
EECS 545 – Machine Learning - Homework #2

David Ke Hong

Due: 11:00pm 02/08/2016

Homework Policy: Working in groups is fine, but each member must submit their own writeup. Please write the members of your group on your solutions. There is no strict limit to the size of the group but we may find it a bit suspicious if there are more than 4 to a team. Questions labelled with **(Challenge)** are not strictly required, but you'll get some participation credit if you have something interesting to add, even if it's only a partial answer. **For coding problems, please append your code to your submission and report your results (values, plots, etc.) in your written solution.** You will lose points if you only include them in your code submissions. Homework will be submitted via Gradescope (<https://gradescope.com/>).

1) **Linear Regression (20 pts).** In this problem, you will implement linear regression (using polynomial features) to reproduce curves similar to those you have seen in lecture, exploring how various choices of model-complexity parameters affect training and test error.

- (a) The training and testing data is provided to you, named **train_graphs_f16_autopilot_cruise.csv** and **test_graphs_f16_autopilot_cruise.csv**. Every row in these csv files correspond to a single datapoint \mathbf{x} , \mathbf{t} . Columns of this graphs are named id, rolling speed, elevation speed, elevation jerk, elevation, roll, elevation acceleration, controller input. Your task is to predict controller input from rolling speed, elevation speed, elevation jerk, elevation, roll, and elevation acceleration. Note that id is a primary key, and is unique to each datapoint, so you'd want to remove it from the feature set. Ponder why this is necessary, and what would happen otherwise. Also, you might need to append an additional feature that is constant for all data points, for example a feature with a value of 1, in order to model the intercept term.
- (i) Fit the data by applying the psuedo-inverse approach of linear regression using x_j^i , $i = 1, \dots, M$ as features (try $M = 1, 2, \dots, 6$), where x_j represents the j^{th} component of vector \mathbf{x} . In other words, you will need to raise every element in vector \mathbf{x} to the power of i , for $i = 1, 2, \dots, M$ to get the set of features. For example, if $\mathbf{x} = [x_1, x_2]'$ and $M = 2$, then you'd have 4 features, $x_1^1, x_1^2, x_2^1, x_2^2$ along with an additional feature 1 for the intercept term. Plot training error and test error as Root Mean Square Error (RMSE) against M , the order of the polynomial features. Name the generated plot 1.png.
- (ii) Fit the data by using psuedo-inverse approach of regularized linear regression. For this, you need to use all polynomial features (i.e. $M = 6$), and choose values of $\ln \lambda$ using a sweep as follows: $-40, -39, \dots, 18, 19, 20$. Plot training error and test error as Root Mean Square Error (RMSE) against $\ln \lambda$, the regularization coefficient. Name the generated plot 2.png.
- (b) Because locally weighted linear regression is computationally more expensive, this problem is provided with a smaller test dataset. The training and testing data is provided to you, named **train_graphs_f16_autopilot_cruise.csv** and **test_locreg_f16_autopilot_cruise.csv**. For this part, you will not create the features by yourself, but rather use the features given in the dataset. Fit the data by applying the psuedo-inverse approach for solving locally weighted linear regression.

Use as weights the similarity between the point of interest and the example point as defined by

$$r_i(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}\|^2}{2\tau^2}\right)$$

and choose values of τ using a geometric sweep as follows: Generate 10 evenly spaced numbers on a logarithmic scale between 2^{-2} and 2^1 (start and end values inclusive). Plot only the test error as Root Mean Square Error (RMSE) against τ , a hyperparameter that affects the performance of the algorithm. Name the generated plot 3.png. Note: You might want to think about formulating efficient matrix multiplications in your code to implement locally weighted linear regression for this problem.

2) **Open Kaggle challenge (15 pts).** You can use any algorithm and design any features to perform regression on the Steel Ultimate Tensile Strength Dataset. All details for this project are included in the Kaggle webpage <https://inclass.kaggle.com/c/steel-ultimate-tensile-strength>. This problem will be graded separately based on your performance on the private leaderboard.

3) **Weighted Linear Regression (15 pts).** Consider a linear regression problem in which we want to weigh different training examples differently. Specifically, suppose we want to minimize

$$E_D(w) = \frac{1}{2} \sum_{i=1}^N r_i (w^\top x_i - t_i)^2. \quad (1)$$

In class, we worked out what happens for the case where all the weights (r_i 's) are one. In this problem, we will generalize some of those ideas to the weighted setting. In other words, we will allow the weights r_i to be different for each of the training examples.

(a) Show that $E_D(w)$ can also be written as

$$E_D(w) = (Xw - t)^\top R(Xw - t) \quad (2)$$

for an appropriate definition of X , t and R , i.e., write the derivation as part of your submitted answer. State these three quantities clearly as part of your answer (i.e., specify the form, matrix, vector, scalar, their dimensions, and how to build them from the x_i 's, t_i 's and r_i 's).

(b) If all the r_i 's are equal to 1, we showed in class that the normal equation is

$$X^\top Xw = X^\top t \quad (3)$$

and that the value of w^* that minimizes $E_D(w)$ is given by $(X^\top X)^{-1} X^\top t$. By finding the derivative $\nabla_w E_D(w)$ and setting that to zero, generalize the normal equation to this weighted setting, and give the new value of w^* that minimizes $E_D(w)$ in closed form as a function of X , R and t .

(c) Suppose we have a training set $\{(x_i, t_i); i = 1, \dots, N\}$ of N independent examples, but in which the t_i 's were observed with different variances. Specifically, suppose that

$$p(t_i | x_i; w) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(t_i - w^\top x_i)^2}{2(\sigma_i)^2}\right). \quad (4)$$

In other words, t_i has mean $w^\top x_i$ and variance $(\sigma_i)^2$ (where the σ_i 's are fixed, **known** constants). Show that finding the maximum likelihood estimate of w reduces to solving a weighted linear regression problem. State clearly what the r_i 's are in terms of the σ_i 's.

4) Naive Bayes Classifier (35 pts).

(a) Download the files `spambase.train` and `spambase.test`.

The file `spambase.train` contains 2000 training data and `spambase.test` has 2601 test data. Both datasets have 58 columns: the first 57 columns are input features, corresponding to different properties of an email, and the last column is an output label indicating spam (1) or non-spam (0). Please fit the Naive Bayes model using the training data.

As a pre-processing step, quantize each variable to one of two values, say 1 and 2, so that values below the median map to 1, and others map to 2. Please aggregate the training and test test to obtain the median value of each variable, and then use the median to quantize both data sets.

(i) Look up definitions of nominal/ordinal/interval/ratio variables. Which one(s) of them are suitable for the pre-processing described above? Look up what features are used in `spambase` data. Are they all suitable? Report briefly.

(ii) Report the test error (misclassification percentage) of Naive Bayes classifier. As a sanity check, what would be the test error if you always predicted the same class, namely, the majority class from the training data?

(b) Open Kaggle challenge: you can design any features to perform Naive Bayes classification on the Naive Bayes Spam Filter Dataset. All details for this project are included in the Kaggle webpage <https://inclass.kaggle.com/c/naive-bayes-spam-filter>. A utility script `extract_feature.py` that performs lemmatization and stopwords removal of email texts is provided. You may adapt it for your feature extraction.

This problem will be graded separately based on your performance on the private leaderboard. Besides submitting your solutions to Kaggle, briefly describe how you extract features from email texts.

5) Softmax Regression (15 pts).

In this problem, you will generalize logistic regression (for binary/2-class classification) to allow more classes (> 2) of labels (for multi-class classification).

In logistic regression, we had a training set $\{(\phi(\mathbf{x}_n), t_n)\}$ of N examples, where $t_n \in \{0, 1\}$, and the likelihood function is

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n} \quad (5)$$

where $y_n = p(C_1|\phi(\mathbf{x}_n)) = \sigma(\mathbf{w}^T \phi(\mathbf{x}_n))$ and $1 - y_n = p(C_0|\phi(\mathbf{x}_n)) = 1 - \sigma(\mathbf{w}^T \phi(\mathbf{x}_n))$.

Now we generalize the logistic regression to multi-class classification, where the label t can take on K different values, rather than only two. We now have a training set $\{(\phi(\mathbf{x}_n), t_n)\}$ of N examples, where $t_n \in \{0, 1, \dots, K-1\}$. We apply a softmax transformation of linear functions of the feature variables ϕ for the posterior probabilities $p(C_k|\phi)$, so that

$$p(C_k|\phi) = \frac{\exp(\mathbf{w}_k^T \phi)}{\sum_{k=0}^{K-1} \exp(\mathbf{w}_k^T \phi)}$$

where \mathbf{w}_k are the parameters of the linear function for class k , and $\mathbf{w} = \{\mathbf{w}_k\}_{k=0}^{K-1}$ are the parameters for

softmax regression. The likelihood function is then given by

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N \prod_{k=0}^{K-1} p(C_k|\phi(\mathbf{x}_n))^{\mathbf{1}(t_n=k)} \quad (6)$$

where $\mathbf{1}(\cdot)$ is an indicator function so that $\mathbf{1}(\text{a true statement})=1$ and $\mathbf{1}(\text{a false statement})=0$. In the likelihood function, $\mathbf{1}(t_n = k)$ returns 1 if t_n equals k and 0 otherwise.

- (a) Gradient descent for softmax regression. We can optimize the softmax regression model in the same way as logistic regression using gradient descent. Please write down the error function for softmax regression $E(\mathbf{w})$ and derive the gradient of the error function with respect to one of the parameter vectors \mathbf{w}_j , i.e. $\nabla_{\mathbf{w}_j} E(\mathbf{w})$. (Show your derivation and final result using the given notations. The final result should not contain any partial derivative or gradient operator.)
- (b) Overparameterized property of softmax regression parameterization and weight decay.

Softmax regression's parameters are overparameterized, which means that for any hypothesis we might fit to the data, there are multiple parameter settings that result in the same mapping from $\phi(\mathbf{x})$ to predictions. For example, if you add a constant to every w , the softmax regression's predictions does not change. There are multiple solutions for this issue. Here we consider one simple solution: modifying the error function by adding a weight decay term, $\sum_{k=0}^{K-1} \mathbf{w}_k^T \mathbf{w}_k$:

$$E^\lambda(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \sum_{k=0}^{K-1} \mathbf{w}_k^T \mathbf{w}_k \quad (7)$$

For any parameter $\lambda > 0$, the error function $E^\lambda(\mathbf{w})$ is strictly convex, and is guaranteed to have a unique solution. Similarly, please derive the gradient of the modified error function $E^\lambda(\mathbf{w})$ with respect to one of the parameter vectors \mathbf{w}_j , i.e. $\nabla_{\mathbf{w}_j} E^\lambda(\mathbf{w})$. (Show your derivation and final result using the given notations. The final results should not contain any partial derivative or gradient operator.)