

MATH 680 Homework 1 Fall 2015

September 9, 2015

This homework is due on Wednesday, September 16 at 11:59pm. Provide both pdf, R files (and \LaTeX file for exercise 3). The naming convention is yourlastname_hw1.pdf, yourlastname_hw1.R and yourlastname_hw1.lyx.

1. (20%) Walk through the following steps:

- Sign up on <https://github.com/>.
- Create a repository (repo).
- Add me as a collaborator (github ID: emeryyi).
- Download and install **SourceTree** <https://www.sourcetreeapp.com/> or alternatively **GitHub Desktop** <https://desktop.github.com/> as the graphical user interface (GUI) client.
- Clone the repository you created on github.com to your local folder using the GUI client.
- Add an R file in the local folder, then stage and commit changes.
- Push the changes to your origin repo on github.com
- Make a new branch and switch to the new branch using “checkout”.
- Make some changes in the new branch, then stage and commit.
- Switch to the master (old) branch using “checkout”.
- Make some changes in the old branch (but make sure they do not conflict with the changes in the new branch), then stage and commit.
- Merge the old branch with the new branch, and then commit.
- Push all the changes to the remote server.
- Provide the hyperlink and a screen shot of your github repo in a pdf file, which reflects all the commits, branching and merging you made in your repo.

2. (80%) In this exercise, we reproduce the “random function generator” (RFG) model by [1] (the paper is included, see page 19 section 6.1) in R. The RFG model generates very complicated data with non-linearity and higher-order interactions. The data generated by RFG can be used for

testing the performance of many fully non-parametric regression-based learning methods such as kernel support vector machine (KSVM), gradient boosting, random forests and others.

The idea is to generate a data frame with $N = 100$ observations of simulation data $\{y_i, \mathbf{x}_i\}_{i=1}^N$ according to

$$y_i = f(\mathbf{x}_i) + \epsilon_i,$$

where ϵ_i s are independent generated from the normal distribution $\epsilon \sim N(0, 1)$. In the data frame y takes up the first column, and the $p = 10$ dimensional vector \mathbf{x} takes up the rest columns. Hence the data frame will be 100×11 . Each row represents one observation. The f functions is randomly generated as a linear combination of functions $\{g_l\}_{l=1}^{20}$:

$$f(\mathbf{x}) = \sum_{l=1}^{20} a_l g_l(\mathbf{z}_l), \quad (1)$$

where coefficients $\{a_l\}_{l=1}^{20}$ are randomly generated from a uniform distribution $a_l \sim U[-1, 1]$. Each $g_l(\mathbf{z}_l)$ is a function of a randomly selected p_l -size subset (sub-vector) of the p -dimensional variable \mathbf{x} , with $p = 10$, and the size of each subset p_l is randomly chosen by $p_l = \min(\lfloor 1.5 + r \rfloor, p)$, and r is generated from an exponential distribution $r \sim \text{Exp}(0.5)$ with mean 2. Each \mathbf{z}_l is a sub-vector of \mathbf{x} defined as

$$\mathbf{z}_l = \{\mathbf{x}_{W_l(j)}\}_{j=1}^{p_l}, \quad (2)$$

where each W_l is an independent permutation of the integers $\{1, \dots, p\}$. Each function $g_l(\mathbf{z}_l)$ is an p_l -dimensional Gaussian function:

$$g_l(\mathbf{z}_l) = \exp \left[-\frac{1}{2} (\mathbf{z}_l - \boldsymbol{\mu}_l)^T \mathbf{V}_l (\mathbf{z}_l - \boldsymbol{\mu}_l) \right], \quad (3)$$

where each of the mean vectors $\{\boldsymbol{\mu}_l\}_{l=1}^{20}$ is randomly generated from the same distribution as that of the input variables \mathbf{x} . The $p_l \times p_l$ covariance matrix \mathbf{V}_l is also randomly generated by

$$\mathbf{V}_l = \mathbf{U}_l \mathbf{D}_l \mathbf{U}_l^T, \quad (4)$$

where \mathbf{U}_l is a $p_l \times p_l$ uniformly distributed random orthonormal matrix and $\mathbf{D}_l = \text{diag}\{d_{1l} \dots d_{p_l l}\}$. The variables d_{jl} are randomly generated from a uniform distribution $\sqrt{d_{jl}} \sim U[0.1, 2.0]$. We generated x from joint normal distribution $\mathbf{x} \sim N(0, \mathbf{I}_p)$ with $p = 10$.

Hints:

- In this exercise, you will probably use the following R functions:

- `rexp`, `runif`, `rnorm`
- `floor`
- `qr.Q`, `qr`
- `diag`
- `sample`

- To generate the uniformly distributed random orthonormal matrix \mathbf{U}_l , please see this post <http://math.stackexchange.com/questions/138512/sampling-q-uniformly-where-qtq-i>. Specifically, you need to generate a Gaussian matrix G , and use QR decomposition to find the corresponding Q matrix, which will be a uniformly distributed random orthonormal matrix. Use R functions `qr.Q` and `qr` to achieve this. Or much easier, you can just use the `genQ` function in the package `bootSVD` <https://cran.r-project.org/web/packages/bootSVD/index.html> to do the same job.
 - If it is your first time using R, learn the above functions using their online manuals before starting the actual coding.
3. (Extra credit 40%) Walk through the following steps:
- Download and install LyX. To make LyX work correctly, Mac users may also need to install MacTeX <https://tug.org/mactex/> while the Windows users may need to install MikTeX <http://miktex.org/>.
 - Reproduce this document using LyX, make sure that all text contents, mathematical formulas, references and typesetting (bullet, bold text, url, title, etc.) are accurately reproduced. However, the layout of the document do not need to be exactly the same. (paper margin, line space, etc.).

References

- [1] J.H. Friedman. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5):1189–1232, 2001.