# Approximate Leave-One-Out with kernel SVM

Wenda Zhou, Peng Xu

March 6, 2019

## 1 $C$-Support Vector Classification

Let $K$ denote a positive-definite kernel matrix (hence invertible), with $K_{i,j} = K(x_i, x_j)$. The objective of a kernel SVC can be expressed in the "loss + penalty" variational form:

$$\min_{\rho,\alpha} \sum_{j=1}^{n} \max\left[0, 1 - y_j h(x_j)\right] + \frac{\lambda}{2}\alpha^\top K\alpha, \qquad h(x_j) = K_{\cdot,j}^\top \alpha + \rho. \tag{1}$$

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

$$\dot{\ell}(K_{\cdot,j}^\top \alpha) = -y_j \cdot \mathbf{1}\{y_j K_{\cdot,j}^\top \alpha < 1\}, \qquad \ddot{\ell}(K_{\cdot,j}^\top \alpha) = 0.$$

Additionally,

$$\nabla R(\alpha) = \lambda K\alpha, \qquad \nabla^2 R(\alpha) = \lambda K.$$

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

$$K_{\cdot,i}^\top \tilde{\alpha}^{\setminus i} = K_{\cdot,i}^\top \hat{\alpha} + a_i g_{\ell,i},$$

where

$$a_i = \begin{cases} \dfrac{1}{\lambda} K_{\cdot,i}^\top \left[ K^{-1} - K^{-1}K_{\cdot,V}\left(K_{\cdot,V}^\top K^{-1}K_{\cdot,V}\right)^{-1}K_{\cdot,V}^\top K^{-1} \right] K_{\cdot,i} & i \in S, \\[2em] \left[ \lambda \left(K_{\cdot,V}^\top K^{-1}K_{\cdot,V}\right)^{-1}_{ii} \right]^{-1} & i \in V, \end{cases}$$

and

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

## 2   $\varepsilon$-Support Vector Regression

The $\varepsilon$-SVR is associated with the $\varepsilon$-insensitive loss function. Its objective function can be written as:

$$\min_{\rho, \alpha} \sum_{j=1}^{n} \max \left[0, |y_j - h(x_j)| - \varepsilon\right] + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{K} \boldsymbol{\alpha}, \qquad h(x_j) = \boldsymbol{K}_{\cdot,j}^\top \boldsymbol{\alpha} + \rho. \tag{2}$$

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

$$\dot{\ell}(\boldsymbol{K}_{\cdot,j}^\top \boldsymbol{\alpha}) = -\operatorname{sgn}\left(\boldsymbol{K}_{\cdot,j}^\top \boldsymbol{\alpha}\right) \cdot \boldsymbol{1}\left\{\left|y_j - \boldsymbol{K}_{\cdot,j}^\top \boldsymbol{\alpha}\right| \geq \varepsilon\right\}, \qquad \ddot{\ell}(\boldsymbol{K}_{\cdot,j}^\top \boldsymbol{\alpha}) = 0.$$

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

$$g_{\ell,S} = -\operatorname{sgn}\left(\boldsymbol{K}_{\cdot,S}^\top \boldsymbol{\alpha}\right) \odot \boldsymbol{1}\left\{\left|y_j - \boldsymbol{K}_{\cdot,S}^\top \boldsymbol{\alpha}\right| \geq \varepsilon\right\},$$

and

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

## 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 $\hat{f}_{(u,v)}(\boldsymbol{x})$, one for each 2-combination $(u, v)$ of $\mathcal{S}$. The prediction for a new observation $\boldsymbol{x}_{\text{new}}$ is then obtained through voting: $\hat{y}_{\text{new}} = \operatorname{mode}\{\hat{f}_{(u,v)}(\boldsymbol{x}_{\text{new}})\}$, 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 $f_{(u',v')}^{\backslash i}(\boldsymbol{x}_j)$ is exactly the full-data prediction $f_{(u',v')}(\boldsymbol{x}_j)$ when $y_i$ is neither $u'$ nor $v'$, for $\boldsymbol{x}_i$ does not contribute to the training of

$f_{(u',v')}(\cdot)$ in anyway. Therefore, to obtain the multiclass ALO prediction $\tilde{y}_j^{\backslash i}$, 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 $K \in \mathbb{R}^{n \times n}$. 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 $K$: consider the Cholesky decomposition $K = LL^\top$, then for an arbitrary index set $E \subseteq \{1, \ldots, n\}$, we have

$$L^{-1}K_{\cdot,E} = L^{-1}LL_{E,\cdot}^\top = L_{E,\cdot}^\top.$$

By further considering the thin QR factorization $L_{V,\cdot}^\top = Q_V R_V$, where $Q_V$ semi-orthogonal and $R_V$ upper-triangular, we may expressed $a_s$, $s \in S$, as a difference of two quadratic form:

$$
\begin{aligned}
\lambda a_s &= K_{,s}^\top K^{-1} K_{,s} - K_{,s}^\top K^{-1} K_{\cdot,V} \left( K_{\cdot,V}^\top K^{-1} K_{\cdot,V} \right)^{-1} K_{\cdot,V}^\top K^{-1} K_{,s} \\
&= K_{,s}^\top (LL^\top)^{-1} K_{,s} - K_{,s}^\top (LL^\top)^{-1} K_{\cdot,V} \left( K_{\cdot,V}^\top (LL^\top)^{-1} K_{\cdot,V} \right)^{-1} K_{\cdot,V}^\top (LL^\top)^{-1} K_{,s} \\
&= (L^{-1}K_{,s})^\top L^{-1}K_{,s} - (L^{-1}K_{,s})^\top (L^{-1}K_{\cdot,V}) \left[ (L^{-1}K_{\cdot,V})^\top (L^{-1}K_{\cdot,V}) \right]^{-1} (L^{-1}K_{\cdot,V})^\top (L^{-1}K_{,s}) \\
&= L_{s,\cdot}L_{s,\cdot}^\top - L_{s,\cdot}L_{V,\cdot}^\top \left( L_{V,\cdot}L_{V,\cdot}^\top \right)^{-1} L_{V,\cdot}L_{s,\cdot}^\top \\
&= L_{s,\cdot}L_{s,\cdot}^\top - (Q_V^\top L_{s,\cdot})^\top (Q_V^\top L_{s,\cdot}).
\end{aligned}
$$

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

$$(\lambda a_v)^{-1} = \left( L_{V,\cdot}L_{V,\cdot}^\top \right)_{v,v}^{-1} = \left( R_V^\top Q_V^\top Q_V R_V \right)_{v,v}^{-1} = R_{v,\cdot}^{-1} R_{\cdot,v}^{-1}.$$

As a result, with $L$, and thus $L_{S,\cdot}$, $L_{V,\cdot}$ 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.