Pattern Information Processing¹⁶⁸ 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 Problem¹⁶⁹

- 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 Problem¹⁷⁰

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

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

 $\lambda \ (\geq 0)$

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

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 $J_{LS}(\alpha)$:

$$J_{0/1}(\boldsymbol{\alpha}) = \sum_{i=1}^{n} I\left(\operatorname{sign}(\hat{f}(\boldsymbol{x}_i)) \neq y_i\right)$$

 $I(a \neq b) = \begin{cases} 0 & (a = b) \\ 1 & (a \neq b) \end{cases}$

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

Hinge-Loss

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

Use hinge-loss as an approximation:

$$J_H(\boldsymbol{\alpha}) = \sum_{i=1}^n \max\left(0, 1 - u_i\right)$$
$$J_{0/1}(\boldsymbol{\alpha}) = \frac{1}{2} \sum_{i=1}^n \left(1 - \operatorname{sign}\left(u_i\right)\right)$$

$$J_{LS}(\boldsymbol{\alpha}) = \sum_{i=1}^{N} (1 - u_i)^2$$

Note
$$:y_i^2 = 1, \ 1/y_i = y_i$$

 $u_i = \hat{f}(\boldsymbol{x}_i) y_i$

172



How to Obtain Solutions

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

$$J_{H}(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \max \left(0, 1 - u_{i} \right)$$

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

Lemma: $\max(0, 1 - u) = \min_{\xi \in \mathbb{R}} \xi \quad \text{subject to } \xi \ge 1 - u$ $\xi \ge 0$

Proof: Constraints are $\xi \ge \max(0, 1 - u)$, so the lemma holds. Q.E.D. How to Obtain Solutions (cont.)⁷⁴ So we have

$$J_H(\boldsymbol{\alpha}) = \min_{\boldsymbol{\xi} \in \mathbb{R}^n} \langle \mathbf{1}_n, \boldsymbol{\xi} \rangle$$
 subject to $\boldsymbol{\xi} \ge \mathbf{1}_n - \boldsymbol{u}$
 $\boldsymbol{\xi} \ge \mathbf{0}_n$

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

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

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

Support Vector Machines

We focus on the following setting:

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

• R = K

 $\boldsymbol{K}_{i,j} = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$

175

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

$$\widehat{\boldsymbol{\alpha}}_{SVM} = \underset{\boldsymbol{\alpha}, \boldsymbol{\xi} \in \mathbb{R}^n}{\operatorname{argmin}} \left[C \langle \mathbf{1}_n, \boldsymbol{\xi} \rangle + \frac{1}{2} \langle \boldsymbol{K} \boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle \right]$$
subject to $\boldsymbol{\xi} \geq \mathbf{1}_n - \boldsymbol{u}$
 $\boldsymbol{\xi} \geq \mathbf{0}_n$

$$u_i = \widehat{f}(\boldsymbol{x}_i) y_i$$

Efficient Formulation

The SVM solution can be obtained by

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

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

subject to $\mathbf{0}_n \leq \boldsymbol{\beta} \leq C \mathbf{1}_n$

- The number of parameters is decreased from 2n to n !
- This corresponds to considering the Wolfe dual (details are omitted).

Examples

177

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

Examples (cont.)





178

Large C

Small C

Examples

. ... ----3 en en en

Original Derivation of SVMs ¹⁸⁰

- The way SVMs were introduced today is quite different from the original derivation.
- Let's briefly follow the original derivation.
 - Hyper-plane classifier
 - VC theory
 - Margin maximization
 - Soft margin
 - Kernel trick

Hyper-plane Classifier

181

Separate sample space by hyper-plane.

$$\widehat{f}(\boldsymbol{x}) = \operatorname{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b)$$

$$w$$

$$(\mathbf{x}, \mathbf{x}) = \mathbf{w}$$

$$(\mathbf{w}, \mathbf{x}) + 1$$

$$(\mathbf{w}, \mathbf{x}) = \mathbf{w}$$

$$(\mathbf{w}, \mathbf{x}) + 1$$

$$(\mathbf{w}, \mathbf{x}) = \mathbf{w}$$

$$(\mathbf{w}, \mathbf{x}) = \mathbf{w}$$

find \boldsymbol{w}, b such that $y_i \widehat{f}(\boldsymbol{x}_i) \ge 1$ for $i = 1, \dots, n$.

Margin

Margin: "Gap" between two classes



Vapnik-Chevonenkis Theory ¹⁸³ Generalization error:

$$R[\widehat{f}] = \iint I(\widehat{f}(\boldsymbol{x}) \neq y) p(\boldsymbol{x}, y) d\boldsymbol{x} dy$$

Empirical error:

$$R_{\text{emp}}[\widehat{f}] = \frac{1}{n} \sum_{i=1}^{n} I(\widehat{f}(\boldsymbol{x}_i) \neq y_i)$$
$$I(a \neq b) = \begin{cases} 0 & (a = b) \\ 1 & (a \neq b) \end{cases}$$

Generalization error bound ("VC bound") $R[\widehat{f}] \le R_{\text{emp}}[\widehat{f}] + \sqrt{\frac{1}{n} \left(h \left(\log \frac{2n}{h} + 1\right) + \log \frac{4}{\delta}\right)}$

h : VC dimension (model complexity)

with probability $1 - \delta$

Vapnik-Chevonenkis Theory (cont?) VC bound:

$$R[\widehat{f}] \le R_{\text{emp}}[\widehat{f}] + \sqrt{\frac{1}{n} \left(h\left(\log\frac{2n}{h}+1\right) + \log\frac{4}{\delta}\right)}$$

Monotone decreasing w.r.t. VC dimension h (h < n)

If samples are linear separable, empirical error is zero. $R_{emp}[\widehat{f}] = 0$



In VC theory, maximum margin classifier is optimal



Soft Margin

If samples are not linearly separable, margin cannot be defined.

Allow small error ξ_i .



Non-linear Extension

187

- Transform samples to a feature space by a non-linear mapping $\phi(x)$.
- Then find the maximum margin hyperplane in the feature space.



Kernel Trick

Compute inner product in the feature space by a kernel function:

$$egin{aligned} &\langle \phi(m{x}_i), \phi(m{x}_j)
angle = K(m{x}_i, m{x}_j) \ & orall m{x}, m{x}', \ \ K(m{x}, m{x}') \geq 0 \ & \end{aligned} \ & \end{aligne$$

188

Any linear algorithm represented by inner product can be non-linearized by kernels

 E.g.: Support vector machine, k-nearest neighbor classifier, principal component analysis, linear discriminant analysis, k-means clustering,

Notification of Final Assignment

189

- 1. Apply supervised learning techniques to your data set and analyze it.
- 2. Write your opinion about this course

 Final report deadline: Aug 6th (Fri.)
 E-mail submission is also accepted! sugi@cs.titech.ac.jp

Mini-Workshop on Data Mining⁹⁰

- On July 20th (final class), we have a mini-workshop on data mining, instead of regular lecture.
- Several students present their data mining results.
- Those who give a talk at the workshop will have very good grades!

Mini-Workshop on Data Mining⁹¹

- Application (just to declare that you want to give a presentation) deadline: June 29th.
- Presentation: 10-15(?) minutes.
 - Specification of your dataset
 - Employed methods
 - Outcome
- OHP or projector may be used.
- Slides should be in English.
- Better to speak in English, but Japanese is also allowed.

Schedule

June 22nd

June 29th

July 6th

July 13th

July 20th July 27th

- : Preparation for workshop (no lecture)
- : Robust method (regular lecture)
- : Neural Networks (regular lecture)
- : Preparation for workshop (no lecture)
- : Mini-workshop
- : Mini-workshop (if necessary)

Homework

- Prepare a toy binary classification problem (say 2-dim input) and test SVM. Then analyze the results by varying experimental conditions (datasets, kernels, regularization parameter C etc.).
 - Software is available from, e.g., http://www.support-vector.net/software.html
 - You may play with Java implementation, e.g., http://svm.dcs.rhbnc.ac.uk/pagesnew/GPat.shtml