# Pattern Information Processing: Sparse Methods

Masashi Sugiyama

(Department of Computer Science)

Contact:   W8E-505

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

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

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

■ Non-linear learning

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

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

# l1-Constrained LS

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

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

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



$$\ell_1 - \text{norm} \quad \|\boldsymbol{\alpha}\|_1 = \sum_{i=1}^{b} |\alpha_i|$$

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

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.



$\ell_1$ constrained LS

Quadratically constrained LS

# How to Obtain A Solution

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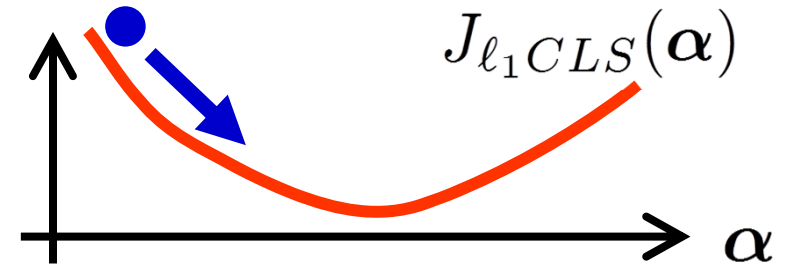
$$\hat{\boldsymbol{\alpha}}_{\ell_1 CLS} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^b}{\arg\min} J_{\ell_1 CLS}(\boldsymbol{\alpha})$$

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

# Gradient Descent



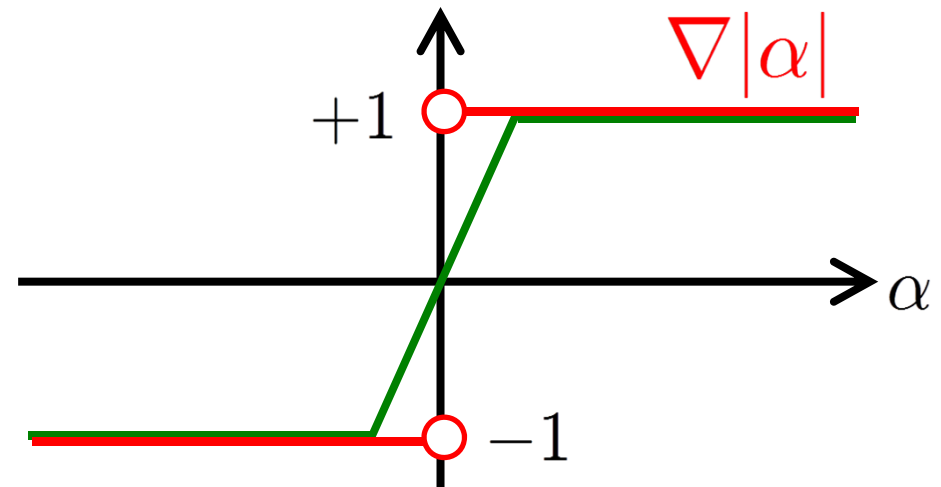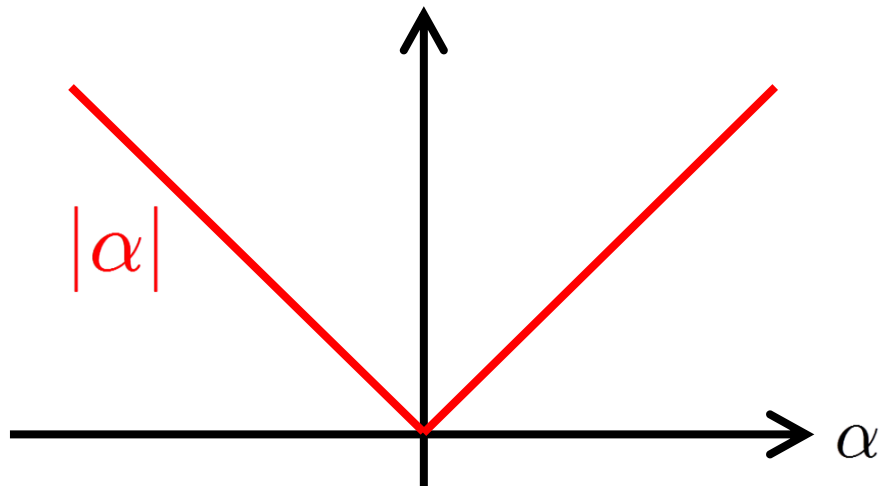- $$\boldsymbol{\alpha} \longleftarrow \boldsymbol{\alpha} - \epsilon \nabla J_{\ell_1 CLS}(\boldsymbol{\alpha})$$

- However, $\ell_1$-norm is not differentiable.
  - Use smooth approximation!



  - You may also use a quasi-Newton method.

# Quadratic Program

■ Use the following expression:

$$|\alpha| = \min_{u \in \mathbb{R}} u \quad \text{subject to} \quad -u \leq \alpha \leq u$$

■ Proof by contradiction:

- Let $\widehat{u} = \underset{u \in \mathbb{R}}{\operatorname{argmin}} \, u \quad \text{subject to} \quad -u \leq \alpha \leq u$.

- The constraint implies $\widehat{u} \geq |\alpha|$.

- Assume $\widehat{u} > |\alpha|$.

- Then such $\widehat{u}$ is not a solution because $\widetilde{u} = |\alpha|$ gives a smaller value.

- This implies that the solution should satisfy $\widehat{u} = |\alpha|$.

# How to Obtain A Solution (cont.)

$$\hat{\boldsymbol{\alpha}}_{\ell_1 CLS} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^b}{\arg\min} \, J_{LS}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|_1$$

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

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

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

Note: Inequality for vectors is component-wise

$$J_{LS}(\boldsymbol{\alpha}) = \sum_{i=1}^{n} (f_{\boldsymbol{\alpha}}(\boldsymbol{x}_i) - y_i)^2 = \|\boldsymbol{X}\boldsymbol{\alpha} - \boldsymbol{y}\|^2$$

# Linearly-Constrained Quadratic Program (QP)

- Standard optimization software can solve QP:

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

$$\text{subject to } \boldsymbol{H}\boldsymbol{\beta} \leq \boldsymbol{h}$$

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

# Transformation into Standard Form

- Let

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

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

- Then

$$\alpha = \Gamma_\alpha \beta$$
$$u = \Gamma_u \beta$$

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

# Standard Form

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

subject to $\boldsymbol{H}\beta \leq \boldsymbol{h}$
$\boldsymbol{G}\beta = \boldsymbol{g}$

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

$$\boldsymbol{Q} = 2\boldsymbol{\Gamma}_{\alpha}^{\top} \boldsymbol{X}^{\top} \boldsymbol{X} \boldsymbol{\Gamma}_{\alpha}$$

$$\boldsymbol{q} = -2\boldsymbol{\Gamma}_{\alpha}^{\top} \boldsymbol{X}^{\top} \boldsymbol{y} + \lambda \boldsymbol{\Gamma}_{u}^{\top} \boldsymbol{1}_b$$

$$\boldsymbol{H} = \begin{pmatrix} -\boldsymbol{\Gamma}_{\alpha} - \boldsymbol{\Gamma}_{u} \\ \boldsymbol{\Gamma}_{\alpha} - \boldsymbol{\Gamma}_{u} \end{pmatrix}$$

$$\boldsymbol{h} = \boldsymbol{0}_{2b}$$

$$\boldsymbol{G} = \boldsymbol{O}_{2b}$$

$$\boldsymbol{g} = \boldsymbol{0}_{2b}$$

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

$$\boldsymbol{\Gamma}_{\alpha} = (\boldsymbol{I}_b, \boldsymbol{O}_b)$$

$$\boldsymbol{\Gamma}_{u} = (\boldsymbol{O}_b, \boldsymbol{I}_b)$$

Proof: Homework!

# Example of Sparse Learning

- ■ Gaussian kernel model:

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

LS $\qquad$ $\ell_2$ -CLS $\qquad$ $\ell_1$ -CLS



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

- ■ 37 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,

$$f_{\boldsymbol{\alpha}}(\boldsymbol{x}) = \boldsymbol{\alpha}^\top \boldsymbol{x} \qquad \boldsymbol{x} = (x^{(1)}, x^{(2)}, \ldots, x^{(d)})^\top$$

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

|  | Sparseness | Model parameter | Parameter learning |
|---|---|---|---|
| Subset LS | Yes | Combina-torial | 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.

■ Final report deadline: Aug 1$^{st}$ (Fri.) 17:00

■ Bring your report to W8E-404.

# Mini-Workshop on Data Mining

- On July 15$^{th}$ and 22$^{nd}$, we will have a <span style="color:red">mini-workshop on data mining</span>.

- Several students present their own data mining results.

- Those who give a talk at the workshop will have <span style="color:red">very good grades!</span>

# Mini-Workshop on Data Mining

- ■ Application (just to declare that you want to give a presentation) deadline: July 1$^{st}$.
- ■ 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.

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

$$\text{subject to } \boldsymbol{H}\boldsymbol{\beta} \leq \boldsymbol{h}$$

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

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

$$\boldsymbol{\Gamma}_{\boldsymbol{\alpha}} = (\boldsymbol{I}_b, \boldsymbol{O}_b)$$

$$\boldsymbol{\Gamma}_{\boldsymbol{u}} = (\boldsymbol{O}_b, \boldsymbol{I}_b)$$

$$\boldsymbol{Q} = 2\boldsymbol{\Gamma}_{\boldsymbol{\alpha}}^{\top} \boldsymbol{X}^{\top} \boldsymbol{X} \boldsymbol{\Gamma}_{\boldsymbol{\alpha}}$$

$$\boldsymbol{q} = -2\boldsymbol{\Gamma}_{\boldsymbol{\alpha}}^{\top} \boldsymbol{X}^{\top} \boldsymbol{y} + \lambda \boldsymbol{\Gamma}_{\boldsymbol{u}}^{\top} \boldsymbol{1}_b$$

$$\boldsymbol{H} = \begin{pmatrix} -\boldsymbol{\Gamma}_{\boldsymbol{\alpha}} - \boldsymbol{\Gamma}_{\boldsymbol{u}} \\ \boldsymbol{\Gamma}_{\boldsymbol{\alpha}} - \boldsymbol{\Gamma}_{\boldsymbol{u}} \end{pmatrix}$$

$$\boldsymbol{h} = \boldsymbol{0}_{2b}$$

$$\boldsymbol{G} = \boldsymbol{O}_{2b}$$

$$\boldsymbol{g} = \boldsymbol{0}_{2b}$$

# 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.

# Solving QP Problems

- **R**: "quadprog"
- **Octave**: "qp"
- **MATLAB**: "quadprog"
  (you need Optimization Toolbox)
  - Various free software seems available, for example, "quadprog2".
    http://www.mathworks.com.au/matlabcentral/fileexchange/7860-quadprog2-convex-qp-solver