

# 用贝叶斯线性回归预测比特币价格变化

财通基金 量化投资部

2018 年 7 月 23 日

## 1 问题概述

## 2 数据预处理

## 3 特征聚类

## 4 拟合系数

## 5 数据预测

## 6 附录

## 1 问题概述

## 2 数据预处理

## 3 特征聚类

## 4 拟合系数

## 5 数据预测

## 6 附录

## 目标与方法

- 基于过去两年的比特币交易信息预测未来价格变化。
- 用贝叶斯线性回归模型预测，建立虚拟账户评估预测表现。

## 数据

- **数据集1**：（目前由于电脑没办法翻墙还没获得）比特币中国交易平台okex.com中2017年1月到2018年6月的价格和定单簿数据

1 问题概述

2 数据预处理

3 特征聚类

4 拟合系数

5 数据预测

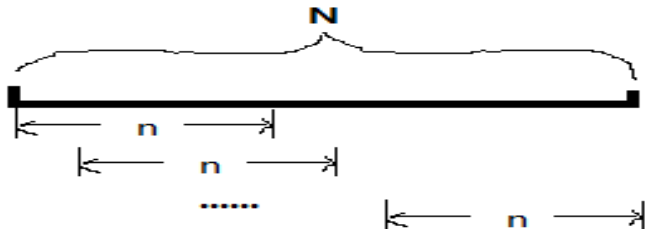
6 附录

## 训练集和测试集划分

- 原始价格数据用 $P$ 表示，包含 $N$ 个数据， $P_i$ 表示 $P$ 中第 $i$ 时刻的价格， $i \in [1, N]$
- 将 $P$ 按时间顺序分为三组长度相同的数据， $P_1 = [P_1 : P_{N/3}]$ ,  $P_2 = [P_{(N/3)+1} : P_{2N/3}]$ ,  $P_3 = [P_{(2N/3)+1} : P_N]$
- 其中 $P_1$ 用于生产时间序列及聚类
- $P_2$ 用于线性回归拟合系数 $w_0, w_1, w_2, w_3, w_4$
- $P_3$ 用于预测比特币价格变动

## 划分 $P_1$ 内的数据

- 用长度为 $n_L$ ,  $1 \leq L \leq 3$ 的移动窗口
- $n_1 = 180, n_2 = 360, n_3 = 720$ 生成 $H^T, 1 \leq T \leq 3$ ,  
 $H$ 有 $N - n + 1$ 行,  $n + 1$ 列
- $j \leq n$ 时 $H^T(i, j) = P_1[i : i + n + 1]$
- $j = n + 1$ 时 $H^T(i, j) = P_1[i + n + 1] - P_1[1 + n]$



1 问题概述

2 数据预处理

3 特征聚类

4 拟合系数

5 数据预测

6 附录



## 聚类方法

- 将 $H^t$ 用KMean的方法划分成 $k$ 个聚类（建议 $k=100$ ）
- 计算每个聚类中最大值与最小值之差 $diff_m$ ，（ $1 \leq m \leq k$ ）  
将 $diff_m$ 排序并挑选 $diff_m$ 最大的 $s$ 个聚类
- 将 $diff_m$ 排序并挑选 $diff_m$ 最大的 $s$ 个聚类(建议 $s=20$ )
- 将每个聚类中心 $c_q$ （ $1 \leq q \leq s$ ）作为矩阵的一行，产生一个大小为 $s * (n + 1)$ 的矩阵 $C$ ， $C$ 的每一列为聚类中心的一个特征值

1 问题概述

2 数据预处理

3 特征聚类

4 拟合系数

5 数据预测

6 附录

## 贝叶斯线性回归公式

$$\blacksquare E_{emp} = \frac{\sum_{i=1}^n y_i * \exp(-1/4 * ||x - x_i||_2^2)}{\sum_{i=1}^n \exp(-1/4 * ||x - x_i||_2^2)} \quad [1]$$

## 参数意义

- $x = P_1[i - n_L + 1 : i - 1]$
- $x_i = C[i + 1, : n_L + 1]$
- $y_i = C[i + 1, n_L + 1]$

## 自变量X与应变变量Y

- 自变量 $X[\Delta p_1, \Delta p_2, \Delta p_3, r]$ , 应变变量 $Y[\Delta p]$
- 分别将 $x^L = P_2[i - n_L : i - 1], 1 \leq L \leq 3$ ,  
 $x_i^L = C_L[i + 1, : n_L + 1], y_i^L = C_L[i + 1, n_L + 1]$ 代入公式[1]得出 $\Delta p_T, 1 \leq T \leq 3$
- $r = \frac{V_{bid} - V_{ask}}{V_{bid} + V_{ask}}$
- $\Delta p = P_2[i + 1] - P_2[i]$

## 拟合X, Y求系数 $w_0 - w_4$

- 用Linear Regression对自变量X和应变变量Y进行fit操作,  $w_0$ 为该线性回归方程与y轴的焦点,  $w_1 - w_3$ 为线性方程回归的系数

1 问题概述

2 数据预处理

3 特征聚类

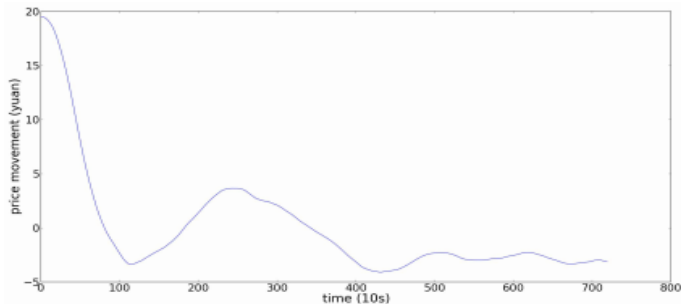
4 拟合系数

5 数据预测

6 附录

根据拟合出的系数  $w_0 - w_4$  得到  $X[\Delta p_1, \Delta p_2, \Delta p_3, r]$  和  $Y[\Delta p]$  之间的线性关系

$$\blacksquare \Delta p = w_0 + \sum_{j=1}^3 w_j * \Delta p_j + w_4 * r$$



- 1 问题概述
- 2 数据预处理
- 3 特征聚类
- 4 拟合系数
- 5 数据预测
- 6 附录

## 公式推导

- $s_1, \dots, s_K \in \mathbb{R}^d$  是  $K$  个不同的隐藏因子, 它们出现的概率为  $\mu_1, \dots, \mu_K$
- $T \in (1, \dots, K)$  是指针, 有  $P(T = k) = \mu_k$ , 其中  $1 \leq k \leq K$ ,  
 $x = s_T + \epsilon$
- $$\begin{aligned} P(y|x) &= \sum_{k=1}^K P(y|x, T = k)P(T = k|x) \\ &\propto \sum_{k=1}^K P(y|x, T = k)P(x|T = k)P(T = k) \\ &= \sum_{k=1}^K P_k(y)P(\epsilon = (x - s_k))\mu_k \\ &= \sum_{k=1}^K P_k(y)\exp(-\frac{1}{2} \|x - s_k\|_2^2)\mu_k \end{aligned}$$
- 用经验数据作为代替来估计给定  $x$  后  $y$  的概率分布
- $$P_{emp}(y|x) = \frac{\sum_{i=1}^n \mathbb{I}(y=y_i)\exp(-\frac{1}{4}\|x-x_i\|_2^2)}{\sum_{i=1}^n \exp(-\frac{1}{4}\|x-x_i\|_2^2)}$$
  

$$\Rightarrow E_{emp}[y|x] = \frac{\sum_{i=1}^n y_i \exp(-\frac{1}{4}\|x-x_i\|_2^2)}{\sum_{i=1}^n \exp(-\frac{1}{4}\|x-x_i\|_2^2)}$$