

Algorithmic perspectives of kernel method

Chanwoo Lee, August 7, 2020

1 A choice of feature mapping Φ

To derive an algorithm, I choose to use Mapping 1 in the previous note for convenience.

$$\begin{aligned}\Phi: \mathbb{R}^{d_1 \times d_2} &\rightarrow \mathcal{H}_r^{d_1} \times \mathcal{H}_c^{d_2} \\ \mathbf{X} &\mapsto (\Phi_r(\mathbf{X})\Phi_c(\mathbf{X})) \stackrel{\text{def}}{=} (\phi_r([\mathbf{X}]_{1:}), \dots, \phi_r([\mathbf{X}]_{d_1:}), \phi_c([\mathbf{X}]_{:1}), \dots, \phi_c([\mathbf{X}]_{:d_2})).\end{aligned}$$

We define decision function,

$$\begin{aligned}f(\mathbf{X}) &= \langle \mathbf{B}, \Phi(\mathbf{X}) \rangle, \text{ where } \mathbf{B} = (\mathbf{B}_r, \mathbf{B}_c) \in \mathcal{H}_r^{d_1} \times \mathcal{H}_c^{d_2} \\ &= \langle \mathbf{B}_r, \Phi_r(\mathbf{X}) \rangle + \langle \mathbf{B}_c, \Phi_c(\mathbf{X}) \rangle \\ &= \sum_{k=1}^n \gamma_k \left(\sum_{i,j \in [d_2]} w_{ij}^{\text{row}} K([\mathbf{X}_k]_{i:}, [\mathbf{X}]_{j:}) + \sum_{i,j \in [d_2]} w_{ij}^{\text{col}} K([\mathbf{X}_k]_{:,i}, [\mathbf{X}]_{:,j}) \right),\end{aligned} \tag{1}$$

where $\mathbf{X}^1, \dots, \mathbf{X}^n$ are sampled matrix features and $\mathbf{W}^{\text{col}}, \mathbf{W}^{\text{row}}$ are some positive semi definite matrices with low rank. We estimate $\mathbf{W}^{\text{col}}, \mathbf{W}^{\text{row}}$, and $\gamma = (\gamma_1, \dots, \gamma_n)$ from the training data set.

2 Algorithm derivation

We solve an optimization problem

$$\begin{aligned}\min_{\mathbf{B}} \quad & \frac{1}{2} \|\mathbf{B}\|_F^2 + c \sum_{i=1}^n \xi_i, \\ \text{subject to } & y_i \langle \mathbf{B}, \Phi(\mathbf{X}_i) \rangle \leq 1 - \xi_i \text{ and } \xi_i \geq 0, i = 1, \dots, n.\end{aligned} \tag{2}$$

where $\|\mathbf{B}\|_F^2 = \|\mathbf{B}_r\|_F^2 + \|\mathbf{B}_c\|_F^2$. From the low rank assumption on \mathbf{B} such that

$$\mathbf{B} = (\mathbf{B}_r, \mathbf{B}_c) = \mathbf{C}\mathbf{P}^T = (\mathbf{C}_r, \mathbf{C}_c)(\mathbf{P}_r, \mathbf{P}_c)^T,$$

where $\mathbf{C} = (\mathbf{C}_r, \mathbf{C}_c) \in \mathcal{H}_r^r \times \mathcal{H}_c^r$ and $\mathbf{P} = (\mathbf{P}_r, \mathbf{P}_c) \in \mathbb{R}^{d_1 \times r} \times \mathbb{R}^{d_2 \times r}$.

1. First we update \mathbf{C} holding \mathbf{P} fixed. The dual problem of Equation (2) is

$$\begin{aligned}\min_{\alpha = (\alpha_1, \dots, \alpha_n)} \quad & - \sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \Phi(\mathbf{X}_i) \mathbf{P} (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T, \Phi(\mathbf{X}_j) \mathbf{P} (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \rangle \\ \text{subject to } & 0 \leq \alpha_i \leq C, i = 1, \dots, n.\end{aligned} \tag{3}$$

Define $\mathbf{K}(i, j) \in \mathbb{R}^{d_1 \times d_1} \times \mathbb{R}^{d_2 \times d_2}$ as

$$\begin{aligned}\mathbf{K}(i, j) &= (\mathbf{K}_r(i, j), \mathbf{K}_c(i, j)) \stackrel{\text{def}}{=} \Phi(\mathbf{X}_i)^T \Phi(\mathbf{X}_j) \\ \text{where } [\mathbf{K}_r(i, j)]_{pq} &= K_r([\mathbf{X}_i]_{p:}, [\mathbf{X}_j]_{q:}), \stackrel{\text{def}}{=} \langle \phi_r([\mathbf{X}_i]_{p:}), \phi_r([\mathbf{X}_j]_{q:}) \rangle,\end{aligned}$$

$$[K_c(i, j)]_{pq} = K_c([X_i]_{:p}, [X_i]_{:q}) \stackrel{\text{def}}{=} \langle \phi_c([X_i]_{:p}), \phi_c([X_j]_{:q}) \rangle.$$

Therefore, we can successfully estimate α with quadratic programming based on K without description of feature mapping ϕ_r, ϕ_c . We update C as

$$C = \sum_{i=1}^n \alpha_i y_i \Phi(X_i) P (P^T P)^{-1} \in \mathcal{H}_r^r \times \mathcal{H}_c^r. \quad (4)$$

2. Second, we update P holding C fixed. The dual problem of Equation (2) is

$$\min_{\alpha=(\alpha_1, \dots, \alpha_n)} - \sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle C ((C^T C)^{-1} C^T \Phi(X_i)), C ((C^T C)^{-1} C^T \Phi(X_j)) \rangle, \quad (5)$$

subject to $0 \leq \alpha_i \leq C, i = 1, \dots, n$,

Notice $C ((C^T C)^{-1} C^T \Phi(X_i)) \in \mathcal{H}^{d_1} \times \mathcal{H}^{d_2}$ is well defined by matrix product: for $A_1 \in \mathcal{H}^r$ and $A_2 \in \mathcal{H}^d$, $A_1^T A_2 = \llbracket a_{ij} \rrbracket \in \mathbb{R}^{r \times d}$, where $a_{ij} = \langle [A_1]_i, [A_2]_j \rangle$. We can find an optimizer of (5) without the feature mapping. To show this, notice that by plugging (4) into (5), we have

$$C^T C = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (P^T P)^{-1} P^T K(i, j) P (P^T P)^{-1} \in \mathbb{R}^{r \times r} \times \mathbb{R}^{r \times r}, \quad (6)$$

$$C^T \Phi(X_i) = \sum_{j=1}^n \alpha_j y_j (P^T P)^{-1} P^T K(i, j) \in \mathbb{R}^{r \times d_1} \times \mathbb{R}^{r \times d_2}.$$

(6) makes inner product in (5) expressed in terms of only P and $\{K(i, j): i, j \in [n]\}$ by the following equation.

$$\langle C ((C^T C)^{-1} C^T \Phi(X_i)), C ((C^T C)^{-1} C^T \Phi(X_j)) \rangle = \text{tr} \left((C^T \Phi(X_i))^T (C^T C)^{-1} (C^T \Phi(X_j)) \right).$$

We update P from an optimal coefficient α of (5) and the formulas in (6).

$$P = \sum_{i=1}^n \alpha_i y_i (C^T C)^{-1} C^T \Phi(X_i).$$

Based on the procedure, we can obtain estimated $\hat{\alpha}$ and $\hat{P} = (\hat{P}_r, \hat{P}_c)$, Therefore, our estimated classifier is

$$\hat{f}(X) = \sum_{k=1}^n \hat{\alpha}_k y_k \left(\sum_{i=1}^{d_1} \sum_{j=1}^{d_1} [\hat{P}_r (\hat{P}_r^T \hat{P}_r)^{-1} \hat{P}_r^T]_{ij} K_r([X_k]_{:i}, [X]_{:j}) + \sum_{i=1}^{d_2} \sum_{j=1}^{d_2} [\hat{P}_c (\hat{P}_c^T \hat{P}_c)^{-1} \hat{P}_c^T]_{ij} K_c([X_k]_{:i}, [X]_{:j}) \right).$$

Notice $\hat{W}_r = \hat{P}_r (\hat{P}_r^T \hat{P}_r)^{-1} \hat{P}_r^T$, $\hat{W}_c = \hat{P}_c (\hat{P}_c^T \hat{P}_c)^{-1} \hat{P}_c^T$ and $\hat{\gamma} = y \circ \hat{\alpha}$.

Remark 1. From derivation of algorithm, we are looking for parameters α, μ, P, C which optimize

$$L_P = \frac{1}{2} \|CP^T\|_F^2 + c \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i \langle CP^T, \Phi(X_i) \rangle - (1 - \xi_i)) - \sum_{i=1}^n \mu_i \xi_i,$$

where α, μ are Lagrange multiplies. So dual problem (3) updates (C, α) and (5) updates (P, α) and $\mu = C - \alpha$. Therefore, we do not have to distinguish coefficients α, β in two dual problems as I specified before.

3 Relation with the previous algorithm symmetric trick

Define symmetric feature matrix $\tilde{X} = \begin{pmatrix} 0_{d_1 \times d_2} & X \\ X^t & 0_{d_2 \times d_1} \end{pmatrix} \in \mathbb{R}^{(d_1+d_2) \times (d_1+d_2)}$. Feature mapping 3 is defined as

$$\begin{aligned} \tilde{\Phi}: \mathbb{R}^{d_1 \times d_2} &\rightarrow \mathcal{H}^{d_1+d_2} \\ X &\mapsto \left(\phi([\tilde{X}]_{1:}), \dots, \phi([\tilde{X}]_{d_1+d_2:}) \right) \end{aligned}$$

where ϕ is induced by kernel $K: \mathbb{R}^{(d_1+d_2) \times (d_1+d_2)} \times \mathbb{R}^{(d_1+d_2) \times (d_1+d_2)} \rightarrow \mathbb{R}$. Since all entries of $\Phi_r(X)$ are corresponding to $[\tilde{\Phi}(X)]_{1:d_1}$ and $\Phi_c(X)$ to $[\tilde{\Phi}(X)]_{d_1+1:d_1+d_2}$, we have an equivalent representation of (1)

$$\begin{aligned} f(X) &= \langle B, \Phi(X) \rangle \\ &= \langle B_r, \Phi_r(X) \rangle + \langle B_c, \Phi_c(X) \rangle \\ &= \langle \tilde{B}_r, [\tilde{\Phi}(X)]_{1:d_1} \rangle + \langle \tilde{B}_c, [\tilde{\Phi}(X)]_{d_1+1:d_1+d_2} \rangle, \text{ where } \tilde{B}_r \in \mathcal{H}^{d_1}, \tilde{B}_c \in \mathcal{H}^{d_2} \\ &= \langle \tilde{B}, \tilde{\Phi}(X) \rangle, \text{ where } \tilde{B} = (\tilde{B}_r, \tilde{B}_c) \in \mathcal{H}^{d_1+d_2}. \end{aligned}$$

Assume that $\tilde{B} = \tilde{C} \tilde{P}^T$ where $\tilde{C} \in \mathcal{H}^r$, $\tilde{P} = (\tilde{P}_r, \tilde{P}_c) \in \mathbb{R}^{(d_1+d_2) \times r}$ and $\tilde{P}_r \in \mathbb{R}^{d_1 \times r}$, $\tilde{P}_c \in \mathbb{R}^{d_2 \times r}$. Let Π_r, Π_c are permutation operators such that

$$\begin{aligned} \text{Proj}_{\mathcal{H}_r} \left(\Pi_r [\tilde{\Phi}(X)]_{1:d_1} \right) &= \Phi_r(X) \\ \text{Proj}_{\mathcal{H}_c} \left(\Pi_c [\tilde{\Phi}(X)]_{d_1+1:d_1+d_2} \right) &= \Phi_c(X). \end{aligned}$$

Here, we denote $\text{Proj}_{\mathcal{H}_c}: \mathcal{H} \rightarrow \mathcal{H}_r$ and $\text{Proj}_{\mathcal{H}_r}: \mathcal{H} \rightarrow \mathcal{H}_c$ as entry-wise projection mappings. Then the following holds

$$\begin{aligned} \langle \tilde{B}, \tilde{\Phi}(X) \rangle &= \langle \tilde{C}(\tilde{P}_r, \tilde{P}_c)^T, \tilde{\Phi}(X) \rangle \\ &= \langle \tilde{C} \tilde{P}_r^T, [\tilde{\Phi}(X)]_{1:d_1} \rangle + \langle \tilde{C} \tilde{P}_c^T, [\tilde{\Phi}(X)]_{d_1+1:d_1+d_2} \rangle \\ &= \langle \Pi_r \tilde{C} \tilde{P}_r^T, \Pi_r [\tilde{\Phi}(X)]_{1:d_1} \rangle + \langle \Pi_c \tilde{C} \tilde{P}_c^T, \Pi_c [\tilde{\Phi}(X)]_{d_1+1:d_1+d_2} \rangle \\ &= \langle \tilde{C}_r \tilde{P}_r^T, \Phi_r(X) \rangle + \langle \tilde{C}_c \tilde{P}_c^T, \Phi_c(X) \rangle, \end{aligned}$$

where $\tilde{C}_r = \text{Proj}_{\mathcal{H}_r}(\Pi_r \tilde{C})$ and $\tilde{C}_c = \text{Proj}_{\mathcal{H}_c}(\Pi_c \tilde{C})$. Therefore, we can conclude that the low rankness of the coefficient on the feature image of $\tilde{\Phi}(X)$ implies the same low rankness of the coefficient of the feature image of $\Phi(X)$. The other direction is also true.