

HW2: Linear Regression and Classification, Kernel Methods

1 Written Exercises

1. Consider the following variant of the sum-of-squared-errors objective for linear regression:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \quad (1)$$

What we have basically done is modulate each of the error terms in the above summation using a variable r_n . Derive the expression for the weight vector \mathbf{w} that minimizes $E(\mathbf{w})$. Also explain what role does r_n play in Equation 1? (**10 marks**)

Distribute the r_n into the terms of errors:

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N r_n (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \\ &= \frac{1}{2} \sum_{n=1}^N (\sqrt{r_n} y_n - \sqrt{r_n} \mathbf{w}^\top \mathbf{x}_n)^2 \\ &= \frac{1}{2} (\sqrt{R} Y - \sqrt{R} X \mathbf{w})^T (\sqrt{R} Y - \sqrt{R} X \mathbf{w}) \end{aligned}$$

Substitute $\sqrt{R} X$ for F and $\sqrt{R} Y$ for G :

$$E(\mathbf{w}) = \frac{1}{2} (G - F \mathbf{w})^T (G - F \mathbf{w})$$

Then, the derivative follows as shown in class (for the now-identical equation $\frac{1}{2} (Y - X \mathbf{w})^T (Y - X \mathbf{w})$):

$$\begin{aligned} \nabla E(\mathbf{w}) &= -F^T (G - F \mathbf{w}) \\ 0 &= -F^T (G - F \mathbf{w}) \\ \Rightarrow &= F^T F \mathbf{w} = F^T G \end{aligned}$$

Now invert for our solution:

$$\begin{aligned} \hat{\mathbf{w}} &= (F^T F)^{-1} F^T G \\ &= ((\sqrt{R} X)^T \sqrt{R} X)^{-1} (\sqrt{R} X)^T \sqrt{R} Y \end{aligned}$$

Note: If it helps your calculations, you may want to rewrite the above objective in matrix notation (e.g., the standard linear regression objective $\sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$ can be written as $(\mathbf{Y} - \mathbf{X}\mathbf{w})^\top (\mathbf{Y} - \mathbf{X}\mathbf{w})$).

2. (6350 only) We have seen examples of regularizers where we penalize large values of the weights by imposing some penalty on the norm (e.g., ℓ_2 or ℓ_1 norm) of the weight vector. It turns out that there is another way to accomplish the same effect: by adding some noise to each of input variables.

Consider the linear regression model $y = \mathbf{w}^\top \mathbf{x}$ and the sum-of-squared-errors loss function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y_n - \mathbf{w}^\top \mathbf{x}_n\}^2$$

Let's assume we added Gaussian noise ϵ_n with zero mean and variance σ^2 to each input \mathbf{x}_n (so the input now becomes $\mathbf{x}_n + \epsilon_n$). Show that, with the noisy inputs, minimizing the *expected* loss function, i.e., $\mathbb{E}[E(\mathbf{w})]$ is equivalent to minimizing the original sum-of-squared-errors error for noise-free inputs *with the addition of a weight decay regularization term*. The expectation \mathbb{E} is w.r.t. the Gaussian noise distribution for which the following properties hold: $\mathbb{E}[\epsilon_n] = 0$ and $\mathbb{E}[\epsilon_m \epsilon_n] = \delta_{mn} \sigma^2$ (where $\delta_{mn} = 1$ if $m = n$, and 0 otherwise). (10 marks)

3. We saw that the Perceptron loss function can be written as $L(\mathbf{w}, b) = \sum_{n=1}^N \max\{0, -y_n(\mathbf{w}^\top \mathbf{x}_n + b)\}$. Using *stochastic gradient descent* (SGD) on the loss function defined above, derive the Perceptron update rule which states that for a mispredicted training example \mathbf{x}_n :

$$\begin{aligned} \mathbf{w} &= \mathbf{w} + y_n \mathbf{x}_n \\ b &= b + y_n \end{aligned}$$

The stochastic gradient descent is just like gradient descent; the only difference is that, in each iteration, it computes the gradient using a single example (instead of using all the examples which standard gradient descent does), updates the parameters using this gradient information (just like standard gradient descent), moves on to the next example and continues. (5 marks)

Our basis is, straightforwardly, $\mathbf{w}^\top \mathbf{x}_n + b$. We are given the update rules $\mathbf{w} = \mathbf{w} + y_n \mathbf{x}_n$ and $b = b + y_n$, and we begin by plugging them directly into this basis equation. This gives us:

$$\begin{aligned} \mathbf{w}_{new}^\top \mathbf{x}_n + b_{new} &= (\mathbf{w} + y_n \mathbf{x}_n)^\top \mathbf{x}_n + b + y_n \\ &= (\mathbf{w}^\top + y_n \mathbf{x}_n^\top) \mathbf{x}_n + b + y_n \\ &= \mathbf{w}^\top \mathbf{x}_n + y_n \mathbf{x}_n^\top \mathbf{x}_n + b + y_n \\ &= (\mathbf{w}^\top \mathbf{x}_n + b) + y_n \mathbf{x}_n^\top \mathbf{x}_n + y_n \end{aligned}$$

This, of course breaks down into two cases. In the first case, we've misclassified the positive example (i.e., $y_n = 1$), in which case we get $(\mathbf{w}^\top \mathbf{x}_n + b) + \mathbf{x}_n^\top \mathbf{x}_n + 1$, which we can trivially see is less negative than the basis $\mathbf{w}^\top \mathbf{x}_n + b$. In the other case, we've misclassified a negative example (i.e., $y_n = -1$), in which case we get $(\mathbf{w}^\top \mathbf{x}_n + b) - \mathbf{x}_n^\top \mathbf{x}_n - 1$, which is trivially see is less positive than the basis. In both cases, we are iteratively correcting our error.

4. The standard Perceptron algorithm makes an update when the incoming training example is misclassified by the current weight vector. For the binary classification case, we write this condition as $\text{sign}(\mathbf{w}^\top \mathbf{x} + b) \neq y_n$ or equivalently $y_n(\mathbf{w}^\top \mathbf{x} + b) < 0$. We know that if the training data is linearly separable, the standard Perceptron will find *some* hyperplane that separates the two classes.

Now suppose we change the misclassification condition to $y_n(\mathbf{w}^T \mathbf{x} + b) < \tau$ where τ is a positive constant. How does the hyperplane found in this case differ from the standard Perceptron? Do you expect the new solution will generalize better or worse than the standard Perceptron on the test data? **(5 marks)**

The old conditions tested only whether the desired classification y_n and the observed result (*e.g.*, the result of the postulated-hyperplane function) have different signs. This new τ -based condition tests both whether they have different signs *and* whether the observed result is $< \tau$. In other words, we will only consider an answer correct if it is *especially* correct. This will probably generalize better, since we are ignoring the answers most likely to be incorrect, and instead fitting for the general model (though of course this depends on the exact value of τ).

5. In the class, we considered the non-separable case of SVM by penalizing the misclassified training examples with a sum of slacks term ($\sum_n \xi_n$). Consider a variation of this problem where we penalize the sum of **squared** slacks (i.e., $\sum_n \xi_n^2$). The optimization problem would be of the form:

$$\begin{aligned} \text{Minimize } f(\mathbf{w}, b) &= \frac{\|\mathbf{w}\|^2}{2} + \frac{C}{2} \sum_{n=1}^N \xi_n^2 \\ \text{subject to } y_n(\mathbf{w}^T \mathbf{x}_n + b) &\geq 1 - \xi_n, \quad \xi_n \geq 0 \quad n = 1, \dots, N \end{aligned}$$

Starting with this objective function, first write down the primal Lagrangian. Then, optimize the primal Lagrangian with respect to \mathbf{w} and b , plug these solutions back in and write the optimization problem just in terms of the dual (Lagrange) variables α . How does this compare to the dual formulation for the standard SVM where we had a sum of slacks term ($\sum_n \xi_n$)? **(10 marks)**

The Primal Lagrangian is straightforwardly given from the objective function noted above:

$$L_p(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{\|\mathbf{w}\|^2}{2} + \frac{C}{2} \sum_{n=1}^N \xi_n^2 + \sum_{n=1}^N \alpha \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) - \xi_n\} - \sum_{n=1}^N \beta_n \xi_n$$

subject to $\alpha_n, \beta_n \geq 0; n = 1, \dots, N$

We begin by finding the partial derivatives of each relevant term:

$$\begin{aligned} \frac{\partial L_p}{\partial \mathbf{w}} &= \sum_N (\alpha - \alpha y_n(\mathbf{w}^T \mathbf{x}_n + b) - \alpha \xi_n) \\ &= \sum_N \alpha y_n(\mathbf{w}^T \mathbf{x}_n + b) \\ \mathbf{w} &= \sum_N \alpha y_n \mathbf{x}_n \end{aligned}$$

$$\begin{aligned} \frac{\partial L_p}{\partial b} &= \sum_N \alpha (1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) - \xi_n) \\ &= \sum_N \alpha y_n b \\ b &= \sum_N \alpha y_n \end{aligned}$$

$$\begin{aligned} \frac{\partial L_p}{\partial \xi_n} &= \frac{C}{2} \sum_N \xi_n^2 - \sum_N \alpha \xi_n - \beta \xi_n \\ &= C \sum_N \xi_n - \sum_N \alpha - \beta \\ C \sum_N \xi_n &= \sum_N \alpha - \beta \\ \xi_n &= \frac{\alpha - \beta}{C} \end{aligned}$$

Then we plug them back into the Primal Lagrangian to get the Dual Lagrangian. Luckily the notes have given us most of this step already. All that remains is to isolate the parts that are different and add those to the end:

$$\begin{aligned} L_D(\mathbf{w}, b, \xi, \alpha, \beta) &= \sum_N \alpha_n - \frac{1}{2} \sum_{m,n} \alpha_m \alpha_n y_m y_n (\mathbf{x}_m^T \mathbf{x}_n) + \sum_n (\alpha \xi - \beta \xi + \frac{C}{2} \xi^2) \\ &= \sum_N \alpha_n - \frac{1}{2} \sum_{m,n} \alpha_m \alpha_n y_m y_n (\mathbf{x}_m^T \mathbf{x}_n) + \sum_n \alpha \left(\frac{\alpha + \beta}{C} - \beta \frac{\alpha + \beta}{C} + \frac{C}{2} \frac{\alpha + \beta}{C} \right) \\ &\text{subject to } \sum_n \alpha_n y_n = 0, 0 \leq \alpha_n \leq C; n \in [1..N] \end{aligned}$$

6. Explain in 50 words (or less) why, intuitively, large margins are good. (5 marks)

It is said that some margin $\gamma = \frac{1}{\|w\|}$, so for a margin to be large, $\|w\|$ must be pretty small (and thus its constituents $w_i \in w$ must also be usually pretty small). Optimizing for simple weights tends to make simpler solutions, and this tends to point to good generalization.

7. Why having a regularizer that prefers most of the entries of the weight vector to be *exactly* zero is the most desirable (call it the “optimal regularizer”)? Why is it difficult to accomplish this in practice? What alternative regularizers are actually used in practice? Rank them (a) in the order of closeness to the optimal regularizer (closest first), (b) ease of optimization (easiest first). Which of these regularizers do you like the most, and why? **(5 marks)**

Generally we want prediction to be based on as small a number of features as possible; as regularizers go, this typically means that we prefer as many weights in w to be 0 as possible, since this leaves prediction to the few weights that are nonzero. Unfortunately, minimalization of the count of nonzero weights in regularizers is NP-hard, so this is usually not a feasible requirement. Alternatives are: norm-based regularizers,

8. Consider a kernel function $k(x, z) = x^\top z + 2(x^\top z)^2$. As we know, each kernel has an associated feature mapping ϕ such that $k(x, z) = \phi(x)^\top \phi(z)$. Assume that each example has two features (so $x = \{x_1, x_2\}, z = \{z_1, z_2\}$). Write down the feature mapping ϕ associated with the kernel defined above. **(5 marks)**

Given $k(x, z) = x^\top z + 2(x^\top z)^2$ for $x = \{x_1, x_2\}, z = \{z_1, z_2\}$:

$$\begin{aligned} k(x, z) &= (x_1 z_1 + x_2 z_2) + 2(x_1 z_1 + x_2 z_2)^2 \\ &= (x_1 z_1 + x_2 z_2) + 2(x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2) \\ &= x_1 z_1 x_2 z_2 + 2x_1^2 z_1^2 + 2x_2^2 z_2^2 + 4x_1 x_2 z_1 z_2 \\ &= (x_1, x_2, \sqrt{2}x_1^2, \sqrt{2}x_2^2, 2x_1 x_2)^T (z_1, z_2, \sqrt{2}z_1^2, \sqrt{2}z_2^2, 2z_1 z_2) \end{aligned}$$

Thus:

$$\begin{aligned} \phi(x) &= (x_1, x_2, \sqrt{2}x_1^2, \sqrt{2}x_2^2, 2x_1 x_2) \\ \phi(z) &= (z_1, z_2, \sqrt{2}z_1^2, \sqrt{2}z_2^2, 2z_1 z_2) \end{aligned}$$

9. The Euclidean distance between two points x and z is defined as $\|x - z\|^2$. It's a linear distance metric in the sense that it measures the shortest direct distance between two points. It turns out that we can turn it to a *nonlinear* distance metric using kernels which may be useful if you want to measure the distance between two points *along* some nonlinear curve/surface. As we discussed while talking about kernel methods, so long as we can express an operation in form of dot (inner) products, we can kernelized (i.e., non-linearize) it by replacing each inner product by an inner product in a high dimensional feature space. Using this fact, show that the Euclidean distance computation can be kernelized as well, given some kernel function k . Give the expression of the nonlinearized version of the distance function in terms of k . Also, can you think of an algorithm (among the ones you have already seen in the class) where you could use this? **(5 marks)**

Let the vectors $\mathbf{x} = \{x_1, x_2\}$ and $\mathbf{z} = \{z_1, z_2\}$. Then, let the kernel function $k(\mathbf{x}, \mathbf{z}) = \{\mathbf{x}^T \mathbf{z}\}$. Then:

$$\begin{aligned} \|\mathbf{x} - \mathbf{z}\|_2 &= \sqrt{(\mathbf{x} - \mathbf{z})^2} \\ &= \sqrt{\mathbf{x}^T \mathbf{x} + \mathbf{z}^T \mathbf{z} - 2(\mathbf{x}^T \mathbf{z})} \\ &= \sqrt{k(\mathbf{x}, \mathbf{x}) + k(\mathbf{z}, \mathbf{z}) - 2k(\mathbf{x}, \mathbf{z})} \end{aligned}$$

2 Programming Exercises

1. Implement the standard Perceptron algorithm. Train it using the provided training data and test it on the test data (ignore the provided development data for this part). The provided `run.m` script calls the Perceptron code (`perceptron.m`) which returns the Perceptron weight vector \mathbf{w} , the bias term b , and the accuracy on the test data. I have provided some skeleton code in `perceptron.m` and your job is to fill in the required pieces to make it work. See the comments in the code. Basically, you have to write code for testing the Perceptron update condition, the Perceptron update equations, and computing the predictions for the test data. Once you are done with these, you can simply call `run.m` which will print the accuracy on the test data. Report the accuracy on the test data. **(10 marks)**

Now, consider the *averaging extension* of the Perceptron. Note that the standard Perceptron produces a sequence of weight vectors (basically, an updated weight vector everytime the Perceptron makes a mistake). However it only uses the final weight vector for making its predictions. The averaging extension of Perceptron, on the other hand, makes use of all these weight vectors by combining them in the following way: note that in the standard Perceptron, each weight vector is “good” until the next mistake is made. Assume that the Perceptron produced K distinct weight vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$ during its training, and they were good for c_1, c_2, \dots, c_K examples respectively. The averaged Perceptron weight vector is given by $\sum_{k=1}^K c_k \mathbf{w}_k$ and the bias is given by $\sum_{k=1}^K c_k b_k$. In a sense, these updates are like a weighted sum.

Implement the averaged Perceptron keeping the above logic in mind. Use the provided skeleton code `averaged_perceptron.m` (it’s basically the same as the standard Perceptron’s skeleton code :-), and complete it to make it work for the averaged Perceptron case. Now, train the averaged version using the provided training data and test it on the test data (again, ignore the development data). The script `run.m` does all this for you once your averaged Perceptron implementation is complete. Report the accuracy on the test data. **(15 marks)**

2. In this part, we will experiment using other people’s software, so the programming aspect of this part will just be to write some short scripts for converting data formats. The software package we will use is libSVM <http://www.csie.ntu.edu.tw/~cjlin/libsvm/> (which, not surprisingly, implements support vector machines). You can either download and install these yourself, or you can use the version I have installed on CADE (in `~cs5350/bin`, the libSVM executables are `svm-train` and `svm-predict`).

The libSVM file format is one-example-per-line. The first column is the class (for binary this should be -1 or 1; for multiclass this should be a number 0, 1, 2, ...). The remaining columns have format $D : N$, where D is a feature number and N is the feature value. I.e., our matlab format might have a row like:

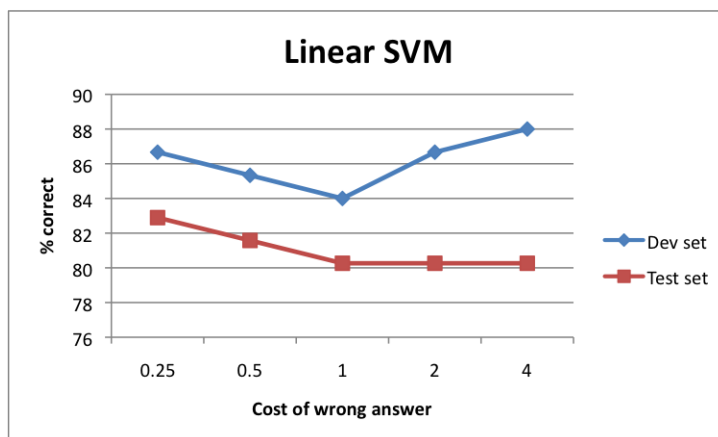
-1 0.2 1.5 -3.2 9.3 0 1.2

The corresponding libSVM line would be something like:

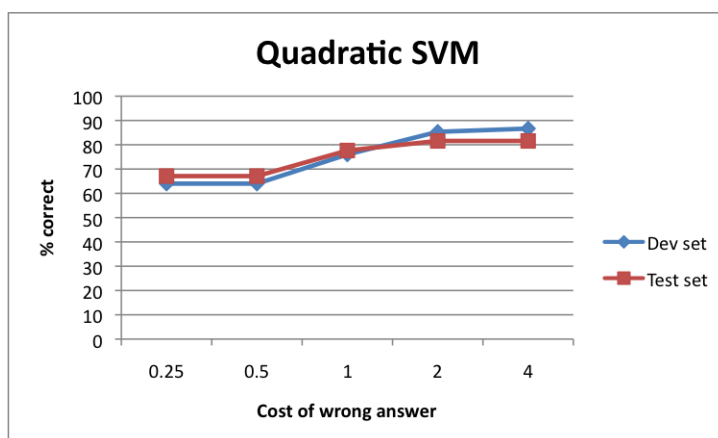
-1 1:0.2 2:1.5 3:-3.2 4:9.3 6:1.2

Note that we could exclude feature 5 because its value was zero. The features must appear sorted (i.e., feature 1 followed by feature 2, and so on). Using the provided datasets, train the following SVMs (using just the training data): a linear SVM, a quadratic SVM, a cubic SVM and an RBF SVM. For each, try different values of C ranging in 0.25, 0.5, 1, 2, 4. For the RBF SVM, try γ values of 0.25, 0.5, 1, 2, 4. Plot error rates on both the development data and test data for the different values of C . How many support vectors are used for each model? Should this increase or decrease with C (why?)? (15 marks)

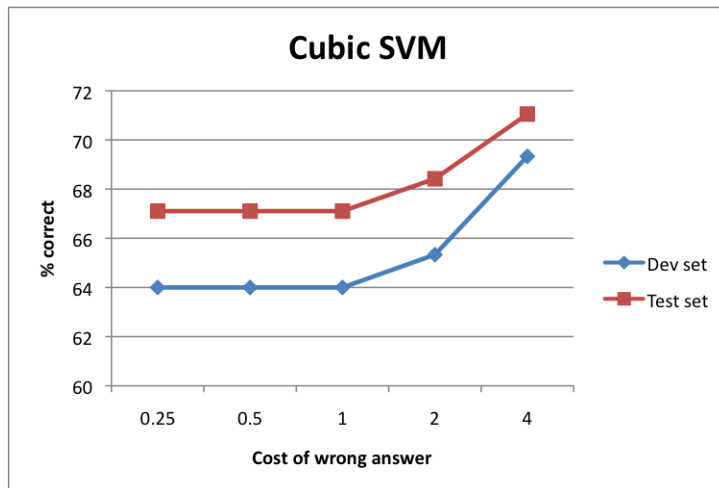
Notice in the following graphs that the number of support vectors steadily decreases for each of them. To answer the last question directly, the reason is because for large C in $C \sum_N \xi_n$, the slack will dominate misclassifications, causing us (all else equal) to have a smaller number of misclassified examples. Unfortunately this also makes our margin smaller; either way, this is the explanation for the phenomenon you will see below.



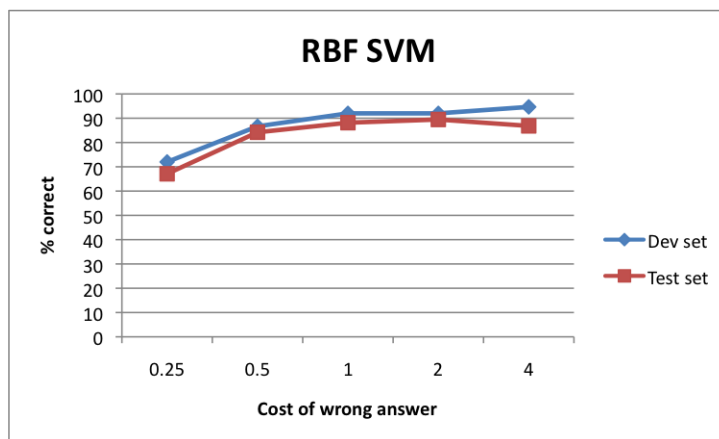
The number of support vectors at each cost point for the linear SVM was, respectively: 85, 76, 73, 66, and 65.



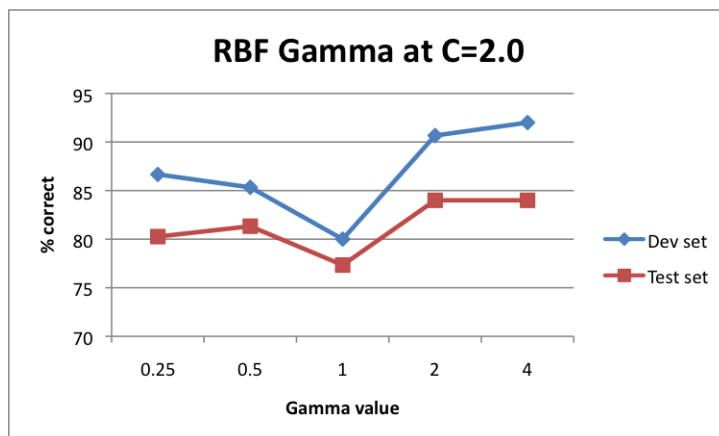
The number of support vectors at each cost point for the quadratic SVM was, respectively: 155, 155, 150, 135, and 120.



The number of support vectors at each cost point for the cubic SVM was, respectively: 152, 152, 152, 152, and 145.



The number of support vectors at each cost point for the RBF SVM was, respectively: 144, 121, 102, 86, and 81.



And, of course, the results of our gamma experiments follow. Rather than look at every single possible experiment, I just ran them on the best-scoring C .