

10-701 Machine Learning: Assignment 2

Due on March 11, 2014 at 11:59

Barnabas Poczos, Aarti Singh

Instructions: Failure to follow these directions may result in loss of points.

- Your solutions for this assignment need to be in a pdf format and should be submitted to the blackboard and the webpage <http://barnabas-cmu-10701.appspot.com> for peer-reviewing.
- For the programming question, your code should be well-documented, so a TA can understand what is happening.
- We are NOT going to use Autolab in this assignment.
- DO NOT include any identification (your name, andrew id, or email) in both the content and filename of your submission.

Q1) Support Vector Machines (Pulkit)

SVM [20 points]

In class we discussed support vector classification, but did not cover the regression problem. In this question you have to derive the dual form of SV regression (SVR). Your training data is $(x_1, y_1), \dots, (x_n, y_n)$, where $x_i \in \mathbb{R}^m$, $y_i \in \mathbb{R}$.

Since the hinge loss that we used in class is only designed for classification we cannot use that for regression. A frequently used loss function for regression is the epsilon sensitive loss:

$$L_\epsilon(x, y, f) = |y - f(x)|_\epsilon = \max(0, |y - f(x)| - \epsilon)$$

Here x is the input, y is the output, and f is the function used for predicting the label.

Using this notation, the SVR cost function is defined as

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n L_\epsilon(x_i, y_i, f),$$

where $f(x) = w^T x$, and $C, \epsilon > 0$ are parameters.

1. Introduce appropriate slack variables, and rewrite this problem as a quadratic problem (i.e. quadratic objective with linear constraints). This form is called the Primal form of support vector regression. (3 points)
2. Write down the Lagrangian function for the above primal form. (2 points)
3. Using the Karush Kunh Tucker conditions derive the dual form. (5 points)
4. Can we use quadratic optimization solvers to solve the dual problem? (1 point)

5. How would you define support vectors in this problem? (2 points)
6. Write down the equation that can be used for predicting the label of an unseen sample X . (2 points)
7. Is it possible to kernelize this algorithm? (1 points)
8. What happens if we change ϵ ? (2 point)
9. What happens if we change C ? (2 point)

Q2) Density Estimation (Prashant)

Kernel Density Estimation (10 points)

Let us explore kernel density estimation with a boxcar kernel. Given data, X_1, \dots, X_n , the empirical pdf is defined as follows

$$\hat{p}_h(x) = \frac{1}{n} \frac{\sum_{j=1}^n I(|x - X_i| \leq h)}{2h} \quad (1)$$

Here, $|y|$ indicates absolute value and $I(y)$ is an indicator function that is 1 if the boolean argument evaluates to true and 0 otherwise.

1. Draw $\hat{p}_h(x)$ for all $x \in [-7, 7]$ where our data consists of the following observations $\{X_i\} = \{-7, -4, -3, -2, 1, 5\}$ and $h = 3$.
2. Let $F(x) = P(X \leq x)$ represent the true cdf over a random variable X . Show that

$$E(\hat{p}_h(x)) = \frac{F(x+h) - F(x-h)}{2h} \quad (2)$$

3. Determine if our kernel density estimator will be an unbiased estimator in the following scenarios
 1. $pdf(x)$ is uniform between 0 and 1. $\hat{p}_h(\frac{1}{2})$ is estimated using $h = \frac{1}{3}$
 2. $pdf(x) = 2x$ for $0 \leq x \leq 1$. $\hat{p}_h(\frac{1}{4})$ is estimated using $h = \frac{1}{5}$
 3. $pdf(x) = \frac{3}{2}x^2$ for $-1 \leq x \leq 1$. $\hat{p}_h(0)$ is estimated using $h = \frac{1}{5}$

Histograms as MLE(10 points)

Recall that histograms can be used for non-parametric density estimation. Let X_1, \dots, X_n be n data points drawn from a random variable X which takes values between 0 and 1. We further divide the range $[0, 1]$ into N different, equal sized, bins, B_1, \dots, B_N each with width h .

Define the empirical histogram p.d.f. is as follows

$$\hat{p}_h = \sum_{j=1}^N \frac{\hat{\theta}_j}{h} I(x \in B_j) \quad (3)$$

$$\hat{\theta}_j = \frac{1}{n} \sum_{i=1}^n I(X_i \in B_j) \quad (4)$$

Now assume that we model the true pdf of X as a piecewise constant function where the constants are the θ_j which is constant over the bins, B_j for j in $1, \dots, N$. Show that the MLE parameters under this model are identical to the histogram based approach for density estimation.

Q3) K-Nearest Neighbors (Jit)

Spam Detection [10 points]

One of the earliest and most successful applications of machine learning has been in developing spam detection algorithms. For this problem, you are going to implement K-Nearest Neighbors algorithm on a dataset containing thousands of instances of spam and non-spam emails.

For this problem, you are going to submit your files to the predictive analytics site Kaggle. We have created a private Kaggle competition in which only users with a *cmu.edu* email can access. To access the competition, please go to the following link: <https://kaggle.com/join/tyranitar>

For this assignment, you will need to write code to complete the following functions:

- `distance(email1, email2)`: Calculate some notion of distance or similarity between two e-mails. As an aside, keep in mind that higher similarities should correspond to smaller distances. (1 points)
- `KNNClassify(testData, trainingData, trainingLabels)`: Classifies each email in the testing dataset using the kNN algorithm. (3 points)

There are many valid distance metrics you can implement, and we encourage you to explore how different distance metrics lead to different results. Furthermore, you are expected to try out different values of k to see which neighborhood leads to the best accuracy on the testing data. Please plot the misclassification rate of at least two different distance metrics and four different values of k . Using your results, comment on the behavior of the kNN algorithm with varying distance metric, and varying values of k . (6 points)

You will be submitting a .csv file to the Kaggle submission system (The format of the CSV file is on the Kaggle website. It needs to contain a header: EmaiID, Prediction, and then for each row, it contains a ID number, and a label 0,1 similar to the training data labels file), and upload your code along it. We will be reviewing your code, so please document your code appropriately. For Python, you are allowed to use the NumPy package, but you cannot use machine learning packages that have implementations of kNN. For this assignment, we ask that you use either Python, Julia, or Matlab, so the TAs will be able to review your code without much difficulty. Please upload your plots and your analysis as part of your solutions that you will submit to the peer review site and Blackboard.

Hidely-Ho, K-Nearest Neighborinos [10 points]

Recall, the KNN algorithm on binary-labeled data works as follows

Data: Documents $\mathbf{X} = [x_1, x_2, \dots, x_n]$, labels $\mathbf{Y} = [y_1, y_2, \dots, y_n]$; $Y_i \in \{-1, 1\}$, integer k , and query x
Result: label y
for $i = 1$ **to** m **do**
 Compute distance $d(x_i, x)$;
end
Compute set I containing indices of the k smallest distances $d(x_i, x)$;
Return majority label of $\{y_i$ where $i \in I\}$;

Algorithm 1: Algorithm: K-Nearest Neighbors Classification

While this algorithm is effective, computing the distances between all training points and the input query x

make take a long time depending on the size of the training points, the dimensionality of the data, and the distance function.

To simplify this, let's precompute $\mu_+ := \frac{1}{n_+} \sum_{y_i=1} x_i$ and $\mu_- := \frac{1}{n_-} \sum_{y_i=-1} x_i$, where n_+ is the number of positive examples, and n_- is the number of negative examples.

Using the fact that for vectors $\mathbf{u}, \mathbf{v} \in \mathbf{R}^n$

$$\|\mathbf{u} - \mathbf{v}\|^2 = \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 - 2 \langle \mathbf{u}, \mathbf{v} \rangle$$

Derive a new classifier $f(x)$ using only μ_+, μ_- , and x that compares each point to the average distance of each class, and express it in the form $f(x) = \text{sgn}(\sum_{i=1}^n \alpha_i \langle x_i, x \rangle + b)$ for some α_i and b , such that $f(x) = 1$ when the predicted label of the new point x is 1, and $f(x) = -1$ when the predicted label of the new point x is -1.

Q4) Kernel Regression (Pengtao)

Kernel Regression (10 points)

Consider local linear regression where the predicted output value of x is $\hat{f}(x) = \hat{\alpha} + \hat{\beta}x$, where

$$\hat{\alpha}, \hat{\beta} = \underset{\alpha, \beta}{\text{argmin}} \sum_{i=1}^n w_i(x) (y_i - \alpha - \beta x_i)^2 \quad (5)$$

where $w_i(x) = K(\frac{x-x_i}{h}) / \sum_{i=1}^n K(\frac{x-x_i}{h})$.

1. Show that the objective function can be re-written as

$$(\mathbf{y} - \mathbf{B}\mathbf{a})^\top \Omega(x) (\mathbf{y} - \mathbf{B}\mathbf{a})$$

where $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_n]^\top$, $\mathbf{a} = [\alpha \ \beta]^\top$, $\mathbf{B} = [1 \ x_1; 1 \ x_2; \dots; 1 \ x_n]$ and $\Omega(x)$ is a diagonal matrix with diagonal $[w_1(x) \ w_2(x) \ \dots \ w_n(x)]$.

2. Show that $\hat{f}(x)$ is a linear combination of $\{y_i\}_{i=1}^n$, namely $\hat{f}(x)$ can be written as $\hat{f}(x) = \sum_{i=1}^n \ell_i(x) y_i = \boldsymbol{\ell}(x)^\top \mathbf{y}$, where $\ell_i(x)$ is some quantity defined over x and $\boldsymbol{\ell}(x) = [\ell_1(x) \ \ell_2(x) \ \dots \ \ell_n(x)]^\top$.

Kernelized Ridge Regression (10 points)

The nonparametric kernel regression does a local fit around the test point. Lets now investigate the use of kernels for regression in another way. Similar to the kernel trick for SVMs, we can apply the kernel trick in regression as follows. (However, note that this produces a global fit to the data.)

Consider the ridge regression problem where we have a set of data points $\{\mathbf{x}_i\}_{i=1}^N$ and corresponding response values $\{y_i\}_{i=1}^N$. To achieve better performance, we use a feature mapping function ϕ to map the original d -dimensional feature vector \mathbf{x} to a new D -dimensional feature vector $\hat{\mathbf{x}} = \phi(\mathbf{x})$, where $D \gg d$. Let $\Phi \in \mathbf{R}^{N \times D}$ denote the design matrix, whose i th row contains the new feature vector $\hat{\mathbf{x}}_i^\top$ of the i th data point. Let \mathbf{y} denote the response value vector whose i th component is the response value y_i of the i th point.

For ridge regression, the objective function is $J(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \Phi \boldsymbol{\beta}\|^2 + \frac{\lambda}{2} \|\boldsymbol{\beta}\|^2$, where λ is the tradeoff parameter. Let $\boldsymbol{\beta}^*$ denote the solution to the ridge regression. In the following steps, we are going to derive the

kernel ridge regression.

(a) First, show that β^* is in the space spanned by rows in Φ , i.e., β^* can be written in the form $\beta^* = \Phi^T \alpha^*$, where $\alpha^* \in \mathbb{R}^{N \times 1}$. (Hint: β can be decomposed into $\beta = \beta_{\parallel} + \beta_{\perp}$, where β_{\perp} is orthogonal to the span of rows in Φ and β_{\parallel} lies in the span of rows in Φ .)

(b) In (a), we have proved that $\beta^* = \Phi^T \alpha^*$. In this step, show that $\alpha^* = (\Phi \Phi^T + \lambda \mathbf{I})^{-1} \mathbf{y}$.

(c) In practice, it is very hard to design the feature mapping function ϕ . Even if we can design it, computing the inner product $\phi(\mathbf{x})^T \phi(\mathbf{x}')$ between two points can be costly. Instead, we can use a kernel function $k(\mathbf{x}, \mathbf{x}')$ to implicitly compute the inner product of high dimensional features, i.e., $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$. Basically, kernelization means using the kernel function to replace inner products. In this step, given the kernel function $k(\mathbf{x}, \mathbf{x}')$, where x and x' are the original feature vectors of dimension d , try to kernelize the ridge regression. You need to consider both training and testing. (Hint: In training phase, you need to replace all inner products in α^* with kernel function. In testing phase, given a new test point \mathbf{x} , you need to compute $\phi(\mathbf{x})^T \beta^*$. Replace all inner products in $\phi(\mathbf{x})^T \beta^*$ with kernel function.)

Q5) Model Selection (Dani)

In this homework, you will perform model selection on a sentiment analysis dataset of music reviews [1]. The dataset consists of reviews from Amazon.com for musics. The ratings have been converted to a binary label, indicating a negative review or a positive review. We will use lasso logistic regression [3] for this problem. The lasso logistic regression objective function to minimize during training is:

$$\mathcal{L}(\beta) = \log(1 + \exp(-y\beta^T \mathbf{x})) + \lambda \|\beta\|_1$$

In lasso logistic regression, we penalize the loss function by an L_1 norm of the feature coefficients. As you have learned in class, penalization with an L_1 norm tends to produce solutions where some coefficients are exactly 0. This makes it attractive for high-dimensional data such as text, because in most cases most words can typically be ignored. Furthermore, since we are often left with only a few nonzero coefficients, the lasso solution is often easy to interpret. The goal of model selection here is to choose λ since each setting of λ implies a different model size (number of non-zero coefficients).

You do not need to implement lasso logistic regression. You can download an implementation from <https://github.com/redpony/creg>, and the dataset can be found here: <http://www.cs.cmu.edu/~dyogata/10701/>. There are three feature files and three response (label) files (all response files end with `.res`). They are already in the format required by the implementation you will use. The files are:

- Training data: `music.train` and `music.train.res`
- Development data: `music.dev` and `music.dev.res`
- Test data: `music.test` and `music.test.res`

Important note: the code outputs **accuracy**, whereas you need to plot **classification error** here. You can simply transform accuracy to error by using $1 - \text{accuracy}$.

Error on development (validation) data [8 points]

In the first part of the problem, we will use error on a development data to choose λ . Run the model with $\lambda = \{10^{-8}, 10^{-7}, 10^{-6}, \dots, 10^{-1}, 1, 10, 100\}$.

1. Plot the error on training data and development data as a function of $\log \lambda$.
2. Plot the model size (number of nonzero coefficients) on development data as a function of $\log \lambda$.
3. Choose λ that gives the lowest error on development data. Run it on the test data and report the test error.

Briefly discuss all the results.

Model Complexity and Bias-Variance Tradeoff [3 points]

Give a high-level explanation on the relation between λ and the bias and variance of parameter estimates $\hat{\beta}$. Does larger λ correspond to higher or lower bias? What about the variance? Does larger λ lead to a more complex or a less complex model?

Resolving a tie [2 points]

If there are more than one λ that minimizes the error on the development data, which one will you pick? Explain your choice.

Random search [5 points]

An alternative way to search λ is by randomly sampling its value from an interval.

1. Sample eleven random values \log uniformly from an interval $[10^{-8}, 100]$ for λ and train a lasso logistic regression model. Plot the error on development data as a function of $\log \lambda$.
2. Choose λ that gives the lowest error on development data. Run it on the test data and report the test error.

Random vs. grid search [2 points]

Which one do you think is a better method for searching values to try for λ ? Why?

References

- [1] John Blitzer, Mark Dredze, and Fernando Pereira. Biographies, Bollywood, Boom-boxes and Blenders: Domain Adaptation for Sentiment Classification. In Proc. of ACL, 1997.
- [2] James Bergstra and Yoshua Bengio. Random Search for Hyper-Parameter Optimization, Journal of Machine Learning Research, 13:282:305, 2012.
- [3] Robert Tibshirani. Regression shrinkage and selection via the lasso, Journal of Royal Statistical Society B, 58(1):267:288, 1996.