Pattern Information Processing¹²⁶ Support Vector Machines

> Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/

(Binary) Classification Problem¹²⁷

- Output values are $y_i = \pm 1$
- We want to predict whether output values of unlearned input points are positive/negative.
- Multi-class problem can be transferred to several binary classification problems:
 - One-versus-rest
 - One-versus-one

(Binary) Classification Problem¹²⁸

In classification, we may still use the same learning methods, e.g., quadraticallyconstrained least-squares:

 $\hat{\boldsymbol{\alpha}}_{QCLS} = \operatorname*{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^b} \left[J_{LS}(\boldsymbol{\alpha}) + \lambda \langle \boldsymbol{R} \boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle \right]$

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

Prediction:

$$\widehat{y} = \operatorname{sign}\left(\widehat{f}(\boldsymbol{x})\right)$$

0/1-Loss

- In classification, only the sign of the learned function is used.
- It is natural to use 0/1-loss instead of squared-loss $J_{LS}(\alpha)$:

$$J_{0/1}(\boldsymbol{\alpha}) = \frac{1}{2} \sum_{i=1}^{n} \left(1 - \operatorname{sign}\left(u_{i}\right)\right)$$
$$u_{i} = \hat{f}(\boldsymbol{x}_{i})y_{i}$$

 $J_{0/1}(\alpha)$ corresponds to the number of misclassified samples (thus natural).

Hinge-Loss

However, $J_{0/1}(\alpha)$ is non-convex so we may not obtain the global minimizer.

Use hinge-loss as an approximation:



Hinge-loss is a tighter upper bound on 0/1loss than squared-loss (thus better?). How to Obtain Solutions $\hat{\alpha}_{SVM} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^{b}}{\operatorname{argmin}} \begin{bmatrix} J_{H}(\boldsymbol{\alpha}) + \lambda \langle \boldsymbol{R}\boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle \end{bmatrix}$ $J_{H}(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \max(0, 1 - u_{i})$ 131

How to deal with "max"? Use following lemma:

Lemma
$$\max(0, 1 - u) = \min_{\xi \in \mathbb{R}} \xi$$

subject to $\xi \ge 1 - u$
 $\xi \ge 0$

Proof: Constraint is $\xi \ge \max(0, 1 - u)$, so $\min_{\xi \in \mathbb{R}} \xi = \max(0, 1 - u)$ Q.E.D. How to Obtain Solutions (cont.)³² So we have $J_H(\alpha) = \min_{\boldsymbol{\xi} \in \mathbb{R}^n} \langle \mathbf{1}_n, \boldsymbol{\xi} \rangle$ subject to $\boldsymbol{\xi} \ge \mathbf{1}_n - \boldsymbol{u}$ $\boldsymbol{\xi} \ge \mathbf{0}_n$

Then $\hat{\alpha}_{SVM}$ is given as

$$\hat{\boldsymbol{\alpha}}_{SVM} = \operatorname*{argmin}_{\boldsymbol{lpha} \in \mathbb{R}^{b}, \boldsymbol{\xi} \in \mathbb{R}^{n}} \begin{bmatrix} \langle \mathbf{1}_{n}, \boldsymbol{\xi} \rangle + \lambda \langle \boldsymbol{R} \boldsymbol{lpha}, \boldsymbol{lpha} \rangle \end{bmatrix}$$

subject to $\boldsymbol{\xi} \geq \mathbf{1}_{n} - \boldsymbol{u}$
 $\boldsymbol{\xi} \geq \mathbf{0}_{n}$

Support Vector Machines

133

We focus on the following setting:

n

•
$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i K(\boldsymbol{x}, \boldsymbol{x}_i) + b$$

• $\boldsymbol{R} = \boldsymbol{K}$ $\boldsymbol{K}_{i,j} = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$ $\exists \boldsymbol{K}^-$

Putting $\lambda = (2C)^{-1}$ (convention), we have

$$\hat{\boldsymbol{\alpha}}_{SVM} = \underset{\boldsymbol{\alpha}, \boldsymbol{\xi} \in \mathbb{R}^{n}, b \in \mathbb{R}}{\operatorname{argmin}} \begin{bmatrix} C \langle \boldsymbol{1}_{n}, \boldsymbol{\xi} \rangle + \frac{1}{2} \langle \boldsymbol{K} \boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle \end{bmatrix}$$
subject to $\boldsymbol{\xi} \geq \boldsymbol{1}_{n} - \boldsymbol{u}$
$$\boldsymbol{\xi} \geq \boldsymbol{0}_{n}$$
$$u_{i} = \hat{f}(\boldsymbol{x}_{i})y_{i}$$

How to Obtain Solutions (cont.)³⁴

- SVM solution can be obtained by solving a linearly constrained quadratic programming problem.
- However, we need to optimize 2n+1 variables, which could be time consuming.
- Consider a dual formulation.

How to Obtain Solutions (cont.)³⁵

Lagrangian:

$$L(\boldsymbol{\alpha}, b, \boldsymbol{\xi}, \boldsymbol{\beta}, \boldsymbol{\gamma}) = C\langle \mathbf{1}_n, \boldsymbol{\xi} \rangle + \frac{1}{2} \langle \boldsymbol{K} \boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle$$
$$-\langle \boldsymbol{\beta}, \boldsymbol{\xi} + \boldsymbol{u} - \mathbf{1}_n \rangle - \langle \boldsymbol{\gamma}, \boldsymbol{\xi} \rangle$$

β, γ :Lagrange multiplier
 Wolfe-dual problem is simpler:

$$\min_{\substack{\boldsymbol{\alpha},\boldsymbol{\xi}\in\mathbb{R}^n,b\in\mathbb{R}\\ \boldsymbol{\alpha},\boldsymbol{\xi}\in\mathbb{R}^n,b\in\mathbb{R}}} \begin{bmatrix} C\langle \mathbf{1}_n,\boldsymbol{\xi}\rangle + \frac{1}{2}\langle \boldsymbol{K}\boldsymbol{\alpha},\boldsymbol{\alpha}\rangle \end{bmatrix} = \max_{\substack{\boldsymbol{\beta},\boldsymbol{\gamma}\in\mathbb{R}^n\\ \boldsymbol{\beta},\boldsymbol{\gamma}\in\mathbb{R}^n\\ \boldsymbol{\beta},\boldsymbol{\gamma$$

How to Obtain Solutions (cont.)³⁶ Constraints yield • $\frac{\partial L}{\partial \boldsymbol{\xi}} = C \mathbf{1}_n - \boldsymbol{\beta} - \boldsymbol{\gamma} = \mathbf{0}_n$ $oldsymbol{\gamma} \geq oldsymbol{0}_n$ $\beta > \mathbf{0}_n$ $v_i = \beta_i y_i$ • $\frac{\partial L}{\partial \alpha} = K\alpha - Kv = \mathbf{0}_n$ $\boldsymbol{\square} \boldsymbol{\square} \boldsymbol{\alpha} = v$

How to Obtain Solutions (cont.)³⁷

Then Wolfe-dual problem is simplified as

$$\widehat{\boldsymbol{\beta}}_{SVM} = \underset{\boldsymbol{\beta} \in \mathbb{R}^n}{\operatorname{argmax}} \left[\sum_{i=1}^n \beta_i - \frac{1}{2} \sum_{i,j=1}^n \beta_i \beta_j y_i y_j \boldsymbol{K}_{i,j} \right]$$

subject to $\mathbf{0}_n \leq \boldsymbol{\beta} \leq C \mathbf{1}_n$ $\langle \boldsymbol{y}, \boldsymbol{\beta} \rangle = 0$

So we have

$$[\hat{\boldsymbol{\alpha}}_{SVM}]_i = [\hat{\boldsymbol{\beta}}_{SVM}]_i y_i$$

How to Obtain Solutions (cont.)³⁸

Standard form of the Wolfe-dual problem is

$$egin{aligned} &\min_{oldsymbol{eta}\in\mathbb{R}^n}\left[rac{1}{2}\langleoldsymbol{Q}eta,oldsymbol{eta}
ight
angle+\langleoldsymbol{eta},oldsymbol{q}
ight
angle\ & ext{ subject to }oldsymbol{V}oldsymbol{eta}\leqoldsymbol{v}\ & ext{ }Goldsymbol{eta}=oldsymbol{g} \end{aligned}$$

$$egin{aligned} m{Q}_{i,j} &= m{K}_{i,j} y_i y_j & m{V} &= egin{pmatrix} -m{I}_n \ m{I}_n \end{pmatrix} & m{G} &= m{y}^{ op} \ m{q} &= -m{1}_n & m{v} &= egin{pmatrix} m{0}_n \ Cm{1}_n \end{pmatrix} & m{g} &= m{0}_n \end{aligned}$$

How to Obtain Solutions (cont.)³⁹

 \hat{b}_{SVM} is computed using any j such that $0 < [\widehat{\boldsymbol{\beta}}_{SVM}]_i < C$ as $\hat{b}_{SVM} = y_j - \sum [\hat{\alpha}_{SVM}]_i K(\boldsymbol{x}_j, \boldsymbol{x}_i)$ i=1In practice, we may take the average: $\hat{b}_{SVM} = \frac{1}{m} \sum_{j:0 < [\widehat{\boldsymbol{\beta}}_{SVM}]_j < C} \left[y_j - \sum_{i=1}^n [\widehat{\boldsymbol{\alpha}}_{SVM}]_i K(\boldsymbol{x}_j, \boldsymbol{x}_i) \right]$

m: Number of non-zero elements in $\hat{\alpha}_{SVM}$

Sparseness



140

 $\beta_i(\xi_i + u_i - 1) = 0$ for all *i* Therefore, some β_i (and thus α_i) could be zero.

$$\alpha_i = \beta_i y_i$$

See e.g.,

- Evgeniou, Pontil & Poggio, Regularization Networks and Support Vector Machines, Advances in Computational Mathematics, 2000, 13(1), 1-50.
- B. Schölkopf and A. Smola.Learning with Kernels. MIT Press, 2002
- http://www.kernel-machines.org

KKT condition:

Examples



Gaussian kernel: $K(x, x') = \exp\left(-\frac{||x - x'||^2}{2c^2}\right)$

Examples (cont.)





142

Small C



Examples

. ... ----.....

Homework

- Originally, SVMs are derived within a totally different framework, i.e., maximum margin principle. Read the following article on SVMs (Sec.1-Sec.4) and write your opinion.
 - B. Schölkopf: Statistical learning and kernel methods.

ftp://ftp.research.microsoft.com/pub/tr/tr-2000-23.pdf

Homework (cont.)

- Implement SVM by yourself or find a suitable software. Then perform simulations under various settings (various data sets, changing kernel, changing C etc.) and analyze the results.
 - Software is available from, e.g., http://www.kernel-machines.org
 - You may play with Java implementation, e.g., http://svm.dcs.rhbnc.ac.uk/pagesnew/GPat.shtml