

CS 536 :What is the Gaussian Kernel even doing?

16:198:536

SVMs are frequently used in an ‘off-the-shelf’ way - take your data, pick your favorite kernel, throw the nearest SVM solver at the problem, and see what happens. A common choice for the kernel, before anything else is known, is the Gaussian Kernel

$$K(\underline{x}, \underline{y}) = \exp\left(-\frac{1}{\sigma^2} \|\underline{x} - \underline{y}\|^2\right). \quad (1)$$

But what is this kernel doing? What does the ‘Gaussian Kernel SVM’ mean? What is the importance of σ^2 ?

To illustrate this by example, consider the usual XOR data:

$$\{A = ((-1, +1), +1), B = ((-1, -1), -1), C = ((+1, -1), +1), D = ((+1, +1), -1)\}.$$

Writing out the dual SVM in this case, we have

$$\begin{aligned} \max_{\alpha_A, \alpha_B, \alpha_C, \alpha_D} \quad & [\alpha_A + \alpha_B + \alpha_C + \alpha_D] - \frac{1}{2} [\alpha_A^2 K(A, A) + \alpha_B^2 K(B, B) + \alpha_C^2 K(C, C) + \alpha_D^2 K(D, D)] \\ & - [-\alpha_A K(A, B) \alpha_B + \alpha_A K(A, C) \alpha_C - \alpha_A K(A, D) \alpha_D] \\ & - [-\alpha_B K(B, C) \alpha_C + \alpha_B K(B, D) \alpha_D - \alpha_C K(C, D) \alpha_D] \\ \text{(s.t.)} \quad & \alpha_A - \alpha_B + \alpha_C - \alpha_D = 0 \\ & \alpha_A, \alpha_B, \alpha_C, \alpha_D \geq 0. \end{aligned} \quad (2)$$

To analyze the solution to the Gaussian Kernel SVM and the effect of σ^2 , we need an explicit, general solution, so it’s worth taking the time to solve out the above explicitly.

To simplify this, note that regardless of σ^2 , we have that $K(A, A) = K(B, B) = K(C, C) = K(D, D) = 1$. Additionally, $K(A, B) = \exp(-(1/\sigma^2)2^2) = \exp(-4/\sigma^2)$. Similarly, $K(A, C) = K(B, D) = \exp(-8/\sigma^2)$. Denote $\rho = \exp(-4/\sigma^2)$ for notational convenience. We wish to solve

$$\begin{aligned} \max_{\alpha_A, \alpha_B, \alpha_C, \alpha_D} \quad & [\alpha_A + \alpha_B + \alpha_C + \alpha_D] - \frac{1}{2} [\alpha_A^2 + \alpha_B^2 + \alpha_C^2 + \alpha_D^2] \\ & - [- (\alpha_A \alpha_B) \rho + (\alpha_A \alpha_C) \rho^2 - (\alpha_A \alpha_D) \rho] \\ & - [- (\alpha_B \alpha_C) \rho + (\alpha_B \alpha_D) \rho^2 - (\alpha_C \alpha_D) \rho] \\ \text{(s.t.)} \quad & \alpha_A - \alpha_B + \alpha_C - \alpha_D = 0 \\ & \alpha_A, \alpha_B, \alpha_C, \alpha_D \geq 0. \end{aligned} \quad (3)$$

This simplifies further still:

$$\begin{aligned} \max_{\alpha_A, \alpha_B, \alpha_C, \alpha_D} \quad & [\alpha_A + \alpha_B + \alpha_C + \alpha_D] - \frac{1}{2} [\alpha_A^2 + \alpha_B^2 + \alpha_C^2 + \alpha_D^2] + [\alpha_A + \alpha_C] [\alpha_B + \alpha_D] \rho - [\alpha_A \alpha_C + \alpha_B \alpha_D] \rho^2 \\ \text{(s.t.)} \quad & \alpha_A + \alpha_C = \alpha_B + \alpha_D \\ & \alpha_A, \alpha_B, \alpha_C, \alpha_D \geq 0. \end{aligned} \quad (4)$$

Noting that $\alpha_A + \alpha_C = \alpha_B + \alpha_D$, it’s convenient to introduce a new variable ϵ to represent this sum, so that $\alpha_C = \epsilon - \alpha_A$, $\alpha_D = \epsilon - \alpha_B$:

$$\begin{aligned} \max_{\alpha_A, \alpha_B, \epsilon} \quad & [2\epsilon] - \frac{1}{2} [\alpha_A^2 + \alpha_B^2 + (\epsilon - \alpha_A)^2 + (\epsilon - \alpha_B)^2] + \epsilon^2 \rho - [\alpha_A(\epsilon - \alpha_A) + \alpha_B(\epsilon - \alpha_B)] \rho^2 \\ \text{(s.t.)} \quad & \alpha_A \leq \epsilon, \alpha_B \leq \epsilon \\ & \alpha_A, \alpha_B, \epsilon \geq 0, \end{aligned} \quad (5)$$

or (grouping like terms)

$$\begin{aligned} \max_{\alpha_A, \alpha_B, \epsilon} & [2\epsilon + \epsilon^2 + \epsilon^2 \rho] - ([\alpha_A^2 - \epsilon \alpha_A] + [\alpha_A \epsilon - \alpha_A^2] \rho^2) - ([\alpha_B^2 - \epsilon \alpha_B] + [\alpha_B \epsilon - \alpha_B^2] \rho^2) \\ \text{(s.t.) } & \alpha_A \leq \epsilon, \alpha_B \leq \epsilon \\ & \alpha_A, \alpha_B, \epsilon \geq 0. \end{aligned} \quad (6)$$

One observation to make here is that the above is completely separable and symmetric in terms of α_A and α_B - we can imagine trying to maximize with respect to these for any given value of ϵ , but because the objective function is separable and symmetric, for any given value of ϵ we'll get the same solution for both, $\alpha_A = \alpha_B = \alpha$ for some α that depends on ϵ . This simplifies the problem even further:

$$\begin{aligned} \max_{\alpha, \epsilon} & [2\epsilon + \epsilon^2 + \epsilon^2 \rho] - 2([\alpha^2 - \epsilon \alpha] + [\alpha \epsilon - \alpha^2] \rho^2) \\ \text{(s.t.) } & 0 \leq \alpha \leq \epsilon. \end{aligned} \quad (7)$$

The (unconstrained) maximum of the above occurs when $\alpha = \epsilon/2$ (why?) so since this is feasible for the constraints, this is the solution α for any given value of ϵ . Since $\alpha_A = \alpha_B = \alpha = \epsilon/2$, and $\alpha_C = \epsilon - \alpha_A, \alpha_D = \epsilon - \alpha_B$, we get that $\alpha_A = \alpha_B = \alpha_C = \alpha_D = \epsilon/2$. But what is the optimal value of ϵ ? Going back to the original problem, for effect, we have

$$\begin{aligned} \max_{\epsilon} & [4(\epsilon/2)] - \frac{1}{2} [4(\epsilon/2)^2] - (\epsilon/2)^2 [-\rho + \rho^2 - \rho] - (\epsilon/2)^2 [-\rho + \rho^2 - \rho] \\ \text{(s.t.) } & \epsilon \geq 0. \end{aligned} \quad (8)$$

or

$$\begin{aligned} \max_{\epsilon} & [2\epsilon] - \frac{1}{2} [\epsilon^2] - 2(\epsilon/2)^2 [\rho^2 - 2\rho] \\ \text{(s.t.) } & \epsilon \geq 0. \end{aligned} \quad (9)$$

Maximizing this in an unconstrained way, the maximum occurs when $\epsilon = 2/(1 - \rho)^2$, which satisfies $\epsilon > 0$, so we get a final solution of

$$\alpha_A = \alpha_B = \alpha_C = \alpha_D = \frac{1}{(1 - e^{-4/\sigma^2})^2}. \quad (10)$$

To reconstruct the primal, note that each vector is a support vector, taking A for instance we want (looking at the kernelized version of $y^i(\underline{w} \cdot \underline{x}^i + b) = 1$,

$$y^A(\alpha_A y^A K(A, A) + \alpha_B y^B K(B, A) + \alpha_C y^C K(C, A) + \alpha_D y^D K(D, A) + b) = 1, \quad (11)$$

or

$$\left(\frac{1}{(1 - e^{-4/\sigma^2})^2} - 2 \frac{1}{(1 - e^{-4/\sigma^2})^2} \exp(-4/\sigma^2) + \frac{1}{(1 - e^{-4/\sigma^2})^2} \exp(-8/\sigma^2) + b \right) = 1, \quad (12)$$

or

$$\frac{1}{(1 - e^{-4/\sigma^2})^2} (1 - 2 \exp(-4/\sigma^2) + \exp(-8/\sigma^2)) + b = 1, \quad (13)$$

and it's relatively straightforward to show that the above solves for $b = 0$.

Having solved for the bias, we are now in a position to build the SVM classifier:

$$\text{classify}(\underline{x}) = \text{sign} \left(\frac{K(A, \underline{x}) - K(B, \underline{x}) + K(C, \underline{x}) - K(D, \underline{x})}{(1 - e^{-4/\sigma^2})^2} \right). \quad (14)$$

However, scaling by a positive constant doesn't actually influence the result of the classifier, so we get a final classifier of

$$\text{classify}(\underline{x}) = \text{sign} \left(e^{-\frac{1}{\sigma^2} \|A - \underline{x}\|^2} - e^{-\frac{1}{\sigma^2} \|B - \underline{x}\|^2} + e^{-\frac{1}{\sigma^2} \|C - \underline{x}\|^2} - e^{-\frac{1}{\sigma^2} \|D - \underline{x}\|^2} \right). \quad (15)$$

At this point we are in a position to ask - what does the value σ^2 do to the classifier? Notice what happens as $\sigma^2 \rightarrow 0$. As we take this limit, we are dividing by a smaller and smaller value, which is effectively ‘blowing up’ the exponents. The exponential term with the smallest exponent will be *the point with the minimal distance to \underline{x}* ! This term is going to ‘dominate’ the overall sum in the limit as $\sigma^2 \rightarrow 0$, which means that the sign of the sum is going to be determined by the point \underline{x} is closest to!

We see then that what the Gaussian classifier is doing in the limit is performing a version of ‘k-Nearest Neighbors’ classification. k-NN Classification classifies a new point by looking at its nearest neighbors among the training data, and classifying the new point based on the majority classes of its neighbors. The Gaussian Kernel Classifier is effectively approximating this, where σ^2 controls how far away you are willing to look for neighbors (note that as $\sigma^2 \rightarrow \infty$, all exponents are roughly the same (near 0) so all points are considered equally regardless of distance/neighborliness), and the sign of the sum is determined by the ‘weighted majority class’ of these nearest neighbors.

The Gaussian Kernel SVM has two advantages potentially over ‘raw’ k-Nearest Neighbors classification. The first is that the model is ‘smooth’ - k-NN has very discontinuous properties based on how many neighbors you choose to classify with; the Gaussian Kernel SVM scales point importance with distance. The second is that while k-Nearest Neighbors requires you to remember *every* training point, the Gaussian Kernel SVM identifies only the potential neighbors that really matter to classification (the support vectors). As $\sigma^2 \rightarrow \infty$, fewer neighbors will matter for classification and it is easy to imagine the number of support vectors should go way down; as $\sigma^2 \rightarrow 0$, you’re essentially looking at 1-NN classification, and it is reasonable that *all* the vectors will be support vectors on some (high dimensional) boundary. But the Gaussian Kernel SVM provides a very convenient framework for both solving a given model and exploring what kinds of models (based on σ^2) you might need for your data.