Pattern Information Processing パターン情報処理

Masashi Sugiyama (Department of Computer Science) 杉山 将(計算工学専攻)

Contact: W8E-505 sugi@cs.titech.ac.jp

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

Learning Methods

Overfitting

- Subset LS learning
- Sparseness
- Regularization
- Small exercise

Over-fitting

- When noise is large, the learned function is over-fitted to noisy examples.
- In order to prevent this phenomenon, model should be restricted.



Moore-Penrose Generalized Inverse

$$egin{array}{rcl} egin{array}{ccc} egin{array}{ccc} AXA &=& A\ XAX &=& X\ (AX)^{ op} &=& AX\ (XA)^{ op} &=& XA \end{array}$$

 ${}^{\exists}A^{-1} \Longrightarrow A^{\dagger} = A^{-1}$ ${}^{\exists}(A^{\top}A)^{-1} \Longrightarrow A^{\dagger} = (A^{\top}A)^{-1}A^{\top}$ ${}^{\exists}(AA^{\top})^{-1} \Longrightarrow A^{\dagger} = A^{\top}(AA^{\top})^{-1}$

Example of SLS



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

6

Sparseness of Solution

- If the subspace is spanned by a subset of basis functions $\{\varphi_i(x)\}_{i=1}^p$, some of the parameters $\{\alpha_i\}_{i=1}^p$ are zero.
- Sparse solution is computationally advantageous in calculating the output values.

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

However, the choice of subspaces is discrete.
Combinatorial explosion! 2^p

Quadratically Constrained LS

8

C > 0

Restrict the search space within a hyper-sphere. $\hat{\alpha}_{QCLS} = \operatorname*{argmin}_{\mathcal{A}} J_{LS}(\boldsymbol{\alpha})$ subject to $\|\boldsymbol{\alpha}\|^2 \leq C$



How to Obtain Solutions

Lagrangian:

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

 $\lambda (\geq 0) : \text{Lagrange multiplier}$ In practice, we start from $\lambda (\geq 0)$ and solve $\hat{\alpha}_{QCLS} = \operatorname*{argmin}_{\boldsymbol{\alpha}} J_{QCLS}(\boldsymbol{\alpha})$

Solution is given by

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

It is often called quadratically regularized LS.

Interpretation of QCLS

QCLS tries to avoid overfitting by adding penalty to the "goodness-of-fit" term.

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

Good- Penalty
ness of fit

Example of QCLS

Gaussian kernel model: $\hat{f}(\boldsymbol{x}) = \sum \alpha_i K(\boldsymbol{x}, \boldsymbol{x}_i)$ $K(x, x') = \exp(-||x - x'||^2/2)$ i=10.5 0.5 0 -0.5 -0.5 -2 -3 _1 0 -2 0 -1 $(\lambda = 1)$ Over-fit can be avoided by properly choosing the regularization factor.

Property of QCLS

Choice of models is continuous: λ

However, solution is not generally sparse.

Generalization

Restrict the search space within a hyper-ellipsoid.

$\hat{\boldsymbol{\alpha}}_{QCLS} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} J_{LS}(\boldsymbol{\alpha})$ subject to $\langle \boldsymbol{R}\boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle \leq C$



 $C \ge 0$ **R** : Positive semidefinite matrix $\forall \alpha, \langle R\alpha, \alpha \rangle \ge 0$

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



Prove that the solution of

$$\hat{\boldsymbol{\alpha}}_{QCLS} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} J_{LS}(\boldsymbol{\alpha})$$

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

is given by

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