
CS589: Machine Learning - Fall 2017

Homework 3: Kernels

Assigned: October 20th, 2017 Due: November 3rd, 2017

Getting Started: In this assignment, you will train and evaluate different kernels for both classification and regression on three datasets. Please install Python 3.6 via Anaconda on your personal machine. For this homework you will only be using numpy, scipy, sklearn and matplotlib packages. Download the homework file HW03.zip via Moodle. Unzipping this folder will create the directory structure shown below,

```
HW03
--- HW03.pdf
--- Data
    |--Synthetic
    |--CreditCard
    |--Tumor
--- Submission
    |--Code
    |--Figures
    |--Predictions
        |--CreditCard
        |--Tumor
```

The data files for each data set are in 'Data' directory respectively. You will write your code under the Submission/Code directory. Make sure to put the deliverables (explained below) into the respective directories.

Deliverables: This assignment has three types of deliverables:

- **Report:** The solution report will give your answers to the homework questions (listed below). Try to keep the maximum length of the report to 5 pages in 11 point font, including all figures and tables. Reports longer than five pages will only be graded up until the first five pages. You can use any software to create your report, but your report must be submitted in PDF format. You can use an additional two pages to answer extra credit questions only.
- **Code:** The second deliverable is the code that you wrote to answer the questions, which will involve implementing kernel methods. Your code must be Python 3.6 (no iPython notebooks or other formats). You may create any additional source files to structure your code. However, you should aim to write your code so that it is possible to re-produce all of your experimental results exactly by running *python run_me.py* file from the Submissions/Code directory.
- **Kaggle Submissions:** We will use Kaggle, a machine learning competition service, to evaluate the performance of your regression models. You will need to register on Kaggle using a *.edu email address to submit to Kaggle (you can use any user name you like). You will generate test prediction files, save them in Kaggle format (helper code provided called *Code/kaggle.py*) and upload them to

Kaggle for scoring. Your scores will be shown on the Kaggle leaderboard. The Kaggle links for each data set are given under respective questions.

Submitting Deliverables: When you complete the assignment, you will upload your report and your code using the Gradescope.com service. Place your final code in Submission/Code, and the Kaggle prediction files for your best-performing submission only for each data set in Submission/Predictions/<Data Set>/best.csv. Naming your files something other than best.csv breaks our grading scripts. Please try to name them best.cav. If you used Python to generate report figures, place them in Submission/Figures. Finally, create a zip file of your submission directory, Submission.zip (NO rar, tar or other formats). Upload this single zip file on Gradescope as your solution to the 'HW02-Kernel-Programming' assignment. Gradescope will run checks to determine if your submission contains the required files in the correct locations. Finally, upload your pdf report to the 'HW03-Kernel-Report' assignment. When you upload your report please make sure to select the correct pages for each question respectively. Failure to select the correct pages will result in point deductions. The submission time for your assignment is considered to be the later of the submission timestamps of your code, report and Kaggle submissions.

Academic Honesty Statement: Copying solutions from external sources (books, internet, etc.) or other students is considered cheating. Sharing your solutions with other students is also considered cheating. Posting your code to public repositories like GitHub, stackoverflow is also considered cheating. Any detected cheating will result in a grade of -100% on the assignment for all students involved, and potentially a grade of F in the course.

Task: In this homework you will experimenting with kernel methods for regression and classification problems using kernel ridge regression and support vector machine (SVM).

1. **Kernel Ridge Regression:** This is a kernel regression method. Let $\{x_i, y_i\}$ be the set of inputs, the problem consists on finding the w that minimizes,

$$\sum_i (x_i \cdot w - y_i)^2 + \lambda \|w\|_2^2 \quad (1)$$

where, the first term corresponds to the squared loss in the estimation, and the second is a regularization term. As seen during the lectures, one can apply a transformation to the samples (which increases their dimensionality) and perform the linear regression in this new higher dimensional space, which will correspond to a non-linear regression in the original space. An extremely important note regarding this is that, in order to estimate the value corresponding to a new sample x_{new} , only the inner product $x_{new} \cdot x_i$ is necessary. Thus, the kernel trick is applicable.

2. **SVM:** This is a classification method that can assign classes to new samples using only the inner product between the new samples and the samples in the training data. This allows us to use several different kernels, which makes SVM an extremely powerful classification method.

Data: You will work with three datasets,

- **Synthetic:** Use this dataset to check equivalence between basis expansion and the use of kernels in regression settings. You are provided four files: data_train.txt, label_train.txt, data_test.txt and label_test.txt.

- **CreditCard:** This dataset has eight attributes and two outputs, all real numbers. We treat this as a regression problem. The attributes correspond to credit card activity of individuals, and the outputs are two measure of risk as established by experts in the bank. You are provided three files: `data_train.txt`, `label_train.txt` and `data_test.txt`.
- **Tumor:** This dataset has nine attributes and a binary output. We treat this as a classification problem. The attributes are different measurements obtained from a medical imaging, and the output corresponds to the presence/absence of tumor. You are provided three files: `data_train.txt`, `label_train.txt` and `data_test.txt`.

Code to load datasets is provided in `run_me.py`. Below is a summary of the datasets, allowed python functions per dataset and performance metric to report in your HW03 pdf report which matches with Kaggle's performance reporting.

Dataset	Python functions	Use	Metric
Synthetic	sklearn.kernel_ridge.* Inbuilt Basis expansions sklearn.linear_model.Ridge	Not allowed Not allowed Allowed	Mean squared error (regression)
Credit Card	sklearn.kernel_ridge.* sklearn.model_selection.*	Allowed Allowed	Mean squared error (regression)
Tumor	sklearn.svm.* sklearn.model_selection.*	Allowed Allowed	Accuracy (classification)

Questions:

1. (75 points) Kernel Ridge Regression:

(11) a. Kernel Ridge Regression was introduced above. The goal is to find w that minimizes

$$\sum_i (x_i \cdot w - y_i)^2 + \lambda \|w\|_2^2 \quad (2)$$

Once the optimal w , w^* , is found, it can be used to estimate the value of new samples as $y_{est} = w^* \cdot x_{new}$. Show that, without using basis expansion nor kernels, the optimal w is,

$$w^* = \left(\sum_i x_i x_i^T + \lambda I \right)^{-1} \sum_i x_i y_i \quad (3)$$

(4) b. Show the solution for the case in which a basis expansion $x \rightarrow \Phi(x)$ is used, what is the expression of w^* ?

(10) c. As mentioned above, for this type of regression the kernel trick is applicable. This means that it is not necessary to know the transformation Φ ; having an expression for $\Phi(x_1) \cdot \Phi(x_2)$ (inner product of two samples in the new space) suffices. Given a new sample x_{new} , derive an expression for y_{new} that depends only on inner products between samples.

(30) d. In this exercise you will perform Kernel Ridge Regression using the synthetic dataset with samples $x_i \in \mathbb{R}$ and labels $y_i \in \mathbb{R}$. The training data is in `data_train.txt`, `label_train.txt`, and the testing data is in `data_test.txt`, `label_test.txt`. The training data is shown in Fig. 1.

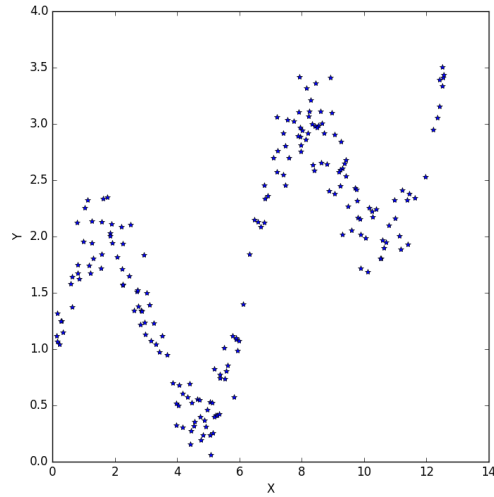


Figure 1: Data

We would like you to implement kernel ridge regression from scratch. Hence **only** for this question when it reads ‘Kernel ridge regression scratch (KRRS)’ we expect you to use your **own handwritten** kernel ridge regression, when it reads ‘Basis Expansion + ridge regression (BERR)’ we expect you to compute the basis expansion again by hand but you are free to use these expanded basis feature vectors with sklearn’s ridge regression or your own implementation of ridge regression. To make it blatantly obvious you are NOT allowed to use ‘`sklearn.kernel_ridge.KernelRidge`’ or any of its variants for KRRS and built in basis expansion functions for BERR.

We would like you to compare the performance of the two approaches on the synthetic dataset to better understand kernel methods. As mentioned in the class, the implementation using a kernel $k(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$ should return the same results as the implementation using the basis expansion Φ followed by ridge regression. You will use the following parameters for KRRS/BERRs (with $\lambda = 0.1$):

Polynomial order i: (for $i \in \{1, 2, 4, 6\}$)

1. Kernel ridge regression scratch (KRRS):

$$k(x_1, x_2) = (1 + x_1 \times x_2)^i \quad (4)$$

2. Basis expansion + ridge regression (BERR):

$$\Phi(x) = [1, x^1, x^2, \dots, x^i] \quad (5)$$

Trigonometric order i: (for $i \in \{3, 6, 10\}$)

1. Kernel ridge regression scratch (KRRS):

$$k(x_1, x_2) = 1 + \sum_{k=1}^i (\sin(k \delta x_1) \times \sin(k \delta x_2) + \cos(k \delta x_1) \times \cos(k \delta x_2)) \quad (6)$$

2. Basis expansion + ridge regression (BERR):

$$\Phi(x) = [1, \sin(\delta x), \cos(\delta x), \sin(2\delta x), \cos(2\delta x), \dots, \sin(i \delta x), \cos(i \delta x)] \quad (7)$$

With $\delta = 0.5$.

1. Train your models using the train data and predict the outputs for the test data. Visually inspect the quality of predictions by plotting the predicted test labels on top of the original test labels as a function of test data. Make these plots corresponding to polynomial kernels of degree 2 and 6, and trigonometric kernels of degree 5 and 10 only. A sample plot is provided in Figure 2. In each plot include the test samples as blue “*” and plot the predictions made on the test set as red circles. Make sure to label the axis, and include a title for each graph with implementation, kernel, degree and lambda used. Arrange your plots as a 4×2 grid (i.e. four rows and 2 columns; see *matplotlib.pyplot.subplot*) where the four rows correspond to the kernels and the two columns correspond to the two implementations. Write in at most three sentences what your observations are about these eight plots.

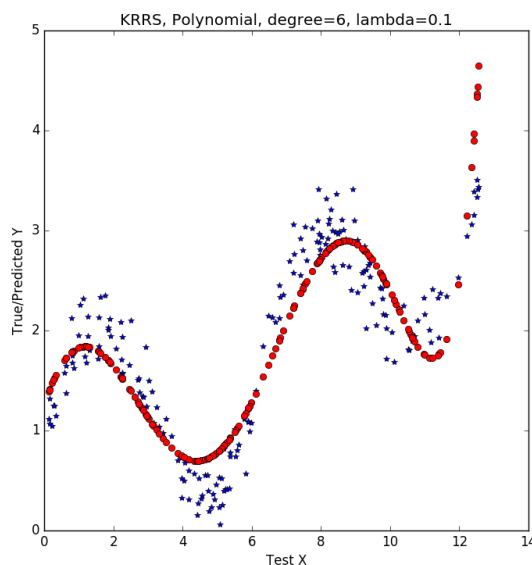


Figure 2: Test samples are plotted in blue “*” and predicted test labels are plotted in red “o”

2. For each of the kernels/basis expansions mentioned above, show the mean squared error over the **test set** using a table like the one below (please make the table exactly as shown below to make grading more efficient),

	Kernel $k(x_1, x_2)$	Basis Expansion $\Phi(x)$
Polynomial degree 1		
Polynomial degree 2		
Polynomial degree 4		
Polynomial degree 6		
Trigonometric degree 3		
Trigonometric degree 5		
Trigonometric degree 10		

where in each empty box you need to write the mean squared error obtained via KRRS or BERR.

- (20) e. In this exercise you will perform Kernel Ridge Regression using the credit card activity dataset. For this question you are allowed to (and encouraged) use sklearn's implementation of Kernel Ridge Regression. You will train several kernels (radial basis functions, polynomial degree 3, linear) with parameters - $\alpha \in \{1, 0.0001\}$ and $\gamma \in \{\text{default}, 1, 0.001\}$ (the default parameter is when you don't specify any γ). In total you will need to train and test 14 models (the γ parameter does not have any effect on linear kernels). For each kernel, pick your choice of model selection method to estimate the out of sample mean squared error. Report the results in a table format as shown below (please make the table exactly as shown below to make grading more efficient),

	RBF	Poly. degree 3	Linear
$\alpha = 1, \gamma = \text{def}$			
$\alpha = 1, \gamma = 1$			—
$\alpha = 1, \gamma = 0.001$			—
$\alpha = 0.0001, \gamma = \text{def}$			
$\alpha = 0.0001, \gamma = 1$			—
$\alpha = 0.0001, \gamma = 0.001$			—

Choose the model with lowest estimated out of sample error, train it using the full training set, make predictions on the test set, kagglize your outputs (helper code given in *Code/kaggle.py*) and finally upload your predictions to Kaggle. The Kaggle link to this question is

<https://www.kaggle.com/t/1776eed5f39a48d4928b1979c14c5f5a>.

Report back the mean squared error that you obtained on the public leader board as well as your Kaggle display name so that we can match your results with your rank on the leaderboard. Make sure to clearly indicate which model provided the best results as well as your choice of model selection method. Finally, consider the results that you obtained on the train set (via model selection) as well as your test results. Compare the two and draw some conclusions on your chosen kernel, specific parameter settings and its relation to the dataset itself.

2. (25 points) SVM:

- (25) a. For this question you will use perform classification using a SVM on the tumor dataset. For this question you are allowed to (and encouraged) use sklearn's implementation of SVM. You will train an SVM classifier using the kernels shown in the below table, and complete the table using estimates

of the out of sample accuracy (NOT error but classification accuracy) obtained using your choice of model selection method.

	RBF	Poly. degree 3	Poly. degree 5	Linear
$C = 1, \gamma = 1$				
$C = 1, \gamma = 0.01$				—
$C = 1, \gamma = 0.001$				—
$C = 0.01, \gamma = 1$				
$C = 0.01, \gamma = 0.01$				—
$C = 0.01, \gamma = 0.001$				—
$C = 0.0001, \gamma = 1$				
$C = 0.0001, \gamma = 0.01$				—
$C = 0.0001, \gamma = 0.001$				—

Again here the γ parameter does not have any effect on linear kernels.

Choose the model with highest estimated out of sample accuracy, train it using the full training set, make predictions on the test set, kagglize your outputs (helper code given in *Code/kaggle.py*) and finally upload your predictions to Kaggle. The Kaggle link to this question is

<https://www.kaggle.com/t/ad816cf2d0a140b78e83bd16770ea3d6>.

Report back the classification accuracy that you obtained on the public leader board as well as your Kaggle display name so that we can match your results with your rank on the leaderboard. Make sure to clearly indicate which model provided the best results and your choice of model selection method. Finally, consider the results that you obtained on the train set (via model selection) as well as your test results. Compare the two and draw some conclusions on your chosen kernel, specific parameter settings and its relation to the dataset itself.

Extra Credit: These questions are deliberately open-ended, leaving you more space for creativity. As a result, you will need to carefully describe exactly what you did for each question. Also, note that these questions carry small point values. To maximize your score with limited time, you should make sure the above questions are done thoroughly and ignore these. We will be very stingy in giving credit for these questions – do them only for the glory, and only at your own risk!

3. (5 points) Extra credit: For the credit card dataset perform regression with your choice of kernels (you are free to even create your own kernels). Experiment with different parameter ranges and/or model selection methods with a single goal of improving performance on the held out test set. For this extra credit question the kaggle URL to submit your predictions is

<https://www.kaggle.com/t/953acf0a51014c04b5f3747b97520370>.

Make sure to clearly describe your approach and list your performance on the public leaderboard. To avoid any confusion make sure to write the outputs to Predictions/CreditCard/best_extra_credit.csv. Note that your grades will be dependent on your ranking in the private leaderboard.

4. (5 points) Extra credit: For the tumor dataset perform classification with your choice of kernels. Experiment with different parameter ranges and/or model selection methods with a single goal of improving performance on the held out test set. For this extra credit question the kaggle URL to submit your predictions is

<https://www.kaggle.com/t/adcd36c33ef4199ad8faef7f88c9e90>.

Make sure to clearly describe your approach and list your performance on the public leaderboard. To avoid any confusion make sure to write the outputs to Predictions/Tumor/best_extra_credit.csv. Note that your grades will be dependent on your ranking in the private leaderboard.

5. (5 points) Code Quality:

- (5) Your code should be sufficiently documented and commented that someone else (in particular, the TAs and graders) can easily understand what each method is doing. Adherence to a particular Python style guide is not required, but if you need a refresher on what well-structured Python should look like, see the Google Python Style Guide: <https://google.github.io/styleguide/pyguide.html>. You will be scored on how well documented and structured your code is.