# L2 wtr

Tianrui Wang

April 2019

# 1 Supporting Vector Machine

**Brief History**
**This is Theoretically Based and Global Minima**

**Optimization Goal**
**(1)Classification problem**
**Training set S=$\{(\mathbf{x}_1, y_1), (x_2, y_2) \ldots\}$**
*Our goal is to learn a function $g(x)$ where $f(x) = sgn(g(x))$ fits the samples best linear function $g(x) = w^T x + b$*
**(2)Destination**
**labels : $\mathbf{y} \in \{-1, 1\}$**
**geometricmargin : $\gamma^{(\mathbf{i})} = \mathbf{y^{(i)}}(\mathbf{w^T x} + \mathbf{b})$**
**in order to simplifiy the calculation , we make $|\mathbf{w}| = 1$**
**$\gamma = \min \gamma^{(\mathbf{i})}$**
**We use $\gamma$ to judge how well the function proforms**
**find $\max_{\gamma, \mathbf{w}, \mathbf{b}} \gamma$**
**s.t.$\mathbf{y^{(i)}}(\mathbf{w^T x^{(i)}} + \mathbf{b}) \geq \gamma$**
**where $||\mathbf{w}|| = 1$**
**(3)Transformation**
**As Unit Circle is not a convex set , we transform the origin question to...**
**find $\max_{\gamma, \mathbf{w}, \mathbf{b}} \ \gamma/||\mathbf{w}||$**
**s.t.$\mathbf{y^{(i)}}(\mathbf{w^T x^{(i)}} + \mathbf{b}) \geq \gamma$**
**We suppose $||\gamma|| = 1$ to simplify the question**
**So the problem is to find the $\min_{\mathbf{w}, \mathbf{b}} \ \mathbf{0.5}||\mathbf{w}||^\mathbf{2}$**
**(4)Lagrangianform**
**$\mathbf{L}(\mathbf{w}, \mathbf{b}, \alpha) = \frac{||\mathbf{w}||^\mathbf{2}}{\mathbf{2}} - \sum_{\mathbf{i=0}}^{\mathbf{N}} \alpha_\mathbf{i}[\mathbf{y^{(i)}}(\mathbf{w^T x^{(i)}} + \mathbf{b}) - \mathbf{1}]$**
**s.t. $\alpha_\mathbf{i} \geq \mathbf{0}$**

**Dual Problem :**
**$\max_\alpha \min_{\mathbf{w}, \mathbf{b}} \mathbf{L}(\mathbf{w}, \mathbf{b}, \mathbf{a}) = \mathbf{L}(\mathbf{w}, \mathbf{b}, \alpha) = \frac{||\mathbf{w}||^\mathbf{2}}{\mathbf{2}} - \sum_{\mathbf{i=0}}^{\mathbf{N}} \alpha_\mathbf{i}[\mathbf{y^{(i)}}(\mathbf{w^T x^{(i)}} + \mathbf{b}) - \mathbf{1}]$**
**With KKT conditions the problem′s answer is the origin problem′s**
**We swap the min and the max to make it easier**
**figure the Partial guidance**

$\mathbf{w} = \sum_{\mathbf{i=1}}^{\mathbf{N}} \alpha_{\mathbf{i}} \mathbf{y^{(i)}} \mathbf{x^{(i)}}$

**Improvement**
**(1)Allow error**
**We allow error $\xi_{\mathbf{i}}$ in classification;**
**it is based on the output of the discriminant function $\mathbf{w^T x + b}$.**
**So we can reach a balance between acuracy and generalization.**
**we replace the origin formula with $\min_{\mathbf{w,b},\xi} \frac{||\mathbf{w}||}{\mathbf{2}} + \mathbf{C} \sum_{\mathbf{i}} = \mathbf{1^N} \xi_{\mathbf{i}}$**
**s.t. $\mathbf{y(w^T x^{(i)} + b)} \geq \mathbf{1} - \xi_{\mathbf{i}} \quad \xi_{\mathbf{i}} \geq \mathbf{0}$**
**(2)Kernel function**
**In many real situations, dividing the samples linearly is not enough, so we use a map to increas**
**in order to get a linear dividing in a high dimensional space**
**define $\mathbf{K(x,z)} = \phi(\mathbf{x})^{\mathbf{T}} \phi(\mathbf{x})$**
**and replace all $\mathbf{(x^{(i)}, x^{(j)})}$ with $(\phi(\mathbf{x^{(i)}}), \phi(\mathbf{x^{(j)}}))$**
**However, it is difficult to find the proper kernel function . . .**
**Here are some commonly used function**
**$\mathbf{K(x,y) = (xy+1)^2}$**
**$\mathbf{K(x,y) = exp(\frac{||x-y||^2}{2\sigma^2})}$**

**Implementation**
**(1)Coordinate descend method**
**change others with one fixed**
**Loop until convergence**
**{**
**For(i = 1 : m)**
**{**
**$\mathbf{a_i = argmax_a iL(a_1, a_2, \ldots, a_{i-1}, a_{i+1}, \ldots, a_m)}$**
**}**
**}**
**it is slow with low efficiency**
**(2)SMO algorithm**
**Repeat until convergence**
**{**
**First, select one pair $\mathbf{a_i}$ and $\mathbf{a_j}$ to update next**
**Second, reoptimize $\mathbf{L(a)}$ with respect to $\mathbf{a_i}$ and $\mathbf{a_j}$ without changing other $\mathbf{a_k}$**
**}**
**(3)Multi − class classification**
**Method 1 one − versus − rest**
**Advantages : low cost**
**Disadvantages : unbalanced data**
**Examples : circle and not circle, triangle and not triangle . . .**
**Method 1 one − versus − one**
**Advantages : make judgment by voting**
**Disadvantages : high cost $\mathbf{(C_n^2}$ models)**
**Examples : circle and tirangle, circle and square, tirangle and sqaure . . .**