# EE6550 Machine Learning HW4 Kernel SVR README

An Extension of HW3 kernel SVM. This is a Support Vector Regression, trained with Sequential Minimal Optimization.

- Author: Yu-Chun (Howard) Lo
- Email: howard.lo@nlplab.cc

## User Manual

### Dev Environment

- Developed under Anaconda 4.3.0 (x86_64).
- Require Numpy for matrix operations.
- Tested on Python 3.6.0.

### File Structure

- `/dataset`: **(Important)** The program reads datasets from this folder and performs training and testing on the specified training and test data files. (See the *Dataset Format* section below)
- `/logs`: **(Important)** Output reports are stored in this folder.
- `/hypothesis`: **(Important)** Output hypotheses are stored in this folder. (There is **timestamp** at the end of the file name, which corresponds to the report)
- `svr.py`: Support vector regression (SVR) model.
- `utils.py`: Some utilities used by this program, such as loading dataset, normalize labels, etc.
- `main.py`: **(Important)** The main program. User should train a SVR by running this program.

### Dataset Format

- Currently, the program only supports reading `.csv` file.
- The true label (real value) of each item should locate at the first column.
- **(Important)** If you want to train your SVR with your own dataset, please be sure that you've followed the required format described above, and have placed your own training and test data files in the `/dataset` folder.

### Getting Started

Train your SVM by running `python main.py` in terminal. Be sure that your terminal is under the same directory as `main.py`.

Note that we've set default values for required input arguments. Run `python main.py --help` to view input arguments information shown below.

```
usage: main.py [-h] [--header_filename HEADER_FILENAME]
```

```
                [--train_filename TRAIN_FILENAME]
                [--test_filename TEST_FILENAME] [--K K] [--C C] [--tol TOL]
                [--epsilon EPSILON] [--kernel_type KERNEL_TYPE]
                [--poly_degree POLY_DEGREE] [--rbf_sigma RBF_SIGMA]
                [--enable_heuristic] [--enable_kernel_cache]
                [--max_iteration MAX_ITERATION]

 Kernel support vector regression.

 optional arguments:
   -h, --help            show this help message and exit
   --header_filename HEADER_FILENAME
                         Training dataset with header csv. This is used to
                         output hypothesis. (Default:
                         "airfoil_self_noise_training_header.csv")
   --train_filename TRAIN_FILENAME
                         Training dataset csv. (Default:
                         "airfoil_self_noise_training.csv")
   --test_filename TEST_FILENAME
                         Training dataset csv. (Default:
                         "airfoil_self_noise_testing.csv")
   --K K                 Denotes for "K"-fold cross-validation for determine
                         the optimal hyper-parameters for SVR. (Default: None)
   --C C                 Parameter for penalty term. (Default: 0.1)
   --tol TOL             Tolerance for KKT conditions. (Default: 1e-2)
   --epsilon EPSILON     Epsilon in the epsilon-SVR model. (Default: 0.1)
   --kernel_type KERNEL_TYPE
                         Kernel type to be used in SVR. Acceptable kernel type:
                         "linear", "poly", "rbf". (Default: None)
   --poly_degree POLY_DEGREE
                         Degree of the polynomial kernel function ("poly").
                         Ignored by all other kernels. (Default: 3)
   --rbf_sigma RBF_SIGMA
                         Sigma term in RBF (guassian). Ignored by all other
                         kernels. (Default: 0.5)
   --enable_heuristic    Whether use Platts heuristics to train SVR. (Defualt:
                         False)
   --enable_kernel_cache
                         Whether precompute kernel results. This can speed up
                         training but need time to initialize when data is
                         large. (Defualt: True)
   --max_iteration MAX_ITERATION
                         Max iteration for SMO training algorithm to avoid not
                         converging. (Defualt: 5000)
```

# For Grading Session

Here we show some guides for different test scenarios:

- Place the training data file(e.g. `xxx_training.csv`), testing data file(e.g. `xxx_testing.csv`) and training header file(e.g. `xxx_training_header.csv`) in the `/dataset` folder before running `main.py` with specified `--train_filename`, `--test_filename` and `--header_filename`.
- For choosing a PDS kernel, for example:
  - For polynomial kernel, specify `--kernel_type="poly"` or maybe along with the free parameter

```
    --poly_degree=3.
```
  ◦ For RBF kernel, specify `--kernel_type="rbf"` or maybe along with the free parameter
    `--rbf_sigma=0.5`.

- For performing K-fold cross-validation, for example, specify `--K=5`.
- All the required output information, such as class label mapping, cross-validation history, optimal hyper-parameters, etc., are stored at the `logs/` folder. Note that the log file name indicates what kernel type and number of K-fold you choosed, and when you run the program. This naming convension aims to help graders to choose which report to check after running the program.
- Note that the csv file name of SVM hypothesis in the `hypothesis/` folder is concatenated with the timestamp. This aims to let the graders know which hypothesis file to choose to test after running the program. (You may remove the timestamp for grading)
- Summing up, you may want to run the following commands in the different test scienarios:

```
(For 5-fold cross-validation on polynomial kernel)
>> python main.py --header_filename="xxx_training_header.csv"
--train_filename="xxx_training.csv" --test_filename="xxx_testing.csv" --K=5
--kernel_type="poly"

(For 5-fold cross-validation on RBF kernel)
>> python main.py --header_filename="xxx_training_header.csv"
--train_filename="xxx_training.csv" --test_filename="xxx_testing.csv" --K=5
--kernel_type="rbf"

(For specifying hyper-parameters when choosing polynomial kernel)
>> python main.py --header_filename="xxx_training_header.csv"
--train_filename="xxx_training.csv" --test_filename="xxx_testing.csv"
--kernel_type="poly" --C=1.0 --poly_degree=3

(For specifying hyper-parameters when choosing RBF kernel)
>> python main.py --header_filename="xxx_training_header.csv"
--train_filename="xxx_training.csv" --test_filename="xxx_testing.csv"
--kernel_type="rbf" --C=1.0 --rbf_sigma=0.5
```

## Report

All required output information are stored at `logs/`.

Snippet (See more in `logs/`):

```
    ----------------------------------------------------------------------------
    ----------------
    [*] Cross validation history:
     -  Parameter: {'C': 1e-07, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
    'poly_degree': 2} | Cross validation error: 4.44724766242486
     -  Parameter: {'C': 1e-09, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
    'poly_degree': 2} | Cross validation error: 4.91084179652826
     -  Parameter: {'C': 1e-15, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
    'poly_degree': 2} | Cross validation error: 5.001990352909237
     -  Parameter: {'C': 1e-15, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
    'poly_degree': 3} | Cross validation error: 21.61332034530309
```

```
- Parameter: {'C': 0.775, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2} | Cross validation error: 65.32086004713993
- Parameter: {'C': 0.325, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2} | Cross validation error: 69.08583664826315
- Parameter: {'C': 0.001, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2} | Cross validation error: 117.62040993662359
- Parameter: {'C': 1.0, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2} | Cross validation error: 439.58008070347495
- Parameter: {'C': 0.55, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2} | Cross validation error: 696.4116476458846
- Parameter: {'C': 0.1, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2} | Cross validation error: 1007.4677956825835
- Parameter: {'C': 0.325, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 3496.125409745657
- Parameter: {'C': 0.1, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 4024.8259614097165
- Parameter: {'C': 0.001, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 4740.655452640667
- Parameter: {'C': 0.55, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 5410.4720372760285
- Parameter: {'C': 1e-09, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 9566.635704103344
- Parameter: {'C': 1.0, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 16123.291955513589
- Parameter: {'C': 0.775, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 39942.169216952716
- Parameter: {'C': 1e-07, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 3} | Cross validation error: 56155.5246502839
- Parameter: {'C': 1e-07, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 99147.07748827015
- Parameter: {'C': 1e-15, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 115320.1163507979
- Parameter: {'C': 0.001, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 132336.50857563427
- Parameter: {'C': 1e-09, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 183195.68687372873
- Parameter: {'C': 1.0, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 255183.7234049135
- Parameter: {'C': 0.775, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 271708.6547973535
- Parameter: {'C': 0.1, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 397473.1726039498
- Parameter: {'C': 0.55, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 404648.8464865573
- Parameter: {'C': 0.325, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 4} | Cross validation error: 626281.2435841758
--------------------------------------------------------------------------------
----------------
[*] Best parameter: {'C': 1e-07, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2}
[*] Best cross validation error: 4.44724766242486
[*] Start to train on full training data and evaluate on test data ...
--------------------------------------------------------------------------------
----------------
[*] Train file path: "dataset/airfoil_self_noise_training.csv"
[*] Test file path: "dataset/airfoil_self_noise_testing.csv"
[*] Datetime: 19:02:58
```

```
[*] Best parameter: {'C': 1e-07, 'kernel_type': 'poly', 'tol': 0.01, 'epsilon': 0.1,
'poly_degree': 2}
[*] Sample mean of bias: 128.653934387
[*] Sample std of bias: 4.34025587068
[*] Performance: Train error: 3.91426335831 | Test error: 3.83562189503
-------------------------------------------------------------------------------
-----------------
[*] Saving SVR hypothesis to "hypothesis/SVR_hypothesis_header-[19:02:58].csv" ...
[*] Output SVR hypothesis to "hypothesis/SVR_hypothesis_header-[19:02:58].csv"
success.
-------------------------------------------------------------------------------
-----------------
```

- For polynomial kernel:
  - For 5-fold cross-validation results, check `logs/svr-poly-5-fold-[HH:MM:SS]`
  - For 10-fold cross-validation results, check `logs/svr-poly-10-fold-[HH:MM:SS]`
- For RBF kernel:
  - For 5-fold cross-validation results, check `logs/svr-rbf-5-fold-[HH:MM:SS]`
  - For 10-fold cross-validation results, check `logs/svr-rbf-10-fold-[HH:MM:SS]`

Note that the corresponding hypothesis csv file is indicated by the `HH:MM:SS` timestamp, which is the time that the program finished training.