Pattern Information Processing<sup>186</sup> 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 Problem<sup>187</sup>

- Output values are  $y_i = \pm 1$ .
- We want to predict whether output values of unlearned input points are positive or negative.
- Multi-class problem can be transferred to several binary classification problems:
  - One-versus-rest (1vs.2&3, 2vs.1&3, 3vs.1&2)
  - One-versus-one (1vs.2, 1vs.3, 2vs.3)

# (Binary) Classification Problem<sup>188</sup>

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

$$\hat{lpha}_{QCLS} = \operatorname*{argmin}_{oldsymbol{lpha}\in\mathbb{R}^b} \left[J_{LS}(oldsymbol{lpha}) + \lambda \langle oldsymbol{R}oldsymbol{lpha},oldsymbol{lpha}
angle
ight] \ \lambda \;(\geq 0) \ J_{LS}(oldsymbol{lpha}) = \sum_{i=1}^n \left(f_{oldsymbol{lpha}}(oldsymbol{x}_i) - y_i
ight)^2$$

$$\widehat{y} = \operatorname{sign}\left(f_{\widehat{oldsymbol{lpha}}}(oldsymbol{x})
ight)$$

### 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  $J_{LS}(\alpha)$ :

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

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

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

### Hinge-Loss

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

Use hinge-loss as an approximation:

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

$$u_{i} = f$$

$$: \text{Samp}$$

$$J_{0/1}(\boldsymbol{\alpha}) = \frac{1}{2} \sum_{i=1}^{n} (1 - \text{sign}(u_{i}))$$

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

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

$$u_i = f_{\alpha}(x_i)y_i$$
  
Sample-wise margin  
$$\int_{a}^{a} \int_{a}^{b} \int_{$$

190

How to Obtain A Solution 191  

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

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

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 A Solution (cont.)<sup>92</sup> So we have

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

$$\hat{oldsymbol{lpha}}_{SVM} = \operatorname*{argmin}_{oldsymbol{lpha} \in \mathbb{R}^b, oldsymbol{\xi} \in \mathbb{R}^n} ig[ \langle oldsymbol{1}_n, oldsymbol{\xi} 
angle + \lambda \langle oldsymbol{R}oldsymbol{lpha}, oldsymbol{lpha} 
angle] \ \mathrm{subject \ to} \ oldsymbol{\xi} \geq oldsymbol{1}_n - oldsymbol{u} \ oldsymbol{\xi} \geq oldsymbol{0}_n \end{cases}$$

## Support Vector Machines

We focus on the following setting:

•
$$f_{oldsymbol{lpha}}(oldsymbol{x}) = \sum_{i=1}^n lpha_i K(oldsymbol{x},oldsymbol{x}_i)$$
  
•  $oldsymbol{R} = oldsymbol{K}$ 

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

193

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

$$egin{aligned} \widehat{oldsymbol{lpha}}_{SVM} = \operatorname*{argmin}_{oldsymbol{lpha},oldsymbol{\xi} \in \mathbb{R}^n} \left[ C \langle oldsymbol{1}_n,oldsymbol{\xi} 
angle + rac{1}{2} \langle oldsymbol{K}oldsymbol{lpha},oldsymbol{lpha} 
angle 
ight] \ \mathrm{subject \ to} \ oldsymbol{\xi} \geq oldsymbol{1}_n - oldsymbol{u} \ oldsymbol{\xi} \geq oldsymbol{0}_n \ oldsymbol{\xi} \geq oldsymbol{0}_n \ oldsymbol{u}_i = f_{oldsymbol{lpha}}(oldsymbol{x}_i) y_i \end{aligned}$$

## **Efficient Formulation**

The SVM solution can be obtained by

 $[\widehat{oldsymbol{lpha}}_{SVM}]_i = [\widehat{oldsymbol{eta}}_{SVM}]_i y_i$  , where

Proof: Homework!

194

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

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

The number of parameters is reduced to n.
QP standard form:

$$\min_{oldsymbol{eta} \in \mathbb{R}^n} egin{bmatrix} rac{1}{2} \langle oldsymbol{Q} oldsymbol{eta}, oldsymbol{eta} 
angle + \langle oldsymbol{eta}, oldsymbol{q} 
angle \end{bmatrix} & oldsymbol{Q}_{i,j} = oldsymbol{K}_{i,j} y_i y_j \quad oldsymbol{q} = -oldsymbol{1}_n \ \mathbf{Subject to} \ oldsymbol{H} oldsymbol{eta} \leq oldsymbol{h} & oldsymbol{H} = \begin{pmatrix} -oldsymbol{I}_n \\ oldsymbol{I}_n \end{pmatrix} oldsymbol{h} = \begin{pmatrix} oldsymbol{0}_n \\ C oldsymbol{1}_n \end{pmatrix} \end{pmatrix}$$

### **Sparseness**

195

#### KKT optimality condition implies β<sub>i</sub>(ξ<sub>i</sub> + u<sub>i</sub> - 1) = 0 for all i u<sub>i</sub> = f̂(x<sub>i</sub>)y<sub>i</sub> Therefore, some β<sub>i</sub> (and thus α<sub>i</sub> = β<sub>i</sub>y<sub>i</sub> also) could be zero.

### Examples



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

#### 197 Examples (cont.)



Large C

## Examples



198

# Original Derivation of SVMs<sup>199</sup>

- 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 200

Separate sample space by hyper-plane.

 $f_{\boldsymbol{w}}(\boldsymbol{x}) = \langle \boldsymbol{w}, \boldsymbol{x} \rangle + b$   $\widehat{y} = \operatorname{sign}(f_{\boldsymbol{w}}(\boldsymbol{x}))$  $-1 \xrightarrow{\times \times}{\times \times} \xrightarrow{\circ \circ}{\times \times} +1$ find  $\boldsymbol{w}, \boldsymbol{b}$ such that  $y_i f_{\boldsymbol{w}}(\boldsymbol{x}_i) \geq 1$  for  $i = 1, \ldots, n$ .

# Margin

### Margin: "Gap" between two classes



201

Vapnik-Chevonenkis Theory <sup>202</sup> Generalization error:

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

with probability  $1 - \delta$ 

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)

### 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 with respect to 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  $\xi_i$ .



### Non-linear Extension 206

- 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:

207

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,



### Homework

1. Prove that the solution of SVM,

$$egin{aligned} \widehat{oldsymbol{lpha}}_{SVM} &= rgmin_{oldsymbol{lpha},oldsymbol{\xi} \in \mathbb{R}^n} \left[ C \langle oldsymbol{1}_n,oldsymbol{\xi} 
angle + rac{1}{2} \langle oldsymbol{K}oldsymbol{lpha},oldsymbol{lpha} 
angle 
ight] \ ext{subject to } oldsymbol{\xi} \geq oldsymbol{1}_n - oldsymbol{u}, \ oldsymbol{\xi} \geq oldsymbol{0}_n \ f_{oldsymbol{lpha}}(oldsymbol{x}) = \sum_{i=1}^n lpha_i K(oldsymbol{x},oldsymbol{x}_i) \ rac{1}{2} \langle oldsymbol{K}oldsymbol{lpha},oldsymbol{lpha} 
angle 
ight] \ oldsymbol{u}_i = oldsymbol{1}_n - oldsymbol{u}, \ oldsymbol{\xi} \geq oldsymbol{0}_n \ oldsymbol{k}_i = oldsymbol{f}_{oldsymbol{lpha}}(oldsymbol{x}_i) y_i \ oldsymbol{K}_{i,j} = K(oldsymbol{x}_i,oldsymbol{x}_j) 
angle \ oldsymbol{K}_{i,j} = K(oldsymbol{x}_i,oldsymbol{x}_j) 
angle \ oldsymbol{k}_i = oldsymbol{L}_i \ oldsymbol{K}_i = oldsymbol{L}_i \ oldsymbol{K}_i, oldsymbol{x}_i = oldsymbol{L}_i \ oldsymbol{L}_i \ oldsymbol{1}_i = oldsymbol{L}_i \ oldsymbol{L}_i \ oldsymbol{K}_i, oldsymbol{x}_i = oldsymbol{L}_i \ ol$$

is given by  $[\widehat{\alpha}_{SVM}]_i = [\widehat{\beta}_{SVM}]_i y_i$ , where  $\widehat{\beta}_{SVM} = \operatorname*{argmax}_{\beta \in \mathbb{R}^n} \left[ \sum_{i=1}^n \beta_i - \frac{1}{2} \sum_{i,j=1}^n \beta_i \beta_j y_i y_j K_{i,j} \right]$ 

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

Hint: Use Wolfe dual

### Homework (cont.)

#### Lagrangian:

$$L(\boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\beta}, \boldsymbol{\gamma}) = C\langle \mathbf{1}_n, \boldsymbol{\xi} \rangle + \frac{1}{2} \langle \boldsymbol{K} \boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle$$
$$-\langle \boldsymbol{\beta}, \boldsymbol{\xi} + \boldsymbol{u} - \mathbf{1}_n \rangle - \langle \boldsymbol{\gamma}, \boldsymbol{\xi} \rangle$$
$$\boldsymbol{\beta}, \boldsymbol{\gamma} : \text{Lagrange multiplier}$$
$$\mathbf{Wolfe duality:}$$

$$egin{aligned} \min_{oldsymbol{lpha},oldsymbol{\xi}\in\mathbb{R}^n} \left[ C\langle oldsymbol{1}_n,oldsymbol{\xi}
angle \ +rac{1}{2}\langle oldsymbol{K}oldsymbol{lpha},oldsymbol{lpha}
ight] &= \max_{oldsymbol{eta},oldsymbol{\gamma}\in\mathbb{R}^n} L(oldsymbol{lpha},oldsymbol{\xi},oldsymbol{\gamma}) \ & ext{ subject to }oldsymbol{\xi}\geq oldsymbol{1}_n-oldsymbol{u} & ext{ subject to }oldsymbol{eta}\geq oldsymbol{0}_n & oldsymbol{\gamma}\geq oldsymbol{0}_n \\ oldsymbol{\xi}\geq oldsymbol{1}_n-oldsymbol{u} & ext{ subject to }oldsymbol{eta}\geq oldsymbol{0}_n & oldsymbol{\gamma}\geq oldsymbol{0}_n \\ oldsymbol{\xi}\geq oldsymbol{0}_n & ext{ }oldsymbol{eta}\in oldsymbol{0}_n \\ oldsymbol{\xi}\geq oldsymbol{0}_n & ext{ }oldsymbol{eta}\in oldsymbol{0}_n \\ oldsymbol{eta}\geq oldsymbol{0}_n & ext{ }oldsymbol{B}=oldsymbol{0}_n & oldsymbol{\gamma}\geq oldsymbol{0}_n \\ oldsymbol{eta}\geq oldsymbol{0}_n & ext{ }oldsymbol{eta}\in oldsymbol{0}_n \\ oldsymbol{eta}\geq oldsymbol{0}_n & ext{ }oldsymbol{eta}\geq oldsymbol{0}_n & oldsymbol{eta}\geq oldsymbol{0}_n \\ oldsymbol{eta}\geq oldsymbol{0}_n & ext{ }oldsymbol{1}_n \ oldsymbol{eta}=oldsymbol{0}_n & ext{ }oldsymbol{eta}\in oldsymbol{0}_n \\ oldsymbol{eta}\geq oldsymbol{0}_n & ext{ }oldsymbol{eta}=oldsymbol{0}_n & ext{ }oldsymbol{eta}=oldsymbol{0}_n \end{array} 
ight.$$

### Homework

211

- 2. 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

# Mini-Workshop on Data Mining<sup>212</sup>

- On July 15<sup>th</sup> and 22<sup>nd</sup>, we will have a miniworkshop on data mining.
- Several students present their own data mining results.
- Those who give a talk at the workshop will have very good grades!

# Mini-Workshop on Data Mining<sup>213</sup>

- Application (just to declare that you want to give a presentation) deadline: July 1<sup>st</sup>.
  - Come to me after the class
- Presentation: 10-15 minutes (?).
  - Specification of your dataset
  - Methods used
  - Outcome
- Slides should be in English.
- Better to speak in English, but Japanese is also allowed.

# Notification of Final Assignment

214

- 1. Apply supervised learning techniques to your data set and analyze it.
- Final report deadline: Aug 1<sup>st</sup> (Fri.) 17:00
   Bring your report to W8E-404.