

Pattern Information Processing:⁴⁵ Constrained Least-Squares

Masashi Sugiyama
(Department of Computer Science)

Contact: W8E-505

sugi@cs.titech.ac.jp

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

Over-fitting

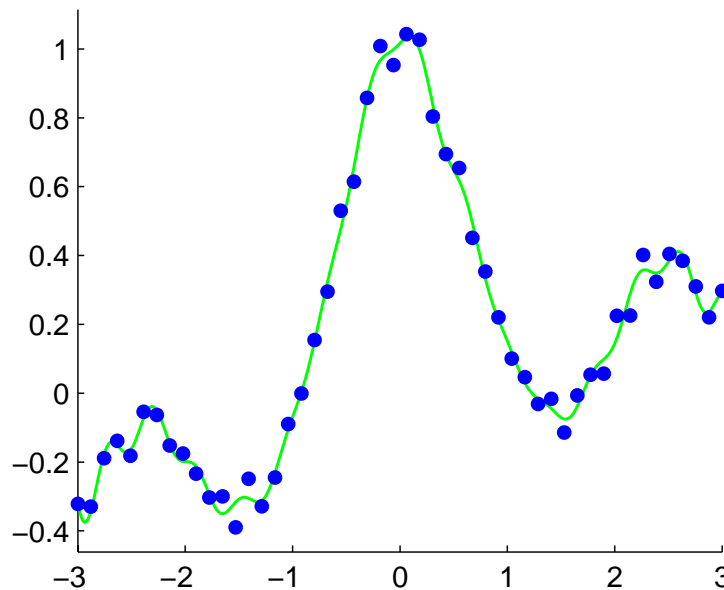
- LS is proved to be a good learning method:
 - Unbiased and BLUE in realizable cases
 - Asymptotically unbiased and asymptotically efficient in unrealizable cases
- However, the learned function can **over-fit** to noisy examples (e.g., when noise variance is large).

Over-fitting

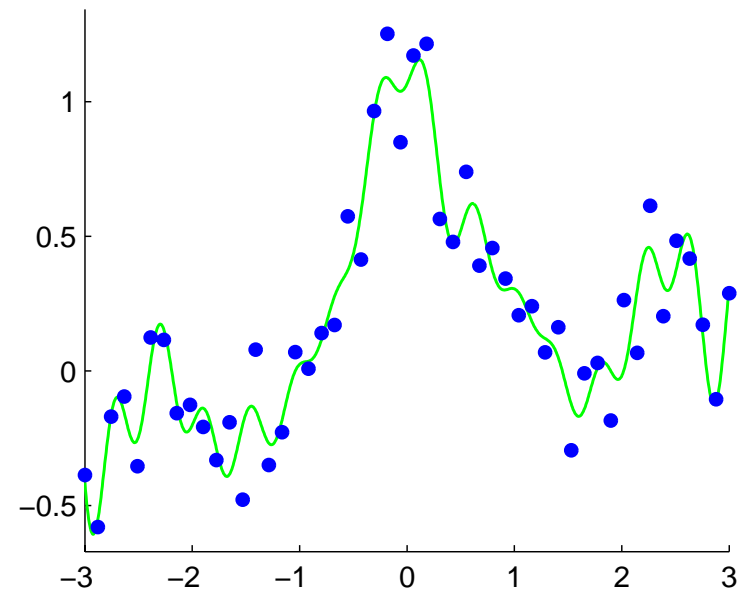
- Trigonometric polynomial model: $\hat{f}(x) = \sum_{i=1}^b \alpha_i \varphi_i(x)$

$$\varphi_i(x) = \{1, \sin x, \cos x, \dots, \sin 15x, \cos 15x\}$$

Small noise



Large noise



- In order to prevent over-fitting, model should be restricted.

Today's Plan

- Two approaches to restrict models:
 - Subspace LS
 - Quadratically constrained LS
- Sparseness and model choice

Subspace LS

- Restrict the search space within a **subspace**

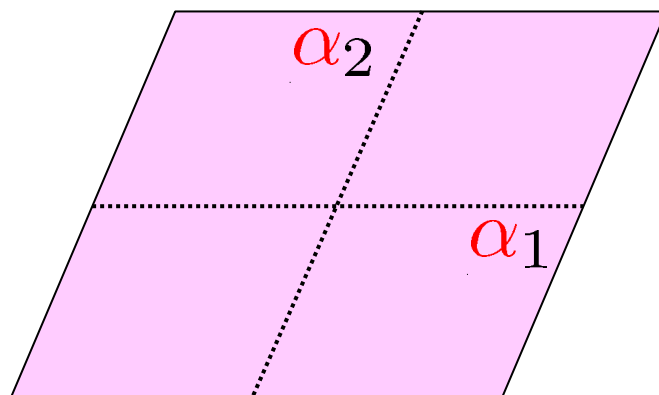
$$\hat{\alpha}_{SLS} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} J_{LS}(\alpha)$$

subject to $P\alpha = \alpha$

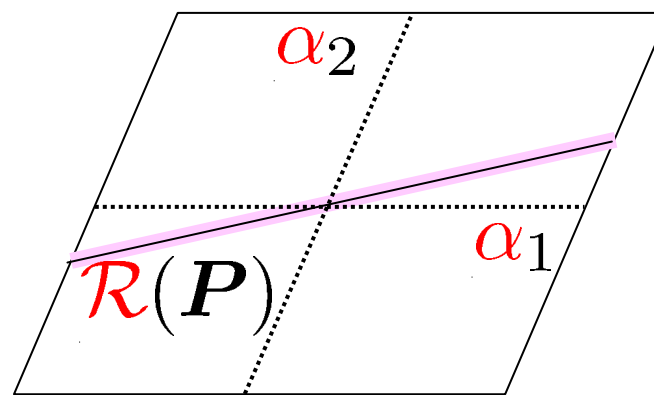
$$J_{LS}(\alpha) = \sum_{i=1}^n \left(\hat{f}(x_i) - y_i \right)^2$$

P : orthogonal projection
onto the subspace

$$P^2 = P \quad P^\top = P$$



Ordinary LS



Subspace LS

How to Obtain Solutions

■ Since

$$J_{LS}(\alpha) = \|X\alpha - y\|^2$$

just replacing X with XP gives a solution:

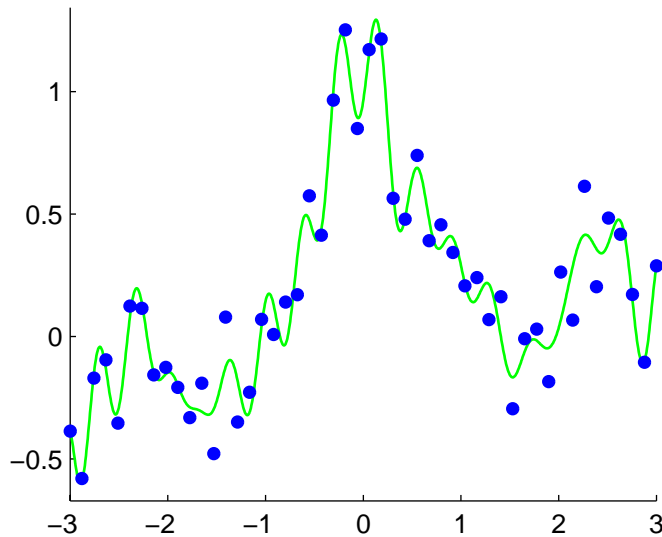
$$\begin{aligned} L_{SLS} &= (PX^\top XP)^\dagger PX^\top \\ &= (XP)^\dagger \end{aligned}$$

†: Moore-Penrose
generalized inverse

$$\begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}^\dagger = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/3 \end{pmatrix}$$

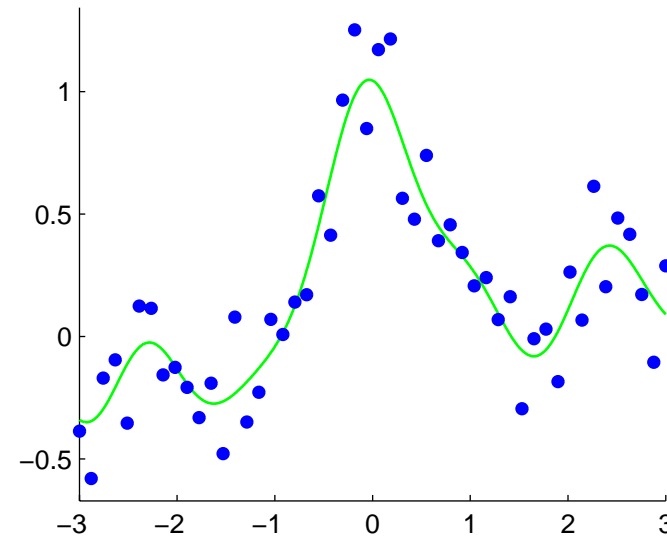
$$\begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix}^\dagger = \begin{pmatrix} 1/2 & 0 \\ 0 & 0 \end{pmatrix}$$

Example of SLS



Full LS

$1, \dots, \sin 15x, \cos 15x$



Subspace LS

$1, \dots, \sin 5x, \cos 5x$

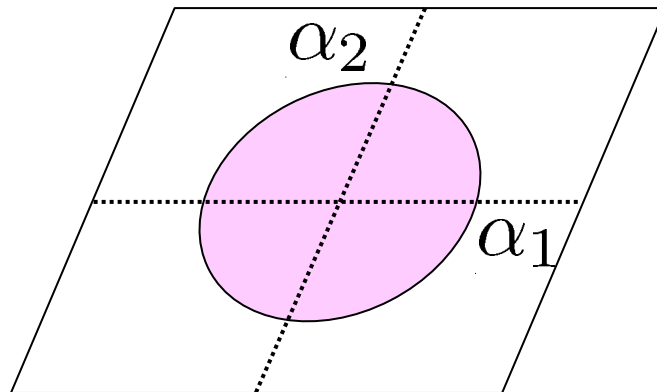
- Over-fit can be avoided by properly choosing the subspace.

Quadratically Constrained LS ⁵²

- Restrict the search space within a hypersphere.

$$\hat{\alpha}_{QCLS} = \underset{\alpha \in \mathbb{R}^b}{\operatorname{argmin}} J_{LS}(\alpha) \quad \text{subject to } \|\alpha\|^2 \leq C$$

$$C \geq 0$$



How to Obtain Solutions

- Lagrangian:

$$J_{QCLS}(\boldsymbol{\alpha}, \lambda) = J_{LS}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|^2$$

- λ : Lagrange multiplier

- Karush-Kuhn-Tucker (KKT) condition: the solution $\hat{\boldsymbol{\alpha}}_{QCLS}$ satisfies with some λ^*

1.
$$\frac{\partial J_{QCLS}(\hat{\boldsymbol{\alpha}}_{QCLS}, \lambda^*)}{\partial \boldsymbol{\alpha}} = \mathbf{0}$$

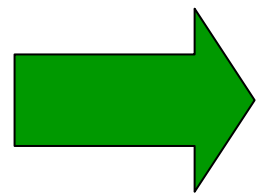
2.
$$\lambda^* \geq 0$$

3.
$$\|\hat{\boldsymbol{\alpha}}_{QCLS}\|^2 - C \leq 0$$

4.
$$\lambda^* (\|\hat{\boldsymbol{\alpha}}_{QCLS}\|^2 - C) = 0$$

How to Obtain Solutions (cont.)⁵⁴

■ $\frac{\partial J_{QCLS}(\hat{\alpha}_{QCLS}, \lambda^*)}{\partial \alpha} = \mathbf{0}$



$$\hat{\alpha}_{QCLS} = L_{QCLS} \mathbf{y}$$

$$L_{QCLS} = (\mathbf{X}^\top \mathbf{X} + \lambda^* \mathbf{I})^{-1} \mathbf{X}^\top$$

■ λ^* is obtained from $\lambda^* (\|\hat{\alpha}_{QCLS}\|^2 - C) = 0$

■ In practice, we start from $\lambda (\geq 0)$ and solve

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

$$J_{QCLS}(\alpha) = J_{LS}(\alpha) + \lambda \|\alpha\|^2$$

Interpretation of QCLS

- QCLS tries to avoid **overfitting** by adding penalty (**regularizer**) to the “goodness-of-fit” term.

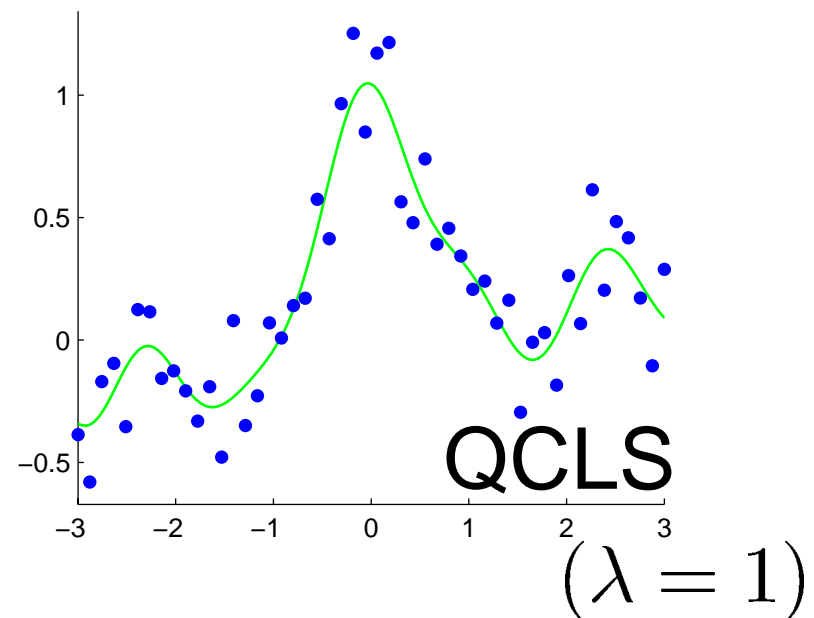
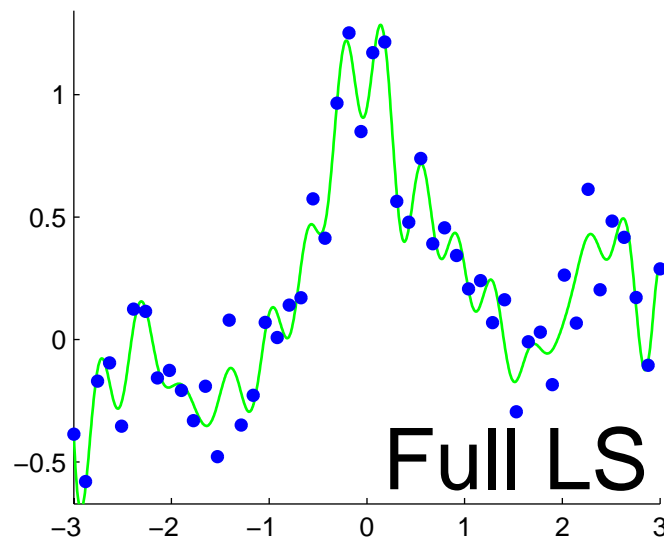
$$J_{QCLS}(\boldsymbol{\alpha}) = \underbrace{J_{LS}(\boldsymbol{\alpha})}_{\text{Goodness of fit}} + \underbrace{\lambda \|\boldsymbol{\alpha}\|^2}_{\text{Penalty (regularizer)}}$$

- For this reason, QCLS is also called **quadratically regularized LS**.
- λ is called the **regularization parameter**.

Example of QCLS

- Gaussian kernel model: $\hat{f}(x) = \sum_{i=1}^n \alpha_i K(x, x_i)$

$$K(x, x') = \exp(-\|x - x'\|^2/2)$$



- Over-fit can be avoided by properly choosing the regularization parameter.

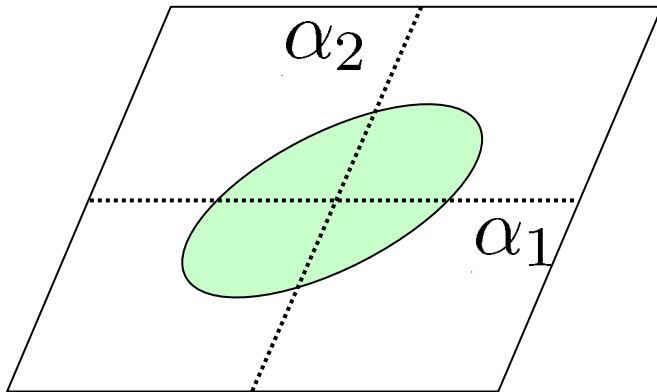
Generalization

- Restrict the search space within a **hyper-ellipsoid**.

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

subject to $\langle \mathbf{R}\alpha, \alpha \rangle \leq C$

$$C \geq 0$$

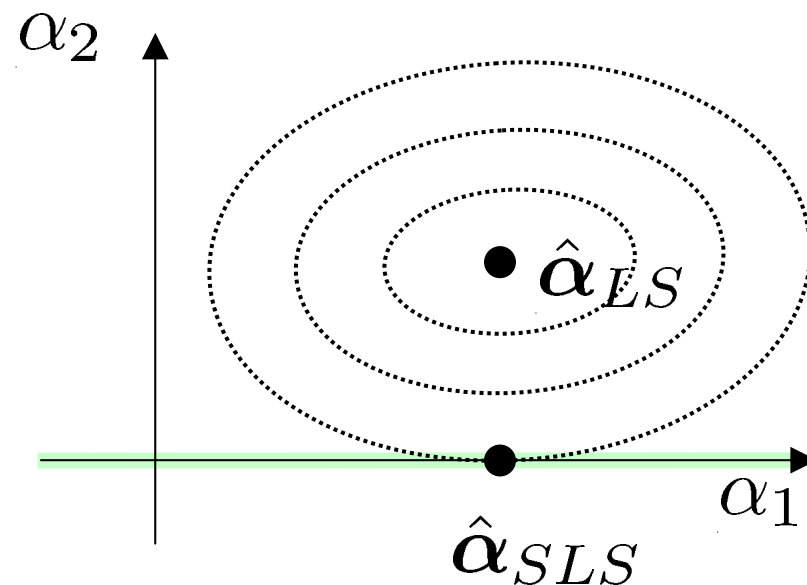


\mathbf{R} : Positive semi-definite matrix
 (“**regularization matrix**”)
 $\forall \alpha, \langle \mathbf{R}\alpha, \alpha \rangle \geq 0$

$$\mathbf{L}_{QCLS} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{R})^{-1} \mathbf{X}^\top$$

Sparseness of Solution

- In SLS, if the subspace is spanned by a subset of basis functions $\{\varphi_i(\mathbf{x})\}_{i=1}^b$, some of the parameters $\{\alpha_i\}_{i=1}^b$ are exactly zero.



Model Choice

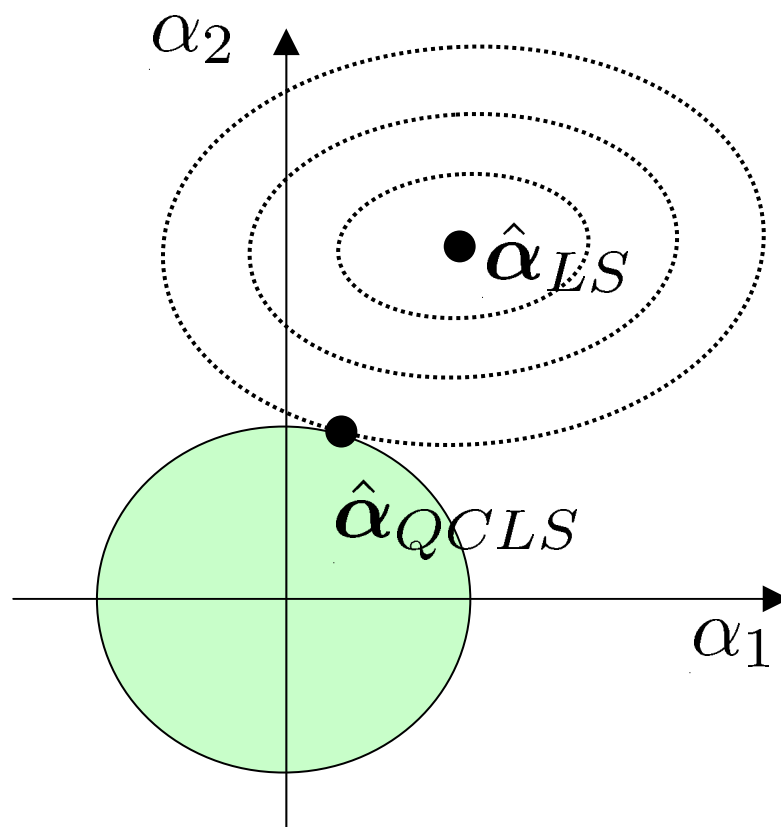
- **Sparse solution** is computationally advantageous in calculating the output values.

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

- However, the possible choices of such subspaces are **combinatorial**: 2^b
- Infeasible to search the best subset.

Property of QCLS

- In QCLS, model choice is continuous: λ
- However, solution is not generally sparse.



Homework

1. Prove that the solution of

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

subject to $\langle \mathbf{R}\boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle \leq C$

is given as

$$\hat{\boldsymbol{\alpha}}_{QCLS} = \mathbf{L}_{QCLS} \mathbf{y}$$

$$\mathbf{L}_{QCLS} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{R})^{-1} \mathbf{X}^\top$$

Homework (cont.)

2. For your own toy 1-dimensional data, perform simulations using

- Gaussian kernel models
- Quadratically-constrained least-squares learning

and analyze the results, e.g., changing

- Target functions
- Number of samples
- Noise level
- Width of Gaussian kernel
- Regularization parameter/matrix

Suggestions

- Please look for software which can solve
 - Linearly constrained quadratic programming

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

subject to $\mathbf{V}\beta \leq \mathbf{v}$ and $\mathbf{G}\beta = \mathbf{g}$

- Linearly constrained linear programming

$$\min_{\beta} \langle \beta, \mathbf{q} \rangle$$

subject to $\mathbf{V}\beta \leq \mathbf{v}$ and $\mathbf{G}\beta = \mathbf{g}$

- The software is not necessarily sophisticated; just elementary one is enough.