

CS 229, Spring 2016

Problem Set #2: Naive Bayes, SVMs, and Theory

Due Wednesday, May 4 at 11:00 pm on Gradescope.

Notes: (1) These questions require thought, but do not require long answers. Please be as concise as possible. (2) If you have a question about this homework, we encourage you to post your question on our Piazza forum, at <https://piazza.com/stanford/spring2016/cs229>. (3) If you missed the first lecture or are unfamiliar with the collaboration or honor code policy, please read the policy on Handout #1 (available from the course website) before starting work. (4) For problems that require programming, please include in your submission a printout of your code (with comments) and any figures that you are asked to plot.

If you are scanning your document by cellphone, please check the Piazza forum for recommended cellphone scanning apps and best practices.

1. [15 points] Constructing kernels

In class, we saw that by choosing a kernel $K(x, z) = \phi(x)^T \phi(z)$, we can implicitly map data to a high dimensional space, and have the SVM algorithm work in that space. One way to generate kernels is to explicitly define the mapping ϕ to a higher dimensional space, and then work out the corresponding K .

However in this question we are interested in direct construction of kernels. I.e., suppose we have a function $K(x, z)$ that we think gives an appropriate similarity measure for our learning problem, and we are considering plugging K into the SVM as the kernel function. However for $K(x, z)$ to be a valid kernel, it must correspond to an inner product in some higher dimensional space resulting from some feature mapping ϕ . Mercer's theorem tells us that $K(x, z)$ is a (Mercer) kernel if and only if for any finite set $\{x^{(1)}, \dots, x^{(m)}\}$, the matrix K is symmetric and positive semidefinite, where the square matrix $K \in \mathbb{R}^{m \times m}$ is given by $K_{ij} = K(x^{(i)}, x^{(j)})$.

Now here comes the question: Let K_1, K_2 be kernels over $\mathbb{R}^n \times \mathbb{R}^n$, let $a \in \mathbb{R}^+$ be a positive real number, let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be a real-valued function, let $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^d$ be a function mapping from \mathbb{R}^n to \mathbb{R}^d , let K_3 be a kernel over $\mathbb{R}^d \times \mathbb{R}^d$, and let $p(x)$ a polynomial over x with *positive* coefficients.

For each of the functions K below, state whether it is necessarily a kernel. If you think it is, prove it; if you think it isn't, give a counter-example.

- (a) [1 points] $K(x, z) = K_1(x, z) + K_2(x, z)$
- (b) [1 points] $K(x, z) = K_1(x, z) - K_2(x, z)$
- (c) [1 points] $K(x, z) = aK_1(x, z)$
- (d) [1 points] $K(x, z) = -aK_1(x, z)$
- (e) [5 points] $K(x, z) = K_1(x, z)K_2(x, z)$
- (f) [2 points] $K(x, z) = f(x)f(z)$
- (g) [2 points] $K(x, z) = K_3(\phi(x), \phi(z))$
- (h) [2 points] $K(x, z) = p(K_1(x, z))$

[Hint: For part (e), the answer is that the K there *is* indeed a kernel. You still have to prove it, though. (This one may be harder than the rest.) This result may also be useful for another part of the problem.]

2. [10 points] Kernelizing the Perceptron

Let there be a binary classification problem with $y \in \{-1, 1\}$. The perceptron uses hypotheses of the form $h_\theta(x) = g(\theta^T x)$, where $g(z) = \text{sign}(z) = 1$ if $z \geq 0$, -1 otherwise. In this problem we will consider a stochastic gradient descent-like implementation of the perceptron algorithm where each update to the parameters θ is made using only one training example. However, unlike stochastic gradient descent, the perceptron algorithm will only make one pass through the entire training set. The update rule for this version of the perceptron algorithm is given by

$$\theta^{(i+1)} := \begin{cases} \theta^{(i)} + \alpha y^{(i+1)} x^{(i+1)} & \text{if } h_{\theta^{(i)}}(x^{(i+1)}) y^{(i+1)} < 0 \\ \theta^{(i)} & \text{otherwise,} \end{cases}$$

where $\theta^{(i)}$ is the value of the parameters after the algorithm has seen the first i training examples. Prior to seeing any training examples, $\theta^{(0)}$ is initialized to $\vec{0}$.

Let K be a Mercer kernel corresponding to some very high-dimensional feature mapping ϕ . Suppose ϕ is so high-dimensional (say, ∞ -dimensional) that it's infeasible to ever represent $\phi(x)$ explicitly. Describe how you would apply the “kernel trick” to the perceptron to make it work in the high-dimensional feature space ϕ , but without ever explicitly computing $\phi(x)$. [Note: You don't have to worry about the intercept term. If you like, think of ϕ as having the property that $\phi_0(x) = 1$ so that this is taken care of.] Your description should specify

- How you will (implicitly) represent the high-dimensional parameter vector $\theta^{(i)}$, including how the initial value $\theta^{(0)} = \vec{0}$ is represented (note that $\theta^{(i)}$ is now a vector whose dimension is the same as the feature vectors $\phi(x)$);
- How you will efficiently make a prediction on a new input $x^{(i+1)}$. I.e., how you will compute $h_{\theta^{(i)}}(x^{(i+1)}) = g(\theta^{(i)T} \phi(x^{(i+1)}))$, using your representation of $\theta^{(i)}$; and
- How you will modify the update rule given above to perform an update to θ on a new training example $(x^{(i+1)}, y^{(i+1)})$; i.e., using the update rule corresponding to the feature mapping ϕ :

$$\theta^{(i+1)} := \theta^{(i)} + \alpha \mathbf{1}\{\theta^{(i)T} \phi(x^{(i+1)}) y^{(i+1)} < 0\} y^{(i+1)} \phi(x^{(i+1)}).$$

[Hint: our discussion of the representer theorem may be useful.]

3. [30 points] Spam classification

In this problem, we will use the naive Bayes algorithm and an SVM to build a spam classifier.

In recent years, spam on electronic newsgroups has been an increasing problem. Here, we'll build a classifier to distinguish between “real” newsgroup messages, and spam messages. For this experiment, we obtained a set of spam emails, and a set of genuine newsgroup messages.¹ Using only the subject line and body of each message, we'll learn to distinguish between the spam and non-spam.

¹Thanks to Christian Shelton for providing the spam email. The non-spam messages are from the 20 newsgroups data at <http://www-2.cs.cmu.edu/afs/cs.cmu.edu/project/theo-20/www/data/news20.html>.

All the files for the problem are in http://cs229.stanford.edu/materials/spam_data.tgz.

Note: Please do not circulate this data outside this class. In order to get the text emails into a form usable by naive Bayes, we've already done some preprocessing on the messages. You can look at two sample spam emails in the files `spam_sample_original*`, and their preprocessed forms in the files `spam_sample_preprocessed*`. The first line in the preprocessed format is just the label and is not part of the message. The preprocessing ensures that only the message body and subject remain in the dataset; email addresses (EMAILADDR), web addresses (HTTPADDR), currency (DOLLAR) and numbers (NUMBER) were also replaced by the special tokens to allow them to be considered properly in the classification process. (In this problem, we'll going to call the features "tokens" rather than "words," since some of the features will correspond to special values like EMAILADDR. You don't have to worry about the distinction.) The files `news_sample_original` and `news_sample_preprocessed` also give an example of a non-spam mail.

The work to extract feature vectors out of the documents has also been done for you, so you can just load in the design matrices (called document-word matrices in text classification) containing all the data. In a document-word matrix, the i^{th} row represents the i^{th} document/email, and the j^{th} column represents the j^{th} distinct token. Thus, the (i, j) -entry of this matrix represents the number of occurrences of the j^{th} token in the i^{th} document.

For this problem, we've chosen as our set of tokens considered (that is, as our vocabulary) only the medium frequency tokens. The intuition is that tokens that occur too often or too rarely do not have much classification value. (Examples tokens that occur very often are words like "the," "and," and "of," which occur in so many emails and are sufficiently content-free that they aren't worth modeling.) Also, words were stemmed using a standard stemming algorithm; basically, this means that "price," "prices" and "priced" have all been replaced with "price," so that they can be treated as the same word. For a list of the tokens used, see the file `TOKENS_LIST`.

Since the document-word matrix is extremely sparse (has lots of zero entries), we have stored it in our own efficient format to save space. You don't have to worry about this format.² The file `readMatrix.m` provides the `readMatrix` function that reads in the document-word matrix and the correct class labels for the various documents. Code in `nb_train.m` and `nb_test.m` shows how `readMatrix` should be called. The documentation at the top of these two files will tell you all you need to know about the setup.

- (a) [11 points] Implement a naive Bayes classifier for spam classification, using the multinomial event model and Laplace smoothing.

You should use the code outline provided in `nb_train.m` to train your parameters, and then use these parameters to classify the test set data by filling in the code in `nb_test.m`. You may assume that any parameters computed in `nb_train.m` are in memory when `nb_test.m` is executed, and do not need to be recomputed (i.e., that `nb_test.m` is executed immediately after `nb_train.m`)³.

Train your parameters using the document-word matrix in `MATRIX.TRAIN`, and then report the test set error on `MATRIX.TEST`.

Remark. If you implement naive Bayes the straightforward way, you'll find that the computed $p(x|y) = \prod_i p(x_i|y)$ often equals zero. This is because $p(x|y)$, which is

²Unless you're not using Matlab/Octave, in which case feel free to ask us about it. We have provided Julia code to read the file in `MatrixReading.jl`.

³Matlab note: If a .m file doesn't begin with a function declaration, the file is a script. Variables in a script are put into the global namespace, unlike with functions.

the product of many numbers less than one, is a very small number. The standard computer representation of real numbers cannot handle numbers that are too small, and instead rounds them off to zero. (This is called “underflow.”) You’ll have to find a way to compute naive Bayes’ predicted class labels without explicitly representing very small numbers such as $p(x|y)$. [Hint: Think about using logarithms.]

- (b) [3 points] Intuitively, some tokens may be particularly indicative of an email being in a particular class. We can try to get an informal sense of how indicative token i is for the SPAM class by looking at:

$$\log \frac{p(x_j = i|y = 1)}{p(x_j = i|y = 0)} = \log \left(\frac{P(\text{token } i|\text{email is SPAM})}{P(\text{token } i|\text{email is NOTSPAM})} \right).$$

Using the parameters fit in part (a), find the 5 tokens that are most indicative of the SPAM class (i.e., have the highest positive value on the measure above). The numbered list of tokens in the file `TOKENS_LIST` should be useful for identifying the words/tokens.

- (c) [3 points] Repeat part (a), but with training sets of size ranging from 50, 100, 200, ..., up to 1400, by using the files `MATRIX.TRAIN.*`. Plot the test error each time (use `MATRIX.TEST` as the test data) to obtain a learning curve (test set error vs. training set size). You may need to change the call to `readMatrix` in `nb_train.m` to read the correct file each time. Which training-set size gives the best test set error?
- (d) [11 points] Train an SVM on this dataset using stochastic gradient descent and the radial basis function (also known as the Gaussian) kernel, which sets

$$K(x, z) = \exp \left(-\frac{1}{2\tau^2} \|x - z\|_2^2 \right).$$

In this case, recall that (as proved in class) the objective with kernel matrix $K = [K^{(1)} \dots K^{(m)}] \in \mathbb{R}^{m \times m}$ is given by

$$J(\alpha) = \frac{1}{m} \sum_{i=1}^m \left[1 - y^{(i)} K^{(i)T} \alpha \right]_+ + \frac{\lambda}{2} \alpha^T K \alpha$$

where $[t]_+ = \max\{t, 0\}$ is the positive part function. In this case, the gradient (actually, this is known as a *subgradient*) of the individual loss terms is

$$\nabla_{\alpha} \left[1 - y^{(i)} K^{(i)} \alpha \right]_+ = \begin{cases} -y^{(i)} K^{(i)} & \text{if } y^{(i)} K^{(i)T} \alpha < 1 \\ 0 & \text{otherwise.} \end{cases}$$

In your SVM training, you should perform stochastic gradient descent, where in each iteration you choose an index $i \in \{1, \dots, m\}$ uniformly at random, for a total of $40 \cdot m$ steps, where m is the training set size, and your kernel should use $\tau = 8$ and regularization multiplier $\lambda = \frac{1}{64m}$. For this part of the problem, you should also replace each training or test point $x^{(i)}$ with a zero-one vector $z^{(i)}$, where $z_j^{(i)} = 1$ if $x_j^{(i)} > 0$ and $z_j^{(i)} = 0$ if $x_j^{(i)} = 0$. Initialize your SGD procedure at $\alpha = 0$.

The output of your training code, which you should implement in `svm_test.m`, should be the α vector that is the *average* of all the α vectors that your iteration updates. At iteration t of stochastic gradient descent you should use stepsize $1/\sqrt{t}$.

Similar to part (c), train an SVM with training set sizes 50, 100, 200, ..., 1400, by using the file `MATRIX.TRAIN.50` and so on. Plot the test error each time, using `MATRIX.TEST` as the test data.

(A few hints for more efficient Matlab code: you should try to vectorize creation of the Kernel matrix, and you should call the method `full` to make the matrix non-sparse, which will make the method faster. In addition, the training data uses labels in $\{0, 1\}$, so you should change the output of the `readMatrix` method to have labels $y \in \{-1, 1\}$.)

- (e) [2 points] How do naive Bayes and Support Vector Machines compare (in terms of generalization error) as a function of the training set size?

4. [20 points] Properties of VC dimension

In this problem, we investigate a few properties of the Vapnik-Chervonenkis dimension, mostly relating to how $VC(H)$ increases as the set H increases. For each part of this problem, you should state whether the given statement is true, and justify your answer with either a formal proof or a counter-example.

- (a) Let two hypothesis classes H_1 and H_2 satisfy $H_1 \subseteq H_2$. Prove or disprove: $VC(H_1) \leq VC(H_2)$.
- (b) Let $H_1 = H_2 \cup \{h_1, \dots, h_k\}$. (I.e., H_1 is the union of H_2 and some set of k additional hypotheses.) Prove or disprove: $VC(H_1) \leq VC(H_2) + k$. [Hint: You might want to start by considering the case of $k = 1$.]
- (c) Let $H_1 = H_2 \cup H_3$. Prove or disprove: $VC(H_1) \leq VC(H_2) + VC(H_3)$.

5. [20 points] Training and testing on different distributions

In the discussion in class about learning theory, a key assumption was that we trained and tested our learning algorithms on the same distribution \mathcal{D} . In this problem, we'll investigate one special case of training and testing on different distributions. Specifically, we will consider what happens when the training labels are *noisy*, but the test labels are not.

Consider a binary classification problem with labels $y \in \{0, 1\}$, and let \mathcal{D} be a distribution over (x, y) , that we'll think of as the original, "clean" or "uncorrupted" distribution. Define \mathcal{D}_τ to be a "corrupted" distribution over (x, y) which is the same as \mathcal{D} , except that the labels y have some probability $0 \leq \tau < 0.5$ of being flipped. Thus, to sample from \mathcal{D}_τ , we would first sample (x, y) from \mathcal{D} , and then with probability τ (independently of the observed x and y) replace y with $1 - y$. Note that $\mathcal{D}_0 = \mathcal{D}$.

The distribution \mathcal{D}_τ models a setting in which an unreliable human (or other source) is labeling your training data for you, and on each example he/she has a probability τ of mislabeling it. Even though our training data is corrupted, we are still interested in evaluating our hypotheses with respect to the original, uncorrupted distribution \mathcal{D} .

We define the generalization error *with respect to* \mathcal{D}_τ to be

$$\varepsilon_\tau(h) = P_{(x,y) \sim \mathcal{D}_\tau}[h(x) \neq y].$$

Note that $\varepsilon_0(h)$ is the generalization error with respect to the "clean" distribution; it is with respect to ε_0 that we wish to evaluate our hypotheses.

- (a) For any hypothesis h , the quantity $\varepsilon_0(h)$ can be calculated as a function of $\varepsilon_\tau(h)$ and τ . Write down a formula for $\varepsilon_0(h)$ in terms of $\varepsilon_\tau(h)$ and τ , and justify your answer.
- (b) Let $|H|$ be finite, and suppose our training set $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$ is obtained by drawing m examples IID from the corrupted distribution \mathcal{D}_τ . Suppose we pick $h \in H$ using empirical risk minimization: $\hat{h} = \arg \min_{h \in H} \hat{\varepsilon}_S(h)$. Also, let $h^* = \arg \min_{h \in H} \varepsilon_0(h)$.

Let any $\delta, \gamma > 0$ be given. Prove that for

$$\varepsilon_0(\hat{h}) \leq \varepsilon_0(h^*) + 2\gamma$$

to hold with probability $1 - \delta$, it suffices that

$$m \geq \frac{1}{2(1-2\tau)^2\gamma^2} \log \frac{2|H|}{\delta}.$$

Remark. This result suggests that, roughly, m examples that have been corrupted at noise level τ are worth about as much as $(1-2\tau)^2 m$ uncorrupted training examples. This is a useful rule-of-thumb to know if you ever need to decide whether/how much to pay for a more reliable source of training data. (If you've taken a class in information theory, you may also have heard that $(1-\mathcal{H}(\tau))m$ is a good estimate of the information in the m corrupted examples, where $\mathcal{H}(\tau) = -(\tau \log_2 \tau + (1-\tau) \log_2 (1-\tau))$ is the “binary entropy” function. And indeed, the functions $(1-2\tau)^2$ and $1-\mathcal{H}(\tau)$ are quite close to each other.)

- (c) Comment **briefly** on what happens as τ approaches 0.5.

6. [19 points] Boosting and high energy physics

In class, we discussed boosting algorithms and decision stumps. In this problem, we explore applications of these ideas to detect particle emissions in a high-energy particle accelerator. In high energy physics, such as at the Large Hadron Collider (LHC), one accelerates small particles to relativistic speeds and smashes them into one another, tracking the emitted particles. The goal in these problems is to detect the emission of certain interesting particles based on other observed particles and energies.⁴ In this problem, we explore the application of boosting to a high energy physics problem, where we use decision stumps applied to 18 low- and high-level physics-based features. All data for the problem is available at http://cs229.stanford.edu/materials/boost_data.tgz.

For the first part of the problem, we explore how decision stumps based on thresholding can provide a weak-learning guarantee. In particular, we show that for real-valued attributes x , there is an edge $\gamma > 0$ that decision stumps guarantee. Recall that thresholding-based decision stumps are functions indexed by a threshold s and sign $+/-$, such that

$$\phi_{s,+}(x) = \begin{cases} 1 & \text{if } x \geq s \\ -1 & \text{if } x < s \end{cases} \quad \text{and} \quad \phi_{s,-}(x) = \begin{cases} -1 & \text{if } x \geq s \\ 1 & \text{if } x < s. \end{cases}$$

That is, $\phi_{s,+}(x) = -\phi_{s,-}(x)$. We assume for simplicity in the theoretical parts of this exercise that our input attribute vectors $x \in \mathbb{R}$, that is, they are one-dimensional. Now, we would like guarantee that there is some $\gamma > 0$ and a threshold s such that, for *any*

⁴For more, see the following paper: Baldi, Sadowski, Whiteson. Searching for Exotic Particles in High-Energy Physics with Deep Learning. *Nature Communications* 5, Article 4308. <http://arxiv.org/abs/1402.4735>.

distribution p on the training set $\{x^{(i)}, y^{(i)}\}_{i=1}^m$ (where $y^{(i)} \in \{-1, +1\}$ and $x^{(i)} \in \mathbb{R}$, and we recall that p is a distribution on the training set if $\sum_{i=1}^m p_i = 1$ and $p_i \geq 0$ for each i) we have

$$\sum_{i=1}^m p_i \mathbf{1}\{y^{(i)} \neq \phi_{s,+}(x^{(i)})\} \leq \frac{1}{2} - \gamma \quad \text{or} \quad \sum_{i=1}^m p_i \mathbf{1}\{y^{(i)} \neq \phi_{s,-}(x^{(i)})\} \leq \frac{1}{2} - \gamma.$$

For simplicity, we assume that all of the $x^{(i)}$ are *distinct*, so that none of them are equal. We also assume (without loss of generality, but this makes the problem notationally simpler) that

$$x^{(1)} > x^{(2)} > \dots > x^{(m)}.$$

- (a) [3 points] Show that for each threshold s , there is some $m_0(s) \in \{0, 1, \dots, m\}$ such that

$$\sum_{i=1}^m p_i \mathbf{1}\{\phi_{s,+}(x^{(i)}) \neq y^{(i)}\} = \frac{1}{2} - \frac{1}{2} \left(\sum_{i=1}^{m_0(s)} y^{(i)} p_i - \sum_{i=m_0(s)+1}^m y^{(i)} p_i \right)$$

and

$$\sum_{i=1}^m p_i \mathbf{1}\{\phi_{s,-}(x^{(i)}) \neq y^{(i)}\} = \frac{1}{2} - \frac{1}{2} \left(\sum_{i=m_0(s)+1}^m y^{(i)} p_i - \sum_{i=1}^{m_0(s)} y^{(i)} p_i \right)$$

Treat sums over empty sets of indices as zero, so that $\sum_{i=1}^0 a_i = 0$ for any a_i , and similarly $\sum_{i=m+1}^m a_i = 0$.

- (b) [3 points] Define, for each $m_0 \in \{0, 1, \dots, m\}$,

$$f(m_0) = \sum_{i=1}^{m_0} y^{(i)} p_i - \sum_{i=m_0+1}^m y^{(i)} p_i.$$

Show that there exists *some* $\gamma > 0$, which may depend on the training set size m (but should not depend on p), such that for any set of probabilities p on the training set, where $p_i \geq 0$ and $\sum_{i=1}^m p_i = 1$, we can find m_0 with

$$|f(m_0)| \geq 2\gamma.$$

What is your γ ?

(*Hint*: Consider the difference $f(m_0) - f(m_0 + 1)$.)

- (c) [2 points] Based on your answer to part (6b), what edge can thresholded decision stumps guarantee on any training set $\{x^{(i)}, y^{(i)}\}_{i=1}^m$, where the raw attributes $x^{(i)} \in \mathbb{R}$ are all distinct? Recall that the edge of a weak classifier $\phi : \mathbb{R} \rightarrow \{-1, 1\}$ is the constant $\gamma \in [0, \frac{1}{2}]$ such that

$$\sum_{i=1}^m p_i \mathbf{1}\{\phi(x^{(i)}) \neq y^{(i)}\} \leq \frac{1}{2} - \gamma.$$

Can you give an upper bound on the number of thresholded decision stumps required to achieve zero error on a given training set?

- (d) [11 points] Now you will implement boosting on data developed from a physics-based simulation of a high-energy particle accelerator. We provide two datasets, `boosting-train.csv` and `boosting-test.csv`, which consist of training data and test data for a binary classification problem on which you will apply boosting techniques. (For those not using `Matlab`, the files are comma-separated files, the first column of which consists of binary ± 1 -labels $y^{(i)}$, the remaining 18 columns are the raw attributes.) The file `load_data.m`, which we provide, loads the datasets into memory, storing training data and labels in appropriate vectors and matrices, and then performs boosting using *your* implemented code, and plots the results.
- i. [5 points] Implement a method that finds the optimal thresholded decision stump for a training set $\{x^{(i)}, y^{(i)}\}_{i=1}^m$ and distribution $p \in \mathbb{R}_+^m$ on the training set. In particular, fill out the code in the method `find_best_threshold.m`. Include your code in your solution.
 - ii. [2 points] Implement boosted decision stumps by filling out the code in the method `stump_booster.m`. Your code should implement the weight updating at each iteration $t = 1, 2, \dots$ to find the optimal value θ_t given the feature index and threshold. Include your code in your solution.
 - iii. [2 points] Implement *random* boosting, where at each step the choice of decision stump is made completely randomly. In particular, at iteration t random boosting chooses a random index $j \in \{1, 2, \dots, m\}$, then chooses a random threshold s from among the data values $\{x_j^{(i)}\}_{i=1}^m$, and then chooses the t th weight θ_t optimally for this (random) classifier $\phi_{s,+}(x) = \text{sign}(x_j - s)$. Implement this by filling out the code in `random_booster.m`.
 - iv. [2 points] Run the method `load_data.m` with your implemented boosting methods. Include the plots this method displays, which show the training and test error for boosting at each iteration $t = 1, 2, \dots$. Which method is better?

[A few notes: we do not expect boosting to get classification accuracy better than approximately 80% for this problem.]