データマイニング

- テキストマイニング
- 自然言語処理
- 機械学習
- ■情報検索

ナイーブ・ベイズ

ナイーブ・ベイズ学習とは...

■データをクラスに分類するための一手法

- 大量にあるデータをクラスに分類する
- 文書分類:大量にある新聞記事を、記事の内容に対応したクラス に分類する(例:「経済」「スポーツ」「娯楽」等)
- メールフィルタリング:受信したEメールを「迷惑メール」と 「普通メール」に分類する
- データにクラスに対応するラベルを付与することと同義
- 今回は、1つのデータに対して、1つのクラスに対応するラベル を付与

クラス分類の例

■ Yahooのトップページ



ナイーブ・ベイズ学習の例

■ 学習用データ

| Section State Publisher | | AMERICAN CONTRACTOR | | |
|-------------------------|-------------|---------------------|---|--------|
| 天気 | 温度 | 湿度 | 風 | ゴルフプレイ |
| 晴 | 暑 | 高 | 無 | × |
| 晴 | 暑 | 高 | 有 | × |
| 晴 曇 | 暑 | 高高高 | 無 | 0 |
| 雨 | 暖 | 高 | 無 | 0 |
| 雨 | 暖 涼 涼 | 普通 | 無 | 0 |
| 雨 | 涼 | 普通 | 有 | × |
| 雨 曇 晴 | 涼 | 普通 | 有 | 0 |
| 晴 | 暖 | 高 | 無 | × |
| 晴 | 暖 涼 | 普通 | 無 | 0 |
| 雨 | 暖 | 普通 | 無 | 0 |
| 晴 | 暖 | 普通 | 有 | 0 |
| 曇 | | 高 | 有 | 0 |
| 雲 | 暖 暑 | 普通 | 無 | 0 |
| 雨 | 暖 | 高 | 有 | × |

■ テスト用データ(新規事例)

| 000 | 天気 | 温度 | 湿度 | 風 | ゴルフプレイ |
|-----|----|----|----|---|--------|
| Š | 晴 | 涼 | 高 | 有 | ? |



新規事例データ→ 晴 涼 高 有

○ 尤度:0.00529100529100527× 尤度:0.0205714285714286

判定:×

素性

「天気」「温度」「湿度」「風」

素性値

┃「晴」「暑」「高」「無」

ナイーブ・ベイズ

■ クラス(ゴルフプレイ=O, ゴルフプレイ=×)に分類するにおいて、素性(天気、温度、湿度、風)が互いに独立であると仮定して学習を行う手法



素性が互いに独立と仮定(ナイーブ): 現実世界ではあまり成立しないが、それなりにうまくいく

素性「天気」「温度」「湿度」「風」

素性値「晴」「暑」「高」「無」

ナイーブ・ベイズ学習式

■ ナイーブ・ベイズ学習式

$$\hat{c} = \underset{c = \{\bigcirc, \times\}}{\operatorname{arg\,max}} P(c \mid E)$$

c: クラス(ゴルフプレイ=O、×)

E: 新規事例データの素性値列(例:天気=晴, 温度=暑,

湿度=高,風=無)

ナイーブ・ベイズ学習式の変形

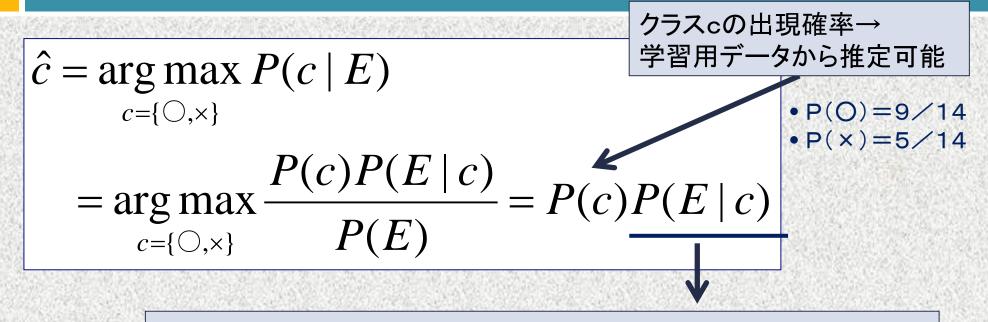
■ 学習データから計算できるようにナイーブ・ベイズ学習式を変形 ※ベイズの定理に変形(証明は黒板に板書)

$$\hat{c} = \underset{c=\{\bigcirc,\times\}}{\operatorname{arg max}} P(c \mid E)$$

$$= \underset{c=\{\bigcirc,\times\}}{\operatorname{arg max}} \frac{P(c)P(E \mid c)}{P(E)} = P(c)P(E \mid c)$$

P(E)はcに依らないので考えなくてよい

ナイーブ・ベイズ学習式の変形



クラスcが与えられたときの素性値列Eの出現確率は、 様々な事例データが存在するため、推定することは困難

c: クラス(ゴルフプレイ=〇、×)

E: 事例データの素性値列(例:天気=晴,温度=暑,湿度=高,風=無)

ナイーブ・ベイズ学習式の変形

■ P(E | c)を計算

$$P(E \mid c) = P(\{e_1, e_2, e_3, e_4\} \mid c)$$

ここで、Eに含まれる各素性値e1,e2,e3,e4が互いに独立に 生起すると仮定すると...

$$|P(E \mid c) = P(\{e_1, e_2, e_3, e_4\} \mid c)$$

$$\approx P(e_1 \mid c)P(e_2 \mid c)P(e_3 \mid c)P(e_4 \mid c) = \prod_{i=1}^4 P(e_i \mid c)$$

ナイーブ・ベイズ学習式のパラメータ推定

■ P(e1|c), P(e2|c), P(e3|c), P(e4|c) を学習用データから推定

$$P(E \mid c) \approx P(e_1 \mid c)P(e_2 \mid c)P(e_3 \mid c)P(e_4 \mid c) = \prod_{i=1}^4 P(e_i \mid c)$$

| 天気 | 温度 | 湿度 高 高 | 風 | ゴルフプレイ |
|--|-------|--------------|---|--------|
| 天気 曇 雨 曇 雨 曇 晴 | 温度暑 | 高 | 無 | 0 |
| 雨 | | 高 | 無 | 0 |
| 雨 | 暖涼涼涼暖 | 普通 | 無 | 0 |
| 曇 | 涼 | 普通 | 有 | 0 |
| 晴 | 涼 | 普通 | 無 | 0 |
| 雨 | 暖 | 普通 | 無 | 0 |
| 晴 | 暖 | 普通 | 有 | 0 |
| 晴 <u>雲</u> 曇 | 暖 | 高 | 有 | 0 |
| 曇 | 暑 | 普通 | 無 | 0 |

P(E | ゴルフプレイ=O)

=P(天気=晴れ | O)×P(温度=涼 | O)

×P(湿度=高 | O)×P(風=有 | O)

 $=(2/9)\times(3/9)\times(3/9)\times(3/9)$

ナイーブ・ベイズ学習式(最終)

■ 学習データからパラメータを推定できるようにナイーブ・ベイズ 学習式を変形

$$\hat{c} = \underset{c = \{\bigcirc, \times\}}{\operatorname{arg max}} P(c \mid E)$$

$$= \underset{c = \{\bigcirc, \times\}}{\operatorname{arg max}} \frac{P(c)P(E \mid c)}{P(E)} \approx P(c) \prod_{i=1}^{4} \frac{P(e_i \mid c)}{P(e_i \mid c)}$$

学習用データから推定可能

課題

■ サンプルデータ(各自、ダウンロードせよ)の学習用データを訓練データとして、テスト用データをナイーブベイズ手法で分類するプログラムを実装せよ。素性として名詞を用いよ。

訓練データの1記事(全体で227記事)

<id> 2012-01-24_1 </id>

<date> 2012/01/24 </date>

<company> エレコム </company>

<class> IT 周辺機器 </class>

<title> エレコム、名刺用紙「なっとく。名刺」シリーズからスーパーファイン用紙など2タイプを発売 </title> しっかりとした厚みで上品に仕上がる「特厚」タイプ!

切り口がすっきりきれいなクリアカットタイプの名刺用紙「なっとく。名刺」の新シリーズを発売

エレコム株式会社(本社:大阪市中央区、取締役社長:葉田順治)は、手持ちのプリンタなどでオリジナルの名刺やメッセージカードが簡単に作成できる名刺用紙「なっとく。名刺」シリーズについて、新たに2タイプ4種類を2月上旬より新発売いたします。 ご家庭のプリンタを使って、簡単・手軽にオリジナルの名刺が作成できる「なっとく。名刺」は、人気の名刺用紙シリーズです。

課題

■ サンプルデータ(各自、ダウンロードせよ)の学習用データを訓練データとして、テスト用データをナイーブベイズ手法で分類するプログラムを実装せよ。素性として名詞を用いよ。

・テスト用データの1記事(全体で12記事)

<id> 2012-04-25_0 </id>

<date> 2012/04/25 </date>

<company> グーグル </company>

推定せよ

<class> </class>

<title>グーグル、全てのファイルを安心して保存・共有できるGoogleドライブを発表 </title> Googleドライブを使って、何でも保存、共有しよう。

Posted by デービッド ウォルツ/Googleドライブ・プロダクトマネージャー

本日、全てのファイルを安心して保存、共有できるGoogleドライブを発表しました。ドライブでは写真、動画、Googleドキュメント、PDFなどさまざまなファイルを一つの場所に保管できるようになります。例えば、顧客とのビジネスプランの相談や、来年度の予算計画を同僚と立てる場合も、共同作業が効率よくでき、プロジェクトを迅速に進められるでしょう。Googleドキュメントが進化したGoogleドライブは、クラウドでの働き方を大きく変えます。

ナイーブ・ベイズ学習式

■ 学習データからパラメータを推定できるようにナイーブ・ベイズ 学習式を変形

$$\hat{c} = \arg\max_{c = \{\bigcirc, \times\}} P(c \mid E)$$

$$= \arg\max_{c = \{\bigcirc, \times\}} \frac{P(c)P(E \mid c)}{P(E)} \approx P(c) \prod_{i=1}^{n} P(e_i \mid c)$$

$$P(c) = \frac{fa(c, S)}{|S|}$$

$$P(e_i | c) = \frac{f(e_i, c)}{f(c)}$$

- fa(c,S): テストデータ記事集合 S中において、クラスcが出現 する記事の頻度
- f(e,c): クラスcの記事集合において、素性の名詞e の出現数
- f(c): クラスcの記事集合において、素性の名詞の総出現数

ナイーブ・ベイズ学習式の実装

■ 確率の積和は、ただでさえ小さい値の積なので、実装においては 問題が生じる→Oに極めて近い値となり、最終的にOになってしまう

$$\hat{c} = \arg\max_{c = \{\bigcirc, \times\}} P(c \mid E)$$

$$= \arg\max_{c = \{\bigcirc, \times\}} \frac{P(c)P(E \mid c)}{P(E)} \approx P(c) \prod_{i=1}^{n} P(e_i \mid c)$$

■ 実装上の解決方法として、確率値のlogをとって、積和を和に変換する → 小さい値も負の整数になって好都合

ナイーブ・ベイズ学習式(最終)

■ 実装上の解決方法として、確率値のlogをとって、積和を和に変換する → 小さい値も負の整数になって好都合

$$\hat{c} = \underset{c=\{\bigcirc,\times\}}{\operatorname{arg max}} P(c \mid E)$$

$$= \underset{c=\{\bigcirc,\times\}}{\operatorname{arg max}} \frac{P(c)P(E \mid c)}{P(E)} \approx P(c) \prod_{i=1}^{n} P(e_i \mid c)$$

$$\hat{c} = \underset{c = \{\bigcirc, \times\}}{\operatorname{arg max}} P(c \mid E)$$

$$= \underset{c = \{\bigcirc, \times\}}{\operatorname{arg max}} \frac{P(c)P(E \mid c)}{P(E)} \approx \log(P(c)) + \sum_{i=1}^{n} \log(P(e_i \mid c))$$