

Topics in Quantitative Finance (FIN 528): Machine Learning for Finance

Jaehyuk Choi

March 9, 2017

Quantitative finance courses in PHBS

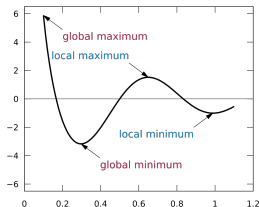
- Y1-M3: Stochastic Finance by Jaehyuk Choi [**required for Qfin MA**]
- Y1-M4: Derivative Pricing by Lei (Jack) Sun
- Y2-M1: Applied Stochastic Processes by Jaehyuk CHOI
Application, Programming, Course project
- Y2-M3: Topics in Quantitative Finance by Jaehyuk CHOI
Machine Learning for Finance (Mon-Thurs 1:30 PM)
- Y2-M3: Numerical Methods and Analysis by Jake ZHAO (Mon-Thurs 3:30 PM)
- Y2-M3: Bayesian Statistics by Qian CHEN (Mon-Thurs 10:30 AM)

AlphaGo vs Humans (LEE Sedol, KE Jie)



We have to think again about joseki / dìngshí (定石/ 定式) .

In Go (围棋), a joseki is the studied sequences of moves for which the result is considered balanced for both black and white sides.



LEE Sedol afterwards ...



- Became more popular and richer
- Perhaps titled as the last human who beat the machine in Go

ML/AI: Rise of the machines?

Probably not!



What and why now?

What is ML?

- Prediction based on data (data into knowledge)
- Extended linear/logistic regression
- Pattern recognition

Why now?

- Abundant Data (Big Data)
- Faster computer (Graphics Processing Unit: GPU)
- Advances in research: Geoffrey Hinton (Google), Yann LeCun (Facebook).

Recent applications of ML

- Automated driving system (Google, Apple, etc)
- Suggestion engine (Amazon, Taobao)
- Cancer diagnosis (IBM Watson)
- Digitizing images (Facebook, Google)
- Shopping without checkout (Amazon Go: big data)

ML in finance?

Prediction

- Asset management / investment
- Trading algorithm (alpha)
- Earnings prediction: e.g., [Prediction Valley](#)

Cost cut / labor reduction

- Automated accounting / tax
- Automated analyst report
- Chat-bot (trading and sales)
- Data analytics: e.g., [Kensho](#)

Softwares to use

Python

- Anaconda (Python distribution + Environment management)
- Python Language Tutorial
- Sci-Kit Learn
- TensorFlow (wrapped by Keras)

Github.com

- Distributed version control system
- Clone or fork a repository to create a local copy
- <https://github.com/PHBS/2016.M3.TQF-ML> (our course)
- <https://github.com/rasbt/python-machine-learning-book> (PML)

Other resources for ML

- Coursera ML course (**CML**) by Andrew Ng (Baidu)
- Stanford CS229 Machine Learning: course notes, student projects, etc
- The Elements of Statistical Learning (**ESL**)
- An Introduction to Statistical Learning (with R) (**ISLR**)
- Pattern Recognition and Machine Learning by Bishop (Microsoft)

Homework

- Install Anaconda (Ver 4.3 Python 3.6)
- Create GitHub account, join [PHBS](#) (organization) and [2016.M3.TQF-ML](#) (team).
- Fork the two repositories (the **PML** book and our course)
- Watch the first lecture of **CML**

Notations and conventions: vector and matrix

General rules (guess from the context)

- Scalar (non-bold): x, y, X, Y
- Vector (lowercase bold): $\mathbf{x} = (x_i), \mathbf{y} = (y_i)$
- Matrix (uppercase bold): $\mathbf{X} = (X_{ij}), \mathbf{Y} = (Y_{ij})$
- The (i, j) component of \mathbf{X} : X_{ij}
- The i -th row vector of \mathbf{X} : $\mathbf{X}_{i*} = (X_{i1}, X_{i2}, \dots, X_{in})^T$
- The j -th column vector of \mathbf{X} : $\mathbf{X}_{*j} = (X_{1j}, X_{2j}, \dots, X_{nj})$

Examples

- Dot product: $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y}$
- Vector norm: $|\mathbf{x}| = \sqrt{\mathbf{x}^T \mathbf{x}}$
- Matrix multiplication: $\mathbf{Z} = \mathbf{X} \mathbf{Y} \rightarrow Z_{ij} = \mathbf{X}_{i*} \mathbf{Y}_{*j}$

Notation and conventions: variables and observations

General rules

- Generic (or representative) variables (uppercase non-bold): X (input), Y (output), G (classification output)
- The predictions: \hat{Y} , \hat{G}
- X (input) may be p -dimensional (features): X_j ($j \leq p$), row vector
- Y (output) may be K -dimensional (responses): Y_k ($k \leq K$), row vector.
- The N observations of X or Y is stacked as rows: \mathbf{X} ($N \times p$), \mathbf{Y} ($N \times K$)
- The i -th observation set: \mathbf{X}_{i*} ($1 \times p$)
- All observation of j -th feature X_j : \mathbf{X}_{*j} ($N \times 1$)
- $\mathbf{X} = (\mathbf{X}_{*1} \cdots \mathbf{X}_{*p})$ (column-wise concatenation)
- The weight vector: β or \mathbf{w} used interchangeably.

Simple Linear Regression (Ordinary Least Square)

For scalar predictor (X) and response (Y),

$$Y \approx \beta_0 + \beta_1 X \quad \longrightarrow \quad \hat{y} = \beta_0 + \beta_1 x.$$

For N observations $(x_1, y_1), \dots, (x_N, y_N)$, the set of $(\hat{\beta}_0, \hat{\beta}_1)$ to minimize the residual sum of squares (RSS):

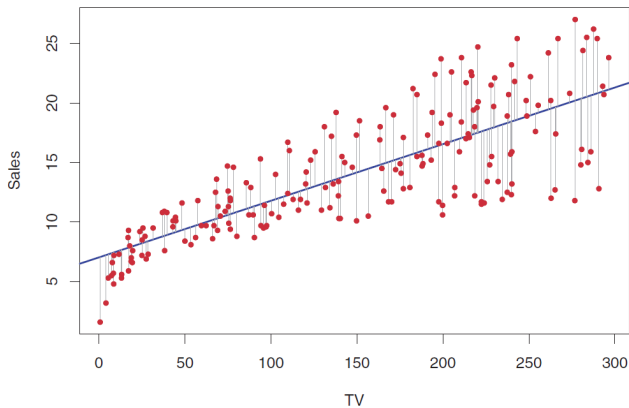
$$\text{RSS}(\beta) = \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_i)^2 = (\mathbf{y} - \beta_0 - \beta_1 \mathbf{x})^T (\mathbf{y} - \beta_0 - \beta_1 \mathbf{x})$$

is given as

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)},$$
$$\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$$

for $\bar{x} = \sum x_i / N$ and $\bar{y} = \sum y_i / N$.

Figure 3.1 (p62) from **ISLR**



Multi-dimensional Linear Regression

For $(p + 1)$ -vector predictor (\mathbf{X}) and scalar response (Y),

$$Y \approx \mathbf{X}\boldsymbol{\beta} \quad \longrightarrow \quad \hat{y} = \mathbf{X}\boldsymbol{\beta},$$

where $X_0 = 1$ ($\mathbf{X}_{*0} = \mathbf{1}$) and $\boldsymbol{\beta}$ is a $(p + 1)$ -column vector.

$$\text{RSS}(\boldsymbol{\beta}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\frac{\partial}{\partial \boldsymbol{\beta}} \text{RSS}(\boldsymbol{\beta}) = -\mathbf{X}^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \quad \Rightarrow \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

For $(p + 1)$ -vector predictor (\mathbf{X}) and K -vector response (\mathbf{Y}), the result is similarly given as

$$\hat{\mathbf{Y}} = \mathbf{X}\mathbf{B} \quad \text{where} \quad \hat{\mathbf{B}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y},$$

which is the independent regressions on Y_j (\mathbf{Y}_{*j}) combined together,

$$\hat{\mathbf{B}}_{*j} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}_{*j}$$

Newton's method

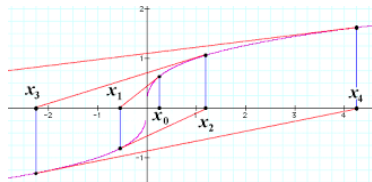
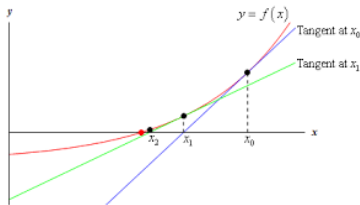
The root x satisfying $f(x) = 0$ can be found by the following iteration:

$$x_{n+1} = x_n - \eta \frac{f(x_n)}{f'(x_n)}$$

In multi-dimensional problems, the gradient is used instead of the differentiation:

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \eta \frac{f(\mathbf{x}_n) \nabla f(\mathbf{x}_n)}{|\nabla f(\mathbf{x}_n)|^2}$$

Typically we use $0 < \eta < 1$ to avoid *overshooting*.



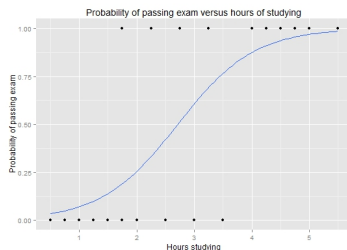
Logistic Regression (Classification)

- Qualitative (categorical) response (binary dependent variable, $Y \in \{0, 1\}$)
- Multiple categories: how to give order?
- Linear regression (quantitative) is not proper
- Logistic (sigmoid) function: $\sigma(\text{logit}) = \text{quantile}$

$$p = \phi(t) = \frac{e^t}{1 + e^t} = \frac{1}{1 + e^{-t}} \quad \text{for } t = X\beta \ (X_0 = 1)$$

- Logit function (the inverse): log odds

$$\phi^{-1}(p) = \log \left(\frac{p}{1 - p} \right) = \log(p) - \log(1 - p)$$



Fitting of logistic regression

Likelihood function

- For a given the prediction model, measures the likelihood of a data set.
- The best prediction model/weight is the one that maximizes the likelihood of the dataset.

For a data set (\mathbf{X}, \mathbf{y}) where $y_i \in \{0, 1\}$,

$$\begin{aligned} L(\beta) &= \prod_i P(y_i = \hat{y}_i) = \prod_{i:y_i=1} \phi(\mathbf{x}_{i*}\beta) \prod_{i:y_i=0} (1 - \phi(\mathbf{x}_{i*}\beta)) \\ &= \prod_i \phi(\mathbf{x}_{i*}\beta)^{y_i} (1 - \phi(\mathbf{x}_{i*}\beta))^{1-y_i} \\ \log L(\beta) &= \sum_i y_i \log \phi(\mathbf{x}_{i*}\beta) + (1 - y_i) \log (1 - \phi(\mathbf{x}_{i*}\beta)) \end{aligned}$$

The cost function (to minimize) is $J(\beta) = -\log L(\beta)$

Updating weights

After some algebra,

$$\mathbf{w} := \mathbf{w} + \Delta \mathbf{w}, \quad \Delta \mathbf{w} = -\eta \nabla J(\mathbf{w}) \quad (1)$$

$$\frac{\partial J(\mathbf{w})}{\partial w_j} = - \sum_i (y_i - \phi(\mathbf{X}_{i*} \mathbf{w})) X_{ij} = -\mathbf{X}_{*j}^T (\mathbf{y} - \phi(\mathbf{X} \mathbf{w})) \quad (2)$$

$$\text{or } \nabla J(\mathbf{w}) = -\mathbf{X}^T (\mathbf{y} - \phi(\mathbf{X} \mathbf{w})). \quad (3)$$

We get a similar weight updating rule as that of linear regression and Adaline! It gives a basis for Perceptron's updating rule.

Regularization

We avoid \mathbf{w} being too big.

$$J(\beta) = -\log L(\mathbf{w}) + \frac{\lambda}{2} |\mathbf{w}|^2 \quad (4)$$

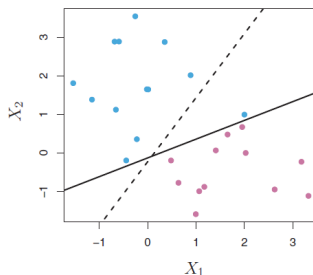
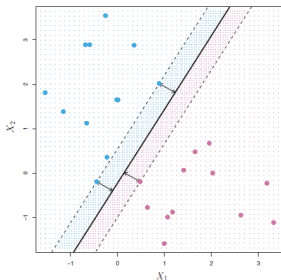
$$J(\beta) = -C \log L(\mathbf{w}) + \frac{1}{2} |\mathbf{w}|^2 \quad (C = \frac{1}{\lambda}, \text{SciKit-Learn}) \quad (5)$$

Maximal margin classifier

For $y_i \in \{-1, 1\}$, maximize the margin of the separating hyperplane M ,

$$y_i(w_0 + \sum_{j=1}^p X_{ij}w_j) = y_i(w_0 + \mathbf{X}_{i*}\mathbf{w}) \geq M > 0 \text{ for all } i, \text{ with } |\mathbf{w}| = 1$$

Maximal margin classifier only works for the separable data set and is sensitive to the change in the *support vectors*.



Support vector classifier

We make maximal margin classifier flexible: maximize the margin of the separating hyperplane M with $\|\mathbf{w}\| = 1$,

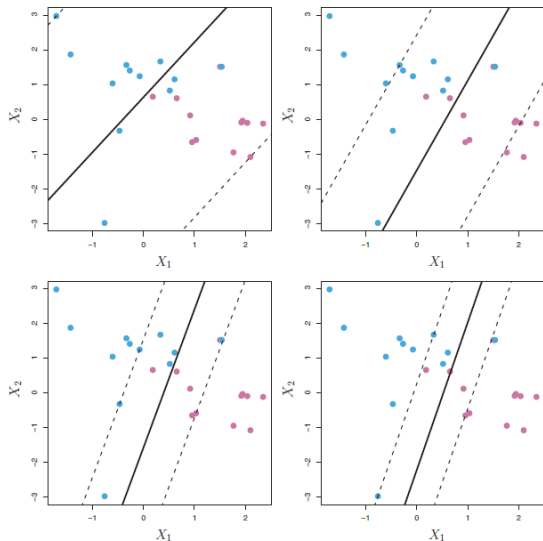
$$y_i(w_0 + \mathbf{X}_{i*}\mathbf{w}) \geq M(1 - \epsilon_i) \text{ for all } i, \quad \sum_{i=1}^n \epsilon_i \leq C,$$

where $\epsilon_i \geq 0$ is *slack variable* indicating the degree of violation ($\epsilon_i = 0$: no violation, $\epsilon_i < 1$: margin violation, $\epsilon_i > 1$: classification violation) and C is a *budget* for the amount of violations by all observations. Alternatively (in PML), we minimize

$$\frac{1}{2}\|\mathbf{w}\|^2 + C' \sum_{i=1}^n \xi_i,$$

where $\{\xi_i \geq 0\}$ satisfies $y_i(w_0 + \mathbf{X}_{i*}\mathbf{w}) \geq 1 - \xi_i$ for all i . The model converges to maximal margin classifier if C' is very large, so the role of C' is opposite to that of C in the original formulation.

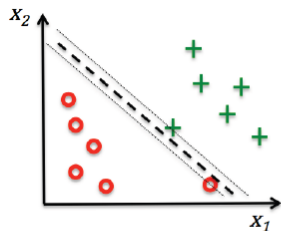
Support vector classifier: role of C



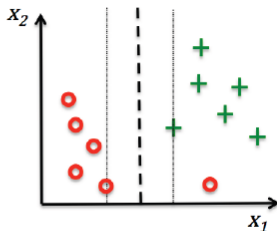
The value of C decreasing from the largest value on top left.

Support vector classifier: role of C'

Large C' (left) vs small C' (right)



Large value for
parameter C



Small value for
parameter C

Support vector machines (SVM)

How can we deal with non-linear decision boundary?

Enlarging feature space

Including high-order terms, $1, X_j, \dots, X_j^2, \dots, X_i X_j, \dots$, can be helpful, but the computation becomes very heavy.

Kernel

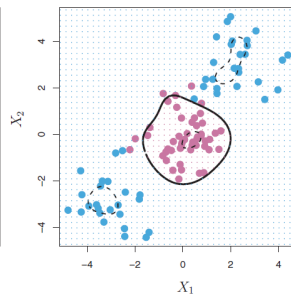
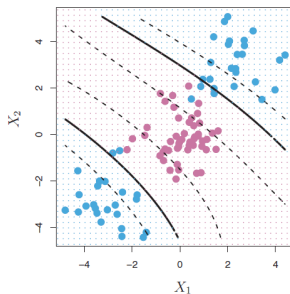
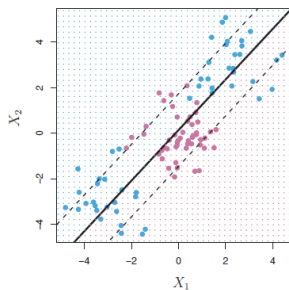
Instead we introduce kernel function, as a generalization of dot product in hyperplane:

- Linear: $K(\mathbf{X}_{i*}, \mathbf{X}_{i'*}) = \mathbf{X}_{i*} \mathbf{X}_{i'*}^T$
- Polynomial: $K(\mathbf{X}_{i*}, \mathbf{X}_{i'*}) = (1 + \mathbf{X}_{i*} \mathbf{X}_{i'*}^T)^d$
- Radial basis: $K(\mathbf{X}_{i*}, \mathbf{X}_{i'*}) = \exp(-\gamma |\mathbf{X}_{i*} - \mathbf{X}_{i'*}|^2)$

Kernel $K(\mathbf{X}_{i*}, \mathbf{X}_{i'*})$ can be understood as a *distance* between two observations: \mathbf{X}_{i*} and $\mathbf{X}_{i'*}$ are similar if the kernel value is high (low) whereas they are different if low (high).

SVM: non-linear decision boundary

SVM classification with linear kernel (left) polynomial kernel of degree 3 (middle) and radial basis kernel (right)



K Nearest Neighbor (KNN)

If $N_K(x)$ is the set of the K nearest neighbors around x ,

$$\text{Regression: } \hat{y} = f(x) = \frac{1}{K} \sum_{x_i \in N_K(x)} y_i$$

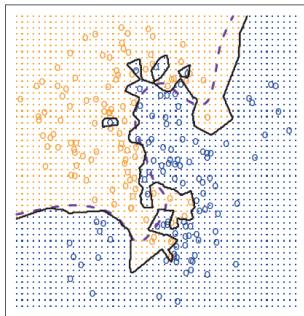
Classifier: $\hat{y} = \text{majority of } \{y_i\} \text{ for } x_i \in N_K(x)$

$$\text{Prob}(y = j|x) = \frac{1}{K} \sum_{x_i \in N_K(x)} I(y_i = j)$$

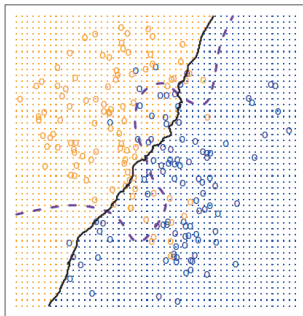
- Parametric vs Non-parametric model
- Learning step is not required, but KNN is intensive in both computation and storage (*memorize* training data set for prediction)

KNN: Classifier example

KNN: $K=1$

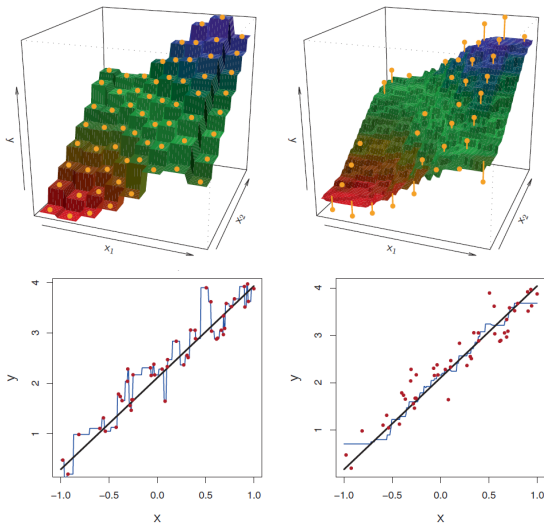


KNN: $K=100$

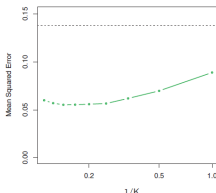
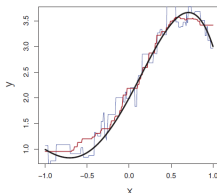
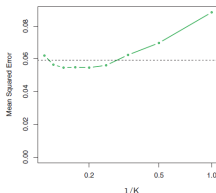
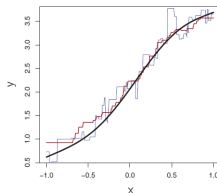
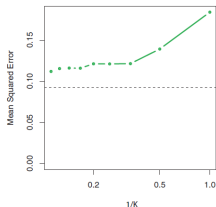
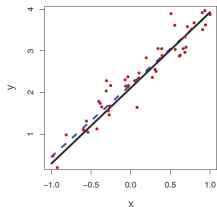


KNN: Regression example

$K = 1$ (left) vs $K = 9$ (right)



KNN: Parametric vs Non-parametric



Non-parametric regression works better as the true function deviates from the basis function (linear function in this example).