

Quadratic Program.

Second Order Cone Program.

Semidefinite Program.

Convex Program.

Nonconvex Quadratic Program.

## Discrete Optimization

Quadratic 0-1 Integer Program.

Linear 0-1 Integer Program.

Linear Integer Program.

Problem name based on Application:

Shortest Path Prob., Travelling Salesman Prob. Knapsack Prob.

# Second-order cone programming

$$\min_{\boldsymbol{x}} \boldsymbol{f}^\top \boldsymbol{x}$$

$$\text{s.t.} \ \|\boldsymbol{A}_i \boldsymbol{x} + \boldsymbol{b}_i\| \le \boldsymbol{c}_i^\top \boldsymbol{x} + d_i, \ \ i = 1, \ldots, m$$

> Euclidean norm
> $$\|\boldsymbol{u}\| = (\boldsymbol{u}^\top \boldsymbol{u})^{1/2}$$

- SOCP can be reformulated as an instance of SDP.
- Convex quadratic programs can also be formulated as SOCPs.
- SOCPs can be solved with great efficiency by interior point methods.

# Optimization Method under Uncertainty

- Robust Optimization

  - ✓ modeling strategies and solution methods for optimization problems that are defined by uncertain inputs

  - ✓ proposed by Ben-Tal & Nemirovski in 1998

- Stochastic Programming

  - ✓ classical framework for modeling optimization problems involving uncertainty (studied since the 1950's).

  - ✓ assuming that probability distributions are known

  - ✓ relation to robust optimization

# Example : Power Generation Planning

T. Electric Company has 2 turbines (Fuel : oil, natural gas).
It wants to determine their production outputs to
**minimize production costs and satisfy electric demands**.

Decision Variable :
$x_i$ : Production Output
[MWh]

Unit Cost （Yen/MWh）

$$\min 135x_1 + 141x_2$$

$$\text{s.t. } x_1 + x_2 \geq 1000$$

Demand

$$L_o \leq x_1 \leq U_o$$

$$L_g \leq x_2 \leq U_g$$

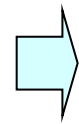Linear Programming: LP
（Simplex Method,
Interior Point Method）

# Formulation of Robust Optimization

Assump.:  uncertain inputs vary within a set (*uncertainty set*).
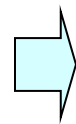**The best decision is done under the *worst-case* scenario.**

uncertainty sets: $\quad \boldsymbol{u}_0 \in \mathcal{U}_0, \ \boldsymbol{u}_i \in \mathcal{U}_i, \forall i$

$$\min_{\boldsymbol{x} \in X} f(\boldsymbol{x}, \boldsymbol{u}_0) \quad \Rightarrow \quad \min_{\boldsymbol{x} \in X} \max_{\boldsymbol{u}_0 \in \mathcal{U}_0} f(\boldsymbol{x}, \boldsymbol{u}_0)$$

$$\text{s.t.} \ g_i(\boldsymbol{x}, \boldsymbol{u}_i) \leq 0, \quad \Rightarrow \quad g_i(\boldsymbol{x}, \boldsymbol{u}_i) \leq 0, \ \forall \boldsymbol{u}_i \in \mathcal{U}_i$$

$$i = 1, \ldots, m \qquad\qquad \longleftrightarrow \ \max_{\boldsymbol{u}_i \in \mathcal{U}_i} g_i(\boldsymbol{x}, \boldsymbol{u}_i) \leq 0$$

# Necessity of Robust Solution

PILOT4 (NETLIB library)

1000 var., 410 const., $\quad x^*$: optimal solution

$a^\top x \equiv$

$-15.79081 x_{826} - 8.598819 x_{827} - 1.88789 x_{828} - 1.362417 \cdots$

$-0.031883 x_{849} - 28.725555 x_{850} - 10.792065 x_{851} - 0.190 \cdots$

$-12.290832 x_{854} + 717.562256 x_{855} - 0.057865 x_{856} - 3.785$

$-122.163055 x_{859} - 6.46609 x_{860} - 0.48371 x_{861} - 0.615264 \cdots$

$-84.644257 x_{864} - 122.459045 x_{865} - 43.15593 x_{866} - 1.712 \cdots$

$+ x_{880} - 0.946049 x_{898} - 0 : 946049 x_{916} \geq \boxed{23.387405} \quad \equiv b$

Change the coeff. $a$ by its 0.1% $\rightarrow \overline{a}$

e.g., $\quad 15.79081 \times 0.001 = 0.0157908$

$x^*$ satisfying $a^\top x^* - b \geq 0$ largely violates the perturbed one:

$$\overline{a}^\top x^* - b < -104.9$$

10

# Applications of Robust Optimization

The obtained solution
- is relatively insensitive to data variations, and
- hedges against catastrophic outcomes.

**Ben-Tal & Nemirovski ['97]**

**Truss topology under the load uncertainties** :
 constructing a building assuming a typical wind load
 → neglecting the possibility of strong wind
 → causing the building to collapse

**Lin, Janak & Floudas ['04]**

**Robust scheduling of chemical processing** :
 scheduling of multiproduct and multipurpose batch plants.
 → neglecting variability of process and environmental data.
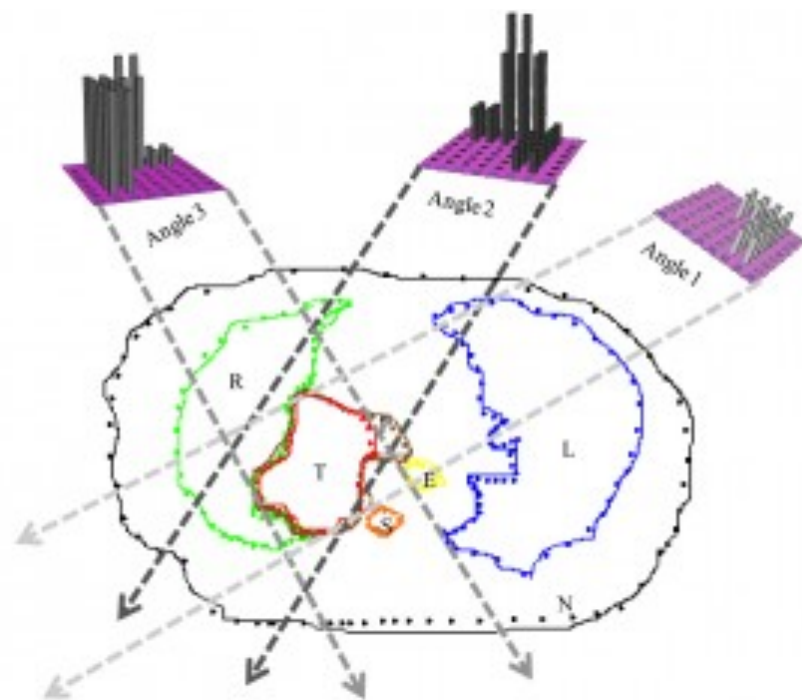 → causing fire and explosion

# Applications to Radio Therapy

**[Radiation Therapy for Cancer Patients]**
T. C. Y. Chan et al. ['06]

Beams of radiation are delivered from different angles around a patient, targeting a tumor in their intersection while trying to spare nearby critical organs.

→ Optimization methods determine the angles of the beams and the intensities of the beamlets, etc.

→ Uncertainty in tumor position (e.g., lung tumors move as the patient breathes during treatment)



http://www.newswise.com/articles/improving-radiation-therapy-for-cancer-patients

# Applications to Solar Energy System
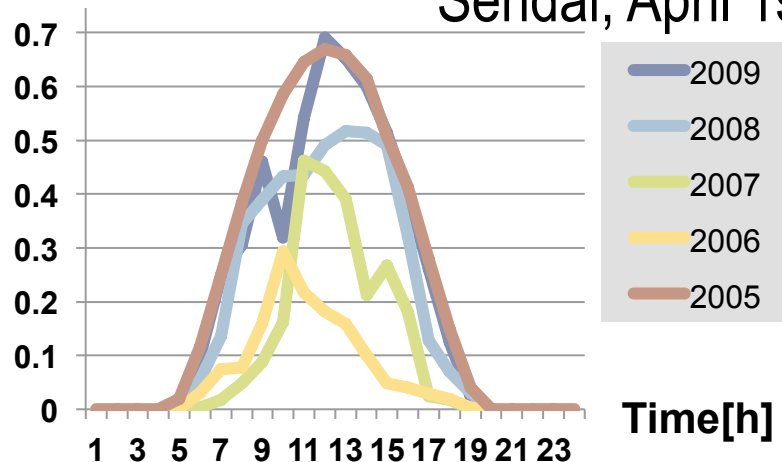
## [Solar Energy System]

Okido & Takeda ['12]

Determining the optimal size of a residential grid-connected solar system to meet a certain CO2 reduction target at a minimum cost.

[project from Japanese local authority]

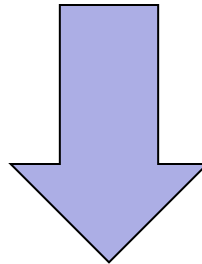→ Useful to determine an amount of subsidy for system owners

→ Taking into consideration uncertainty in the level of solar irradiation (or solar energy) due to weather conditions



Solar energy [kwh/kw] Sendai, April 1st

Legend: 2009, 2008, 2007, 2006, 2005

Time[h]



21

# What is Robust Optimization?

When the data differs from the assumed nominal values,
the generated optimal solution may violate
critical constraints and perform poorly.

**Want to find a solution
immune to data uncertainty.**

Robust optimization:
modeling strategies and solution methods for uncertain problems.

It optimizes against the *worst* instance
that might arise due to uncertain inputs.

# Other Method: Stochastic Programming

Uncertain Optimization Problem:

$$\min_{\boldsymbol{x} \in X} f(\boldsymbol{x}, \boldsymbol{u}_0) \ \text{s.t.} \ g(\boldsymbol{x}, \boldsymbol{u}_1) \le 0$$

$\boldsymbol{u}_0, \boldsymbol{u}_1$ : uncertain data

● Stochastic Programming      **Dantzig ['55], Beale ['55]**

<span style="color:red">Assump.</span> :

Prob. distributions of $\boldsymbol{u}_0, \boldsymbol{u}_1$ are known.

$p(\boldsymbol{u}_0)$

density func.

ex.1) $\displaystyle \min_{\boldsymbol{x} \in X} E_{\boldsymbol{u}_0}[f(\boldsymbol{x}, \boldsymbol{u}_0)]$

  s.t. $E_{\boldsymbol{u}_1}[g(\boldsymbol{x}, \boldsymbol{u}_1)] \le 0$

$\boldsymbol{u}_0$

ex.2)  Chance Const.（Probabilistic Const.）    **Charnes & Cooper ['59]**

$$\Pr_{\boldsymbol{u}_1}(g(\boldsymbol{x}, \boldsymbol{u}_1) \le 0) \ge 1 - \epsilon$$

# Other Method: Sensitivity Analysis

Uncertain Optimization Problem :

$$\min_{\boldsymbol{x} \in X} f(\boldsymbol{x}, \boldsymbol{u}_0) \text{ s.t. } g(\boldsymbol{x}, \boldsymbol{u}_1) \leq 0$$

$\boldsymbol{u}_0, \boldsymbol{u}_1$ : uncertain data

- Post-optimal analysis after obtaining an optimal solution for some $\boldsymbol{u}_0, \boldsymbol{u}_1$.
- It shows whether the optimal solution changes for the data perturbation.

$\boldsymbol{u}_0^{\top} \boldsymbol{x}$

Restrictions: data of objective func. & RHS of LP can be uncertain

# History of Robust Optimization

Robust Optimization:
$$\min_{\boldsymbol{x} \in X} \ \max_{\boldsymbol{u}_0 \in \mathcal{U}_0} f(\boldsymbol{x}, \boldsymbol{u}_0)$$
$$\text{s.t.} \ \max_{\boldsymbol{u}_i \in \mathcal{U}_i} g_i(\boldsymbol{x}, \boldsymbol{u}_i) \leq 0, \ \forall i$$

■ In 1973, A.L.Soyster proposed "inexact LP" using rectangular $\mathcal{U}$ .

Almost no progress (two papers†)

†) **reported by Ben-Tal, El Ghaoui & Nemirovski ['09]**

$u_2$  $\mathcal{U}$ : Rect.

$u_1$

■ In 1998, Ben-Tal & Nemirovski proposed "robust optimization" using ellipsoidal $\mathcal{U}$ .

$u_2$  $\mathcal{U}$ : Ellips.

■ Studies on robust optimization are going on …

$u_1$

# Why robust optimization became popular?

① Inexact LP (=Robust LP with rectangle $\mathcal{U}$ ) only assumes extreme situations. This drawback was solved by ellipsoidal $\mathcal{U}$.

② Resulting in a second-order cone programming (SOCP), semidefinite programming (SDP).

| Inexact LP **Soyster ['73]** | Robust LP **Ben-Tal & Nemirovski ['98]** |

$\mathcal{U}$ is a rectangle $\longrightarrow$    $\mathcal{U}$ is an ellipsoid, etc….

**Extreme situations**

# Various Research Directions

Original Form of Robust Opt.
**Soyster ['73]**

第2部

Stochastic Approach
**Calafiore & Campi ['05,'06]**

第3部

Establishment of Robust Opt.
**Ben-Tal & Nemirovski ['98,'99]**

Application to Finance,
Machine learning,
Energy System, etc.

Conditions for Tractable Robust Opt.
**Goldfarb & Iyengar ['03]**

Extension to Multi-period Model
**Ben-Tal, et.al. ['04]**

第1部：残り時間で．．

19

# Difficult to Be Solved in General

$$\min_{\boldsymbol{x}} \quad -x_1 + x_2$$

infinite number of constraints

$$\text{s.t.} \quad -1 \leq \quad u_1 x_1 + u_2 x_2 \leq 1$$

$$-1 \leq -u_2 x_1 + u_1 x_2 \leq 1$$

$$\forall \boldsymbol{u} \in \mathcal{U} = \left\{ (u_1, u_2) \mid u_1^2 + u_2^2 = 1 \right\}$$

Feasible region at $(u_1, u_2) = (1, 0)$

The optimal value of robust optimization problem

Objective function

min

The optimal value of the deterministic problem with $(u_1, u_2) = (1, 0)$

One research direction:

**Want to define $\mathcal{U}$ so that the RO problem is tractable.**

# Standard Form for Robust Optimization

$$\min_{\boldsymbol{x} \in X} \quad \boldsymbol{c}^\top \boldsymbol{x} \quad \text{s.t.} \quad f_i(\boldsymbol{x}, \boldsymbol{u}_i) \leq 0, \quad \forall \boldsymbol{u}_i \in \mathcal{U}_i,$$
$$i = 1, \ldots, m$$

- <u>Constraint-wise uncertainty</u> is assumed.
- $f_i(\boldsymbol{x}, \boldsymbol{u}_i)$ : convex in $\boldsymbol{x}$ $\quad$ ( $\forall \boldsymbol{u}_i \in \mathcal{U}_i$ )
- $X$ : closed convex set, $\quad$ $\mathcal{U}_i$ : bounded closed set

- **When the objective function is uncertain**

$$\min_{\boldsymbol{x} \in X} \max_{\boldsymbol{u}_0 \in \mathcal{U}_0} f_0(\boldsymbol{x}, \boldsymbol{u}_0)$$

$$\Longrightarrow \quad \min_{\boldsymbol{x} \in X, t} t \quad \text{s.t.} \quad f_0(\boldsymbol{x}, \boldsymbol{u}_0) \leq t, \ \forall \boldsymbol{u}_0 \in \mathcal{U}_0$$

# Tractable Robust LP (Ellipsoidal Case)

$$\min_{x} \ c^\top x \ \text{ s.t. } \ a(u)^\top x \le b, \quad \forall u \in \mathcal{U}$$

**Ellipsoidal uncertainty set:**

$$a(u) = a_0 + Au \quad \mathcal{U} = \{ \ u : \|u\|_2 \le 1\}$$

$$\min_{x} \ c^\top x \ \text{ s.t. } \quad a_0^\top x + x^\top Au \le b, \ \|u\|_2 \le 1,$$

$$a_0^\top x + \left( \max_{u:\|u\|_2 \le 1} x^\top Au \right) \le b$$

$$u^* = \frac{A^\top x}{\|A^\top x\|_2}$$

$$\min_{x} \ c^\top x \ \text{ s.t. } \quad a_0^\top x + \|A^\top x\|_2 \le b$$

**Second order cone programming (SOCP)**

22

# Tractable Robust LP (Rectangle Case)

$$\min_{x} \ c^\top x \ \text{ s.t. } \ a^\top x \le b, \quad \forall a \in \mathcal{U}$$

**Rectangle :**

$$\mathcal{U} = \{u : a_0 - \bar{a} \le u \le a_0 + \bar{a}\} \subset R^n$$
where $\bar{a} \ge 0$

$$\max_{a \in \mathcal{U}} a^\top x = a_0^\top x + \bar{a}^\top |x| \le b$$

A vector constructed by taking absolute values for each element of $x$

**Linear Programming Problem**

$$\min_{x,y} \ c^\top x$$
$$\text{s.t. } \ a_0^\top x + \bar{a}^\top y \le b, \quad -y \le x \le y, \quad y \ge 0$$

23

# Conditions for Tractable Robust Optimization

$$\min_{\boldsymbol{x} \in X} \quad \boldsymbol{c}^\top \boldsymbol{x} \quad \text{s.t.} \quad f_i(\boldsymbol{x}, \boldsymbol{u}_i) \leq 0, \quad \forall \boldsymbol{u}_i \in \mathcal{U}_i,$$

$$i = 1, \ldots, m$$

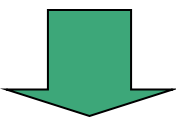➡ Want to transform it to a tractable convex prob.

**Ben-Tal & Nemirovski ['98], Goldfarb & Iyengar ['03]**

**Three Assumptions**

(1) $f(\boldsymbol{x}, \boldsymbol{u})$ is convex quadratic in terms of $\boldsymbol{x}$.

$$f(\boldsymbol{x}, \boldsymbol{u}) = \boldsymbol{x}^\top \boldsymbol{Q}(\boldsymbol{u}) \boldsymbol{x} + \boldsymbol{q}(\boldsymbol{u})^\top \boldsymbol{x} + \gamma(\boldsymbol{u})$$

(2) Uncertain data is linear w.r.t $\boldsymbol{u}$ .

$$\boldsymbol{Q}(\boldsymbol{u}) = \boldsymbol{Q}_0 + \sum_i \boldsymbol{Q}_i u_i$$

$$\boldsymbol{q}(\boldsymbol{u}) = \boldsymbol{q}_0 + \sum_i \boldsymbol{q}_i u_i$$

(3) $\mathcal{U}_i$ is a finite set, its convex hull or ellipsoid.

# Difficulty of Solving Problems

Assump.



✓ $\mathcal{U}$ is an ellipsoidal uncertainty set

✓ Uncertain data is linear with respect to $\boldsymbol{u} \in \mathcal{U}$

$$a(\boldsymbol{u}) = \boldsymbol{a}_0 + \sum_i \boldsymbol{a}_i u_i, \quad \boldsymbol{F}(\boldsymbol{u}) = \boldsymbol{F}_0 + \sum_i \boldsymbol{F}_i u_i$$

● Robust LP → Second-order Cone Programming（SOCP）

● Robust SOCP → Semidefinite Programming（SDP）

● Robust SDP → ×

Approximately solved by SDP

# Tips on Formulation of Robust Optimization

With robust optimization ... ..

✓ How to express uncertainty data is important!

✓ There is a great limitation on its expression

- Uncertainty data is linear w.r.t $u$ .

- The range for $u$ is an ellipse, etc.

If these conditions are satisfied, a RO problem can be converted to a tractable problem.

In the case where the condition is not satisfied

⇒ stochastic approach by sampling a finite number of constraints among infinitely many constraints

26

# Contents

●Robust Optimization

  ✓modeling strategies and solution methods for optimization problems that are defined by uncertain inputs

  ✓proposed by Ben-Tal & Nemirovski in 1998


●Stochastic Programming

  ✓classical framework for modeling optimization problems involving uncertainty (studied since the 1950's).

  ✓assuming that probability distributions are known

  ✓relation to robust optimization

# Stochastic Programming

Uncertain Optimization Problem:

$$\min_{\boldsymbol{x} \in X} f(\boldsymbol{x}, \boldsymbol{u}_0) \text{ s.t. } g(\boldsymbol{x}, \boldsymbol{u}_1) \leq 0$$

$\boldsymbol{u}_0, \boldsymbol{u}_1$ : uncertain data

● Stochastic Programming **Dantzig ['55], Beale ['55]**

Assump. :

Prob. distributions of $\boldsymbol{u}_0, \boldsymbol{u}_1$ are known.

ex.1)
$$\min_{\boldsymbol{x} \in X} E_{\boldsymbol{u}_0}[f(\boldsymbol{x}, \boldsymbol{u}_0)]$$
$$\text{s.t. } E_{\boldsymbol{u}_1}[g(\boldsymbol{x}, \boldsymbol{u}_1)] \leq 0$$

$p(\boldsymbol{u}_0)$

density func.

$\boldsymbol{u}_0$

ex.2)  Chance Const.（Probabilistic Const.） **Charnes & Cooper ['59]**

$$\Pr_{\boldsymbol{u}_1}(g(\boldsymbol{x}, \boldsymbol{u}_1) \leq 0) \geq 1 - \epsilon$$

28

# Examples of Another Risk Measure

Instead of "Expectation", risk measure "CVaR" is often used.

$$\min_{\boldsymbol{x}\in X} E_{\boldsymbol{u}}[f(\boldsymbol{x}, \boldsymbol{u})] \quad \Longrightarrow \quad \min_{\boldsymbol{x}\in X} \phi_\beta(\boldsymbol{x}) \qquad \beta \in (0,1)$$

**CVaR** (**C**onditional **V**alue-**a**t-**R**isk) : $\phi_\beta(\boldsymbol{x})$

Conditional expectation of $f(\boldsymbol{x}, \boldsymbol{u})$ exceeding $\beta$-quantile $\alpha_\beta(\boldsymbol{x})$



density function

**High Risk**

$p(\boldsymbol{u})$  mean

cdf

$\beta = 0.8$

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

$\beta$-quantile (VaR): $\alpha_\beta(\boldsymbol{x})$ $\phi_\beta(\boldsymbol{x})$

$\beta$-quantile (VaR): $\alpha_\beta(\boldsymbol{x})$

29

# Definition of Conditional Value-at-Risk (CVaR)

## Rockafellar & Uryasev ['02]

$\beta \in (0, 1)$

$\alpha_\beta(\boldsymbol{x})$ : $\beta$-VaR (= $\beta$-quantile)

$\phi_\beta(\boldsymbol{x})$ : $\beta$-CVaR

of the loss $f(\boldsymbol{x}, \boldsymbol{u})$ associated with a decision $\boldsymbol{x}$

random vec.

Conditional expectation of $f(\boldsymbol{x}, \boldsymbol{u})$ exceeding $\beta$-quantile $\alpha_\beta(\boldsymbol{x})$

$$\phi_\beta(\boldsymbol{x}) = \frac{1}{1-\beta} \int_{f(\boldsymbol{x},\boldsymbol{u}) \geq \alpha_\beta(\boldsymbol{x})} f(\boldsymbol{x}, \boldsymbol{u}) p(\boldsymbol{u}) \, \mathrm{d}\boldsymbol{u}$$

density function

$$\alpha_\beta(\boldsymbol{x}) \in \arg \min_\alpha F_\beta(\boldsymbol{x}, \alpha)$$

$$\phi_\beta(\boldsymbol{x}) = \min_\alpha F_\beta(\boldsymbol{x}, \alpha)$$

$$F_\beta(\boldsymbol{x}, \alpha) := \alpha + \frac{1}{1-\beta} \int_{\boldsymbol{u}} [f(\boldsymbol{x}, \boldsymbol{u}) - \alpha]^+ p(\boldsymbol{u}) \, \mathrm{d}\boldsymbol{u}$$



$p(\boldsymbol{u})$

mean

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

$\alpha_\beta(\boldsymbol{x})$  $\phi_\beta(\boldsymbol{x})$

# CVaR for Discrete Distribution

When random variables follow a discrete dist. or normal dist., CVaR minimization can be tractable.

ex.) For some $\beta \in (0,1)$ and $\boldsymbol{x}$,

**Rockafellar & Uryasev ['02]**

$$\phi_\beta(\boldsymbol{x}) = \min_\alpha \; \alpha + \frac{1}{1-\beta} \sum_{i=1}^{N} p_i [f(\boldsymbol{x}, \boldsymbol{u}_i) - \alpha]^+$$

opt.sol: $\alpha^* \approx \alpha_\beta(\boldsymbol{x})$

Histogram of $f(\boldsymbol{x}, \boldsymbol{u}_i)$, $i = 1, 2, \ldots N$

frequency

**High Risk**

$\beta$

$\phi_\beta(\boldsymbol{x})$

For the finite support:
$$\mathcal{U} = \{\boldsymbol{u}_1, \ldots, \boldsymbol{u}_N\}$$
$$\Pr(\boldsymbol{u} = \boldsymbol{u}_i) = p_i$$

$\beta$-quantile (VaR): $\alpha_\beta(\boldsymbol{x})$

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

# Tractable Form for CVaR Minimization

$$\min_{\boldsymbol{x}\in X} \phi_\beta(\boldsymbol{x})$$

$$\Rightarrow \min_{\boldsymbol{x}\in X,\alpha} \alpha + \frac{1}{1-\beta}\sum_{i=1}^{N} p_i[f(\boldsymbol{x},\boldsymbol{u}_i)-\alpha]^+$$

$$\Rightarrow \min_{\boldsymbol{z},\boldsymbol{x},\alpha} \alpha + \frac{1}{1-\beta}\sum_{i=1}^{N} p_i z_i$$

$$\text{s.t. } f(\boldsymbol{x},\boldsymbol{u}_i)-\alpha-z_i \le 0, \ \forall i$$

$$\boldsymbol{z}\ge \boldsymbol{0}, \ \boldsymbol{x}\in X$$

If $f(\boldsymbol{x},\boldsymbol{u}_i)$ is convex in $\boldsymbol{x}$ and $X$ is a convex set, this is a convex optimization prob.

# Parameter $\beta$ of CVaR

$\beta \in (0, 1)$

**CVaR** (**C**onditional **V**alue-**a**t-**R**isk) : $\phi_\beta(\boldsymbol{x})$

Conditional expectation of $f(\boldsymbol{x}, \boldsymbol{u})$ exceeding $\beta$-quantile $\alpha_\beta(\boldsymbol{x})$

$$\min_{\boldsymbol{x} \in X} \phi_\beta(\boldsymbol{x})$$

$\beta \to 0$

$$\min_{\boldsymbol{x} \in X} E_{\boldsymbol{u}}[f(\boldsymbol{x}, \boldsymbol{u})]$$

traditional stochastic program.

Histogram of $f(\boldsymbol{x}, \boldsymbol{u}_i)$, $i = 1, 2, \ldots N$

frequency

**High Risk**

$\beta$

$\phi_\beta(\boldsymbol{x})$

$\beta$-quantile: $\alpha_\beta(\boldsymbol{x})$

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

$\beta \to 1$

$$\min_{\boldsymbol{x} \in X} \max_{i=1,\ldots,N} f(x, u_i)$$

robust optimization

33

# CVaR for Normal Distribution

$\beta \in (0, 1)$

**CVaR** (**C**onditional **V**alue-**at**-**R**isk) : $\phi_\beta(\boldsymbol{x})$

Conditional expectation of $f(\boldsymbol{x}, \boldsymbol{u})$ exceeding $\beta$-quantile $\alpha_\beta(\boldsymbol{x})$

Random variable: $\boldsymbol{u} \sim \mathcal{N}_n(\bar{\boldsymbol{u}}, \Sigma)$

$$L = \boldsymbol{u}^\top \boldsymbol{x} \sim \mathcal{N}(\underbrace{\bar{\boldsymbol{u}}^\top \boldsymbol{x}}_{\mu}, \underbrace{\boldsymbol{x}^\top \Sigma \boldsymbol{x}}_{\sigma^2})$$

Probability density of the normal dist. :

$$p(L) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(L-\mu)^2}{2\sigma^2}\right)$$

density func.

$p(L)$

**High Risk**

mean

$L$

$\beta$-VaR: $\alpha_\beta(\boldsymbol{x})$ $\phi_\beta(\boldsymbol{x})$

CVaR is defined as $\phi_\beta(\boldsymbol{x}) = \dfrac{1}{1-\beta} \displaystyle\int_{\alpha_\beta}^{\infty} L \cdot p(L)\, \mathrm{d}L$

$$= \bar{\boldsymbol{u}}^\top \boldsymbol{x} + C\sqrt{\boldsymbol{x}^\top \Sigma \boldsymbol{x}}$$

# Probabilistic Constraint

$$\Pr_{\boldsymbol{u}}(g(\boldsymbol{x}, \boldsymbol{u}) \leq 0) \geq 1 - \epsilon$$

ex.)

$$\Pr_{\boldsymbol{u}}(\boldsymbol{u}^\top \boldsymbol{x} \leq b) \geq \eta \quad (\text{ただし}, \eta \geq 0.5)$$

Under the assump.: $\boldsymbol{u} \sim \mathcal{N}_n(\bar{\boldsymbol{u}}, \Sigma)$

$$\Longleftrightarrow \quad \Pr_{\boldsymbol{u}}\left( \underbrace{\frac{\boldsymbol{u}^\top \boldsymbol{x} - \bar{\boldsymbol{u}}^\top \boldsymbol{x}}{\sqrt{\boldsymbol{x}^\top \Sigma \boldsymbol{x}}}}_{\sim \mathcal{N}(0,1)} \leq \frac{b - \bar{\boldsymbol{u}}^\top \boldsymbol{x}}{\sqrt{\boldsymbol{x}^\top \Sigma \boldsymbol{x}}} \right) \geq \eta$$

$\Phi(z) = \eta$

$z$

$$\Longleftrightarrow \quad \frac{b - \bar{\boldsymbol{u}}^\top \boldsymbol{x}}{\sqrt{\boldsymbol{x}^\top \Sigma \boldsymbol{x}}} \geq \Phi^{-1}(\eta)$$

$: \eta\text{-quantile}$

$\Phi(z)$ : cumulative dist.
func. (cdf) of $\mathcal{N}(0,1)$

second-order cone constr.

$$\Longleftrightarrow \quad \bar{\boldsymbol{u}}^\top \boldsymbol{x} + \Phi^{-1}(\eta) \|\Sigma^{1/2} \boldsymbol{x}\| \leq b$$

35

# Relation to Robust Constraint

Probabilistic Const.

Assump.: $\boldsymbol{u} \sim \mathcal{N}_n(\bar{\boldsymbol{u}}, \Sigma)$

$$\text{Pr}_{\boldsymbol{u}}(\boldsymbol{u}^\top \boldsymbol{x} \leq b) \geq \eta \quad \Longleftrightarrow \quad \bar{\boldsymbol{u}}^\top \boldsymbol{x} + \Phi^{-1}(\eta)\|\Sigma^{1/2}\boldsymbol{x}\| \leq b$$

Robust Const.

Assump.: $\boldsymbol{u} \in \mathcal{U} := \{\bar{\boldsymbol{u}} + \Sigma^{1/2}\boldsymbol{v} : \|\boldsymbol{v}\| \leq \Phi^{-1}(\eta)\}$

$$\max_{\boldsymbol{u} \in \mathcal{U}} \boldsymbol{u}^\top \boldsymbol{x} \leq b$$

$$\Longleftrightarrow \quad \bar{\boldsymbol{u}}^\top \boldsymbol{x} + \max_{\boldsymbol{v} : \|\boldsymbol{v}\| \leq \Phi^{-1}(\eta)} \boldsymbol{x}^\top \Sigma^{1/2} \boldsymbol{v} \leq b$$

$$= \bar{\boldsymbol{u}}^\top \boldsymbol{x} + \Phi^{-1}(\eta)\|\Sigma^{1/2}\boldsymbol{x}\|$$

# Stochastic Interpretation for Uncertainty Set

Assump.: $\boldsymbol{u} \sim \mathcal{N}_n(\bar{\boldsymbol{u}}, \Sigma)$

$\mathrm{Pr}_{\boldsymbol{u}}(\boldsymbol{u}^{\top}\boldsymbol{x} \leq b) \geq \eta$

Relation of two Asumptions?

Assump.: $\boldsymbol{u} \in \mathcal{U} := \{\bar{\boldsymbol{u}} + \Sigma^{1/2}\boldsymbol{v} : \|\boldsymbol{v}\| \leq \Phi^{-1}(\eta)\}$

$\max_{\boldsymbol{u} \in \mathcal{U}} \boldsymbol{u}^{\top}\boldsymbol{x} \leq b$

$\mathrm{Pr}(\boldsymbol{u} \in \mathcal{U}) = \mathcal{F}_n((\Phi^{-1}(\eta))^2)$

chi-squared distribution with $n$ degrees of freedom

$n$=2

density

$\bar{\boldsymbol{u}}$

$u^2$

$u^1$

$100\mathcal{F}_n((\Phi^{-1}(\eta))^2)$%
data are covered

support for truncated normal dist.

37

# Two Optimization Methods under Uncertainty

$$\min_{\boldsymbol{x} \in X} f(\boldsymbol{x}, \boldsymbol{u}_0) \text{ s.t. } g(\boldsymbol{x}, \boldsymbol{u}_1) \leq 0$$

$\boldsymbol{u}_0, \boldsymbol{u}_1$

: uncertain data

Probabilistic Const.

Assump.: $\boldsymbol{u} \sim \mathcal{N}_n(\bar{u}, \Sigma)$

Robust Const.

Assump.: $\boldsymbol{u} \in \mathcal{U} := \{\bar{u} + \Sigma^{1/2}\boldsymbol{v} : \|\boldsymbol{v}\| \leq \Phi^{-1}(\eta)\}$

Boundary between two methods is getting blurred.

Recently, studies on robust optimization using "probability" are increased e.g. for setting the uncertainty set $\mathcal{U}$.

# Stochastic Approach for Robust Optimization

Among three assumptions for tractable robust optimization,
  (2) Uncertain data is linear w.r.t $\mathcal{u}$
  (3) $\mathcal{U}$ is a finite set, its convex hull or ellipsoid
can be removed.

$\boldsymbol{u}_1, \cdots, \boldsymbol{u}_N$ : randomly generated following the distribution on $\mathcal{U}$

Solve a relaxation problem having a finite number of const.

**Calafiore & Campi ['05]**

Want to estimate the sample size $N$ to obtain
a relaxed solution with theoretical guarantee.

# How to determine the sample size *N*

$$\boldsymbol{u}_1, \ldots, \boldsymbol{u}_N \overset{\text{i.i.d.}}{\sim} P \quad \text{(Assume the probability distribution on } \mathcal{U})$$

Randomly generated relaxation problem (SCP$_N$) :

$$\min_{\boldsymbol{x} \in X} \boldsymbol{c}^\top \boldsymbol{x} \quad \text{s.t.} \quad f(\boldsymbol{x}, \boldsymbol{u}_i) \leq 0, \quad i = 1, \ldots, N$$

Optimal sol. of (SCP$_N$) : $\widehat{\boldsymbol{x}}_N$

feasible set of robust opt.

Criteria for deciding $N$ :

- Allow $\widehat{\boldsymbol{x}}_N$ to violate some ratio of constraints:

$$V(\widehat{\boldsymbol{x}}_N) = P\{\boldsymbol{u} \in \mathcal{U} : f(\widehat{\boldsymbol{x}}_N, \boldsymbol{u}) > 0\} \leq \epsilon_1$$

min

**Calafiore & Campi ['05, '06]**

- Allow some amount of constraint violation for $\widehat{\boldsymbol{x}}_N$ :

$$\max_{\boldsymbol{u} \in \mathcal{U}} f(\widehat{\boldsymbol{x}}_N, \boldsymbol{u}) \leq \epsilon_2$$

**Kanamori & Takeda ['12]**

# Evaluation for Sample Size

$$N(\epsilon, \eta) := \frac{2}{\epsilon} \log \frac{1}{\eta} + 2n + \frac{2n}{\epsilon} \log \frac{2}{\epsilon}$$

**Calafiore & Campi ['06]**

$$N(\epsilon, \eta) := \min \left\{ N \in \mathbb{N} \mid \sum_{i=0}^{n-1} \binom{N}{i} \epsilon^i (1-\epsilon)^{N-i} \leq \eta \right\}$$

**Campi & Garatti ['08]**

Theo. （**Calafiore & Campi ['05,'06], Campi & Garatti ['08]**）

Let $\epsilon \in (0, 1), \ \eta \in (0, 1)$.

The optimal solution $\widehat{\boldsymbol{x}}_N \in R^n$ of (SCP$_N$) generated with $N \geq N(\epsilon, \eta)$ samples satisfies $V(\widehat{\boldsymbol{x}}_N) \leq \epsilon$ with the probability at least $1 - \eta$, that is,

$$P^N \{ \ V(\widehat{\boldsymbol{x}}_N) \leq \epsilon \ \} \geq 1 - \eta$$

$\epsilon \to 0, \eta \to 0 \quad \Rightarrow \quad N(\epsilon, \eta) \to \infty$

Violation probability: $V(\widehat{\boldsymbol{x}}_N) = P\{\boldsymbol{u} \in \mathcal{U} : f(\widehat{\boldsymbol{x}}_N, \boldsymbol{u}) > 0\}$

# A-priori / A-posteriori Evaluations

( A-priori Evaluation )

**Takeda, Taguchi & Tanaka ['10]**

(construction of function *q is* a key idea)

$$\epsilon \in (0, q(B)), \ \eta \in (0, 1), \ N \geq N(\epsilon, \eta)$$

$\Longrightarrow$ Optimal sol. $\widehat{\boldsymbol{x}}_N$ of (SCP$_N$) satisfies

$$P^N \{ \ V(\widehat{\boldsymbol{x}}_N) \leq \epsilon, \ \max_{\boldsymbol{u} \in \mathcal{U}} f(\widehat{\boldsymbol{x}}_N, \boldsymbol{u}) \leq q^{-1}(\epsilon) \} \geq 1 - \eta$$

independent from $\widehat{\boldsymbol{x}}_N$

( A-posteriori Evaluation )

Optimal sol. $\widehat{\boldsymbol{x}}_N$ of (SCP$_N$), $N > 0$, satisfies

$$\delta \in (0, B], \ \eta \in (0, 1), \ M \geq M(\delta, \eta) := \frac{\ln \eta}{\ln(1 - q(\delta))},$$

$$\tilde{\boldsymbol{u}}_1, \ldots, \tilde{\boldsymbol{u}}_M \ \sim \ P,$$

depending on $\widehat{\boldsymbol{x}}_N$

$$P^M \{ \max_{\boldsymbol{u} \in \mathcal{U}} f(\widehat{\boldsymbol{x}}_N, \boldsymbol{u}) < \max_{i=1,\ldots,M} f(\widehat{\boldsymbol{x}}_N, \tilde{\boldsymbol{u}}_i) + \delta \} \geq 1 - \eta$$

# Various Research Directions



43

# ロバスト最適化や確率計画法の機械学習問題への適用

統計数理研究所 / 理化学研究所AIP センター

武田朗子

# Optimization Techniques in ML

● There are trends in optimization techniques used in ML
- ✓ semidefinite program
- ✓ submodular optimization
- ✓ first-order methods such as APG, ADMM, etc.

● Stochastic Program. and Robust Optimization are not popular in ML
- ✓ but they are implicitly used.

# Contents

● Provide a view based on Robust Optimization for various
   Binary Classification Models including
   ✓ Support Vector Machine (SVM),
     Minimax Probability Machine (MPM) and
     Fisher Discriminant Analysis (FDA), etc.

● Provide a view based on Stochastic Programming
   ✓ $\nu$-SVM & E$\nu$-SVM ➜ Generalization Bound
   ✓ Minimum Margin MPM

# Application of Robust Optimization to ML

✓ Introducing the work of Xu, Caramanis and Mannor [2009]

✓ Showing a unified view for various ML models such as SVM MPM, FDA, logistic regression.

We use robust optimization techniques in a different problem setting

# Binary Classification Problem

extendable to nonlinear one using kernel

Find a decision function $f(\boldsymbol{x}) = \widehat{\boldsymbol{w}}^{\top}\boldsymbol{x} + \widehat{b}$

based on given training samples $(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_m, y_m)$

to correctly classify new samples.

EX.) diagnosis of diabetes

$\boldsymbol{x}_i \in R^n$ ⟵ medical examination

$y_i \in \{\pm 1\}$ ⟵ tested positive/negative

$i \in M := \{1, 2, \ldots, m\}$

insulin

$(n = 2)$

$\widehat{\boldsymbol{w}}^{\top}\boldsymbol{x}_i + \widehat{b} < 0$

$y = 1$

$\boldsymbol{x}_j$

$\boldsymbol{x}_i$

**Label??**

$y = -1$

$\widehat{\boldsymbol{w}}^{\top}\boldsymbol{x}_i + \widehat{b} > 0$

blood pressure

# Hard margin SVM (support vector machine)

**Boser, Guyon & Vapnik ['92]**

**Linearly Separable**

$$\boldsymbol{w}^\top \boldsymbol{x}_i + b < 0$$

$$\boldsymbol{w}^\top \boldsymbol{x}_i + b > 0$$

$\boldsymbol{x}_2$

$\boldsymbol{x}_1$

$\boldsymbol{x}_4$

$\boldsymbol{x}_3$

$$\boldsymbol{w}^\top \boldsymbol{x} + b = 0$$

$y = 1$

$y = -1$

Maximize the minimum distance to the hyperplane

$= 1$

$$\max_{w \neq 0, b} \quad \min_{i=1,\ldots,m} \frac{y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b)}{\|\boldsymbol{w}\|}$$

regularization penalty

$$\min_{w,b} \quad \frac{1}{2}\|\boldsymbol{w}\|^2$$

$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq 1$$

$$i = 1, \ldots, m$$

Minimizing a regularization penalty enhances generalization performance (prediction accuracy for test dataset)

# *C*-SVM

$$\min_{\boldsymbol{w},b,\boldsymbol{z}} \quad \frac{1}{2}\|\boldsymbol{w}\|^2 + C\sum_{i=1}^{m} z_i$$

$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq 1 - z_i, \quad (i \in M)$$

$$\boldsymbol{z} \geq \boldsymbol{0}$$

penalized samples

$u\boldsymbol{w}^\top \boldsymbol{x}_i + b < -1$

$y = -1$

$y = 1$

$\boldsymbol{w}^\top \boldsymbol{x}_i + b > 1$

**Margin = 1**

Two conflicting goals

{ ● minimizing training error

● minimizing a regularization penalty

- the trade-off between these goals
  is controlled by *C*

# $\nu$-SVM

**Scholkopf, Smola, Williamson & Bartlett ['00]**

penalized samples



$\boldsymbol{w}^\top \boldsymbol{x}_i + b < -\rho$

$y = -1$

SVs

$y = 1$

$\boldsymbol{w}^\top \boldsymbol{x}_i + b > \rho$

**Margin =** $\rho^*$

$C$ is replaced by an intuitive parameter $\nu$

$$\min_{\boldsymbol{w},b,\boldsymbol{z},\rho} \quad \frac{1}{2}\|\boldsymbol{w}\|^2 - \nu\rho + \frac{1}{m}\sum_{i=1}^{m} z_i$$

$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq \rho - z_i \quad (i \in M)$$

$$\boldsymbol{z} \geq \boldsymbol{0}$$

> C-SVM with $C = \dfrac{1}{m\rho^*} \longleftrightarrow \nu$-SVM

> margin is nonnegative : $\rho^* \geq 0$

> admissible values of $\nu$ are limited
  $$\left( \nu \in (\ \nu_{\min}, \nu_{\max}\ ] \subseteq (0,1] \right)$$

> **0** opt. solution for small $\nu$

# Extended $v$-SVM (E$v$-SVM)

**Perez-Cruz, Weston, Hermann & Scholkopf ['03]**

$$\min_{\boldsymbol{w}, b, \boldsymbol{z}, \rho} \quad -\nu\rho + \frac{1}{m}\sum_{i=1}^{m} z_i$$

$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq \rho - z_i, \quad (i \in M)$$

$$\boldsymbol{z} \geq \boldsymbol{0}, \quad \boldsymbol{w}^\top \boldsymbol{w} = 1$$

**Nonconvex optimization**

- The margin $\rho^*$ is negative for $\nu \in (0, \nu_{\min}]$.

- A non-trivial solution is obtained even for the range.

- The same optimal sol. with $\nu$-SVM for $\nu \in (\nu_{\min}, \nu_{\max}]$

- An iterative algorithm was proposed for a local solution.

# Advantage of Extended Range of $\nu$

# Uncertainty in Dataset

Bi & Zhang ('04), Shivaswamy et al. ('06), Trafalis & Gilbert ('06), etc. applied robust optimization to handle uncertainty in observations.

$$\boldsymbol{x}_i^o \rightarrow \boldsymbol{x}_i^o + \Delta \boldsymbol{x}_i$$
$$\Delta \boldsymbol{x}_i \in \mathcal{U}_i := \{\Delta \boldsymbol{x}_i : \|\Delta \boldsymbol{x}_i\| \leq \delta_i\}$$

$y = -1$

$\widehat{\boldsymbol{w}}^\top \boldsymbol{x}_i + \hat{b} < 0$

$\boldsymbol{x}_i^o$

$\boldsymbol{x}_j^o$

$y = 1$

$\widehat{\boldsymbol{w}}^\top \boldsymbol{x}_i + \hat{b} > 0$

Instead of the deterministic constraint:
$$y_i(\boldsymbol{w}^\top \boldsymbol{x}_i^o + b) \geq 1 - z_i$$

Robust *C*-SVM model

$$\min_{\boldsymbol{w}, b, \boldsymbol{z}} \quad \frac{1}{2}\|\boldsymbol{w}\|^2 + C \sum_{i=1}^{m} z_i$$
$$\text{s.t.} \quad \min_{\Delta \boldsymbol{x}_i \in \mathcal{U}_i} \quad y_i(\boldsymbol{w}^\top (\boldsymbol{x}_i^o + \Delta \boldsymbol{x}_i) + b) \geq 1 - z_i,$$
$$z_i \geq 0, \quad i = 1, \ldots, m \quad \rightarrow \textbf{Second-order cone program}$$

11

# Regularization = Robustness

Regularization penalty

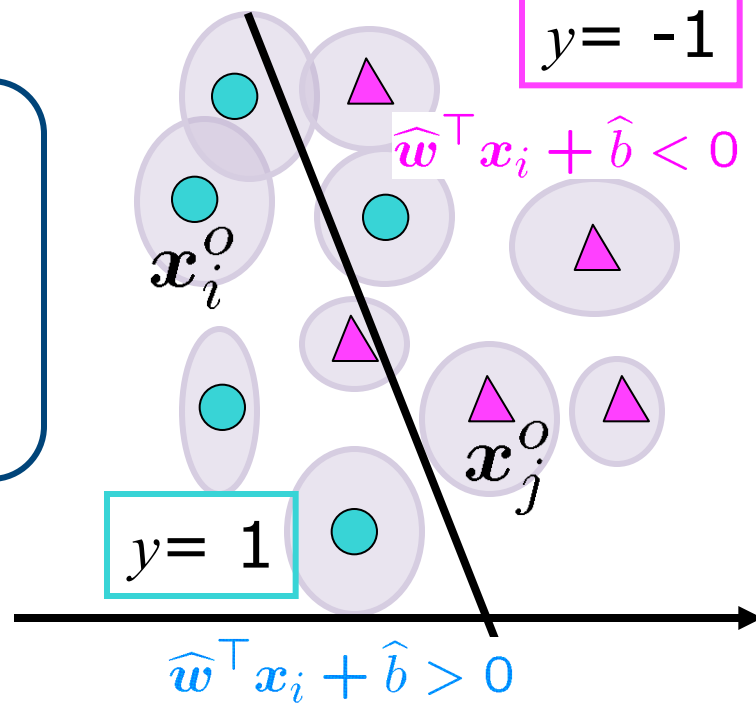$$\min_{\boldsymbol{w},b,\boldsymbol{z}} \quad \delta\|\boldsymbol{w}\| + \sum_{i=1}^{m} z_i$$
$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \boldsymbol{x}_i^o + b) \geq 1 - z_i,$$
$$z_i \geq 0, \quad i = 1,\dots,m$$

$y = -1$

$\widehat{\boldsymbol{w}}^\top \boldsymbol{x}_i + \widehat{b} < 0$

$\boldsymbol{x}_i^o$

$\boldsymbol{x}_j^o$

Equivalent

Remove "regularization"

$y = 1$

$$\min_{\boldsymbol{w},b} \quad \sum_{i=1}^{m} [1 - y_i(\boldsymbol{w}^\top \boldsymbol{x}_i^o + b)]^+$$

$\widehat{\boldsymbol{w}}^\top \boldsymbol{x}_i + \widehat{b} > 0$

Consider "robustness"

$$\min_{\boldsymbol{w},b} \max_{(\triangle \boldsymbol{x}_1,\dots,\triangle \boldsymbol{x}_m)\in \mathcal{U}} \sum_{i=1}^{m} [1 - y_i\{\boldsymbol{w}^\top(\boldsymbol{x}_i^o + \triangle \boldsymbol{x}_i) + b\}]^+$$
$$\mathcal{U} = \{(\triangle \boldsymbol{x}_1,\dots,\triangle \boldsymbol{x}_m) : \sum_{i=1}^{m} \|\triangle \boldsymbol{x}_i\| \leq \delta\}$$

12

# Robust Classification Model (RCM)

Max-min form. finds a robust solution with
the best worst-case performance.

RCM:   $$\max_{\|\boldsymbol{w}\|=1} \quad \min_{\boldsymbol{x}_+ \in \mathcal{U}_+, \boldsymbol{x}_- \in \mathcal{U}_-} (\boldsymbol{x}_+ - \boldsymbol{x}_-)^\top \boldsymbol{w}$$

**Uncertain Inputs**

✓ $\boldsymbol{x}_+, \boldsymbol{x}_-$ : representative points (or means) of each class.

✓ $\mathcal{U}_+$ (resp. $\mathcal{U}_-$) : set of possible points $\boldsymbol{x}_+$ (resp. $\boldsymbol{x}_-$) for each class, called uncertainty set.

✓ $\boldsymbol{w}$ is optimized under the worst-case vectors $\boldsymbol{x}_+^*, \boldsymbol{x}_-^*$ .

✓ $b$ is determined by using $\boldsymbol{x}_+^*$ and $\boldsymbol{x}_-^*$ ;
  e.g., so as to go though in the middle of $\boldsymbol{x}_+^*$ and $\boldsymbol{x}_-^*$ .

# Examples of Uncertainty Sets

$\mathcal{U}_+$ and $\mathcal{U}_-$ are defined with training samples in each class.



Reduced convex hull (RCH) with param. $\kappa$ :

$$\kappa \in \left[\frac{1}{m_+}, 1\right]$$

$$\mathcal{U}_+ = \left\{ \sum_{i \in M_+} \lambda_i \boldsymbol{x}_i : \begin{array}{l} \boldsymbol{e}^\top \boldsymbol{\lambda} = 1, \\ \boldsymbol{0} \leq \boldsymbol{\lambda} \leq \kappa \boldsymbol{e} \end{array} \right\}$$

a set of discrete distributions

$M_+$ : index set of samples with label +1

Ellipsoid with param. $\kappa$ :

$$\mathcal{U}_+ = \left\{ \bar{\boldsymbol{x}}_+ + \Sigma_+^{1/2} \boldsymbol{u} : \|\boldsymbol{u}\| \leq \kappa \right\}$$
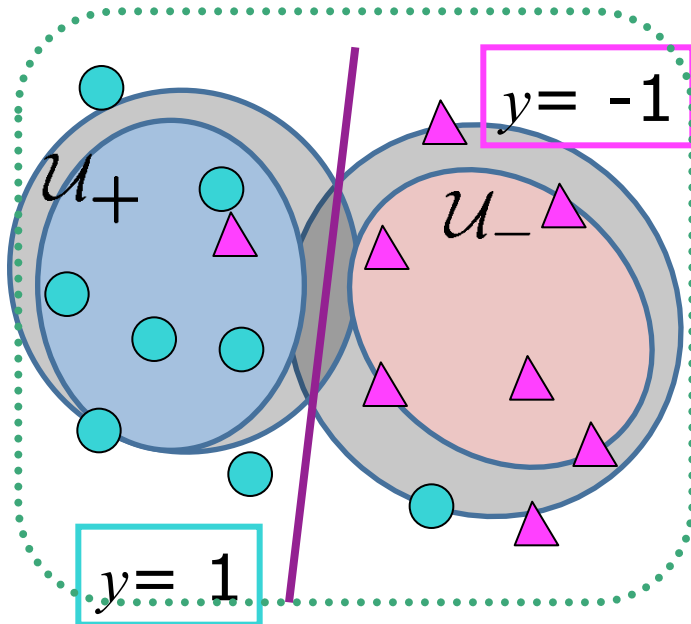
using sample mean : $\bar{\boldsymbol{x}}_+, \ \bar{\boldsymbol{x}}_-$

sample covariance : $\Sigma_+, \Sigma_-$

of samples in each class.

14

# Intersecting or Non-intersecting Uncertainty Set

RCM: $\quad \max\limits_{\|\boldsymbol{w}\|=1} \min\limits_{\boldsymbol{x}_+ \in \mathcal{U}_+, \boldsymbol{x}_- \in \mathcal{U}_-} (\boldsymbol{x}_+ - \boldsymbol{x}_-)^\top \boldsymbol{w}$



$\mathcal{U}_+$

$\mathcal{U}_-$

$y = -1$

$y = 1$

Two uncertainty sets do not intersect.

➡ $\|\boldsymbol{w}\| = 1$ is replaced by $\|\boldsymbol{w}\| \leq 1$.

➡ $\min\limits_{\boldsymbol{x}_+ \in \mathcal{U}_+, \boldsymbol{x}_- \in \mathcal{U}_-} \|\boldsymbol{x}_+ - \boldsymbol{x}_-\|$

Optimal solution: $\boldsymbol{w} = \boldsymbol{x}_+^* - \boldsymbol{x}_-^*$

⟶ RCM is a convex problem.
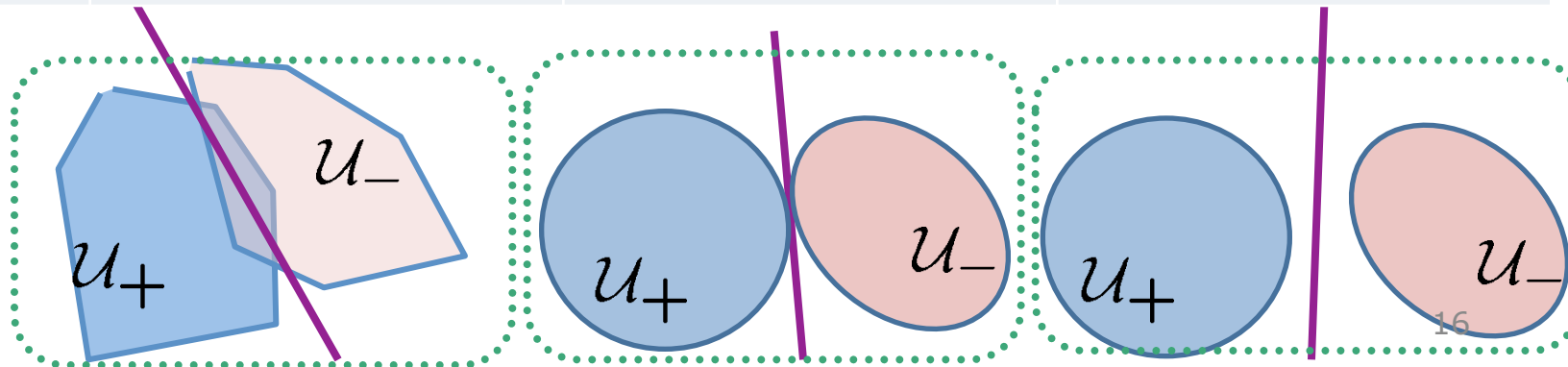
Two uncertainty sets intersect.

➡ $\|\boldsymbol{w}\| = 1$ is replaced by $\|\boldsymbol{w}\| \geq 1$.

⟶ RCM is a non-convex problem.

RCMs with specific sets $\mathcal{U}_\pm$ are reduced to well-known models. 15

# Correspondence to Existing Classifiers

| Uncertainty sets | Intersecting | They touch externally | Non-intersecting |
|---|---|---|---|
| **Ellipsoid 1 :** | **No corresponding model** | **Minimax Probability Machine (MPM)** Lanckriet et al. ('02) | **Minimum Margin-MPM** Nath & Bhattacharyya ('07) |
| **Ellipsoid 2 :** | **No corresponding model** | **Fisher Discriminant Analysis (FDA)** Fukunaga ('90) | **Sparse Feature Selection** Bhattacharyya ('04) |
| **Reduced convex hull :** | E$\nu$-SVM Perez-Cruz et al. ('03) | $\nu_{\min}$ Crisp & Burges ('00) | $\nu$-SVM ( = *C*-SVM ) Scholkopf et al. ('00) |
| **Convex hull :** $\nu \rightarrow \infty$ | ---- | ---- | **Hard Margin SVM** Boser et al. ('92) |

# What Can We Achieve from Robust-Opt View?

We could give an unified interpretation as robust optimization for some existing classification models.

✓ Main difference of those models is **in the definition of their uncertainty sets** for the mean of each class.

✓ New models can be available by defining new uncertainty sets.

✓ The parameter range can be extended so that the intersection of two sets are allowed.

✓ **Unified solution method based on APG** is applicable to convex models (nonintersecting cases).

# Correspondence to Existing Classifiers

| Uncertainty sets | Intersecting | They touch externally | Non-intersecting |
|---|---|---|---|
| Ellipsoid 1 : | No corresponding model | Minimax Probability Machine (MPM) Lanckriet et al. ('02) | Minimum Margin-MPM Nath & Bhattacharyya ('07) |
| Ellipsoid 2 : | No corresponding model | Fisher Discriminant Analysis (FDA) Fukunaga ('90) | Sparse Feature Selection Bhattacharyya ('04) |
| Reduced convex hull : | Eν-SVM Perez-Cruz et al. ('03) | $\nu_{min}$ Crisp & Burges ('00) | ν-SVM ( = C-SVM ) Scholkopf et al. ('00) |
| Convex hull : | | | Hard Margin SVM Boser et al. ('92) |

Analyze these models
by stochastic programming approach



$\mathcal{U}_+$   $\mathcal{U}_-$   $\mathcal{U}_+$   $\mathcal{U}_-$   $\mathcal{U}_+$   $\mathcal{U}_-$
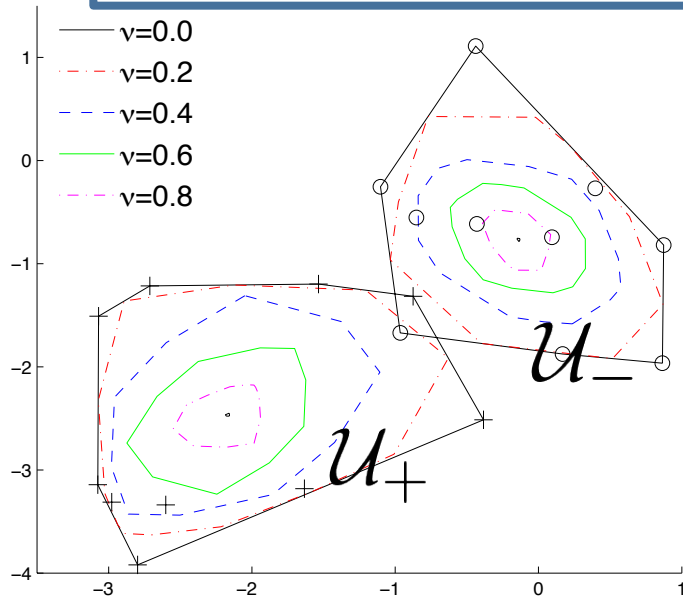
18

# Contents

● Provide a view based on Robust Optimization for various Binary Classification Models including
  - ✓ Support Vector Machine (SVM),
    Minimax Probability Machine (MPM) and
    Fisher Discriminant Analysis (FDA), etc.

● Provide a view based on Stochastic Programming
  - ✓ $\nu$-SVM & E$\nu$-SVM    → Generalization Bound
  - ✓ Minimum Margin MPM

# ν-SVM & Eν-SVM (dual form.)

Robust Classification Model

$$\max_{\|\boldsymbol{w}\|=1} \min_{\boldsymbol{x}_+ \in \mathcal{U}_+, \boldsymbol{x}_- \in \mathcal{U}_-} (\boldsymbol{x}_+ - \boldsymbol{x}_-)^\top \boldsymbol{w}$$

$\Big\{$ If two RCHs do not intersect (with large ν) ➡ **ν-SVM**

If two RCHs intersect       (with small ν) ➡ **Eν-SVM**



Reduced convex hull (RCH) with param. ν :

$$\mathcal{U}_+ = \left\{ \sum_{i \in M_+} \lambda_i \boldsymbol{x}_i : \begin{array}{l} e^\top \boldsymbol{\lambda} = 1, \\ 0 \le \boldsymbol{\lambda} \le \frac{2}{\nu m} e \end{array} \right\}$$

$\kappa$

Shrunk polytopes toward the centers by increasing ν.

$$\nu \in \left(0, 2\frac{\min(m_+, m_-)}{m}\right]$$

20

# ν-SVM & Eν-SVM (primal form.)

Robust Classification Model

$$\max_{\|\boldsymbol{w}\|=1} \min_{\boldsymbol{x}_+ \in \mathcal{U}_+, \boldsymbol{x}_- \in \mathcal{U}_-} (\boldsymbol{x}_+ - \boldsymbol{x}_-)^\top \boldsymbol{w}$$

$$\|\boldsymbol{w}\| \geq 1 \qquad\qquad \|\boldsymbol{w}\| \leq 1$$

Two RCHs intersect.  Two RCHs do not intersect.

Nonconvex Program  Convex Program

$\boldsymbol{\nu} = 0$

## CVaR Minimization??

$= 2\frac{\min(m_+, m_-)}{m}$

**(Eν-SVM )** Perez-Cruz, Weston, Hermann & Schoelkopf ('03)

$$\min_{\boldsymbol{w}, b, \boldsymbol{z}, \rho} \quad -\boldsymbol{\nu}\rho + \frac{1}{m}\sum_{i=1}^m z_i$$
$$\text{s.t.} \quad z_i + y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) - \rho \geq 0,$$
$$i = 1, \ldots, m,$$
$$\boldsymbol{z} \geq \boldsymbol{0}, \quad \boldsymbol{w}^\top \boldsymbol{w} = 1$$

**(ν-SVM )** Schoelkopf, Smola, Williamson & Bartlett ('00)

$$\min_{\boldsymbol{w}, b, \boldsymbol{z}, \rho} \quad \frac{1}{2}\|\boldsymbol{w}\|^2 - \boldsymbol{\nu}\rho + \frac{1}{m}\sum_{i=1}^m z_i$$
$$\text{s.t.} \quad z_i + y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) - \rho \geq 0,$$
$$i = 1, \ldots, m,$$
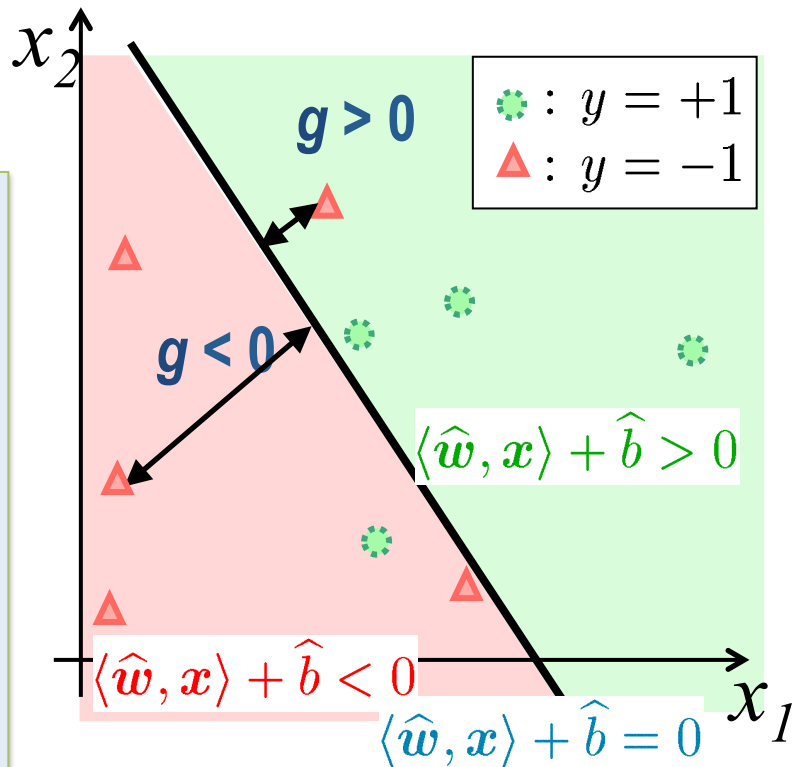$$\boldsymbol{z} \geq \boldsymbol{0}, \quad \cancel{\rho \geq 0}$$

21

# CVaR of Distance



For a hyperplane:
$$\boldsymbol{w}^\top \boldsymbol{x} + b = 0$$
compute the **signed distance (score)** from a point $\boldsymbol{x}_i$ to the hyperplane for all training samples by

$$g(\boldsymbol{w}, b; \boldsymbol{x}_i, y_i) = -\frac{y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b)}{\|\boldsymbol{w}\|}$$

In the figure:
- $\cdot$ : $y = +1$
- $\triangle$ : $y = -1$
- $g > 0$
- $g < 0$
- $\langle \widehat{\boldsymbol{w}}, \boldsymbol{x} \rangle + \widehat{b} > 0$
- $\langle \widehat{\boldsymbol{w}}, \boldsymbol{x} \rangle + \widehat{b} < 0$
- $\langle \widehat{\boldsymbol{w}}, \boldsymbol{x} \rangle + \widehat{b} = 0$

$g < 0$   correctly classified,     $g > 0$   misclassified

$$y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) > 0$$

Minimize CVaR $\phi_\beta(\boldsymbol{w}, b)$ with $\beta = 1 - \nu$ using $g(\boldsymbol{w}, b; \boldsymbol{x}_i, y_i),\ i = 1, \ldots, m$

⟹ hyperplane of **(E)ν-SVM**

22

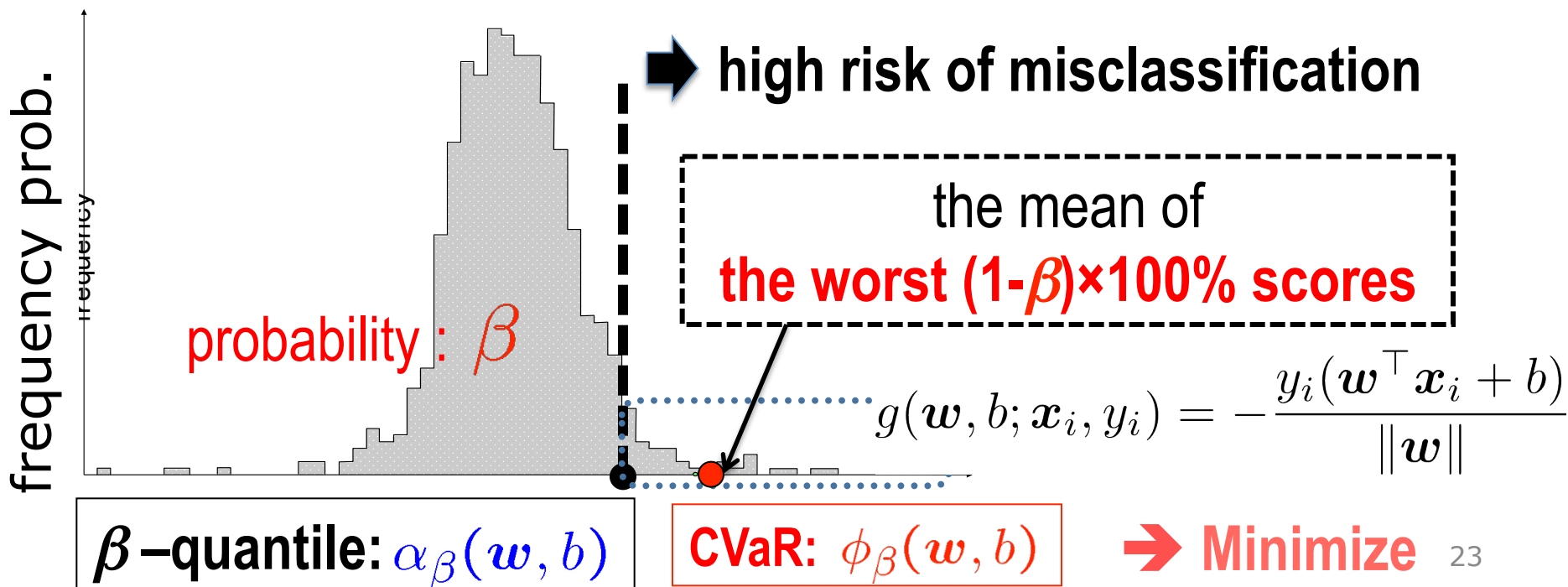# CVaR Minimization for Classification

✓ Minimize CVaR $\phi_\beta(\boldsymbol{w}, b)$ with $\beta = 1 - \nu$

and $\Pr((\boldsymbol{x}, y) = (\boldsymbol{x}_i, y_i)) = \dfrac{1}{m}$

by

$$\min_{\boldsymbol{x}, b, \alpha} \alpha + \frac{1}{m\nu} \sum_{i=1}^{m} [g(\boldsymbol{w}, b; \boldsymbol{x}_i, y_i) - \alpha]^+$$



➡️ **high risk of misclassification**

**frequency prob.**

probability : $\beta$

the mean of
**the worst (1-$\beta$)×100% scores**

$$g(\boldsymbol{w}, b; \boldsymbol{x}_i, y_i) = -\frac{y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b)}{\|\boldsymbol{w}\|}$$

**$\beta$ –quantile:** $\alpha_\beta(\boldsymbol{w}, b)$ | **CVaR:** $\phi_\beta(\boldsymbol{w}, b)$ | ➡️ **Minimize**

# New interpretation for Eν-SVC

$$\min_{\boldsymbol{x},b,\alpha} \alpha + \frac{1}{m\nu}\sum_{i=1}^{m}\left[-\frac{y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b)}{\|\boldsymbol{w}\|} - \alpha\right]^+$$

If $\phi_{1-\nu} > 0$

variable: $\rho = -\alpha$

If $\phi_{1-\nu} \leqq 0$

**(Eν-SVM )** **Perez-Cruz, Weston, Hermann & Schoelkopf ('03)**

$$\min_{\boldsymbol{w},b,\boldsymbol{z},\rho} \quad -\nu\rho + \frac{1}{m}\sum_{i=1}^{m} z_i$$

$$\text{s.t.} \quad z_i + y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) - \rho \geq 0,$$
$$i = 1,\dots,m,$$
$$\boldsymbol{z} \geq \boldsymbol{0}, \quad \boldsymbol{w}^\top \boldsymbol{w} = 1$$

**(ν-SVM )** **Schoelkopf, Smola, Williamson & Bartlett ('00)**

$$\min_{\boldsymbol{w},b,\boldsymbol{z},\rho} \quad \frac{1}{2}\|\boldsymbol{w}\|^2 - \nu\rho + \frac{1}{m}\sum_{i=1}^{m} z_i$$

$$\text{s.t.} \quad z_i + y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) - \rho \geq 0$$
$$i = 1,\dots,m,$$
$$\boldsymbol{z} \geq \boldsymbol{0}, \quad \rho \geq 0 \quad \boldsymbol{w}^\top \boldsymbol{w} \leq 1$$

$$\rho^* = -\alpha^* \approx -\alpha_{1-\nu}$$
: margin of Eν-SVC
negative margin $\leftrightarrow \alpha_{1-\nu} > 0$

**misclassification**

bad samples =SVs

$\beta = 1-\nu$

$0$

$\alpha_{1-\nu}$ $\phi_{1-\nu}$

24

# Three Cases depending on $\nu$

If $\phi_{1-\nu} > 0$                                    If $\phi_{1-\nu} \leqq 0$
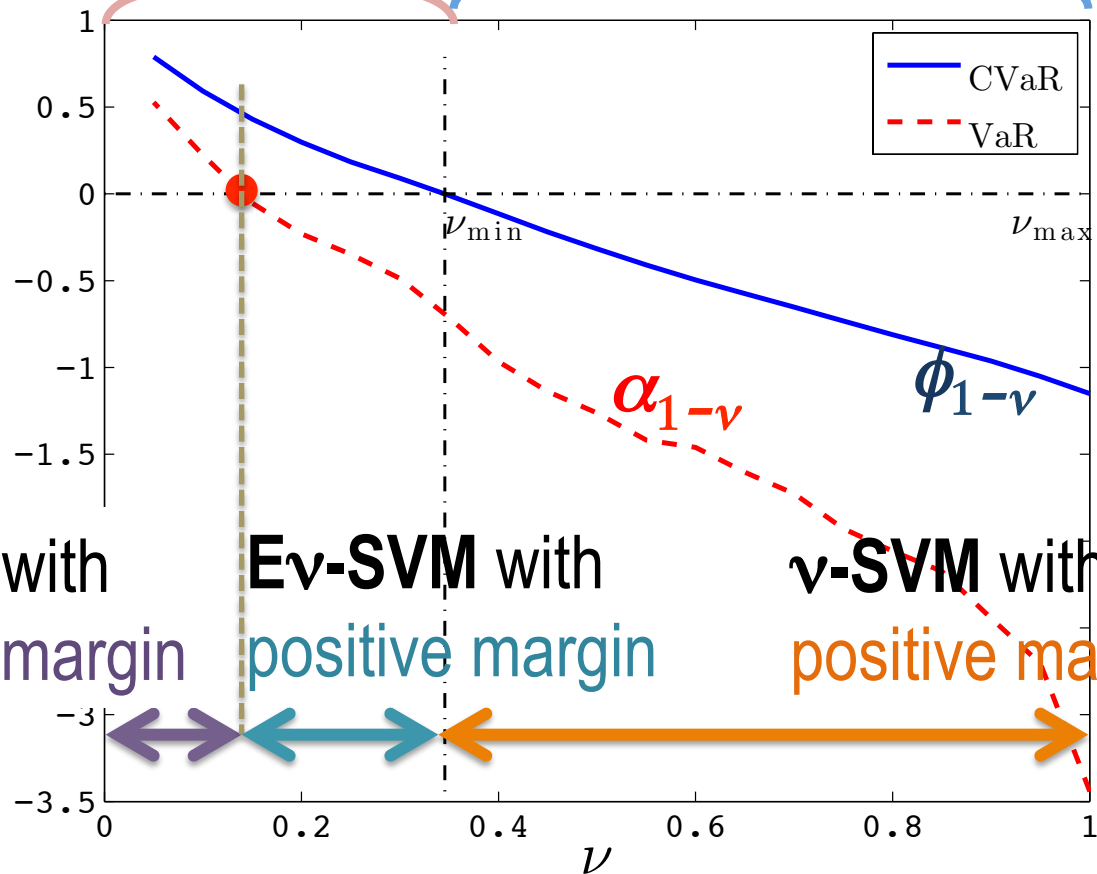
(E$\nu$-SVM) **Perez-Cruz, Weston, Hermann & Schoelkopf ('03)** **Nonconvex Problem**

($\nu$-SVM) **Schoelkopf, Smola, Williamson & Bartlett ('00)** **Convex Problem**

Margin:
$$\rho^* = -\alpha_{1-\nu}$$



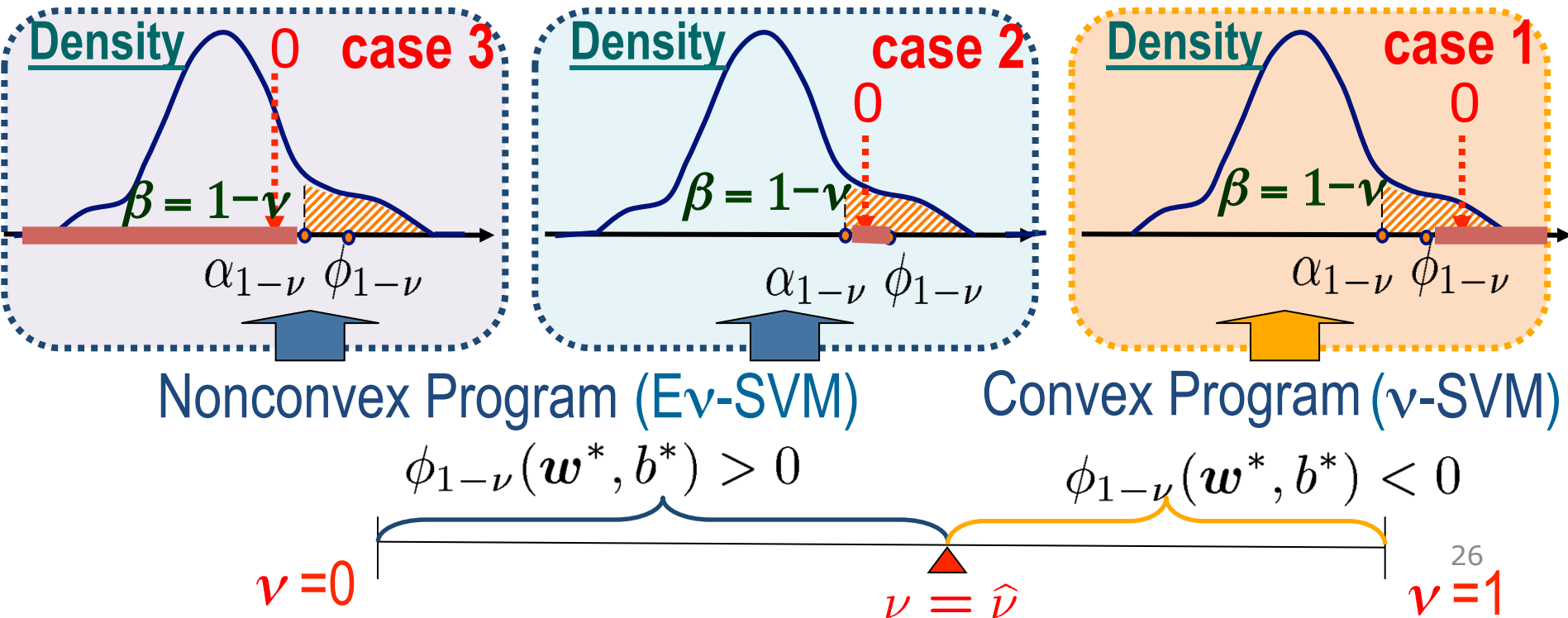**E$\nu$-SVM** with negative margin    **E$\nu$-SVM** with positive margin    **$\nu$-SVM** with positive margin

25

# Generalization Error Bounds

New <u>generalization error</u> bounds of Eν-SVM include
the CVaR risk measure ⌐ error rates for test (new) samples

➜ Minimizing the CVaR lowers the bound

➜ It justifies the use of Eν-SVM & ν-SVM



Nonconvex Program (Eν-SVM)          Convex Program (ν-SVM)

$$\phi_{1-\nu}(\boldsymbol{w}^*, b^*) > 0 \qquad \phi_{1-\nu}(\boldsymbol{w}^*, b^*) < 0$$

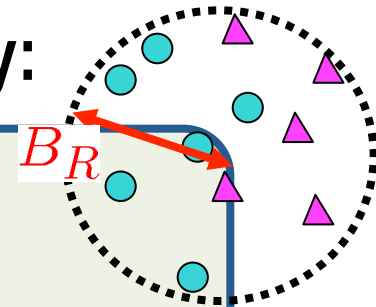$\nu = 0$        $\nu = \hat{\nu}$        $\nu = 1$

# Generalization Error Bound (case 1)

**Takeda-Sugiyama [ '08]**

**Theorem : (case 1 )**

**For a feasible sol.** $(w, b)$ **of ($\nu$-SVM), the inequality:**

( generalization error with $f(x) = w^\top x + b$ )

$$\leq \nu + G(\alpha_{1-\nu}(w, b)) \leq \nu + G(\phi_{1-\nu}(w, b))$$

$< 0$

$$G(\gamma) := \sqrt{\frac{2}{m}\left(\frac{4c^2(1+B_R^2)^2}{\gamma^2}\log_2(2m) - 1 + \log\left(\frac{2}{\delta}\right)\right)}$$

$G(\rho^*)$ is used for (**$\nu$-SVM**) in **Schoelkopf, Smola, Williamson & Bartlett ('00)**

**holds with probability at least** $1 - \delta$

➔ CVaR min. gives an opt. solution which minimizes the bound.

➔ $\nu$-SVM is reasonable.
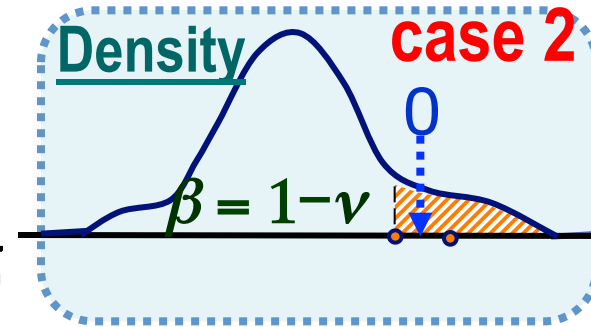
# Generalization Error Bound (cases 2&3)

**For a feasible sol.** $(\boldsymbol{w}, b)$ **of (E$\nu$-SVM)**

(generalization error with $f(\boldsymbol{x}) = \boldsymbol{w}^\top \boldsymbol{x} + b$)

$\leqq \quad \nu + G(\boxed{\alpha_{1-\nu}(\boldsymbol{w}, b)})$

holds with probability at least $1 - \delta$

$$G(\gamma) := \sqrt{\frac{2}{m}\left(\frac{4c^2(1+B_R^2)^2}{\gamma^2}\log_2(2m) - 1 + \log\left(\frac{2}{\delta}\right)\right)}$$

**Density** **case 2**
0
$\boldsymbol{\beta} = 1 - \boldsymbol{\nu}$

**For a feasible sol.** $(\boldsymbol{w}, b)$ **of (E$\nu$-SVM)**

(generalization error with $f(\boldsymbol{x}) = \boldsymbol{w}^\top \boldsymbol{x} + b$)

$\geqq \quad \nu - G(\boxed{\alpha_{1-\nu}(\boldsymbol{w}, b)})$

This bound is upper -bounded as

$\nu - G(\alpha_{1-\nu}(\boldsymbol{w}, b)) \leq \nu - G(\boxed{\phi_{1-\nu}(\boldsymbol{w}, b)})$

**Density** 0 **case 3**
$\boldsymbol{\beta} = 1 - \boldsymbol{\nu}$

# (E)ν-SVM (classification method)

CVaR Min.:

$$\min_{\boldsymbol{w},b,\rho} \; -\rho + \frac{1}{\nu m} \sum_{i \in M} [g(\boldsymbol{w}, b; \boldsymbol{x}_i, y_i) + \rho]^+$$

$$g(\boldsymbol{w}, b; \boldsymbol{x}_i, y_i) = -\frac{y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b)}{\|\boldsymbol{w}\|}$$

Stochastic Programming

(E)ν-SVM:

$$\min_{\boldsymbol{w},b,\boldsymbol{z},\rho} \quad -\nu\rho + \frac{1}{m} \sum_{i=1}^{m} z_i$$

$$\text{s.t.} \quad z_i + y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) - \rho \geq 0, \quad i \in M,$$

$$\boldsymbol{z} \geq \boldsymbol{0}, \quad \boldsymbol{w}^\top \boldsymbol{w} = 1 \quad (\text{or } \boldsymbol{w}^\top \boldsymbol{w} \leq 1)$$

**Perez-Cruz, Weston, Hermann & Schoelkopf ('03)**

Robust Optimization

by taking dual w.r.t. $b, \boldsymbol{z}, \rho$

RCM:

$$\max_{\|\boldsymbol{w}\|=1} \; \min_{\boldsymbol{x}_+ \in \mathcal{U}_+, \boldsymbol{x}_- \in \mathcal{U}_-} (\boldsymbol{x}_+ - \boldsymbol{x}_-)^\top \boldsymbol{w}$$

$$\mathcal{U}_\pm = \left\{ \sum_{i \in M_\pm} \lambda_i \boldsymbol{x}_i : \boldsymbol{e}^\top \boldsymbol{\lambda} = 1, \; 0 \leq \boldsymbol{\lambda} \leq \frac{2}{\nu m} \boldsymbol{e} \right\}$$

29

# Ellipsoidal Uncertainty Sets

Robust Classification Model

$$\max_{\|\boldsymbol{w}\|=1} \quad \min_{\boldsymbol{x}_+\in\mathcal{U}_+,\boldsymbol{x}_-\in\mathcal{U}_-} (\boldsymbol{x}_+ - \boldsymbol{x}_-)^\top \boldsymbol{w}$$
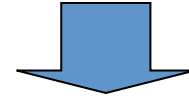


Using sample mean : $\bar{x}_+, \ \bar{x}_-$

sample covariance : $\Sigma_+, \Sigma_-$

of samples in each class, let

$$\mathcal{U}_+ = \left\{ \bar{x}_+ + \Sigma_+^{1/2}\boldsymbol{u} : \|\boldsymbol{u}\| \leq \kappa \right\}$$

$$\mathcal{U}_- = \left\{ \bar{x}_- + \Sigma_-^{1/2}\boldsymbol{v} : \|\boldsymbol{v}\| \leq \kappa \right\}.$$

$$\min_{\|\boldsymbol{w}\|=1} \kappa\|\Sigma_+^{1/2}\boldsymbol{w}\| + \kappa\|\Sigma_-^{1/2}\boldsymbol{w}\| - (\bar{x}_+ - \bar{x}_-)^\top \boldsymbol{w}$$

$\|\boldsymbol{w}\| = 1$ can be replaced by $\|\boldsymbol{w}\| \leq 1$ when $\mathcal{U}_+ \cap \mathcal{U}_- = \emptyset$ .

# Equivalence to Maximum-Margin MPM

**Robust Classification Model (non-intersecting case)**

$$\min_{\|\boldsymbol{w}\|\leq 1} \kappa\|\Sigma_+^{1/2}\boldsymbol{w}\| + \kappa\|\Sigma_-^{1/2}\boldsymbol{w}\| - (\bar{\boldsymbol{x}}_+ - \bar{\boldsymbol{x}}_-)^\top\boldsymbol{w}$$

$$\kappa = \sqrt{\frac{1-\eta}{\eta}}$$

## Maximum-Margin MPM

$$\min_{\boldsymbol{w},b} \frac{1}{2}\|\boldsymbol{w}\|^2$$

**Nath & Bhattacharyya ('07)**

Worst-case misclassified probabilities

$$\text{s.t.} \quad \sup_{\boldsymbol{x}_+\sim(\bar{\boldsymbol{x}}_+,\Sigma_+)} \text{Pr}\{\boldsymbol{x}_+^\top\boldsymbol{w} + b < 1\} \leq \eta$$

Using generalized Chebyshev-Cantelli inequality,

$$\bar{\boldsymbol{x}}_+^\top\boldsymbol{w} + b \geq 1 + \sqrt{\frac{1-\eta}{\eta}}\|\Sigma_+^{1/2}\boldsymbol{w}\|$$

$\boldsymbol{x}_+$, $\boldsymbol{x}_-$ :  random vectors from each of two classes with means and covariance matrices given by $(\bar{\boldsymbol{x}}_+, \Sigma_+)$ and $(\bar{\boldsymbol{x}}_-, \Sigma_-)$.

# Stochastic Problem under Normal Distribution

Robust Classification Model with $\mathcal{U}_{\pm} = \{\bar{\boldsymbol{x}}_{\pm} + \Sigma_{\pm}^{1/2}\boldsymbol{u} : \|\boldsymbol{u}\| \leq \kappa\}$

$$\max_{\|\boldsymbol{w}\| \leq 1} \quad \min_{\boldsymbol{x}_+ \in \mathcal{U}_+, \boldsymbol{x}_- \in \mathcal{U}_-} (\boldsymbol{x}_+ - \boldsymbol{x}_-)^\top \boldsymbol{w}$$

$$\kappa = \sqrt{\frac{1-\eta}{\eta}}$$

$$\min_{\boldsymbol{w},b} \frac{1}{2}\|\boldsymbol{w}\|^2$$

$$\bar{\boldsymbol{x}}_+^\top \boldsymbol{w} + b \geq 1 + \sqrt{\frac{1-\eta}{\eta}}\|\Sigma_+^{1/2}\boldsymbol{w}\|$$

$$\text{s.t.} \quad \sup_{\boldsymbol{x}_+ \sim (\bar{\boldsymbol{x}}_+, \Sigma_+)} \Pr\{\boldsymbol{x}_+^\top \boldsymbol{w} + b < 1\} \leq \eta$$

$$\sup_{\boldsymbol{x}_- \sim (\bar{\boldsymbol{x}}_-, \Sigma_-)} \Pr\{\boldsymbol{x}_-^\top \boldsymbol{w} + b > -1\} \leq \eta$$

The worst-case prob. distribution is considered in **Nath & Bhattacharyya ('07)**

$$\kappa = \Phi^{-1}(1-\eta)$$

Under the assump: $\boldsymbol{x}_+ \sim \mathcal{N}_{m_+}(\bar{\boldsymbol{x}}_+, \Sigma_+)$

$$\Pr\{\boldsymbol{x}_+^\top \boldsymbol{w} + b < 0\} \leq \eta$$

$\Phi(z)$: cumulative dist. func. (cdf) of $\mathcal{N}(0,1)$

$$\bar{\boldsymbol{x}}_+^\top \boldsymbol{w} + b \geq 1 + \Phi^{-1}(1-\eta)\|\Sigma_+^{1/2}\boldsymbol{w}\|$$

# Conclusions

➢ We provided new views based on <span style="color:red">Robust Optimization / Stochastic Programming</span> for existing machine learning classification models (SVM, MPM, FDA and their variants).

➢ We could evaluate <span style="color:red">generalization bounds</span> from the viewpoint of SP and propose an <span style="color:red">efficient algorithm</span> from the viewpoint of RO.

# Summary

● The first textbook on Robust Optimization appears in 2009.

**Ben-Tal, El Ghaoui & Nemirovski ['09]**

● Robust optimization techniques are used in various research areas.

✓ The preface of the book briefly mentions the relation to
Robust Control (H$_\infty$ Control), Robust Statistics,
Machine learning（SVM）, etc.

● Recently, studies on robust optimization using "probability"
are increased. The robust optimization research is still developing.

# 参考文献 -1-

❑ E. M. L. Beale. On minimizing a convex function subject to linear inequalities. J. Roy. Statist. Soc. Ser. B. 17 (1955), 173–184.

❑ A. Ben-Tal, L. El Ghaoui and A. Nemirovski. Robust optimization. Princeton University Press, 2009.

❑A. Ben-Tal, A. Goryashko, E. Guslitzer and A. Nemirovski. Adjustable robust solutions of uncertain linear programs. *Math. Progr.* 99 (2004), 351-376.

❑A. Ben-Tal and A. Nemirovski. Robust convex optimization. *Math. of Oper. Res.* 23:4 (1998), 769-805.

❑A. Ben-Tal and A. Nemirovski. Robust solutions of uncertain linear programs. *OR Letters* 25 (1999), 1-13.

❑Ben-Tal and A. Nemirovski. Robust solutions of linear programming problems contaminated with uncertain data. *Math. Progr.* 88 (2000), 411-424.

❑B. E. Boser, I. M. Guyon and V. N. Vapnik. A training algorithm for optimal margin classiers. COLT (pp. 144-152). ACM Press, 1992.

❑G. Calafiore and M.C. Campi. Uncertain convex programs: Randomized solutions and confidence levels. *Math. Progr. 102:1 (2005), 25–46.*

❑G. Calafiore and M.C. Campi. The scenario approach to robust control design. IEEE Transactions on Automatic Control, 51:5 (2006), 742–753.

❑T.C.Y. Chan, T. Bortfeld and J.N. Tsitsiklis. A robust approach to IMRT optimization. Phys. Med. Biol. 51 (2006), 2567–2583.

# 参考文献 -2-

❑A. Charnes and W. W. Cooper. Uncertain convex programs: Randomize solutions and confidence level. *Management Sci.* 6 (1959), 73–79.

❑C. Cortes and V. Vapnik. Support-vector networks. Machine Learning, 20 (1995), 273-297.

❑G. B. Dantzig. Linear programming under uncertainty. Management Sci., 1 (1955), 197-206.

❑L. El Ghaoui and H. Lebret. Robust solution to least-squares problems with uncertain data. *SIAM J. of Matrix Anal. Appl.* 18 (1997), 1035-1064.

❑D. Goldfarb and G. Iyengar. Robust convex quadratically constrained programs. *Math. Progr.* 97 (2003), 495-515.

❑J. Gotoh and A. Takeda. A linear Classification Model Based on Conditional Geometric Score. Pacific Journal of Optimization 1 (2005), 277-296.

❑P. Kouvelis and G. Yu. Robust discrete optimization and its applications. Kluwer Academic, Dordrecht (1997).

❑X. Lin, S. L. Janak and C. A. Floudas. A new robust optimization approach for scheduling under uncertainty: I. Bounded uncertainty. Computers and Chemical Engineering 28 (2004), 1069–1085.

❑F. Perez-Cruz , J. Weston, D.J.L. Hermann, B. Schölkopf. Extension of the *v*-SVM range for classification. Advances in Learning Theory: Methods, Models and Applications 190 (2003), 179–196.

# 参考文献 -3-

- R.T. Rockafellar and S. Uryasev. Conditional value-at-risk for general loss distributions. J. Bank. Finance 26 (2002) , 1443–1471.
- B. Schoelkopf, A. J. Smola, R. C. Williamson and P. L. Bartlett. New support vector algorithms. Neural Computation 12 (2000), 1207–1245.
- A.L. Soyster. Convex programming with set-inclusive constraints and applications to inexact linear programming. *Oper. Res. 21* (1973), 1154-1157.
- A. Takeda and M. Sugiyama. *ν*-support vector machine as conditional value-at-risk minimization. ICML 2008, 2008.
- A. Takeda, H. Mitsugi and T. Kanamori. A unified robust classification model. ICML2012, 2012.
- H. Xu, C. Caramanis and S. Mannor. Robustness and regularization of support vector machines. Journal of Machine Learning Research, 10 (2009), 1485-1510.
- G. Yu and J. Yang. On the robust shortest path problem. Comput. Oper. Res. 25 (1998), 457–468.

- 寒野 善博, ロバスト性を考慮した設計. 2008年度日本建築学会大会(中国)・構造部門(応用力学)パネルディスカッション資料『建築構造設計における冗長性と頑強性の役割—リダンダンシーとロバスト性とは—』, 14-23.
- 土谷隆、笹川卓, 二次錐計画問題による磁気シールドのロバスト最適化、統計数理53：2（2005）, 297-315.