

リコメンテーションアルゴリズムの研究

研究対象

オンラインDVDレンタルサービス会社のNetflix社は、膨大な情報の中からユーザの嗜好に合った情報を見つけ出し、その情報をユーザに能動的に提供する推薦技術の研究を行っている。また、自社の推薦システムの予測精度を高めるために、自社の顧客の嗜好データを公開し、多くの研究者と協調しながら研究を進めている。提供されたデータを研究対象とする。

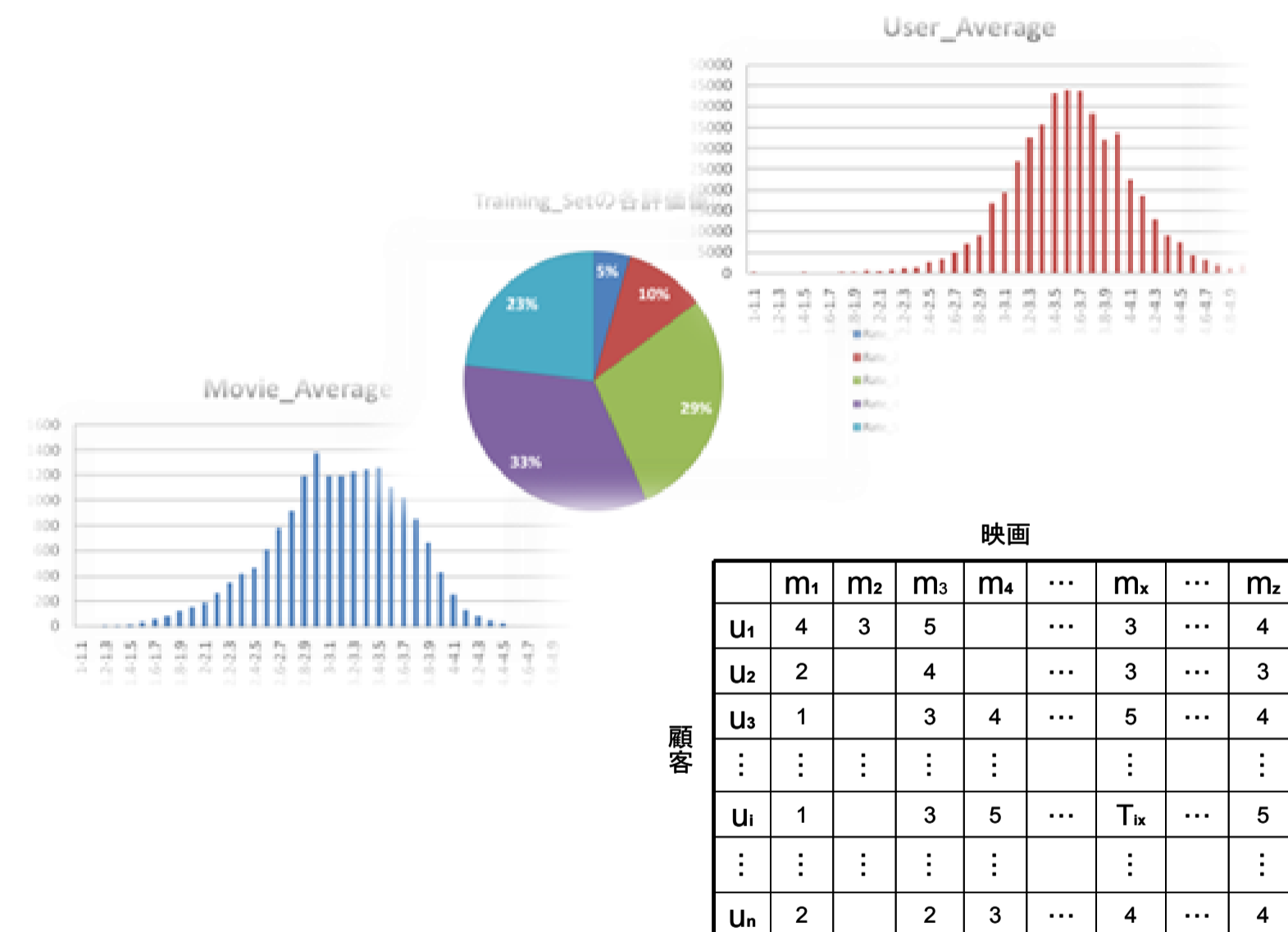


© 1997-2006 Netflix, Inc. All rights reserved. <http://www.netflixprize.com/>

予測手法

データ

Netflixデータは478,615人の顧客の17,170本の映画に対する5段階評価データである。評価されている項目は約1億件である。このデータを分析し、特徴をつかむ。



協調フィルタリング

協調フィルタリング(Collaborative Filtering)とは、多くの顧客の嗜好情報を記録し、ある顧客と類似した嗜好を持つ他の顧客の嗜好情報から、当該の顧客の嗜好を推測する方法論である。

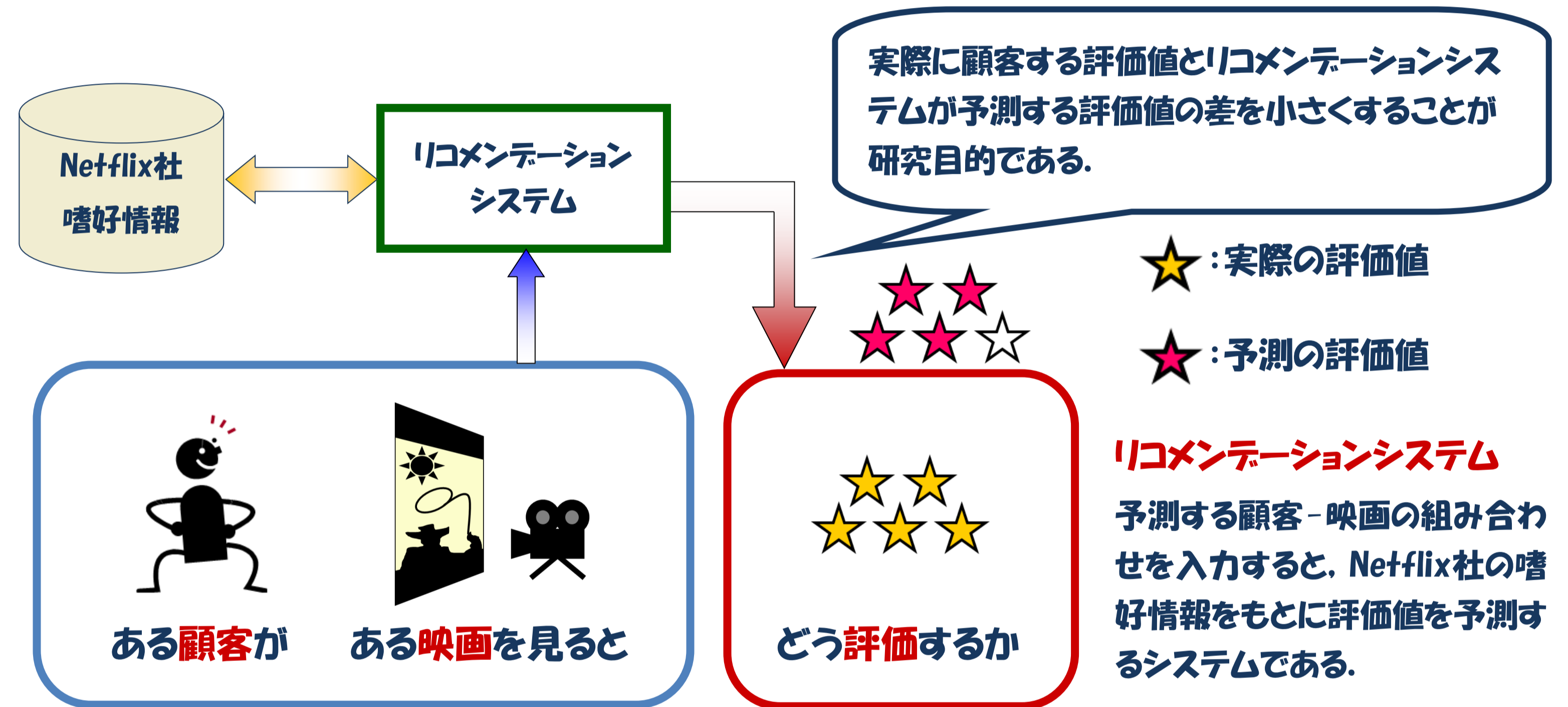


予測対象顧客 u_i が見た映画の平均評価値 T_i を求め、参考顧客 u_j の平均評価 T_j も求める。次に、参考顧客 u_j との相関係数 C_{ij} を以下の式で求める。

$$C_{ij} = \frac{\sum_k (T_{ik} - \bar{T}_i)(T_{jk} - \bar{T}_j)}{\sqrt{\sum_k (T_{ik} - \bar{T}_i)^2 \sum_k (T_{jk} - \bar{T}_j)^2}}$$

研究目的

Netflix社に蓄積された顧客の嗜好情報をもとに、ある顧客がある映画を見るとどう評価するのかを予測する。この予測の精度を高めることが研究目的である。



C_{ij} を用いて未知の評価値 T_{ix} を以下の式から求める。

$$T_{ix} = \bar{T}_i + \frac{\sum_j C_{ij} (T_{jx} - \bar{T}_j)}{\sum_j |C_{ij}|}$$

リコメンテーションアルゴリズムの予測精度の評価

リコメンテーションアルゴリズムの予測精度の評価はRMSEを用いて行う。予測の精度をあらわす指標である。値が0に近いほど正確で、精度の高い予測である。 P_i を予測値、 T_i を正解値、 n を予測回数とする。

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (P_i - T_i)^2}$$

