

**Soft Sum**  
 $\min_{\alpha} \sum_{i=1}^n \alpha_i \log \alpha_i$   
 with  $\sum_{i=1}^n \alpha_i = 1$   
 $\alpha_i \in [0, 1]$   
 A distribution  $(\alpha_i)$  is a probability distribution if  $\sum \alpha_i = 1$  and  $\alpha_i \geq 0$   
 Example:  $\alpha_1 = 0.7, \alpha_2 = 0.3$   
 Theorem: Let  $D$  be an  $n$ -dimensional simplex with vertices  $e_1, \dots, e_n$  and let  $f$  be a convex function. Then the minimum of  $f$  over  $D$  is attained at one of the vertices of  $D$ .  
 $\sqrt{\frac{1}{n+1}}, \sqrt{\frac{2}{n+1}}, \dots, \sqrt{\frac{n}{n+1}}$   
**SOFT SUM**  
 group  $\mathbb{R}^n \rightarrow \mathbb{R}^n$   
 $\psi(x) = \frac{1}{\sum_{i=1}^n e^{x_i}}$   
 $\psi(x) = \frac{1}{\sum_{i=1}^n e^{x_i}}$   
 $\psi(x) = \frac{1}{\sum_{i=1}^n e^{x_i}}$

**Kernel Methods**  
 Def: The  $\mathcal{H}$ -kernel Trick  
 Def: A map  $K: X \times X \rightarrow \mathbb{R}$  is called a kernel of feature space  $\mathcal{H}$  if  $\langle \Psi(x), \Psi(y) \rangle_{\mathcal{H}} = K(x, y)$   
 and  $\langle \Psi(x), \Psi(x) \rangle_{\mathcal{H}} = K(x, x) \geq 0$

Consider the following optimization problem  
 (a)  $\min_{\alpha} \left( \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j) - \sum_{i=1}^n \alpha_i y_i \right)$   
 where  $\mathcal{H} \rightarrow \mathbb{R}$ ,  $n$ -dim,  $\mathbb{R}$ ,  $\mathbb{R}_+ \rightarrow \mathbb{R}_+$   
 feature map  
 $\Rightarrow$  Soft-SVM  
 $R(\alpha) = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j)$   
 $f(\alpha_1, \dots, \alpha_n) = \frac{1}{n} \sum_{i=1}^n \max\{0, 1 - y_i \alpha_i\}$   
 $\Rightarrow$  Hard-SVM  
 $R(\alpha) = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j)$   
 $f(\alpha_1, \dots, \alpha_n) = \begin{cases} 0 & \text{if } \forall i, y_i \alpha_i \geq 1 \\ \infty & \text{otherwise} \end{cases}$

**Representer Theorem**  
 Assume  $\Psi: X \rightarrow \mathcal{H}$  is a feature map. Then the optimal solution  $w^*$  of (\*) lies in  $\text{span}\{\Psi(x_1), \dots, \Psi(x_n)\}$   
 $\Rightarrow \exists \alpha_1, \dots, \alpha_n \in \mathbb{R} : w^* = \sum_{i=1}^n \alpha_i \Psi(x_i)$   
 Proof: omitted

Corollary: (\*) is equivalent to the optimization problem  
 $\min_{\alpha} \left( \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j) - \sum_{i=1}^n \alpha_i y_i \right)$   
 with support  $w^* = \sum_{i=1}^n \alpha_i \Psi(x_i)$ . Thus only  $n$  values on channel  $\langle \Psi(x_i), \Psi(x_j) \rangle = K(x_i, x_j)$   
 Proof: Let  $w = \sum_{i=1}^n \alpha_i \Psi(x_i)$ . Then  
 $\langle w, \Psi(x_i) \rangle = \sum_{j=1}^n \alpha_j \langle \Psi(x_j), \Psi(x_i) \rangle = \sum_{j=1}^n \alpha_j K(x_j, x_i)$   
 and  
 $\|w\|^2 = \langle w, w \rangle = \left\langle \sum_{i=1}^n \alpha_i \Psi(x_i), \sum_{j=1}^n \alpha_j \Psi(x_j) \right\rangle$   
 $= \sum_{i,j=1}^n \alpha_i \alpha_j \langle \Psi(x_i), \Psi(x_j) \rangle = \sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j) \quad \square$

**S2. Implementing Soft-SVM with kernel**  
 We tackle the optimization problem in the feature space  
 $\min_w \left( \frac{1}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max\{0, 1 - y_i \langle w, \Psi(x_i) \rangle\} \right)$

for the sample  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$   
 Algorithm: SBT for solving Soft-SVM  
 Input:  $T$  iterations  
 $\beta^0 = 0$   
 for  $t = 1 \dots T$   
 $\alpha^t \leftarrow \frac{1}{T} \beta^t$   
 $i \leftarrow \arg \min_{i \in \{1, \dots, n\}} \max\{0, 1 - y_i \langle \beta^t, \Psi(x_i) \rangle\}$   
 $j \leftarrow \arg \min_{j \in \{1, \dots, n\}} \max\{0, 1 - y_j \langle \beta^t, \Psi(x_j) \rangle\}$   
 $\beta^{t+1} \leftarrow \beta^t + \alpha^t (\Psi(x_i) - \Psi(x_j))$   
 else:  
 $\beta^{t+1} \leftarrow \beta^t$   
 Output:  $\bar{w} = \sum_{i=1}^n \bar{\alpha}_i \Psi(x_i)$  where  $\bar{\alpha} = \frac{1}{T} \sum_{t=1}^T \alpha^t$

**General procedure**  
 1) Given domain  $X$ , choose  $\Psi: X \rightarrow \mathcal{H}$   
 feature map  
 2) Given labelled training examples  
 $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$   
 compute  $\tilde{S} = \{(\Psi(x_1), y_1), \dots, (\Psi(x_n), y_n)\}$   
 3) Train a (linear) predictor  $h$  over  $\tilde{S}$   
 4) Predict  $l(\hat{x}) = h(\Psi(\hat{x}))$  for a training set  $\hat{S}$

**Expressive power**: Let  $p: \mathbb{R}^n \rightarrow \mathbb{R}$  be a degree  $d$  multivariate polynomial.  
 $\Rightarrow p(x) = \sum_{j_1, \dots, j_d} w_j \prod_{i=1}^n x_i^{j_i}$   
 Then  $p$  can be expressed with the hypothesis  
 $p(x) = \langle w, \Psi(x) \rangle$   
 where  $\Psi(x)_j = \prod_{i=1}^n x_i^{j_i}$  and  $w = (w_j)_j$ .

$E_{x_1, \dots, x_n}$   
 $E_{y_1, \dots, y_n}$   
 $E_{x_1, \dots, x_n, y_1, \dots, y_n}$   
 $k \left( \sum_{i=1}^n k(x_i, x_i) \sum_{j=1}^n k(x_j, x_j) \right)$