

# Approximate Leave-One-Out with kernel SVM

Wenda Zhou, Peng Xu

March 13, 2019

## 1 C-Support Vector Classification

Let  $\mathbf{K}$  denote a positive-definite kernel matrix (hence invertible), with  $\mathbf{K}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$ . The objective of a kernel SVC can be expressed in the “loss + penalty” variational form (notice the optimal  $\gamma^*$  satisfies  $\gamma^* = \mathbf{y} \odot \alpha$ , with  $\alpha$  being the dual solution):

$$\min_{\rho, \gamma} \sum_{j=1}^n \max [0, 1 - y_j h(x_j)] + \frac{\lambda}{2} \gamma^\top \mathbf{K} \gamma, \quad h(x_j) = \mathbf{K}_{:,j}^\top \gamma + \rho. \quad (1)$$

For simplicity we ignore the offset  $\rho$  for now. Let  $S$  and  $V$  be the smooth set and the set of singularities, respectively. For the  $j$ -th observation,  $j \in S$ , we have

$$\dot{\ell}(\mathbf{K}_{:,j}^\top \gamma) = -y_j \cdot \mathbf{1}\{y_j \mathbf{K}_{:,j}^\top \gamma < 1\}, \quad \ddot{\ell}(\mathbf{K}_{:,j}^\top \gamma) = 0.$$

Additionally,

$$\nabla R(\gamma) = \lambda \mathbf{K} \gamma, \quad \nabla^2 R(\gamma) = \lambda \mathbf{K}.$$

Substitute corresponding terms in Thm. 4.1, we deduce the ALO formula for kernel SVC:

$$\mathbf{K}_{:,i}^\top \hat{\gamma}^{\setminus i} = \mathbf{K}_{:,i}^\top \hat{\gamma} + a_i g_{\ell,i},$$

where

$$a_i = \begin{cases} \frac{1}{\lambda} \mathbf{K}_{:,i}^\top \left[ \mathbf{K}^{-1} - \mathbf{K}^{-1} \mathbf{K}_{:,V} \left( \mathbf{K}_{:,V}^\top \mathbf{K}^{-1} \mathbf{K}_{:,V} \right)^{-1} \mathbf{K}_{:,V}^\top \mathbf{K}^{-1} \right] \mathbf{K}_{:,i} & i \in S, \\ \left[ \lambda \left( \mathbf{K}_{:,V}^\top \mathbf{K}^{-1} \mathbf{K}_{:,V} \right)^{-1} \right]^{-1}_{ii} & i \in V, \end{cases}$$

and

$$g_{\ell,S} = -y_S \odot \mathbf{1} \left\{ y_S \mathbf{K}_{:,S}^\top \gamma < 1 \right\}, \quad g_{\ell,V} = \left( \mathbf{K}_{:,V}^\top \mathbf{K}_{:,V} \right)^{-1} \mathbf{K}_{:,V}^\top \left[ \sum_{j \in S: y_j \mathbf{K}_{:,j}^\top \gamma < 1} y_j \mathbf{K}_{:,j} - \lambda \mathbf{K} \alpha \right].$$

## 2 $\varepsilon$ -Support Vector Regression

The  $\varepsilon$ -SVR is associated with the  $\varepsilon$ -insensitive loss function. Its objective function can be written as (notice the optimal  $\gamma^*$  satisfies  $\gamma^* = \alpha^* - \alpha$ , where  $\alpha$  and  $\alpha^*$  are the dual solutions):

$$\min_{\rho, \gamma} \sum_{j=1}^n \max [0, |y_j - h(x_j)| - \varepsilon] + \frac{\lambda}{2} \gamma^\top K \gamma, \quad h(x_j) = K_{:,j}^\top \gamma + \rho. \quad (2)$$

So for the  $j$ -th observation,  $j \in S$ , we have

$$\dot{\ell}(K_{:,j}^\top \gamma) = -\text{sgn}(K_{:,j}^\top \gamma) \cdot \mathbf{1} \left\{ |y_j - K_{:,j}^\top \gamma| \geq \varepsilon \right\}, \quad \ddot{\ell}(K_{:,j}^\top \gamma) = 0.$$

Our ALO recipe will then be exactly the same as in C-SVC, except

$$g_{\ell, S} = -\text{sgn}(K_{:,S}^\top \gamma) \odot \mathbf{1} \left\{ |y_j - K_{:,S}^\top \gamma| \geq \varepsilon \right\},$$

and

$$g_{\ell, V} = \left( K_{:,V}^\top K_{:,V} \right)^{-1} K_{:,V}^\top \left[ \sum_{j \in S: |y_j - K_{:,j}^\top \gamma| \geq \varepsilon} \text{sgn}(K_{:,j}^\top \gamma) K_{:,j} - \lambda K \alpha \right].$$

## 3 $\nu$ -SVC and $\nu$ -SVR

$\nu$ -SVC and  $\nu$ -SVR are introduced in [Sch+00], with  $\nu \in (0, 1]$  a new parameter that can control the number of support vectors. Notably,  $\nu$  replaces  $C$  in C-SVC and  $\varepsilon$  in  $\varepsilon$ -SVR, and in the latter case  $\varepsilon$  itself becomes an optimization argument.

Because of the following equivalencies, there is no need to figure out the variational problem in order to derive ALO for  $\nu$ -SVC and  $\nu$ -SVR:

- Let  $\alpha^*$  be the dual solution of  $\nu$ -SVC with parameters  $(C, \nu)$ , then  $\alpha^*/\rho^*$  is an optimal solution of C-SVC with  $C = 1/(n\rho^*)$ . See [CL01].
- Let  $\alpha^*$  be the dual solution of  $\nu$ -SVR with parameters  $(C, \nu)$ , then  $\alpha^*$  is also the solution of  $\varepsilon$ -SVR with parameters  $(C/n, \varepsilon^*)$ . See [CL02].

Here  $\rho^*$  and  $\varepsilon^*$  are the corresponding primal solutions, and can be worked out quickly from  $\alpha^*$  through equations derived from KKT conditions. See [CL11, Section 4.2].

## 4 Multiclass-classifier

To handle the multinomial classification, LIBSVM adopts the one-versus-one strategy: for a dataset with  $y \in \mathcal{S}$ ,  $|\mathcal{S}| = K$ , LIBSVM will train  $K(K-1)/2$  classifiers  $\hat{f}_{(u,v)}(x)$ , one for each 2-combination  $(u, v)$  of  $\mathcal{S}$ . The prediction for a new observation  $x_{\text{new}}$  is then obtained through voting:  $\hat{y}_{\text{new}} = \text{mode}\{\hat{f}_{(u,v)}(x_{\text{new}})\}$ , and the label with the lowest index will be chosen in case of ties.

Note that for a specific pair  $(u', v')$ , the leave- $i$ -out prediction  $f_{(u', v')}^{\setminus i}(\mathbf{x}_j)$  is exactly the full-data prediction  $f_{(u', v')}(\mathbf{x}_j)$  when  $y_i$  is neither  $u'$  nor  $v'$ , for  $\mathbf{x}_i$  does not contribute to the training of  $f_{(u', v')}(\cdot)$  in anyway. Therefore, to obtain the multiclass ALO prediction  $\tilde{y}_j^{\setminus i}$ , we first construct all  $K(K-1)/2$  full-data predictions, then replace the  $K-1$  ones whose training are affected by  $y_i$  with the corresponding binary ALO predictions, and finally ensemble the result through popular voting.

## 5 Computation

To compute ALO for SVM, we must perform a variety of operations with the kernel matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$ . In practice  $n$  can often be quite large, therefore, naïvely implement the ALO formula will result in severe loss of efficiency. This issue can be, however, alleviated by exploiting the symmetry of  $\mathbf{K}$ : consider the Cholesky decomposition  $\mathbf{K} = \mathbf{L}\mathbf{L}^\top$ , then for an arbitrary index set  $E \subseteq \{1, \dots, n\}$ , we have

$$\mathbf{L}^{-1}\mathbf{K}_{:,E} = \mathbf{L}^{-1}\mathbf{L}\mathbf{L}_{E,:}^\top = \mathbf{L}_{E,:}^\top.$$

By further considering the thin QR factorization  $\mathbf{L}_{V,:}^\top = \mathbf{Q}_V \mathbf{R}_V$ , where  $\mathbf{Q}_V$  semi-orthogonal and  $\mathbf{R}_V$  upper-triangular, we may expressed  $a_s, s \in S$ , as a difference between two quadratic forms:

$$\begin{aligned} \lambda a_s &= \mathbf{K}_{:,s}^\top \mathbf{K}^{-1} \mathbf{K}_{:,s} - \mathbf{K}_{:,s}^\top \mathbf{K}^{-1} \mathbf{K}_{:,V} \left( \mathbf{K}_{:,V}^\top \mathbf{K}^{-1} \mathbf{K}_{:,V} \right)^{-1} \mathbf{K}_{:,V}^\top \mathbf{K}^{-1} \mathbf{K}_{:,s} \\ &= \mathbf{K}_{:,s}^\top (\mathbf{L}\mathbf{L}^\top)^{-1} \mathbf{K}_{:,s} - \mathbf{K}_{:,s}^\top (\mathbf{L}\mathbf{L}^\top)^{-1} \mathbf{K}_{:,V} \left( \mathbf{K}_{:,V}^\top (\mathbf{L}\mathbf{L}^\top)^{-1} \mathbf{K}_{:,V} \right)^{-1} \mathbf{K}_{:,V}^\top (\mathbf{L}\mathbf{L}^\top)^{-1} \mathbf{K}_{:,s} \\ &= (\mathbf{L}^{-1} \mathbf{K}_{:,s})^\top \mathbf{L}^{-1} \mathbf{K}_{:,s} - (\mathbf{L}^{-1} \mathbf{K}_{:,s})^\top (\mathbf{L}^{-1} \mathbf{K}_{:,V}) \left[ (\mathbf{L}^{-1} \mathbf{K}_{:,V})^\top (\mathbf{L}^{-1} \mathbf{K}_{:,V}) \right]^{-1} (\mathbf{L}^{-1} \mathbf{K}_{:,V})^\top (\mathbf{L}^{-1} \mathbf{K}_{:,s}) \\ &= \mathbf{L}_{s,:} \mathbf{L}_{s,:}^\top - \mathbf{L}_{s,:} \mathbf{L}_{V,:}^\top \left( \mathbf{L}_{V,:} \mathbf{L}_{V,:}^\top \right)^{-1} \mathbf{L}_{V,:} \mathbf{L}_{s,:}^\top \\ &= \mathbf{L}_{s,:} \mathbf{L}_{s,:}^\top - (\mathbf{Q}_V^\top \mathbf{L}_{s,:})^\top (\mathbf{Q}_V^\top \mathbf{L}_{s,:}). \end{aligned}$$

Similarly,  $a_v$  with  $v \in V$  can be evaluated through the inner product of two vectors:

$$(\lambda a_v)^{-1} = \left( \mathbf{L}_{V,:} \mathbf{L}_{V,:}^\top \right)_{v,v}^{-1} = (\mathbf{R}_V^\top \mathbf{Q}_V^\top \mathbf{Q}_V \mathbf{R}_V)_{v,v}^{-1} = \mathbf{R}_{v,:}^{-1} \mathbf{R}_{:,v}^{-1}.$$

As a result, with  $\mathbf{L}$ , and thus  $\mathbf{L}_{S,:}, \mathbf{L}_{V,:}$  computed beforehand, we avoided direct multiplication of large matrices and reduced computational costs.

## References

- [CL01] Chih-Chung Chang and Chih-Jen Lin. "Training v-Support Vector Classifiers: Theory and Algorithms". In: *Neural Computation* 13.9 (2001), pp. 2119–2147. DOI: 10.1162/089976601750399335.
- [CL02] Chih-Chung Chang and Chih-Jen Lin. "Training v-Support Vector Regression: Theory and Algorithms". In: *Neural Computation* 14.8 (2002), pp. 1959–1977. DOI: 10.1162/089976602760128081.

- [CL11] Chih-Chung Chang and Chih-Jen Lin. “LIBSVM: A Library for Support Vector Machines”. In: *ACM Trans. Intell. Syst. Technol.* 2.3 (May 2011), 27:1–27:27. ISSN: 2157-6904. DOI: 10.1145/1961189.1961199. URL: <https://www.csie.ntu.edu.tw/~cjlin/papers/libsvm.pdf>.
- [Sch+00] Bernhard Schölkopf et al. “New Support Vector Algorithms”. In: *Neural Computation* 12.5 (2000), pp. 1207–1245. DOI: 10.1162/089976600300015565.