

Pattern Information Processing:⁵⁰ 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

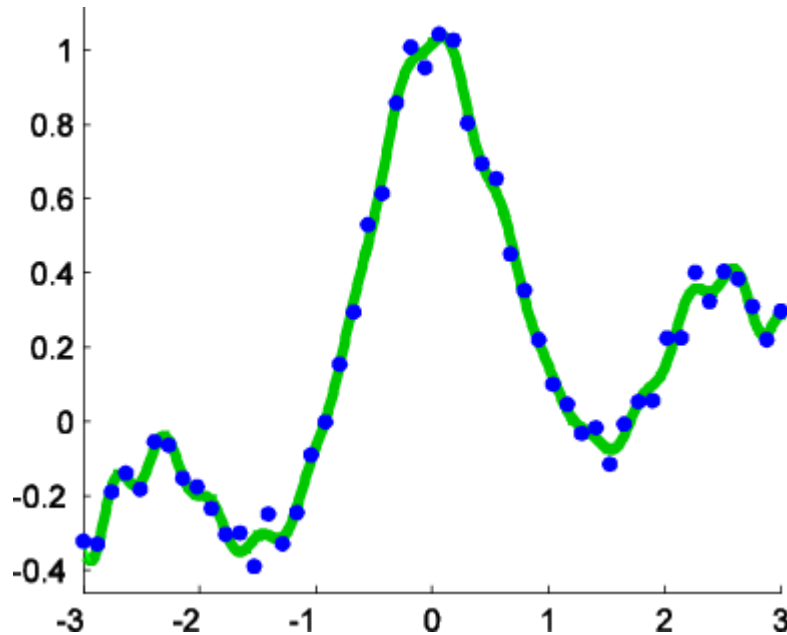
- LS was proved to be a good learning method:
 - Unbiased and BLUE in realizable cases.
 - Asymptotically unbiased and asymptotically efficient in unrealizable cases.
- However, a learned function can **over-fit** to noisy examples (e.g., when the noise level is high).

Over-fitting

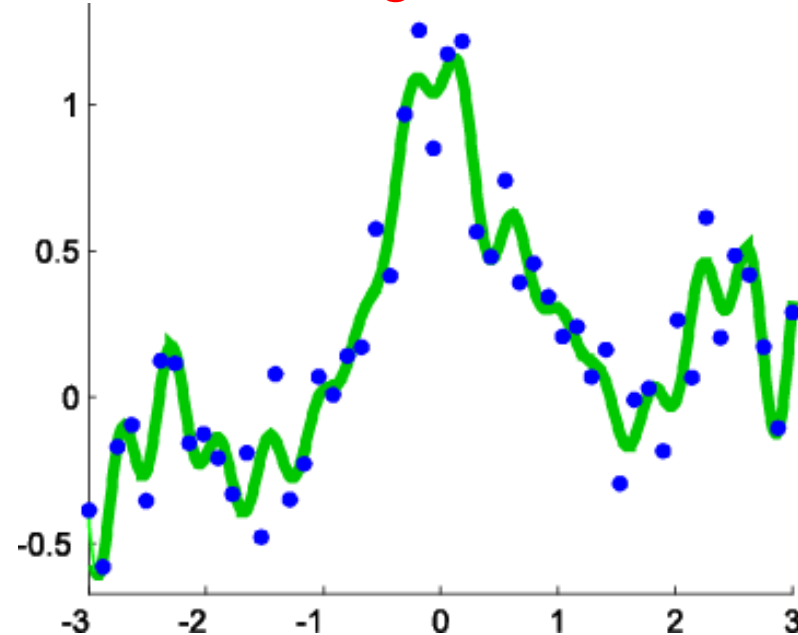
■ Trigonometric polynomial model: $f_{\alpha}(x) = \sum_{i=1}^b \alpha_i \varphi_i(x)$

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

Small noise



Large noise



■ In order to prevent over-fitting, model (search space) should be restricted appropriately.

Today's Plan

- Two approaches to restricting models:
 - Subspace LS
 - Quadratically constrained LS
- Sparseness and model choice.
- We focus on linear/kernel models.

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

Subspace LS

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

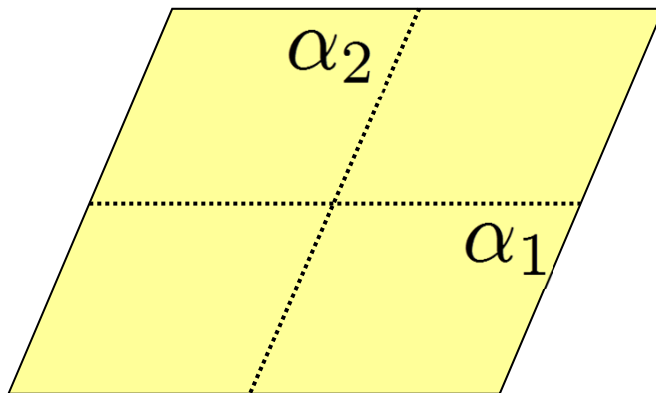
$$\hat{\alpha}_{SLS} = \underset{\alpha \in \mathbb{R}^b}{\operatorname{argmin}} J_{LS}(\alpha)$$

subject to $P\alpha = \alpha$

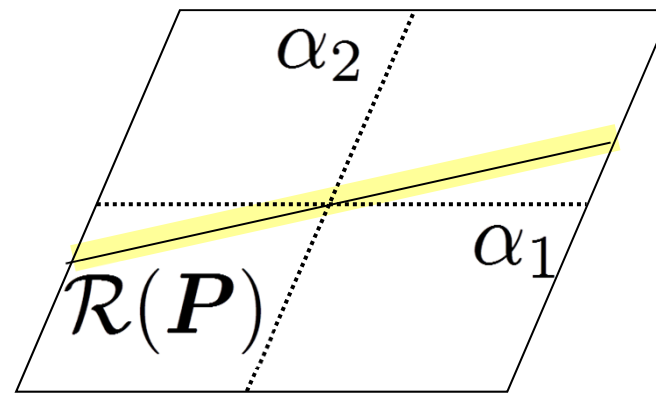
P : orthogonal projection
onto a subspace

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

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



Ordinary LS



Subspace LS

How to Obtain A Solution

■ Since

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

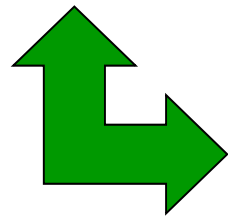
just replacing X with XP gives a solution:

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

$$X_{i,j} = \varphi_j(\mathbf{x}_i)$$

■ \dagger : Moore-Penrose generalized inverse

$$B = A^\dagger$$

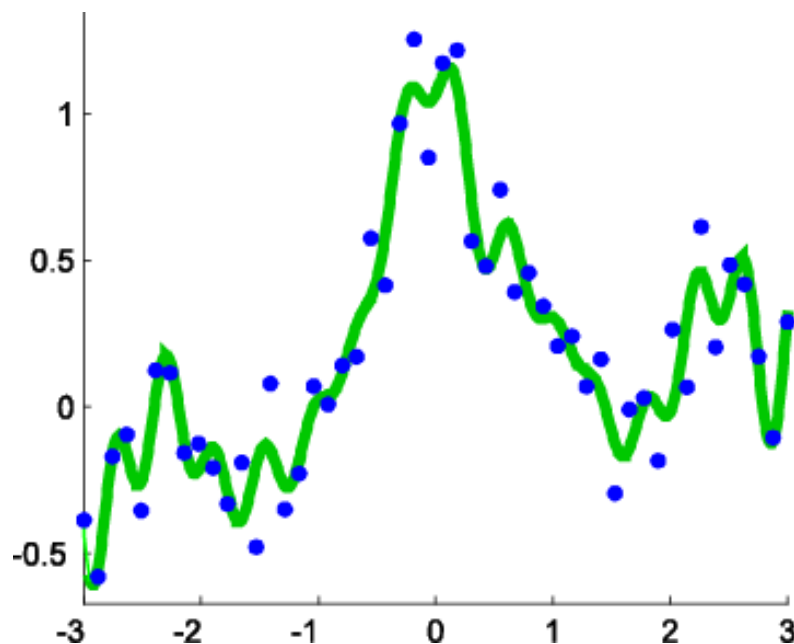


$$\left\{ \begin{array}{l} ABA = A \\ BAB = B \\ (AB)^\top = AB \\ (BA)^\top = BA \end{array} \right.$$

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

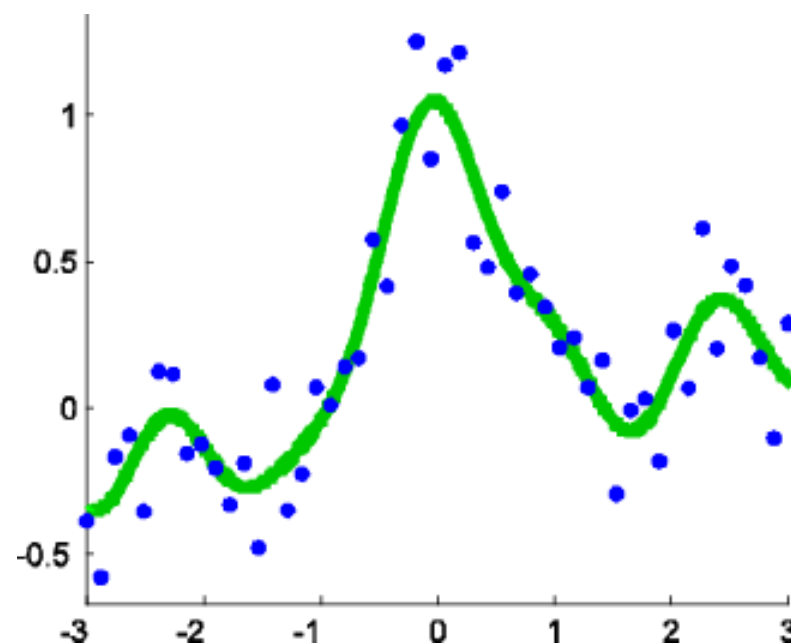
$$\begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix}^\dagger = \begin{pmatrix} 1/2 & 0 \\ 0 & 0 \end{pmatrix}$$

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 a subspace.

Principal Component Regression⁵⁷

- Choose the **maximum-variance subspace**:

$$P = \sum_{k=1}^m \phi_k \phi_k^\top$$

- Eigendecomposition of covariance matrix:

$$X^\top H H X \phi = \lambda \phi$$

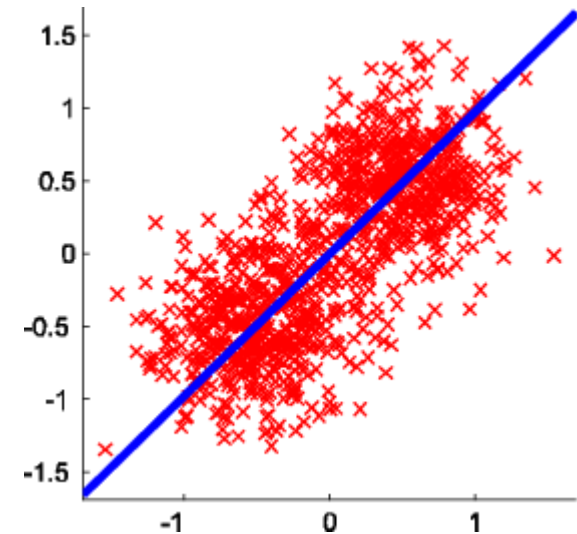
$$H = I_n - \frac{1}{n} \mathbf{1}_{n \times n}$$

- Eigenvalues: $\lambda_1 \geq \dots \geq \lambda_b$

- Eigenvectors: ϕ_1, \dots, ϕ_b

I_n : n -dimensional identity matrix

$\mathbf{1}_{n \times n}$: $n \times n$ matrix with all ones

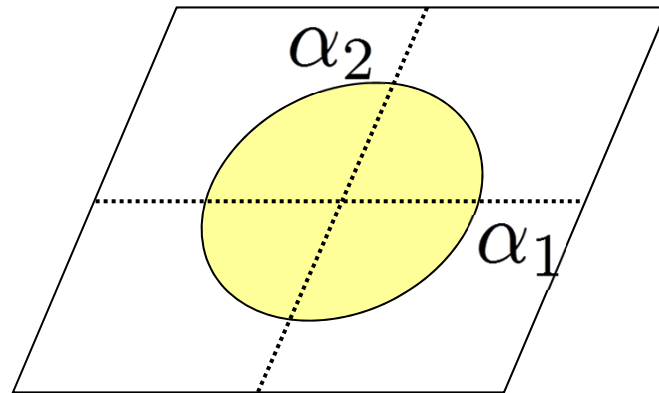


Quadratically Constrained LS ⁵⁸

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

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

$$C \geq 0$$



How to Obtain A Solution

■ Lagrangian:

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

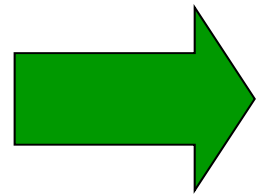
- λ : Lagrange multiplier

■ Karush-Kuhn-Tucker (KKT) condition: for some $\lambda^* \geq 0$, the solution $\hat{\boldsymbol{\alpha}}_{QCLS}$ satisfies

- $\frac{\partial L(\hat{\boldsymbol{\alpha}}_{QCLS}, \lambda^*)}{\partial \boldsymbol{\alpha}} = \mathbf{0}$
- $\|\hat{\boldsymbol{\alpha}}_{QCLS}\|^2 - C \leq 0$
- $\lambda^* (\|\hat{\boldsymbol{\alpha}}_{QCLS}\|^2 - C) = 0$

How to Obtain A Solution (cont.)⁶⁰

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



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

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

- We still need to determine λ^* , but this is not straightforward.
- In practice, we start from setting $\lambda (\geq 0)$ 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 a penalty (**regularizer**) to the “goodness-of-fit” term.

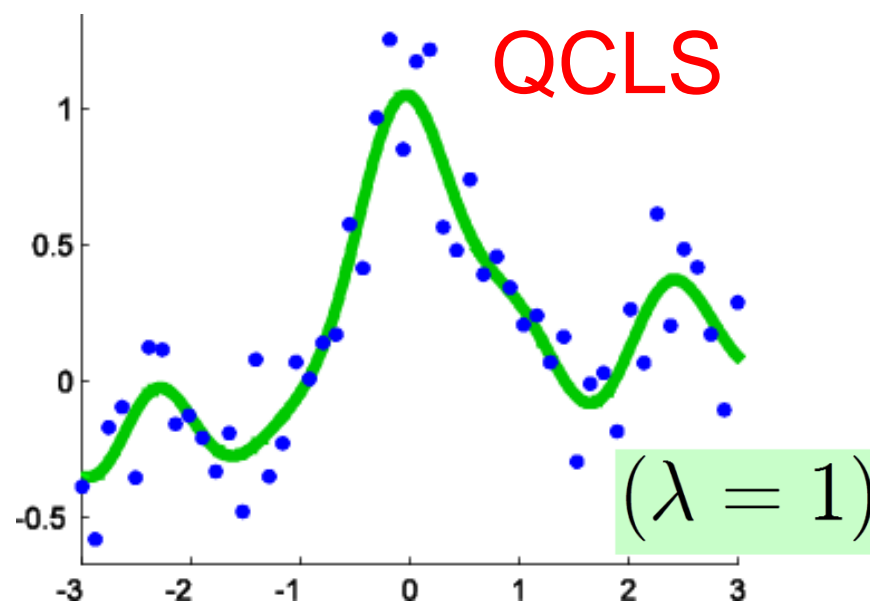
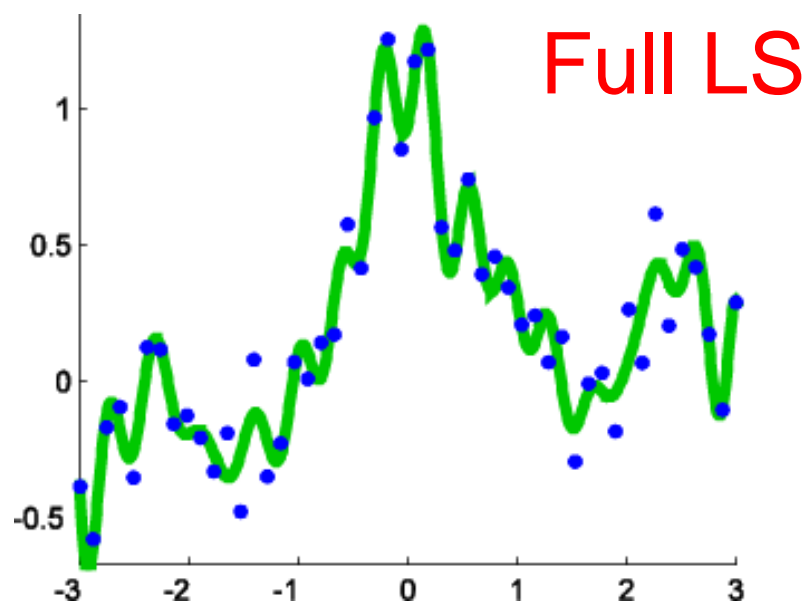
$$J_{QCLS}(\boldsymbol{\alpha}) = \underbrace{J_{LS}(\boldsymbol{\alpha})}_{\text{Goodness of fit}} + \underbrace{\lambda \|\boldsymbol{\alpha}\|^2}_{\text{Penalty (regularizer)}}$$

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

Example of QCLS

■ Gaussian kernel model:

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



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

Generalization

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

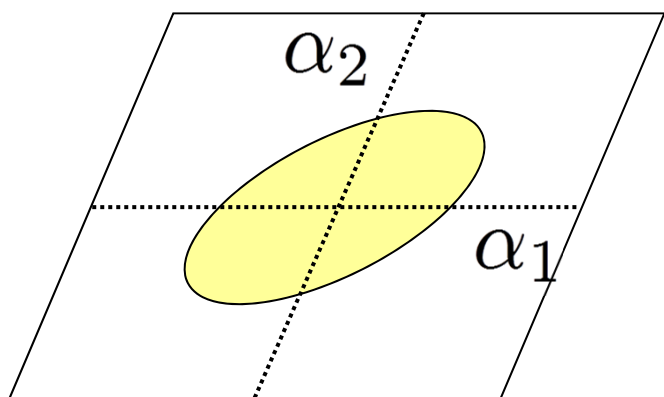
$$\hat{\alpha}_{QCLS} = \underset{\alpha \in \mathbb{R}^b}{\operatorname{argmin}} J_{LS}(\alpha)$$

subject to $\langle R\alpha, \alpha \rangle \leq C$

$$C \geq 0$$

R : Positive semi-definite matrix
 (“**regularization matrix**”)

$$\forall \alpha, \langle R\alpha, \alpha \rangle \geq 0$$

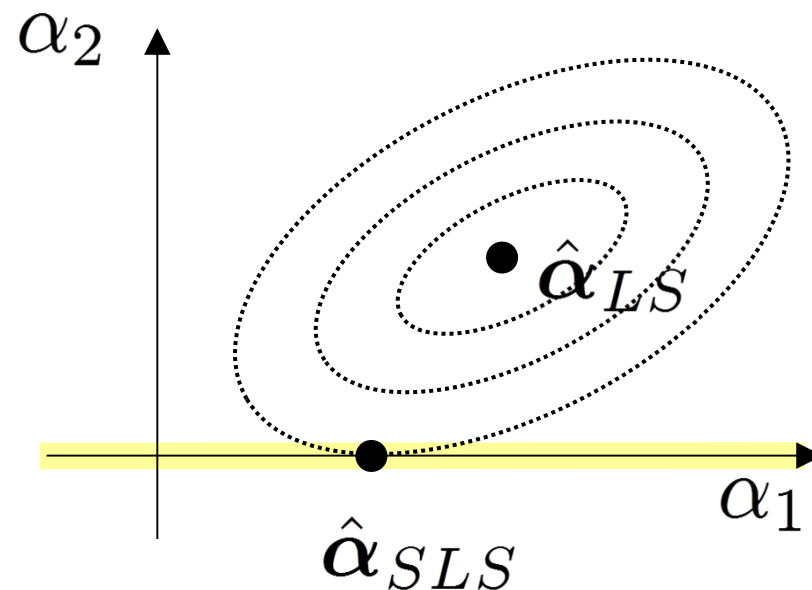


- Solution:** (proof is homework!)

$$L_{QCLS} = (X^T X + \lambda R)^{-1} X^T$$

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 $\{\hat{\alpha}_i\}_{i=1}^b$ are **exactly zero**.



Sparseness and Model Choice⁶⁵

- Having **sparsity** is computationally attractive:

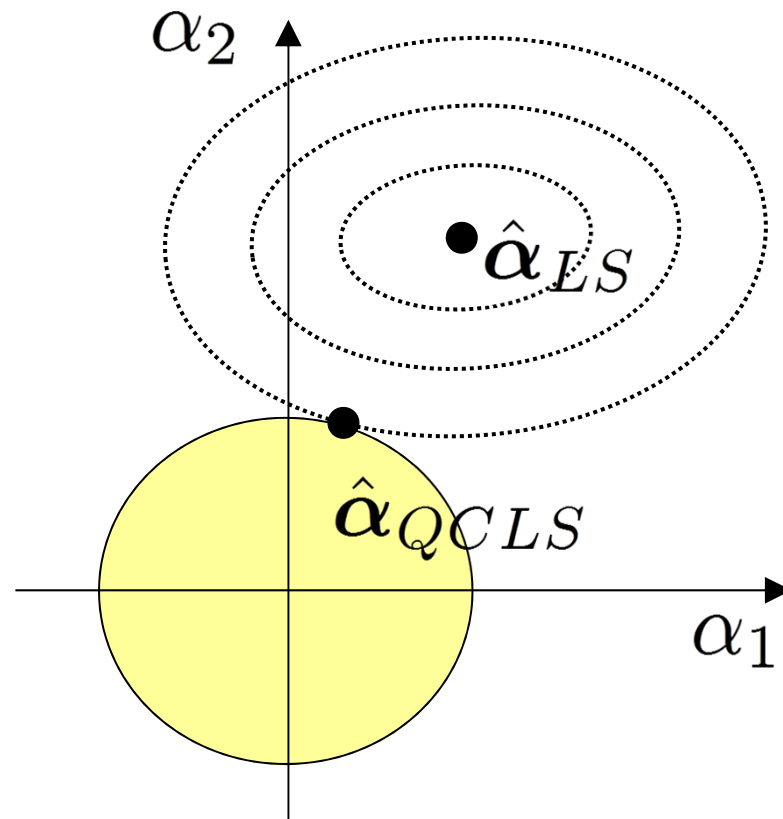
- Calculating output values is easier.

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

- Computing the solution can potentially be easier.
- However, the number of possible subsets is **combinatorial**, 2^b .
 - It is computationally infeasible to find the best subset if b is large.

Sparseness and Model Choice⁶⁶ (cont.)

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



Homework

1. Prove that the solution of

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

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

is given by

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

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

Homework (cont.)

68

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