Pattern Information Processing<sup>184</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>185</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>186</sup>

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

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

 $\lambda \ (\geq 0)$ 

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

n

Prediction:

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

### 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:

$$J_H(\boldsymbol{\alpha}) = \sum_{i=1}^n \max\left(0, 1 - u_i\right)$$
$$I_H(\boldsymbol{\alpha}) = \frac{1}{n} \sum_{i=1}^n \left(1 - a_i\right)$$

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

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

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

$$u_i = f_{\alpha}(x_i)y_i$$
  
: Sample-wise margin

188

How to Obtain A Solution 189  

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

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

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 A Solution (cont<sup>190</sup> So we have

$$J_{H}(\boldsymbol{\alpha}) = \min_{\boldsymbol{\xi} \in \mathbb{R}^{n}} \langle \mathbf{1}_{n}, \boldsymbol{\xi} \rangle \text{ subject to } \boldsymbol{\xi} \geq \mathbf{1}_{n} - \boldsymbol{u}$$
$$\boldsymbol{\xi} \geq \mathbf{0}_{n}$$

#### 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:

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

• R = K

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

191

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!** 

192

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

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

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

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^n} \left[ \frac{1}{2} \langle \boldsymbol{Q} \boldsymbol{\beta}, \boldsymbol{\beta} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{q} \rangle \right] \qquad \boldsymbol{Q}_{i,j} = \boldsymbol{K}_{i,j} y_i y_j \quad \boldsymbol{q} = -\boldsymbol{1}_n \\ \text{subject to } \boldsymbol{H} \boldsymbol{\beta} \leq \boldsymbol{h} \qquad \boldsymbol{H} = \begin{pmatrix} -\boldsymbol{I}_n \\ \boldsymbol{I}_n \end{pmatrix} \quad \boldsymbol{h} = \begin{pmatrix} \boldsymbol{0}_n \\ C\boldsymbol{1}_n \end{pmatrix}$$

### **Sparseness**

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

# Examples (cont.) <sup>195</sup>



Large C

#### Small C

# Examples



196

# Original Derivation of SVMs <sup>197</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 <sup>198</sup>

Separate sample space by hyper-plane.

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



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

# Margin

#### Margin: "Gap" between two classes



### Vapnik-Chevonenkis Theory <sup>200</sup> Generalization error:

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

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 with respect to VC dimension h (h < n)

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

The larger margin is, the smaller VC dim is.



In VC theory, maximum margin classifier is optimal



# Soft Margin

203

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

Allow small error  $\xi_i$ .



## Non-linear Extension 204

- 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:

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

205

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,

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

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.)

208

#### Lagrangian:

$$egin{aligned} L(oldsymbol{lpha},oldsymbol{\xi},oldsymbol{eta},oldsymbol{\gamma}) &= C\langle oldsymbol{1}_n,oldsymbol{\xi}
angle + rac{1}{2}\langle oldsymbol{K}oldsymbol{lpha},oldsymbol{lpha}
angle \ &-\langleoldsymbol{eta},oldsymbol{\xi}+oldsymbol{u}-oldsymbol{1}_n
angle - \langleoldsymbol{\gamma},oldsymbol{\xi}
angle \end{aligned}$$

$$\beta, \gamma$$
 :Lagrange multiplier
Wolfe duality:

$$\min_{\substack{\boldsymbol{\alpha},\boldsymbol{\xi}\in\mathbb{R}^{n}}} \begin{bmatrix} C\langle \mathbf{1}_{n},\boldsymbol{\xi}\rangle & +\frac{1}{2}\langle \boldsymbol{K}\boldsymbol{\alpha},\boldsymbol{\alpha}\rangle \end{bmatrix} = \max_{\substack{\boldsymbol{\beta},\boldsymbol{\gamma}\in\mathbb{R}^{n}}} L(\boldsymbol{\alpha},\boldsymbol{\xi},\boldsymbol{\beta},\boldsymbol{\gamma})$$
  
subject to  $\boldsymbol{\xi}\geq\mathbf{1}_{n}-\boldsymbol{u}$  subject to  $\boldsymbol{\beta}\geq\mathbf{0}_{n}$   $\boldsymbol{\gamma}\geq\mathbf{0}_{n}$   
 $\boldsymbol{\xi}\geq\mathbf{0}_{n}$   $\frac{\partial L}{\partial\boldsymbol{\alpha}}=\mathbf{0}_{n}$   $\frac{\partial L}{\partial\boldsymbol{\xi}}=\mathbf{0}_{n}$ 

# Homework

209

- 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