

ノイズあり属性統計量からの学習

Learning from Attribute Statistics with Noise

吉川 友也 *1

Yuya Yoshikawa

*1 千葉工業大学 人工知能・ソフトウェア技術研究センター

Software Technology and Artificial Intelligence Research Laboratory, Chiba Institute of Technology

Estimating the attributes of SNS users is an important technique for marketing and advertisement delivery systems. The attributes are typically estimated in a supervised learning manner. However, labeling the users with their attributes manually is difficult and expensive. In this paper, we address the problem of learning a classifier that predicts users' attributes only from the attribute statistics of websites that the users have visited. To tackle this problem, we propose a probabilistic generative model for the attribute statistics, in which can capture the users' attributes as hidden variables. Moreover, since the attribute statistics of the SNS and the websites may be different, the proposed model is modeled so as to treat the difference as noise variables. In the experiment, we show that the proposed model outperforms the existing methods on synthetic datasets.

1. はじめに

個人の属性推定は、性別、年齢、職業、人種、居住地等その人を表す特徴量から推定するタスクである。SNS等のインターネットサービスでは、個人の属性情報を使い、マーケティングや広告配信システムを作成する。しかし、そのような属性情報が得られないケースが多いため、属性推定技術の開発が必要である。

機械学習による典型的な属性推定の方法は、正解となる属性情報を教師データとして作成し、そのデータから属性識別器を学習する方法である。しかし、この方法では、個人に属性ラベルを手で付ける作業が発生するが、この作業は簡単ではない。例えば、SNS上のユーザに対して属性ラベルを付けることを考える。性別は、ユーザの顔写真や発言内容等から比較的簡単にラベル付けできる。しかし、年齢は30代と40代の違いを正しくラベル付けすることは難しい。また、属性識別器の性能を上げるためには、より多くの属性ラベル付きデータが必要であるが、人手での作業のためコストがかかる。

図1は本論文で想定するシチュエーションを示す。今、SNS上の個人の属性を推定したいとし、個人を表す特徴量とどの企業アカウントにリンクを張っているか(フォローしているか)が既知であるとする。また、企業アカウントはWebサイトを持っており、そのWebサイトに訪れるユーザの属性統計量(例えば、訪問ユーザの性別比が男性60%、女性40%)は、Quantcast*1等によって提供されるとする。

このシチュエーションの下で、本稿では、個人の属性情報を教師データとして使わずに、個人の属性を推定する手法を提案する。観測値として、人手でラベル付けしたユーザ属性情報の代わりに、そのユーザが訪れたWebサイトとそのWebサイトの属性統計量が与えられる。提案法では、個人の特徴量から属性が生成され、Webサイトの属性統計量はそのWebサイトに訪問したことがある個人の属性の統計量として生成されると仮定する。ただし、Webサイトの属性統計量と、我々が知りたい個人の属性の統計量は異なることが想定される。すなわち、

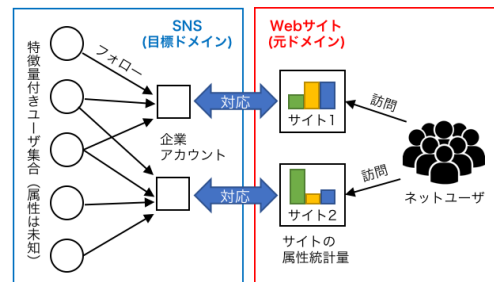


図1: 想定するシチュエーション

Webサイトに訪れる集団と、属性を知りたい個人の集団は異なるため、属性統計量にズレが生じる。提案法では、このズレを考慮するようにモデル化を行う。実験では、提案法が図1のシチュエーションにおいて有効であることを人工データを用いて示す。

2. 関連研究

これまで、インターネット上の個人の属性推定は多数行われてきた。推定する属性は、職業 [Li 14], 性別・年齢 [Rao 10], 民族 [Mislove 11], 居住地 [Cheng 10] 等がある。個人を表す特徴としては、SNS上の発言 [Rao 10], 顔画像 [Kumar 09], 友人関係 [Dong 14] 等が用いられる。提案法は、特定の属性や特徴に特化しない汎用的なモデルであり、このモデルを拡張することにより、様々な属性や特徴に特化したモデルを作ることができる。また、これらの研究では、個人に対して属性ラベルが付けられ、それが教師データとして与えられると仮定する。提案法は、そのような作成にコストの掛かる教師データは使わずに、属性分類器を構築することができる。

提案法は、[Culotta 15]で提案された属性推定のアイデアから着想を得ている。彼らの論文では、図1と同等のシチュエーションを想定し、個人の特徴からWebサイトの属性統計量を回帰するモデルに基づき、属性分類器を構築する。そのようなやり方では、企業アカウントをフォローする個人の属性統計量

連絡先: 吉川 友也, yoshikawa@stair.center

<https://sites.google.com/site/yuyay4ml/>*1 <https://www.quantcast.com/>

と、その企業の Web サイトを訪問する個人の属性統計量が一致する前提で、属性分類器が構築される。提案法は、SNS と Web サイトのドメイン間の属性統計量の違いが考慮できるようにモデル化している点異なる。

3. 提案法

3.1 モデル

提案法を一般的に説明するために、属性を予測したい個人がいる世界を「目標ドメイン」、サイトの属性統計量が観測できる世界を「元ドメイン」と呼ぶ。また、目標ドメインと元ドメインのサイトは一対一対応が取れているとする。

観測値として、目標ドメインにおける U 人の特徴量 $X = [x_1, x_2, \dots, x_U]^T \in \mathbb{R}^{U \times D}$ と、元ドメインにおけるサイトの属性統計量 $S = [s_1, s_2, \dots, s_N]^T \in \mathbb{R}_+^{N \times M}$ が与えられる。ここで、 D は特徴量の次元数、 M は属性のクラス数を表す。また、各サイトの属性統計量の総和は必ず 1 となる。すなわち、 $\sum_{m=1}^M s_{im} = 1$ である。 E_i は、目標ドメインにおけるサイト i を訪問したことがあるユーザ集合を表す。簡単のため、 $E = \{E_1, E_2, \dots, E_N\}$ とする。

提案モデルでは、各個人 $u = 1, 2, \dots, U$ の属性 $y_u \in \{1, 2, \dots, M\}$ の確率は、その個人の特徴量 x_u と各属性の重みベクトル $w_m \in \mathbb{R}^D$ の内積 $w_m^T x_u$ を入力とするソフトマックス関数によって与えるとする。これは多クラスロジスティック回帰と同じアイデアであるが、提案法においては、 y_u は潜在変数である。また、重みベクトル w_m は等分散ガウス分布 $\mathcal{N}(0, \alpha_w^{-1} I_D)$ から生成される。ここで、 I_D は D 次元の単位行列である。

目標ドメインと元ドメインのサイトでは、ユーザの属性分布が異なると仮定する。この違いを表現するために、混同行列 $C = [c_{11}, c_{12}, \dots, c_{M1}, c_{M2}, \dots, c_{MM}]^T$ を導入する。 $c_m = [c_{m1}, c_{m2}, \dots, c_{mM}]$ は、目標ドメインにおいて個人の属性が m の時、元ドメインではどの属性に割り当てられるかの確率を表す。従って、個人 u の属性ラベルが y_u のとき、元ドメインのサイト i における属性 t_{ui} はカテゴリカル分布 $\text{Cat}(c_{y_u})$ から生成される。なお、各属性 m について、 $\sum_{l=1}^M c_{ml} = 1$ となる必要があるため、 c_m はディリクレ分布から生成することが適当である。ただし、多くの場合、目標ドメインと元ドメインの間で属性は変化しないと思われるため、混同行列 C の対角成分の値は比較的大きくなりやすいと考えられる。この事前知識を導入するため、各属性 m ごとにディリクレパラメータ $\beta_m \in \mathbb{R}_+^M$ を超パラメータとして与える。ここで、

$$\beta_{ml} = \begin{cases} \alpha_{c_0} & (m \neq l) \\ \alpha_{c_1} & (m = l). \end{cases} \quad (1)$$

$\alpha_{c_1} > \alpha_{c_0}$ とすることにより、混同行列 C の対角成分は大きな値になりやすい。

t_i を個人 $u \in E_i$ の t_{ui} の統計量とする。すなわち、 $t_i = \frac{1}{|E_i|} \sum_{u \in E_i} v(t_{ui})$ とする。ここで、 $v(t)$ は t 次元の値は 1 でそれ以外は 0 の M 次元ベクトルである。最終的に、元ドメインのサイト i の属性統計量 s_i は、 t_i を平均とする等分散ガウス分布 $s_i \sim \mathcal{N}(t_i, \alpha_s^{-1} I_M)$ から生成される。

まとめると、提案モデルの生成過程は以下ようになる。

1. 各クラス $m = 1, 2, \dots, M$ に対して

(a) 重み行列を生成: $w_m \sim \mathcal{N}(0, \alpha_w^{-1} I_D)$

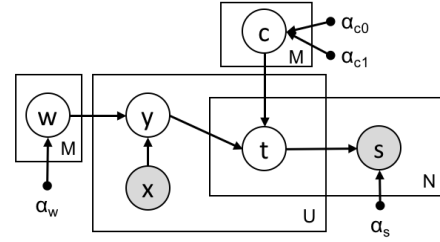


図 2: 提案法のグラフィカルモデル

(b) 混同行列を生成: $c_m \sim \text{Dir}(\beta_m)$

2. 各ユーザ $u = 1, 2, \dots, U$ に対して

(a) 属性ラベルを生成: $y_u \sim \text{Softmax}([w_m^T x_u]_{m=1}^M)$

(b) 各サイト $i = 1, 2, \dots, N$ に対して

i. サイト依存属性ラベルを生成: $t_{ui} \sim \text{Cat}(c_{y_u})$

3. 各サイト $i = 1, 2, \dots, N$ に対して

(a) 属性統計量を生成: $s_i \sim \mathcal{N}(t_i, \alpha_s^{-1} I_M)$.

ここで、 $t_i = \frac{1}{|E_i|} \sum_{u \in E_i} v(t_{ui})$.

図 2 は、提案モデルのグラフィカルモデルを示す。白ノードは潜在変数、灰色ノードは観測値、小点は超パラメータを表す。超パラメータの集合は $\alpha = \{\alpha_w, \alpha_{c_0}, \alpha_{c_1}, \alpha_s\}$ と表記する。

3.2 変分ベイズによる推論

図 2 から分かるように、提案モデルは階層ベイズモデルの一種であるため、変分ベイズ等の近似推論法や MCMC を使うことで効率的な推論が可能になる。本稿では、変分ベイズによる推論法を説明する。

まず、 T, Y を周辺化したパラメータ W, C の対数周辺事後確率は、

$$\begin{aligned} \log p(W, C | X, S, \alpha) & \\ \propto \log \sum_{T, Y} p(S | T, \alpha_s) p(T | Y, C) p(Y | W, X) & \\ + \log p(W | \alpha_w) + \log p(C | \alpha_{c_0}, \alpha_{c_1}) & \end{aligned} \quad (2)$$

と書ける。提案モデルの生成過程より、

$$p(S | T, \alpha_s) = \prod_{i=1}^N \mathcal{N}(s_i | t_i, \alpha_s^{-1} I_M), \quad (3)$$

$$p(T | Y, C) = \prod_{i=1}^N \prod_{u \in E_i} c_{y_u, t_{ui}}, \quad (4)$$

$$p(Y | W, X) = \prod_{u=1}^U \frac{\exp(w_{y_u}^T x_u)}{\sum_{m=1}^M \exp(w_m^T x_u)}, \quad (5)$$

$$p(W | \alpha_w) = \prod_{m=1}^M \mathcal{N}(w_m | 0, \alpha_w^{-1} I_D), \quad (6)$$

$$p(C | \alpha_{c_0}, \alpha_{c_1}) = \prod_{m=1}^M \frac{\Gamma(\sum_{l=1}^M \beta_{ml})}{\prod_{l=1}^M \Gamma(\beta_{ml})} \prod_{l=1}^M c_{ml}^{\beta_{ml}-1}, \quad (7)$$

となる．式 (2) が最大となる W, C を求めることが目的であるが， T, Y のとり得る値の全ての組み合わせを考慮し計算しなければならず，計算量的に困難である．従って，イエンセンの不等式により，式 (2) に対する変分下限 $\mathcal{L}(\Theta)$ を以下のように導出し，これを目的関数として最適化する．

$$\begin{aligned} & \log p(W, C | X, S, \alpha) \\ & \geq \sum_{T, Y} q(T, Y) \log \frac{p(S|T, \alpha_s) p(T|Y, C) p(Y|W, X)}{q(T, Y)} \\ & \quad + \log p(W | \alpha_w) + \log p(C | \alpha_{c_0}, \alpha_{c_1}) \\ & = \mathcal{L}(\Theta). \end{aligned} \quad (8)$$

変分分布 $q(T, Y)$ は， $q(T, Y) = q(Y|\zeta)q(T|\eta)$ の因子分解を仮定する．ここで，

$$q(Y|\zeta) = \prod_{u=1}^U \zeta_{uy_u}, \quad q(T|\eta) = \prod_{i=1}^N \prod_{u=1}^U \eta_{iut_{ui}}, \quad (9)$$

である．また， $\Theta = \{W, C, \zeta, \eta\}$ は最適化するパラメータ集合を表す．

次に， W, C, ζ, η の更新則を導出する．

W の更新: W は閉形式で更新できないため，以下のように勾配を計算し，準ニュートン法等の最適化法に使い，更新する．

$$\frac{\partial \mathcal{L}(\Theta)}{\partial w_m} = \sum_{u=1}^U \left(\zeta_{um} - \frac{\exp(w_m^\top x_u)}{\sum_{l=1}^M \exp(w_l^\top x_u)} \sum_{l=1}^M \zeta_{ul} \right) x_u - \alpha_w w_m. \quad (10)$$

C の更新: $\sum_{l=1}^M c_{ml} = 1$ の制約があるので，ラグランジュの未定乗数法によって，以下の更新則を導出する．

$$c_{ml} = \frac{\sum_{i=1}^N \sum_{u=1}^U \zeta_{um} \eta_{ium} + \beta_{ml} - 1}{\sum_{m'=1}^M \sum_{i=1}^N \sum_{u=1}^U \zeta_{um'} \eta_{ium'} + \sum_{l'=1}^M \beta_{ml'} - M}. \quad (11)$$

ζ の更新: $\sum_{m=1}^M \zeta_{um} = 1$ の制約があるので，ラグランジュの未定乗数法によって，以下の更新則を導出する．

$$\zeta_{um} \propto \exp \left\{ \sum_{i=1}^N \sum_{l=1}^M \eta_{iul} \log c_{ml} + a_{um} - \log \sum_{m'=1}^M \exp(a_{um'}) \right\}. \quad (12)$$

ここで， $a_{um} = x_u^\top w_m$ である．式 (12) の値を計算した後， $\sum_{m=1}^M \zeta_{um} = 1$ となるように，値を正規化する．

η の更新: $\sum_{m=1}^M \eta_{ium} = 1$ の制約があるので，ラグランジュ未定乗数法で更新する．ただし，閉形式で更新はできないため， η_{ium} と λ_{iu} の勾配を以下のように計算し，準ニュートン法等により交互に更新する．

$$\begin{aligned} \frac{\partial \mathcal{L}(\Theta)}{\partial \eta_{ium}} &= -\frac{\alpha_s}{|E_i|} (s_{im} - \mathbb{E}[t_i]_m) + \sum_{l=1}^M \zeta_{ul} \log c_{lm} \\ &\quad - \log \eta_{ium} - 1 + \lambda_{iu}, \end{aligned} \quad (13)$$

$$\frac{\partial \mathcal{L}(\Theta)}{\partial \lambda_{iu}} = \sum_{m=1}^M \eta_{ium} - 1. \quad (14)$$

ここで， $\mathbb{E}[t_i]_m = \frac{1}{|E_i|} \sum_{u \in E_i} \eta_{ium}$ である．

上記のようにパラメータを逐次的に更新し，式 (8) の値が収束したら学習を終了する．超パラメータ α は，クロスバリデーションによって適切な値を求める．

表 1: 人工データによる属性推定精度

	$\alpha_{c_1} = 1$	$\alpha_{c_1} = 10$	$\alpha_{c_1} = 100$
提案法	0.43	0.52	0.45
MTEN [Culotta 15]	0.32	0.51	0.31
Ridge [Culotta 15]	0.24	0.48	0.24

4. 実験

図 1 のシチュエーションにおいて，提案法が有効であることを確認するために，人工データを使った実験を行った．

今回は $M = 4$ 種類の属性の分類問題を解く設定とする．人工データは，3.1 節のモデルの生成過程で説明した手順に従う．ただし，

1. 個人の特徴量を平均 0，分散 1 の $D = 200$ 次元ガウス分布から生成する．
2. サイトの数を $N = 1,000$ ，個人数を $U = 100$ とする．各サイトはランダムに選ばれた個人 30 人からフォローされる．
3. 式 (1) に従って β を設定する．ここで， $\alpha_{c_0} = 1$ ， $\alpha_{c_1} \in \{1, 10, 100\}$ とする．

比較手法として，[Culotta 15] で提案された属性統計量を回帰するモデルに基づいて個人の属性を当てる方法を用いる．この方法では，まず，個人の特徴量からその個人と関係のあるサイトの属性統計量を回帰するモデルを学習する．ここで，属性統計量は一般に多変量になるため，[Culotta 15] は各変量間の関係を捉えられる Multi-Task Elastic Net (MTEN) と，変量間で独立な回帰モデルを構築する Ridge 回帰のケースを考えた．予測の際は，この回帰モデルによって属性統計量を予測し，その値が最も大きい属性をその個人の属性として出力する．

表 1 は， α_{c_1} を変えて人工データを生成した際の，個人の属性推定の精度を示す． α_{c_1} は小さい時，元ドメインと目標ドメインの属性統計量の差が大きくなり，元ドメインの属性統計量の値を回帰するようにモデル化した既存手法では，ランダムに予測した場合とほとんど変わらない精度となった．一方，提案法は，その差異を混同行列 C で捉えることができるため，比較的良い精度で属性を推定可能であった．また， α_{c_1} が大きくなった場合でも，提案法の予測精度が最も良かった．

5. おわりに

本稿では，個人の属性情報を教師データとして使わずに，個人の属性を推定する手法を提案した．具体的には，提案法は，別のドメインから比較的簡単に得られる属性統計量に基づいて，個人の属性情報を推定する．その際，提案法は，我々が属性を知りたい個人の存在するドメインと属性統計量のドメインの間に発生する属性統計量の違いを考慮する．

今後の研究では，実データでの実験を行い，提案法の有効性を検証していく．

参考文献

- [Cheng 10] Cheng, Z., Caverlee, J., and Lee, K.: You Are Where You Tweet: A Content-Based Approach to Geo-locating Twitter Users, in *Proceedings of the 19th ACM International Conference on Information and Knowledge Management*, pp. 759–768 (2010)

-
- [Culotta 15] Culotta, A., Ravi, N. K., and Cutler, J.: Predicting the Demographics of Twitter Users from Website Traffic Data, in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pp. 72–78 (2015)
- [Dong 14] Dong, Y., Yang, Y., Tang, J., Yang, Y., and Chawla, N. V.: Inferring User Demographics and Social Strategies in Mobile Social Networks, *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 15–24 (2014)
- [Kumar 09] Kumar, N., Berg, A. C., Belhumeur, P. N., and Nayar, S. K.: Attribute and Simile Classifiers for Face Verification, in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 365–372 (2009)
- [Li 14] Li, J., Ritter, A., and Hovy, E.: Weakly Supervised User Profile Extraction from Twitter, *ACL*, pp. 165–174 (2014)
- [Mislove 11] Mislove, A., Lehmann, S., Ahn, Y.-y., Onnela, J.-p., and Rosenquist, J. N.: Understanding the Demographics of Twitter Users, in *Fifth International AAAI Conference on Weblogs and Social Media*, pp. 554–557 (2011)
- [Rao 10] Rao, D., Yarowsky, D., Shreevats, A., and Gupta, M.: Classifying Latent User Attributes in Twitter, in *Proceedings of the 2nd international workshop on Search and mining user-generated contents*, p. 37 (2010)