

CSC411 Homework 7Question 1: Representer Theorem**Part 1A**

Formula for reference (expanded  $\|w\|^2 \rightarrow w^T w$ )

$$J(w) = \frac{1}{N} \sum_{i=1}^N L(y^{(i)}, t^{(i)}) + \frac{\lambda}{2} w^T w$$

Expanding  $y^{(i)}$

$$= \frac{1}{N} \sum_{i=1}^N L(g(w^T \psi(x)^{(i)}), t^{(i)}) + \frac{\lambda}{2} w^T w$$

Decomposing according to hint:  $w = w_\Psi + w_\perp$ , where  $w_\Psi$  is the projection of  $w$  onto  $\Psi$  and  $w_\perp$  is orthogonal to  $w_\Psi$

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N L(g((w_\Psi + w_\perp)^T \psi(x)^{(i)}), t^{(i)}) + \frac{\lambda}{2} (w_\Psi + w_\perp)^T (w_\Psi + w_\perp) \\ &= \frac{1}{N} \sum_{i=1}^N L(g(w_\Psi^T \psi(x)^{(i)} + w_\perp^T \psi(x)^{(i)}), t^{(i)}) + \frac{\lambda}{2} (w_\Psi^T w_\Psi + 2w_\Psi^T w_\perp + w_\perp^T w_\perp) \end{aligned}$$

Since  $w_\perp$  is orthogonal to  $w_\Psi$  and every row vector in  $\Psi$ , so their dot products reduce to 0

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N L(g(w_\Psi^T \psi(x)^{(i)}), t^{(i)}) + \frac{\lambda}{2} (w_\Psi^T w_\Psi + w_\perp^T w_\perp) \\ &\geq \frac{1}{N} \sum_{i=1}^N L(g(w_\Psi^T \psi(x)^{(i)}), t^{(i)}) + \frac{\lambda}{2} (w_\Psi^T w_\Psi) = J(w_\Psi) \end{aligned}$$

Thus, we have  $J(w) \geq J(w_\Psi)$  where  $w_\Psi$  is the projection of  $w$  onto  $\Psi$ .

Therefore,  $w_\Psi$  is a linear combination of  $\psi(x^{(1)})^T \dots \psi(x^{(N)})^T$ , which are the rows of  $\Psi$ .

Therefore,  $w_\Psi$  must belong to the row space of  $\Psi$ .

**Part 1B****Final Answer:**

$$\alpha_{opt} = (KK + N\lambda K)^{-1}(K^T t)$$

**Proof:**

Formula for reference:

$$J(w) = \frac{1}{2N} \|t - \Psi w\|^2 + \frac{\lambda}{2} \|w\|^2$$

Rephrasing using definition of squared magnitude

$$\|a\|^2 = a^T a$$

$$\|a - b\|^2 = a^T a - 2a^T b + b^T b$$

$$\Rightarrow J(w) = \frac{1}{2N} (t^T t - 2t^T \Psi w + (\Psi w)^T (\Psi w)) + \frac{\lambda}{2} w^T w$$

$$\Rightarrow J(w) = \frac{1}{2N} (t^T t - 2t^T \Psi w + w^T \Psi^T \Psi w) + \frac{\lambda}{2} w^T w$$

Subbing in ( $w = \Psi^T \alpha$ ,  $w^T = \alpha^T \Psi$ )

$$\Rightarrow J(\Psi^T \alpha) = \frac{1}{2N} (t^T t - 2t^T \Psi \Psi^T \alpha + \alpha^T \Psi \Psi^T \Psi \Psi^T \alpha) + \frac{\lambda}{2} \alpha^T \Psi \Psi^T \alpha$$

Simplifying using gram matrix ( $K = \Psi \Psi^T$ )

$$\Rightarrow J(\alpha) = \frac{1}{2N} (t^T t - 2t^T K \alpha + \alpha^T K K \alpha) + \frac{\lambda}{2} \alpha^T K \alpha$$

Removing fraction for simplicity (won't affect finding minimizing  $\alpha$ )

$$\Rightarrow 2N * J(\alpha) = t^T t - 2t^T K \alpha + \alpha^T K K \alpha + N\lambda \alpha^T K \alpha$$

Rearranging addition

$$\Rightarrow 2N * J(\alpha) = \alpha^T K K \alpha + N\lambda \alpha^T K \alpha - 2t^T K \alpha + t^T t$$

Factoring out  $\alpha$ ,  $\alpha^T$  for first two terms:

$$\Rightarrow 2N * J(\alpha) = (\alpha^T K K + N\lambda \alpha^T K) \alpha - 2t^T K \alpha + t^T t$$

$$\Rightarrow 2N * J(\alpha) = \alpha^T (K K + N\lambda K) \alpha - 2t^T K \alpha + t^T t$$

Dividing entire equation by 2

$$\Rightarrow N * J(\alpha) = \frac{1}{2} (\alpha^T (K K + N\lambda K) \alpha) - t^T K \alpha + \frac{1}{2} (t^T t)$$

Using hint from handout:

$$\operatorname{argmin}_{\alpha} \left( \frac{1}{2} (\alpha^T A \alpha) + b^T \alpha + c \right) = -A^{-1} b$$

Where:

$$A = (KK + N\lambda K)$$

$$[b^T = -t^T K] \Rightarrow [b = -K^T t]$$

$$c = \frac{1}{2}(t^T t)$$

Then:

$$\Rightarrow \operatorname{argmin}_{\alpha} \left( \frac{1}{2} (\alpha^T (KK + N\lambda K) \alpha) - t^T K \alpha + \frac{1}{2} (t^T t) \right) = (KK + N\lambda K)^{-1} (K^T t)$$

**Question 2: Compositional Kernels****Part 2A****Final Answer**

$$\psi_S(x) = \begin{bmatrix} \psi_1(x)_1 \\ \vdots \\ \psi_1(x)_A \\ \psi_2(x)_1 \\ \vdots \\ \psi_2(x)_B \end{bmatrix}$$

Where: A, B are the respective lengths of feature maps  $\psi_1, \psi_2$

**Proof**

Formulae for reference

$$k_1(x, x') = \psi_1(x) \cdot \psi_1(x')$$

$$k_2(x, x') = \psi_2(x) \cdot \psi_2(x')$$

$$k_S(x, x') = k_1(x, x') + k_2(x, x')$$

Want:  $\psi_S$  such that  $k_S(x, x') = \psi_S(x) \cdot \psi_S(x')$

$$\psi_S(x) \cdot \psi_S(x') = \psi_1(x) \cdot \psi_1(x') + \psi_2(x) \cdot \psi_2(x')$$

Expanding  $\psi_1, \psi_2$  (with respective lengths A, B)

$$\psi_S(x) \cdot \psi_S(x') = \sum_{i=1}^A \psi_1(x)_i \times \psi_1(x')_i + \sum_{i=1}^B \psi_2(x)_i \times \psi_2(x')_i$$

Express the above as the dot product of the concatenation of the feature vectors

$$\psi_S(x) \cdot \psi_S(x') = \begin{bmatrix} \psi_1(x)_1 \\ \vdots \\ \psi_1(x)_A \\ \psi_2(x)_1 \\ \vdots \\ \psi_2(x)_B \end{bmatrix} \cdot \begin{bmatrix} \psi_1(x')_1 \\ \vdots \\ \psi_1(x')_A \\ \psi_2(x')_1 \\ \vdots \\ \psi_2(x')_B \end{bmatrix}$$

$$\Rightarrow \psi_S(x) = \begin{bmatrix} \psi_1(x)_1 \\ \vdots \\ \psi_1(x)_A \\ \psi_2(x)_1 \\ \vdots \\ \psi_2(x)_B \end{bmatrix}$$

**Part 2B****Final Answer**

$$\psi_P(x) = \text{vec} \left( \begin{bmatrix} \psi_1(x)_1 \\ \vdots \\ \psi_1(x)_A \end{bmatrix} \times [\psi_2(x)_1 \quad \cdots \quad \psi_2(x)_B] \right)$$

$$\psi_P(x)_{(i,j)} = (\psi_1(x)_i \times \psi_2(x)_j)$$

**Proof**

Formulae for reference

$$k_1(x, x') = \psi_1(x) \cdot \psi_1(x')$$

$$k_2(x, x') = \psi_2(x) \cdot \psi_2(x')$$

$$k_P(x, x') = k_1(x, x') \times k_2(x, x')$$

Want:  $\psi_P$  such that  $k_P(x, x') = \psi_P(x) \cdot \psi_P(x')$

$$\psi_P(x) \cdot \psi_P(x') = \psi_1(x) \cdot \psi_1(x') \times \psi_2(x) \cdot \psi_2(x')$$

Expanding  $\psi_1, \psi_2$  (with respective lengths  $A, B$ )

$$\Rightarrow \psi_P(x) \cdot \psi_P(x') = \left( \sum_{i=1}^A \psi_1(x)_i \times \psi_1(x')_i \right) \times \left( \sum_{i=1}^B \psi_2(x)_i \times \psi_2(x')_i \right)$$

Merging summations and rearranging multiplication

$$\Rightarrow \psi_P(x) \cdot \psi_P(x') = \left( \sum_{i=1}^A \sum_{j=1}^B \psi_1(x)_i \times \psi_1(x')_i \times \psi_2(x)_j \times \psi_2(x')_j \right)$$

$$\Rightarrow \psi_P(x) \cdot \psi_P(x') = \left( \sum_{i=1}^A \sum_{j=1}^B (\psi_1(x)_i \times \psi_2(x)_j) \times (\psi_1(x')_i \times \psi_2(x')_j) \right)$$

Note that the terms associated with  $x$  are the same as the terms associated with  $x'$  (multiplication of a  $\psi_1$  element then a  $\psi_2$  element).

Clearly, the  $(i, j)$ -th element of the desired kernel (represented as a matrix of size  $A \times B$ ) must be:

$$\psi_P(x)_{(i,j)} = (\psi_1(x)_i \times \psi_2(x)_j)$$

We can represent the whole kernel as a vectorization (**flattening**) of the feature mappings' matrix multiplication:

$$\psi_P(x) = \text{vec} \left( \begin{bmatrix} \psi_1(x)_1 \\ \vdots \\ \psi_1(x)_A \end{bmatrix} \times [\psi_2(x)_1 \quad \cdots \quad \psi_2(x)_B] \right)$$