

Machine Learning for Finance (FIN 570)

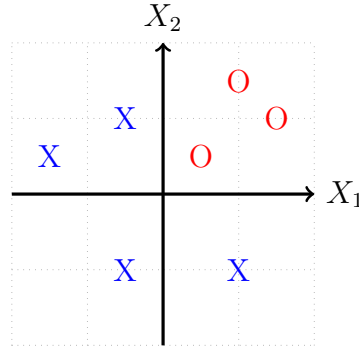
Midterm Exam

Instructor: Jaehyuk Choi

2018-19 Module 1 (2018. 10. 19)

1. (4 points) **(Bias vs variance trade-off)** For each of the following learning method, comment whether bias or variance (or both) is high.
 - (a) Unpruned decision tree (i.e. decision tree with unlimited number of leaves)
 - (b) Support vector machine with high allowance of violation
 - (c) Linear regression for 10 samples in training dataset with 9 features and 1 response.
 - (d) K -NN with large K .
2. (3 points) **(Probability)** Suppose that you are an equity analyst and you arrange two baskets of stocks for your clients. The clients pick a stock to invest in the following way: first, she randomly choose one basket and, then, randomly choose one stock in the basket. (The clients can not see the stocks in each basket.)

Suppose that 101 stocks are currently traded in Shenzhen Stock Exchange (SZSE), and you know that only 51 of the stocks will bring positive return. (You know the exact list of those stocks.) How can you arrange the two baskets in order for your client to maximize the probability of selecting a stock with positive return? What is the maximized probability? The only restriction is that there should be at least one stock in each basket.
3. (2 points) **(Logistic Regression)** Suppose we collect data for a group of students in the MLF class with variables X_1 = hours studied, X_2 = undergrad GPA, and Y = receive an A grade. We fit a logistic regression and produce estimated coefficient, $w_0 = -6$, $w_1 = 0.05$, $w_2 = 1$.
 - (a) Estimate the probability that a student who studies for 40 hours and has an undergrad GPA of 3.5 gets an A in the class.
 - (b) How many hours would the student in part (a) need to study to have a 50% chance of getting an A in the class?
4. (2 points) **(Test accuracy)** Suppose that we take a dataset, divide it into equally-sized training and test sets, and then try out two different classification procedures. First we use logistic regression and get an error rate of 20% on the training data and 30% on the test data. Next we use 1-nearest neighbors (i.e. $K = 1$) and get an average error rate (averaged over both test and training datasets) of 18%. Based on these results, which method should we prefer to use for classification of new observations? Why?
5. (3 points) **(Decision Tree)** Suppose that we want to apply decision tree to the following classification problem (O vs X):



- (a) Compute the Gini impurity of the raw dataset before splitting?
- (b) Let us consider the two criteria: (a) $X_1 > 0$ and (b) $X_2 > 0$. For the first split, which criteria should be selected? Compute the information gain for each.
6. (3 points) **(PCA analysis)** Assume that the covariance matrix of 3-dimensional input (X_1, X_2, X_3) is given as

$$\Sigma = \begin{pmatrix} 7 & 4 & 0 \\ 4 & 9 & 4 \\ 0 & 4 & 11 \end{pmatrix}.$$

- (a) What is the total variance, i.e., $\sum_{k=1}^3 \text{Var}(X_k)$?
- (b) Assume that we obtained the following three PCA component **without** correct order:

$$u_1 = \frac{1}{3}(2, -2, 1), \quad u_2 = \frac{1}{3}(1, 2, 2), \quad u_3 = \frac{1}{3}(2, 1, -2).$$

For example, the projection on the first PCA component is given as $P_1 = \frac{2}{3}X_1 - \frac{2}{3}X_2 + \frac{1}{3}X_3$. Calculate the variance of the projected features, $\text{Var}(P_k)$ for $k = 1, 2, 3$. Order the PCA components correctly according to the computed variance.

- (c) If we use the first two (correctly ordered) PCA components only, what is the ratio of explained variance by those two components over the total variance in (a)?
7. (3 points) **(XOR gate, K-NN, LOOCV)** The table below provides a training dataset containing four observations: two predictors and one classification response variable.

Obs.	X_1	X_2	Y
1	1	1	FALSE
2	1	-1	TRUE
3	-1	1	TRUE
4	-1	-1	FALSE

Suppose that we wish to use this dataset to make a prediction for Y using 1-NN (K -NN with $K = 1$). Here, we use the L_1 norm to compute the distance instead of the L_2 (Euclidean) norm. So, the distance between the two observations, (x_1, x_2) and (y_1, y_2) , is given as

$$D = |x_1 - y_1| + |x_2 - y_2|.$$

- (a) What is the prediction for Y when a new sample point is given as $(X_1, X_2) = (2, 1)$?
- (b) Draw the region of classification on the (X_1, X_2) coordinate.
- (c) For the training set given in the table above, what is the test accuracy if LOOCV method is applied? How does your answer change if you apply 2-NN or 3-NN?