



中文維基百科 用戶分析






Snowyowls@zh.wikipedia.org









簡介

-  分析用戶的編輯歷史
-  瞭解不同用戶參與維基百科的特點
-  運用數學工具對用戶的編輯次數數據進行深入分析










數據來源

-  中文維基百科註冊用戶超過8萬
-  選取其中一小部分進行分析
-  將用戶分為5類
-  機器人、管理員、用戶、破壞者、辯論者






數據來源

-  機器人：中文維基百科全部閤法機器人
-  管理員：中文維基百科管理員
-  用 戶：統計頁面中編輯量Top50的註 冊用戶
-  破壞者：當前的破壞中列名的註冊用戶
-  辯論者：在南京大學和法輪功兩個條目的討論頁發言超過5次的用戶








數據來源

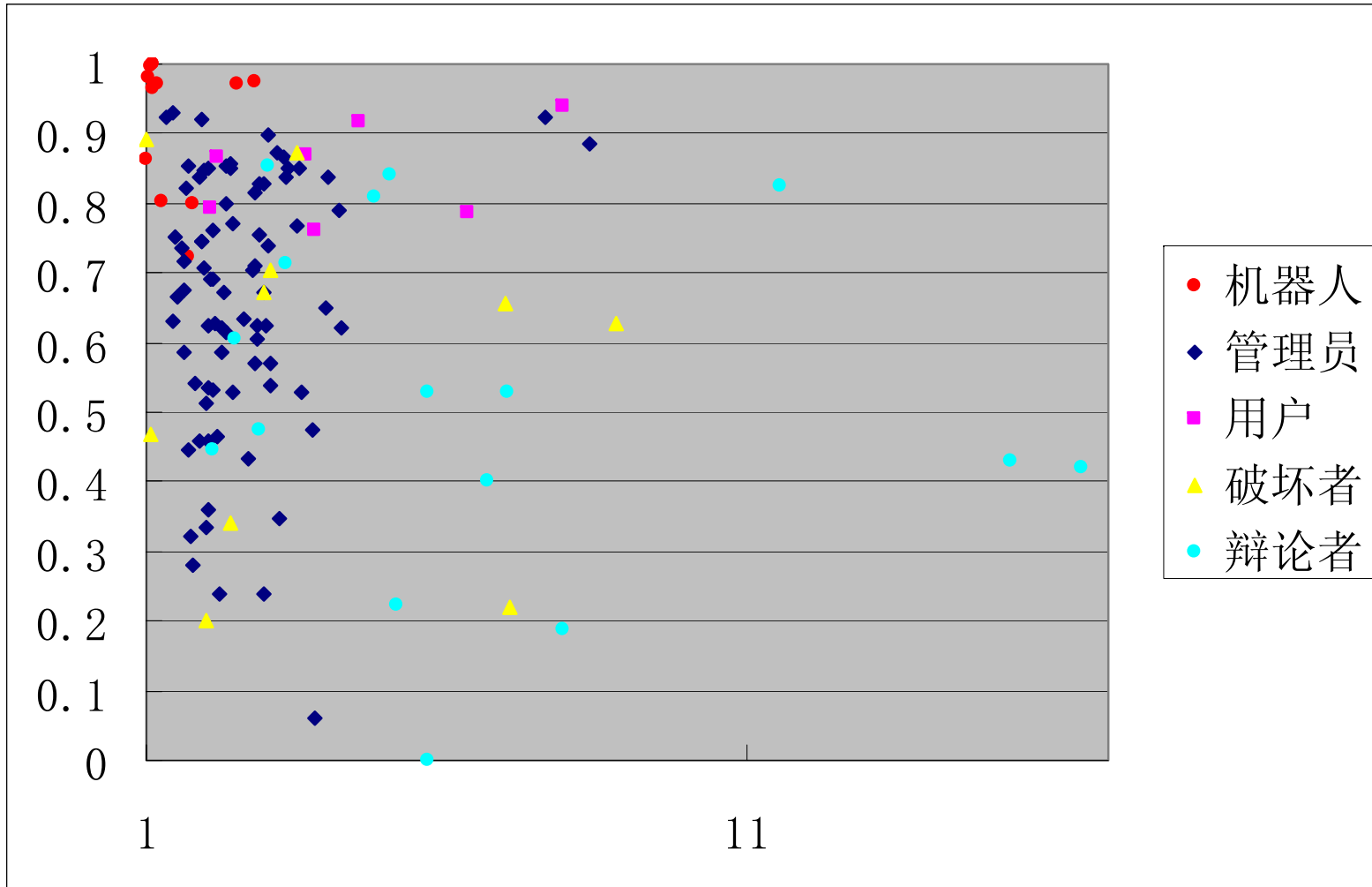
-  用戶的編輯數據用kate的工具從數據庫中提取
-  數據包括：Image uploads、cur and old images、Distinct pages edited、Deleted edits、Edits/pages、Articles、Talk、User、User talk、Project、Project talk、Image、Image talk、MediaWiki、MediaWiki talk、Template、Template talk、Help、Help talk、Category、Category talk、Portal、Portal talk、move



方法與結果




-  用戶的平均頁面編輯次數對內容編輯傾嚮作圖
-  內容編輯傾嚮是用戶的Articles、Template、Portal三項編輯次數之和佔總編輯次數的比例
-  辯論者常具有極高的平均頁面編輯次數

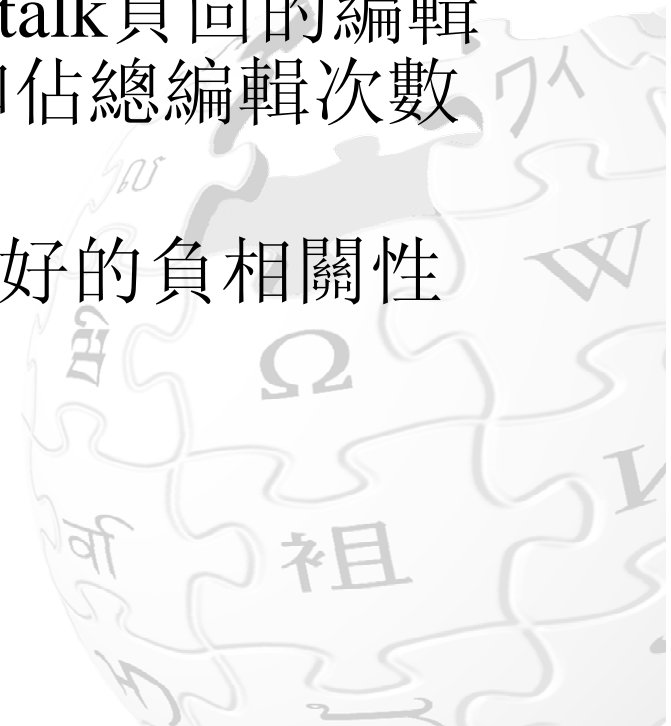


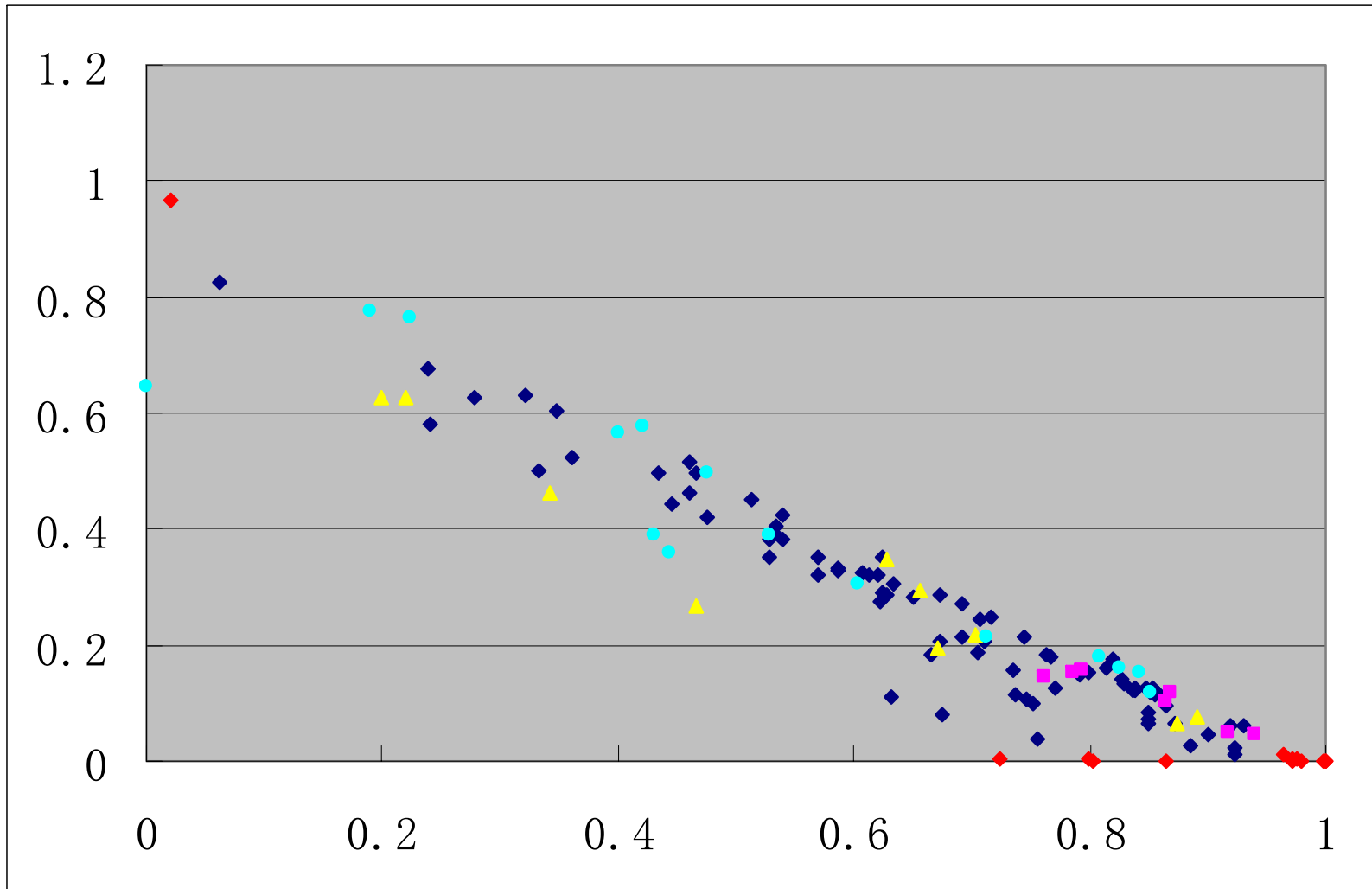




方法與結果




-  用戶的內容編輯傾嚮相對討論傾嚮作圖
-  用戶的討論傾嚮是用戶在所有talk頁面的編輯次數以及在Project編輯次數之和佔總編輯次數的比例
-  圖錶顯示這兩項數據之間有良好的負相關性



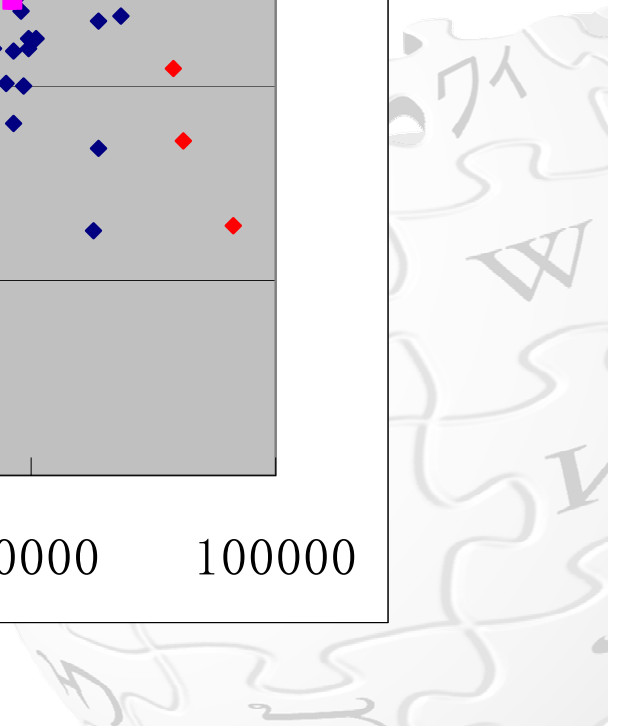
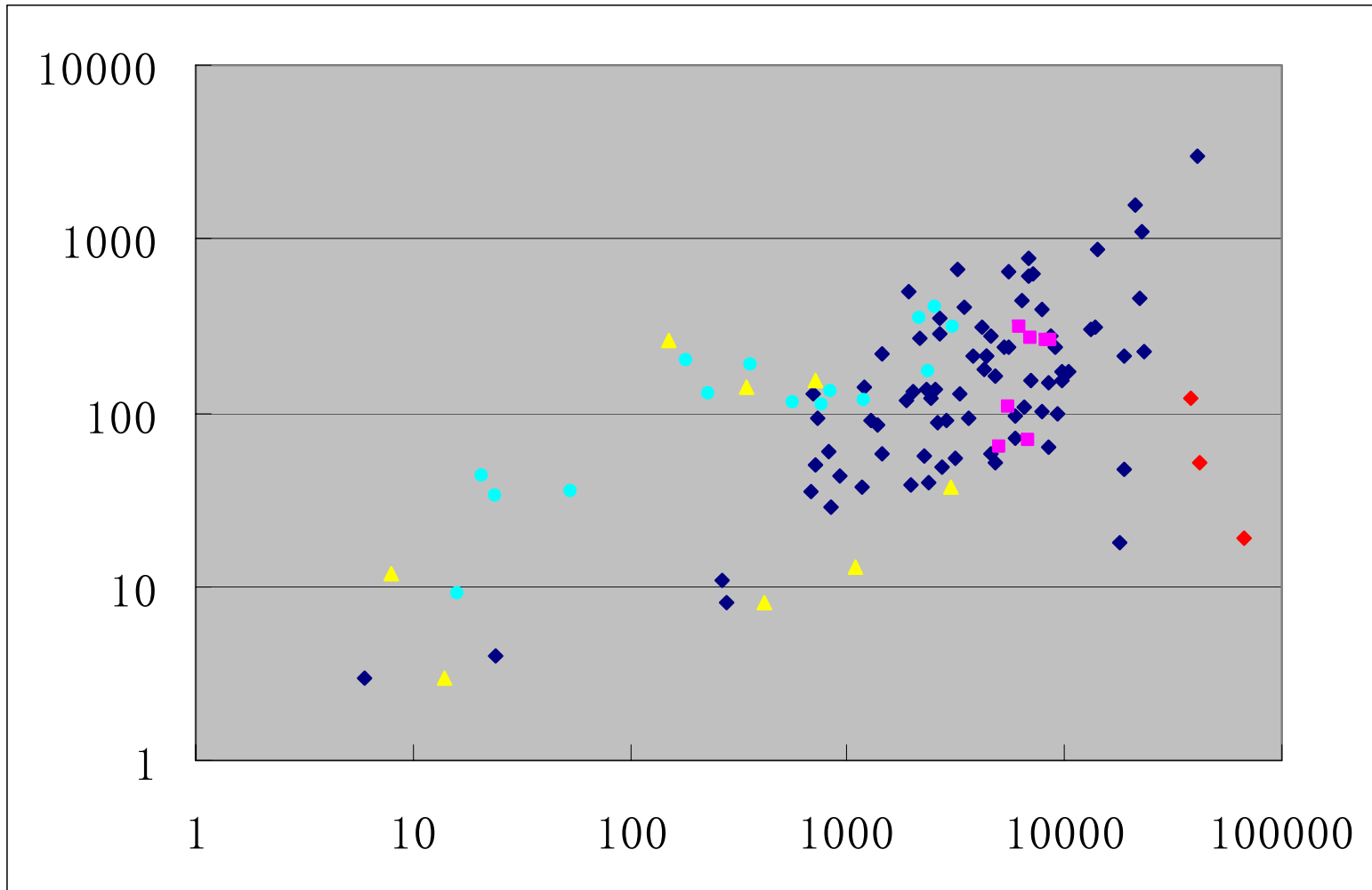




方法與結果




-  用戶的條目編輯次數相對討論頁編輯次數作圖
-  坐標均採用對數坐標
-  一部分破壞者和機器人的討論頁編輯次數為0，在對數坐標中沒有體現

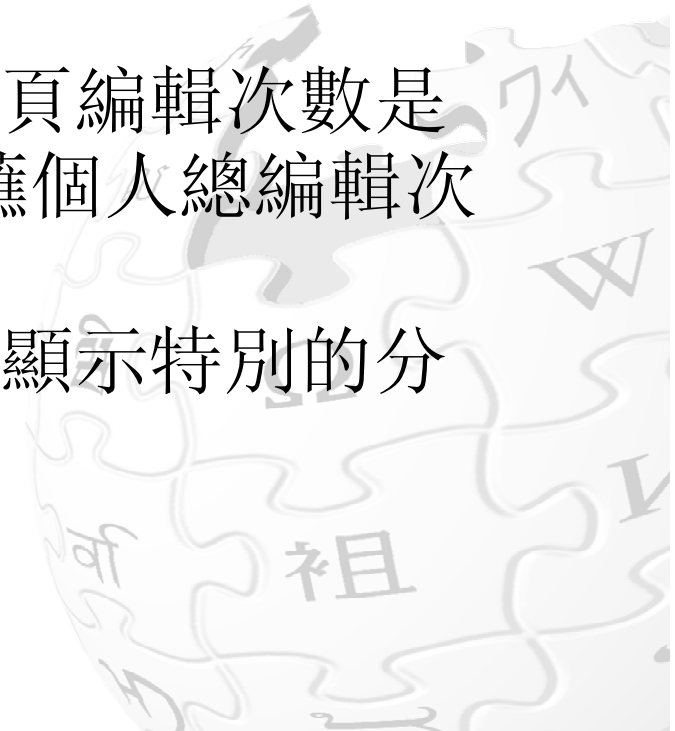


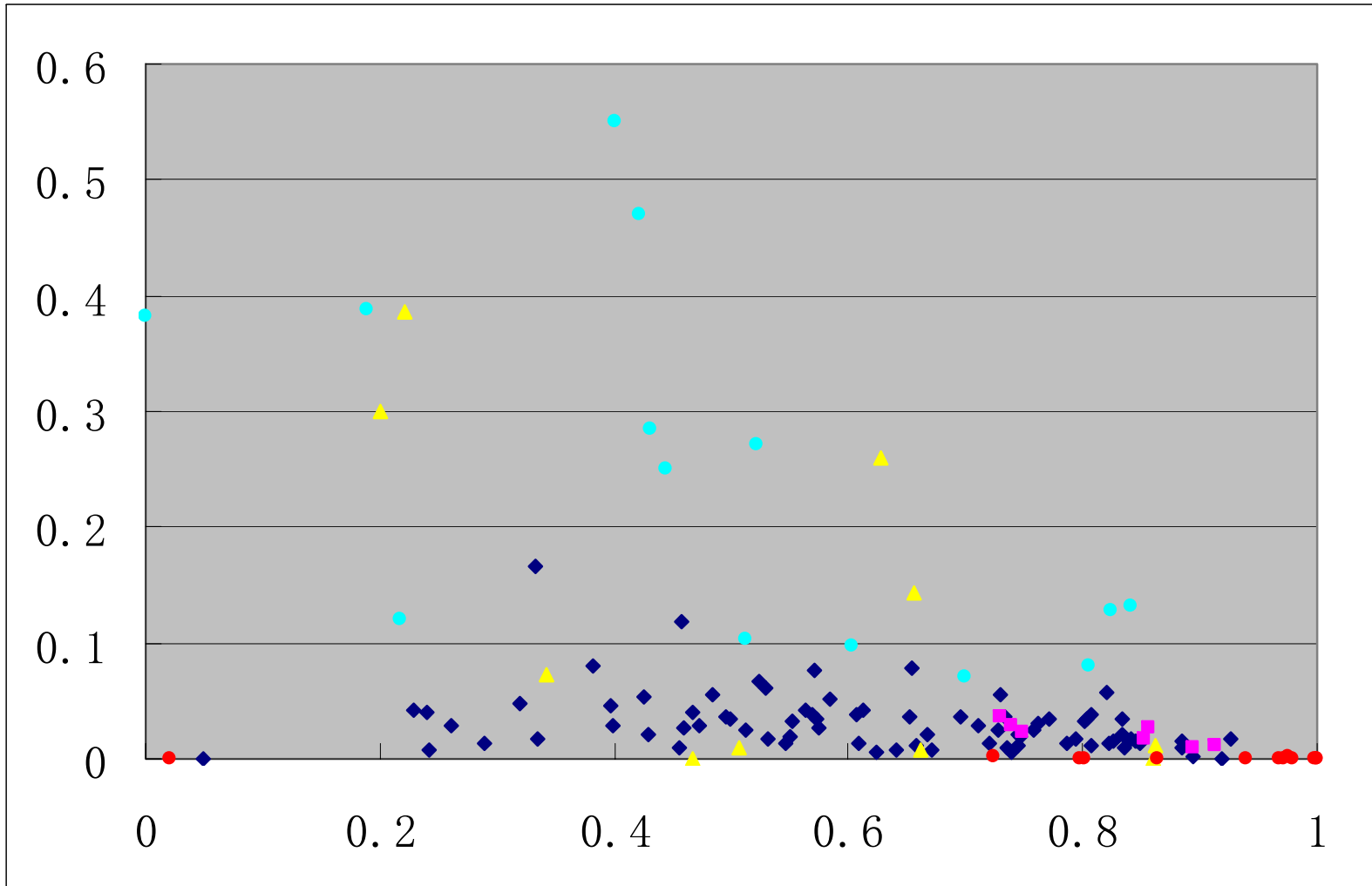




方法與結果

-  用戶的相對條目編輯次數對相對討論頁編輯次數作圖
-  相對條目編輯次數和相對討論頁編輯次數是用戶的條目和討論頁編輯次數蘸個人總編輯次數之比例
-  在這副圖中，辯論者和破壞者顯示特別的分佈

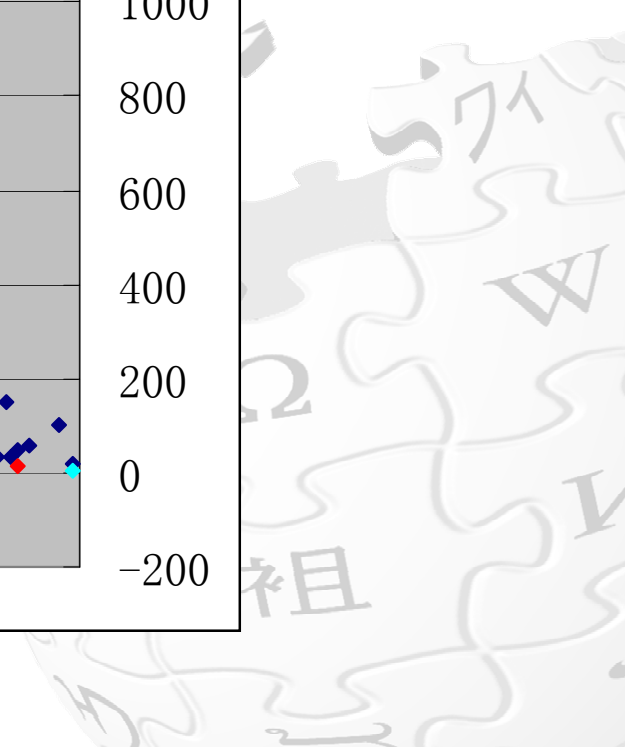
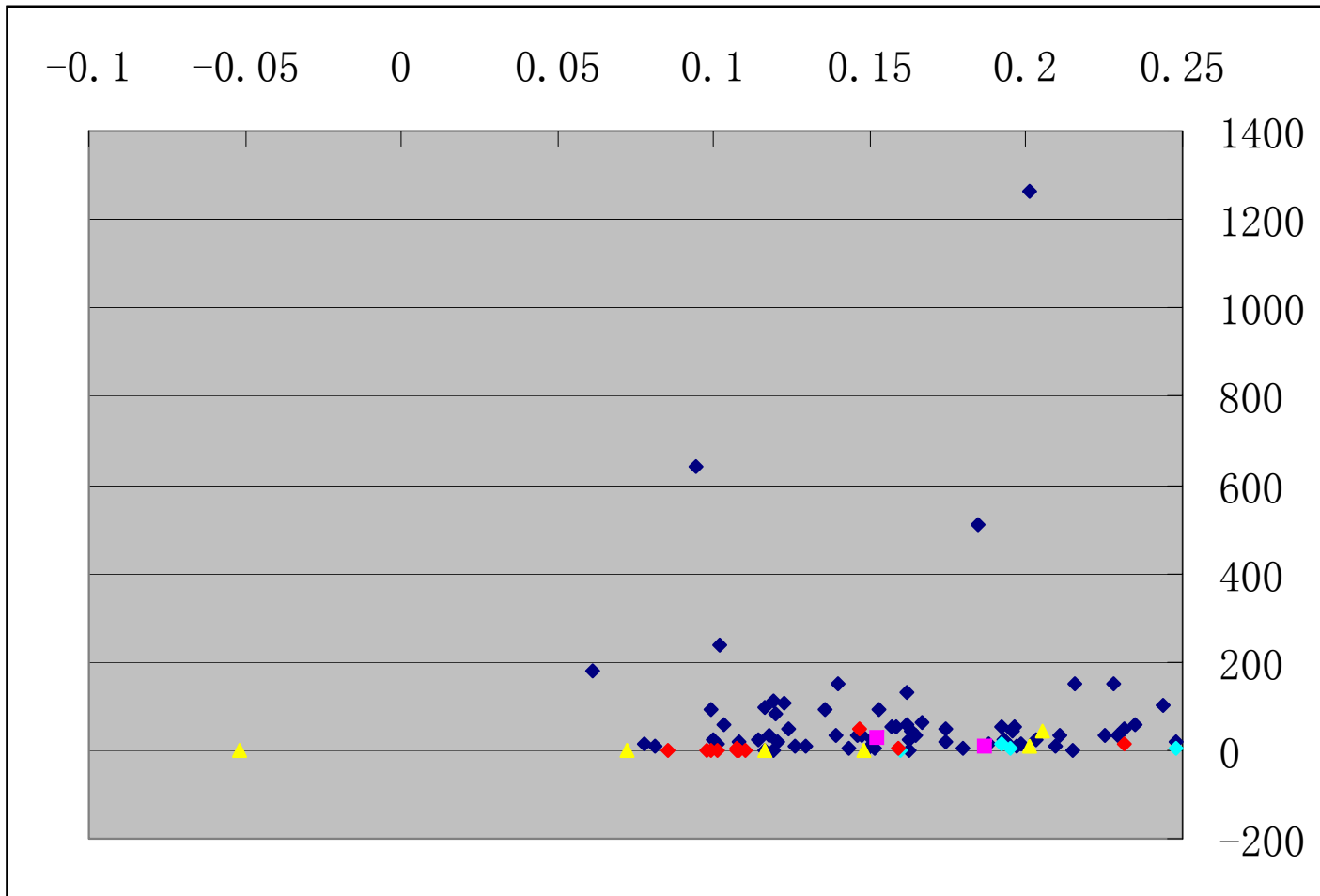






方法與結果

- 對用戶進行主成分分析后第一第二第三主成分組成的綜閣主成分相對第四主成分作圖
- 在統計學中，主成分分析是一種簡化數據集的技術。它把數據變換到一個新的坐標系統中，使得任何數據投影的第一大方差在第一個坐標(稱為第一主成分)上，第二大方差在第二個坐標(第二主成分)上，依次類推。主成分分析經常用減少數據集的維數，同時保持數據集的對方差貢獻最大的特徵

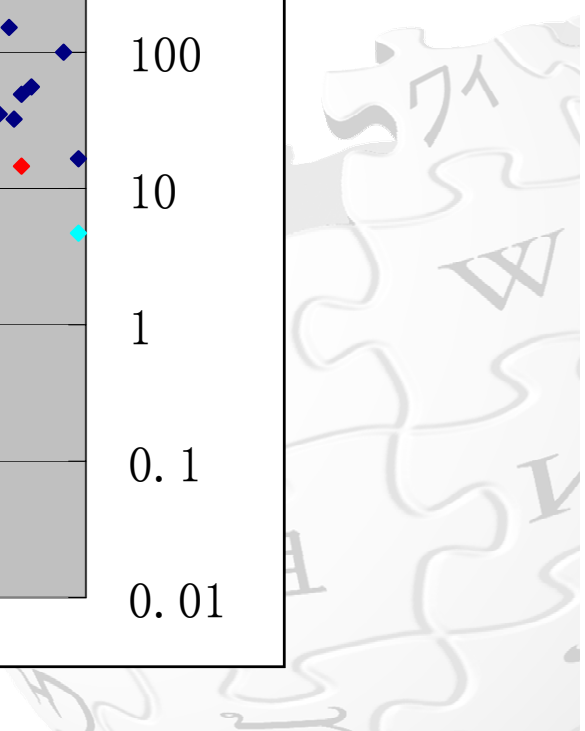
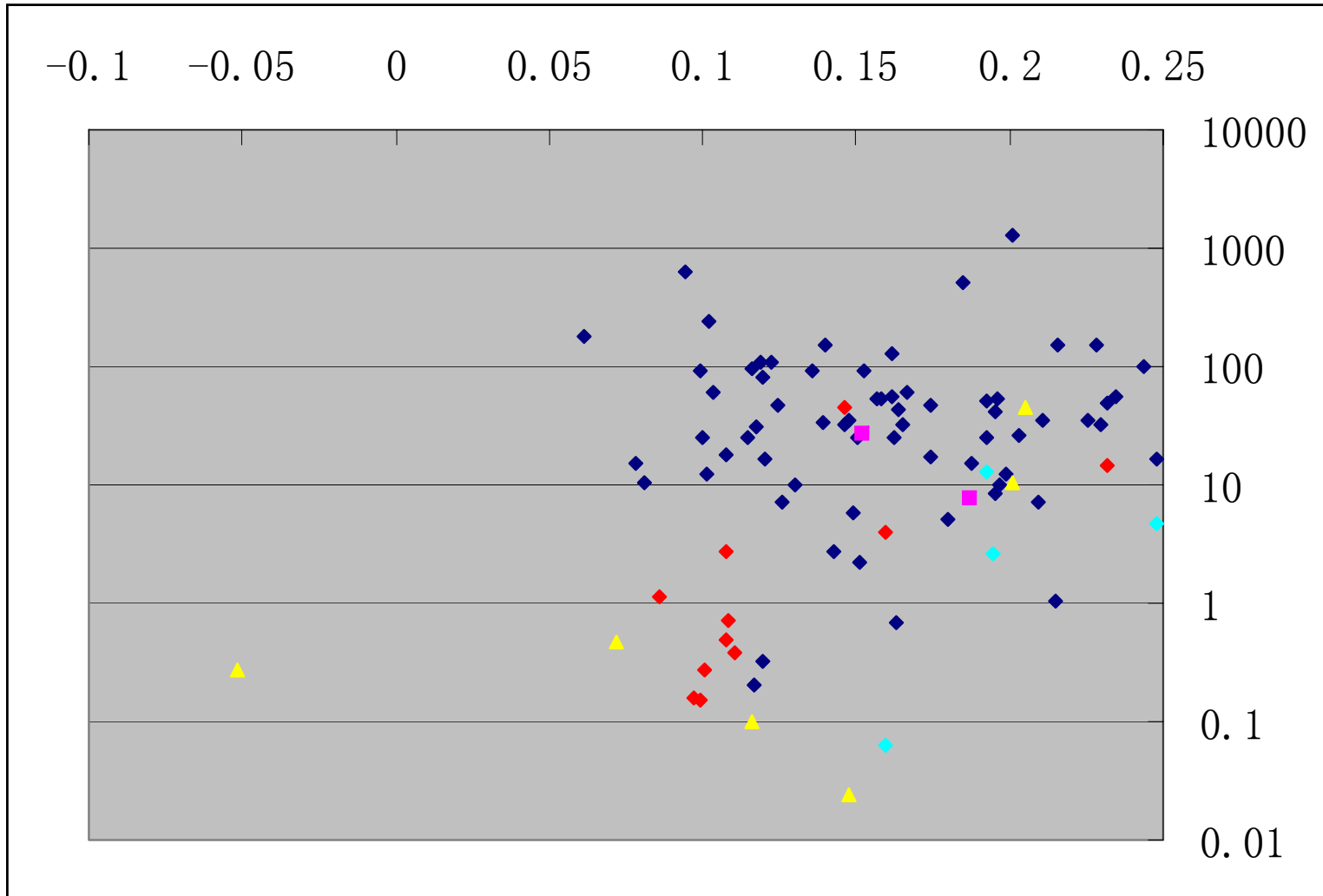


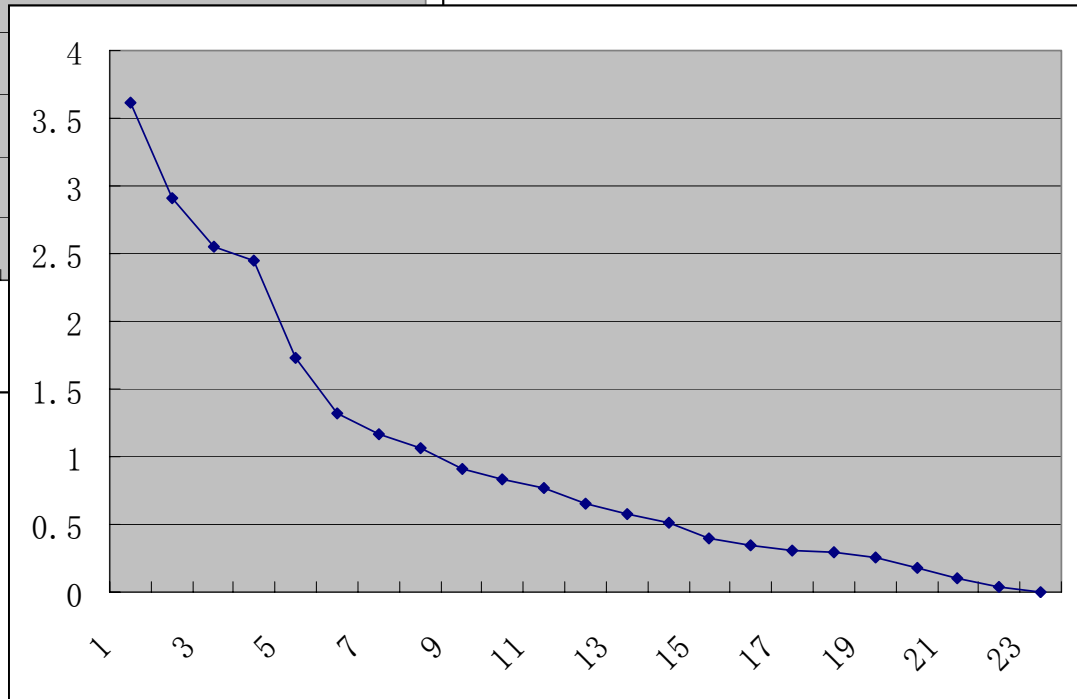
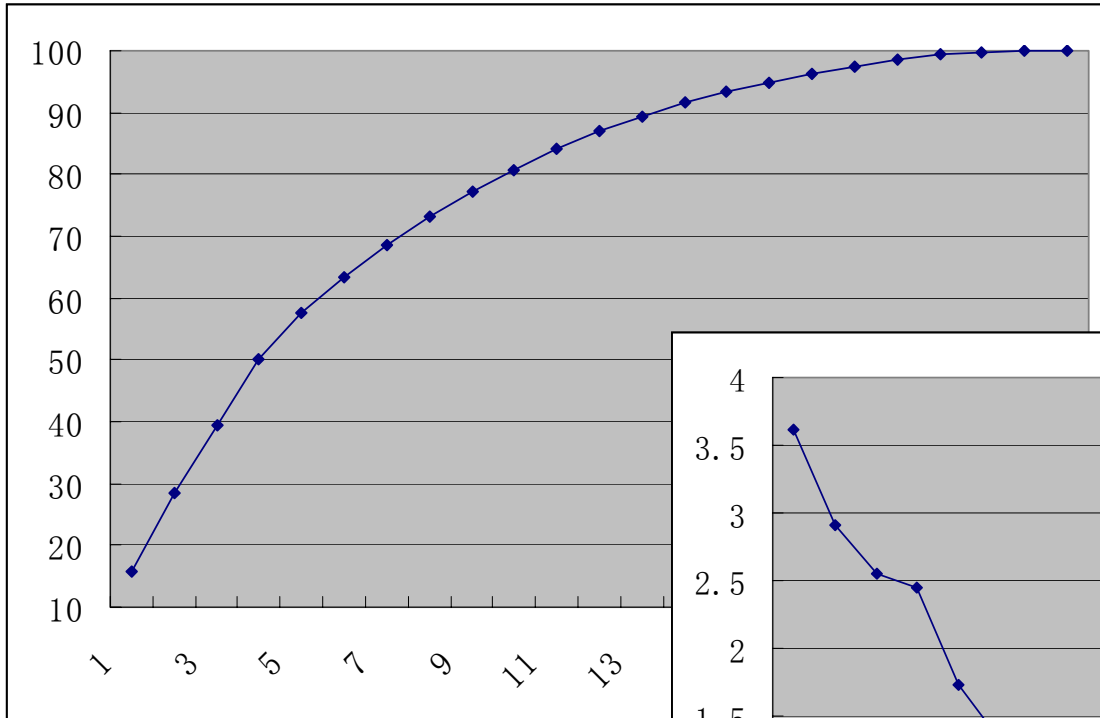


方法與結果

 這是將縱軸變換為對數坐標軸之后的圖像














一點小小的討論

-  不同類型的用戶，編輯行為會有所差異
-  可以通過挖掘這些數據獲得關於用戶的更多信息
-  辯論者和機器人是編輯行為最具有特點的群體，破壞者有不同的類型
-  數據分析的結果顯示管理員並非一個具有共同特征的群體
-  用主成分分析方法處理的結果並不理想，可以嘗試其他多元統計分析方法



謝謝

