

# Pattern Information Processing:<sup>97</sup> Sparse Methods

Masashi Sugiyama  
(Department of Computer Science)

Contact: W8E-505

[sugi@cs.titech.ac.jp](mailto:sugi@cs.titech.ac.jp)

<http://sugiyama-www.cs.titech.ac.jp/~sugi/>

# Sparseness and Continuous Model Choice

- Two approaches for avoiding over-fitting:

	Sparseness	Model parameter
Subset LS	Yes	Combinatorial
Quadratically constrained LS	No	Continuous

- We want to have **sparseness** and **continuous** model choice at the same time.

# Today's Plan

- Sparse learning method
- How to deal with absolute values in optimization
- Approximate gradient descent
- Standard form of quadratic programs

# Non-Linear Learning for Linear / Kernel Models

## ■ Linear / kernel models

$$f_{\alpha}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

$$f_{\alpha}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

## ■ Non-linear learning

$$\hat{\alpha} = \mathbf{L}(\mathbf{y})$$

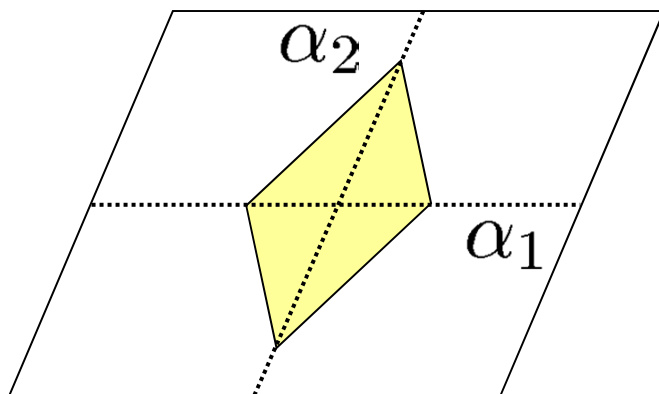
$\mathbf{L}(\cdot)$  : Non-linear function

# $\ell_1$ -Constrained LS

- Restrict the search space within an  $\ell_1$ -ball.

$$\hat{\alpha}_{\ell_1 CLS} = \underset{\alpha \in \mathbb{R}^b}{\operatorname{argmin}} J_{LS}(\alpha)$$

$$\text{subject to } \|\alpha\|_1 \leq C$$



$\ell_1$  - norm

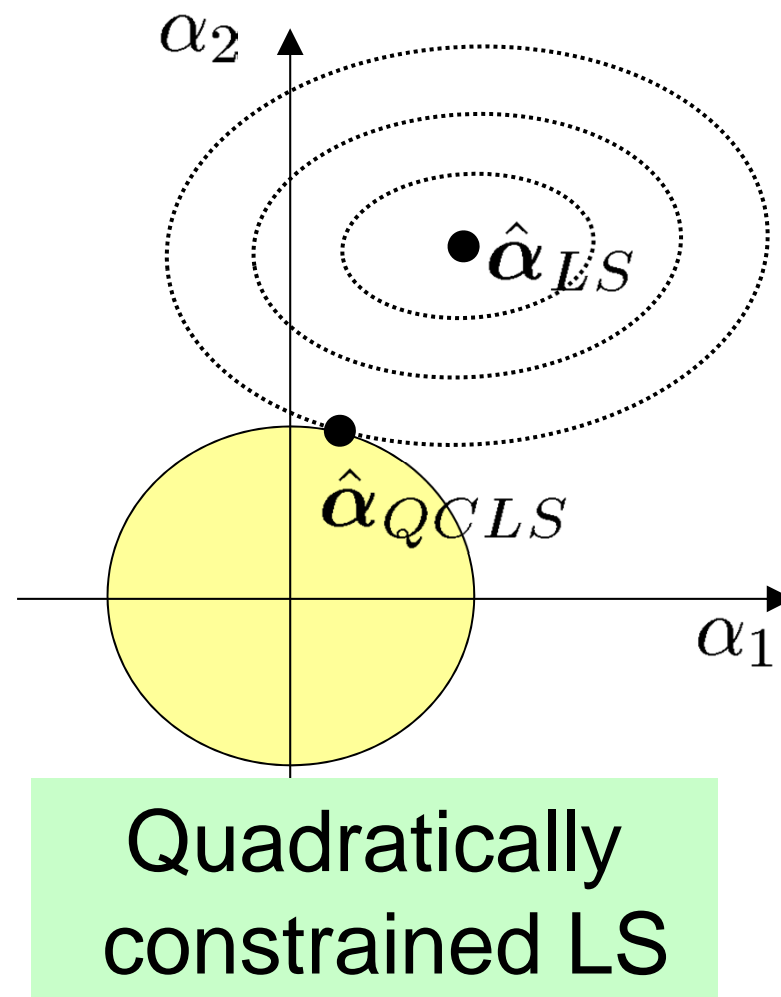
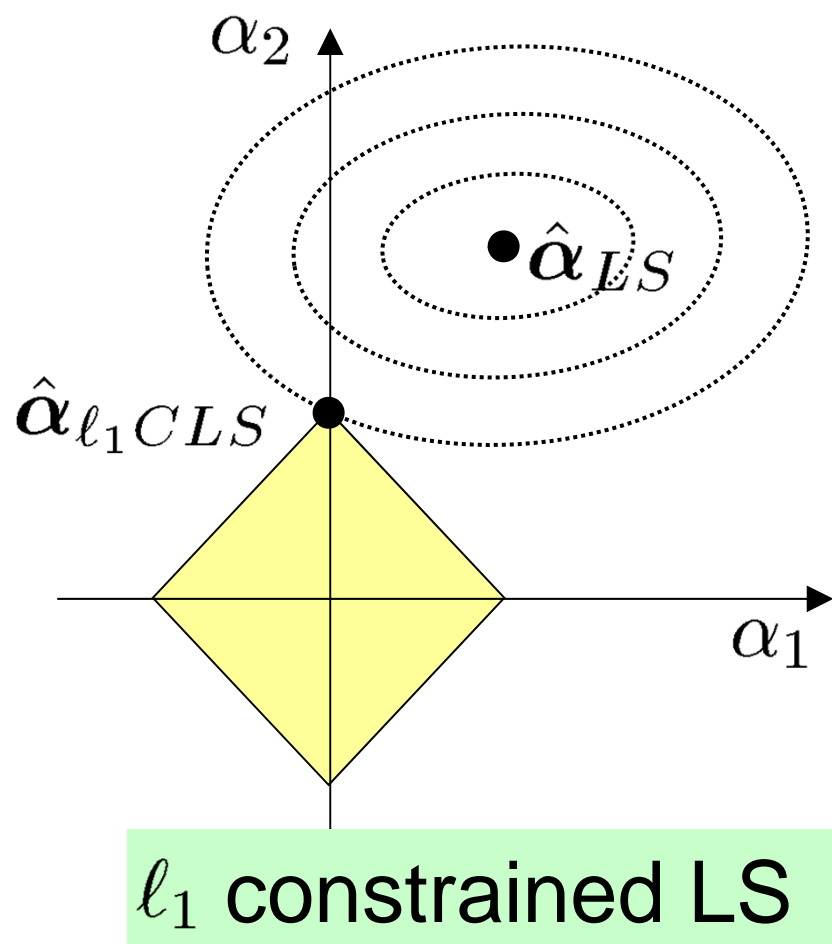
$$\|\alpha\|_1 = \sum_{i=1}^b |\alpha_i|$$

$$J_{LS}(\alpha) = \sum_{i=1}^n (f_{\alpha}(x_i) - y_i)^2$$

Tibshirani, Regression shrinkage and selection via the lasso,  
Journal of the Royal Statistical Society, Series B, 58(1), 267-288, 1996.

# Why Sparse?

- The solution is often exactly on an axis.



# How to Obtain A Solution

103

- Lagrangian:

$$J_{\ell_1 CLS}(\boldsymbol{\alpha}) = J_{LS}(\boldsymbol{\alpha}) + \lambda(\|\boldsymbol{\alpha}\|_1 - C)$$

- $\lambda$ : Lagrange multiplier

- Similarly to QCLS, we practically start from  $\lambda (\geq 0)$  and solve

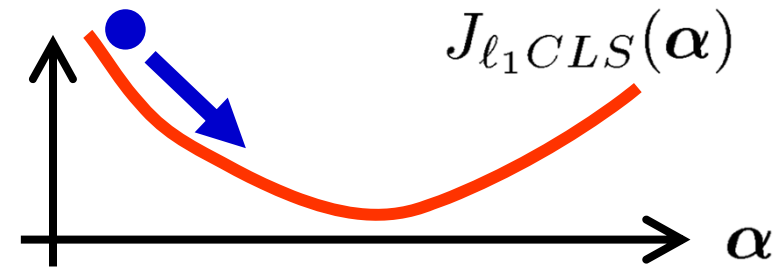
$$\hat{\boldsymbol{\alpha}}_{\ell_1 CLS} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^b}{\operatorname{argmin}} J_{\ell_1 CLS}(\boldsymbol{\alpha})$$

- It is often called  $\ell_1$ -regularized LS.

# Gradient Descent

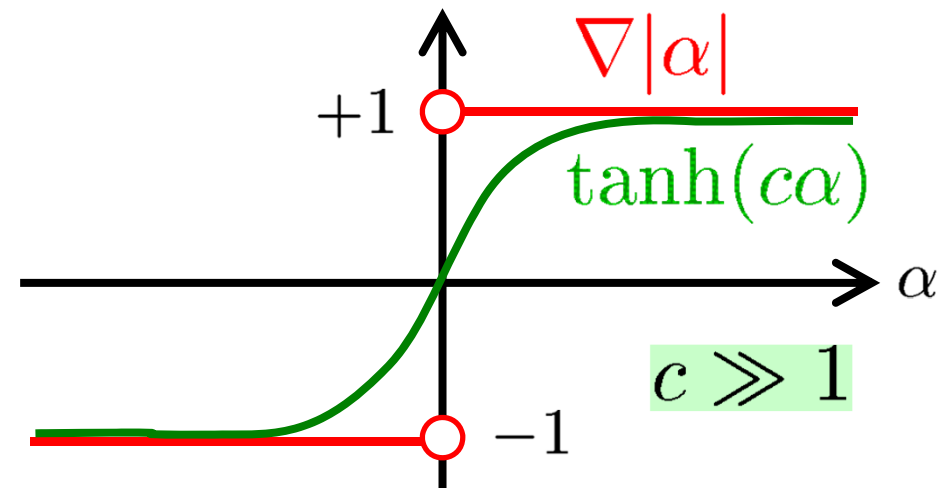
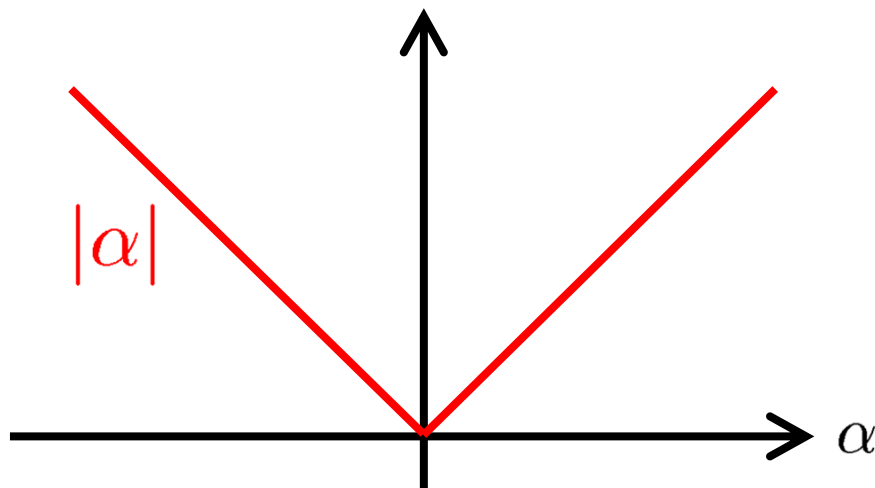
104

■  $\alpha \leftarrow \alpha - \epsilon \nabla J_{\ell_1 CLS}(\alpha)$



■ However,  $\ell_1$ -norm is not differentiable.

- Use smooth approximation!



- You may also use a quasi-Newton method.



# Quadratic Program

- Use the following lemma:

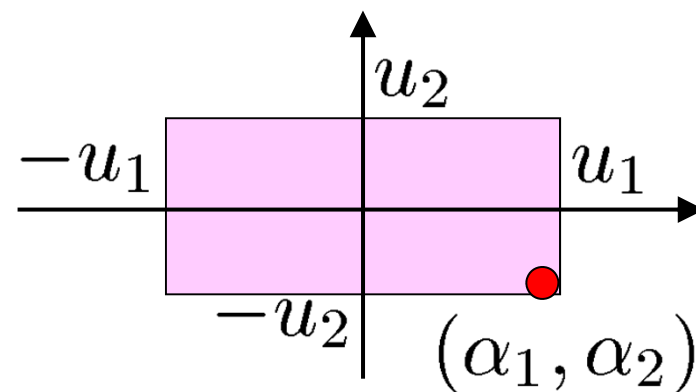
**Lemma**

$$\|\alpha\|_1 = \min_{u \in \mathbb{R}^b} \sum_{i=1}^b u_i$$

subject to  $-u \leq \alpha \leq u$

**Note:** Inequality for vectors is component-wise

**Intuition:** Obtain the smallest box that includes  $\alpha$



# Proof

■ Let  $\hat{u} = \operatorname{argmin}_{u \in \mathbb{R}^b} \sum_{i=1}^b u_i$  subject to  $-u \leq \alpha \leq u$ ,

The constraint implies  $\hat{u}_i \geq |\alpha_i|$ .

Suppose  $\hat{u}_i > |\alpha_i|$ . Then such  $\hat{u}_i$  is not a solution since  $\tilde{u}_i = |\alpha_i|$  gives a smaller value:

$$\sum_{i=1}^b \tilde{u}_i < \sum_{i=1}^b \hat{u}_i$$

This implies that the solution satisfies  $\hat{u}_i = |\alpha_i|$ , which yields

$$\sum_{i=1}^b \hat{u}_i = \sum_{i=1}^b |\alpha_i| = \|\alpha\|_1$$

(Q.E.D.)

# How to Obtain A Solution (cont.)<sup>107</sup>

$$\hat{\alpha}_{\ell_1 CLS} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} J_{\ell_1 CLS}(\alpha)$$

$$J_{\ell_1 CLS}(\alpha) = J_{LS}(\alpha) + \lambda \|\alpha\|_1$$

- $\hat{\alpha}_{\ell_1 CLS}$  is given as the solution of

$$\min_{\alpha, u \in \mathbb{R}^b} \left[ J_{LS}(\alpha) + \lambda \sum_{i=1}^b u_i \right]$$

$$\text{subject to } -u \leq \alpha \leq u,$$

$$J_{LS}(\alpha) = \sum_{i=1}^n (f_{\alpha}(x_i) - y_i)^2$$

$$= \|X\alpha - y\|^2$$

# Linearly Constrained Quadratic Program

- Standard optimization software can solve linearly constrained quadratic programs.

$$\min_{\beta} \left[ \frac{1}{2} \langle \mathbf{Q}\beta, \beta \rangle + \langle \beta, \mathbf{q} \rangle \right]$$

$$\text{subject to } \mathbf{V}\beta \leq \mathbf{v}$$

$$\mathbf{G}\beta = \mathbf{g}$$

# Transformation into Standard Form

■ Let

$$\beta = \begin{pmatrix} \alpha \\ u \end{pmatrix}$$

$$\begin{aligned} \Gamma_{\alpha} &= (I_b, O_b) \\ \Gamma_u &= (O_b, I_b) \end{aligned}$$

■ Then

$$\begin{aligned} \alpha &= \Gamma_{\alpha} \beta \\ u &= \Gamma_u \beta \end{aligned}$$

■ Use these expressions and replace all  $\alpha, u$  with  $\beta$  .

# Standard Form

$$\min_{\beta} \left[ \frac{1}{2} \langle Q\beta, \beta \rangle + \langle \beta, q \rangle \right] \quad \text{subject to } V\beta \leq v$$

$$G\beta = g$$

■  $\ell_1$ -constrained LS can be expressed as

$$\begin{aligned} Q &= 2\Gamma_{\alpha}^{\top} X^{\top} X \Gamma_{\alpha} \\ q &= -2\Gamma_{\alpha}^{\top} X^{\top} y + \lambda \Gamma_u^{\top} \mathbf{1}_b \\ V &= \begin{pmatrix} -\Gamma_{\alpha} & -\Gamma_u \\ \Gamma_{\alpha} & -\Gamma_u \end{pmatrix} \\ v &= \mathbf{0}_{2b} \\ G &= O_{2b} \\ g &= \mathbf{0}_{2b} \end{aligned}$$

$$\beta = \begin{pmatrix} \alpha \\ u \end{pmatrix}$$

$$\begin{aligned} \Gamma_{\alpha} &= (I_b, O_b) \\ \Gamma_u &= (O_b, I_b) \end{aligned}$$

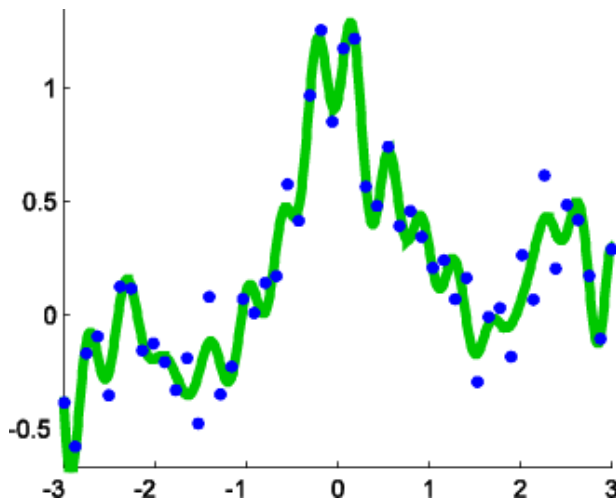
Proof: Homework!

# Example of Sparse Learning <sup>111</sup>

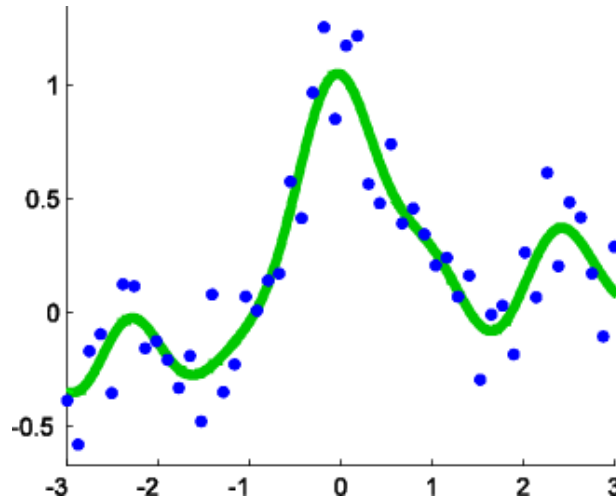
## ■ Gaussian kernel model:

$$f_{\alpha}(\mathbf{x}) = \sum_{i=1}^n \alpha_i \exp \left( -\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2} \right)$$

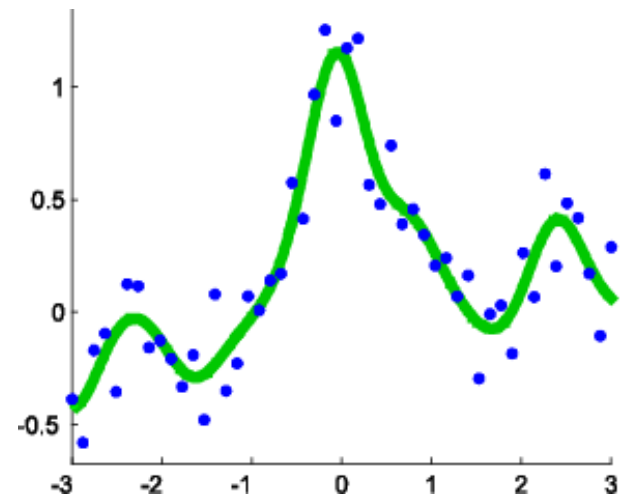
LS



$\ell_2$ -CLS



$\ell_1$ -CLS



■  $\ell_2$ -CLS and  $\ell_1$ -CLS give similar results.

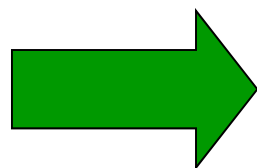
■ 27 out of 50 parameters are exactly zero in  $\ell_1$ .

# Feature Selection

- If  $\ell_1$ -CLS is combined with **linear model with respect to input**,

$$f_{\alpha}(x) = \alpha^{\top} x \quad x = (x^{(1)}, x^{(2)}, \dots, x^{(d)})^{\top}$$

some input variables are not used for prediction.



**Important features  
are automatically selected**

- **Example:** Gene selection
- Generally,  $2^d$  combinations need to be compared for feature selection (cf. subset LS).
- On the other hand,  $\ell_1$ -CLS only involves a continuous model parameter  $\lambda$ .



# Constrained LS

113

	Sparseness	Model parameter	Parameter learning
Subset LS	Yes	Combinatorial	Analytic (Linear)
Quadratically constrained LS	No	Continuous	Analytic (Linear)
$\ell_1$ constrained LS	Yes	Continuous	Iterative (Non-linear)

# Notification of Final Assignment

1. Apply supervised learning techniques to your data set and analyze it.
  2. Write your opinion about this course.
- Final report deadline: Aug 3<sup>rd</sup> (Fri.)
  - E-mail submission is also accepted!  
*sugi@cs.titech.ac.jp*

# Mini-Workshop on Data Mining<sup>115</sup>

- On July 10<sup>th</sup> and 24<sup>th</sup>, we will have a **mini-workshop on data mining**.
- Several students present their own data mining results.
- Those who give a talk at the workshop will have **very good grades!**

# Mini-Workshop on Data Mining<sup>116</sup>

- Application (just to declare that you want to give a presentation) deadline: **June 19<sup>th</sup>**.
- Presentation: **10-15 minutes (?)**.
  - Specification of your dataset
  - Methods used
  - Outcome
- Slides should be in English.
- Better to speak in English, but Japanese is also allowed.

# Homework

1. Derive the standard quadratic programming form of  $\ell_1$ -constrained LS.

$$\min_{\beta} \left[ \frac{1}{2} \langle Q\beta, \beta \rangle + \langle \beta, q \rangle \right]$$

subject to  $V\beta \leq v$

$G\beta = g$

$$\beta = \begin{pmatrix} \alpha \\ u \end{pmatrix}$$

$$\Gamma_{\alpha} = (I_b, O_b)$$

$$\Gamma_u = (O_b, I_b)$$

$$\begin{aligned} Q &= 2\Gamma_{\alpha}^{\top} X^{\top} X \Gamma_{\alpha} \\ q &= -2\Gamma_{\alpha}^{\top} X^{\top} y + \lambda \Gamma_u^{\top} \mathbf{1}_b \\ V &= \begin{pmatrix} -\Gamma_{\alpha} - \Gamma_u \\ \Gamma_{\alpha} - \Gamma_u \end{pmatrix} \\ v &= \mathbf{0}_{2b} \\ G &= O_{2b} \\ g &= \mathbf{0}_{2b} \end{aligned}$$

## Homework (cont.)

2. For your own toy 1-dimensional data, perform simulations using
- Gaussian kernel models
  - $\ell_1$ -constraint least-squares learning
- and analyze the results, e.g., by changing

- Target functions
- Number of samples
- Noise level

Use 5-fold cross-validation for choosing

- Width of Gaussian kernel
- Regularization parameter

Compare the results of QCLS and  $\ell_1$ CLS, e.g., in terms of sparseness and accuracy.