

# 導入

通信路モデル

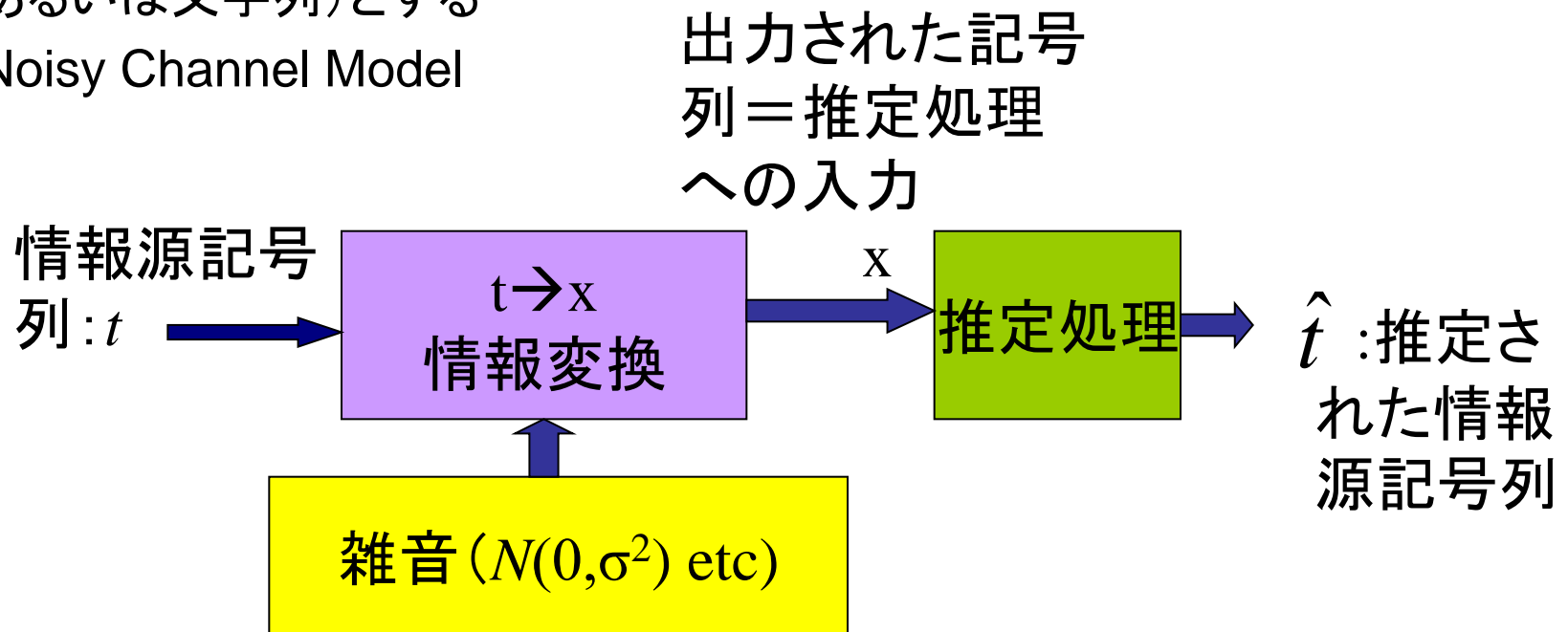
Bayes統計

最尤推定とMAP推定

データの性質

## 機械学習の先史時代 --情報の変換過程のモデル化--

- 情報源を記号列(例えば単語列あるいは文字列)とする
- Noisy Channel Model



出力された記号列＝推定処理への入力データ $x$ から  
情報源記号列 $t$ を推定し  $\hat{t}$  を計算する

# Bayes統計の意義

## ➤ Bayesの定理

$$P(t | x) = \frac{P(x | t)P(t)}{P(x)}$$


- $P(t/x)$ は、新たな出力記号列 $x$ が得られたときの情報源から出力された記号列  $t$  を推定する式で、これを最大化する  $t$  すなわち  $\hat{t} = \arg \max_t P(t | x)$  を求めるのが目標。
- ところが、このままでは、既に得られている情報を使えないので、Bayesの定理で変換する。
- すると、既知の情報源状態と出力記号列のペアに関する条件付き確率 $P(x|t)$ （＝教師データ）
- 情報源についての事前知識 $P(t)$ が使える形になる。

# Bayes統計とは

1. 常にBayesの定理を用いる
2. 用いられる確率は主観確率(＝確信度)
3. 事前情報を利用する
4. 未知量(確率分布のパラメター)は確率的に変動
5. 観測されたデータは絶対的
6. 推測は常に条件付
7. アドホックな手続きを認めない

# Bayes統計を用いた情報変換過程のモデルにおける出力データからの情報源の推定方法

- 通信路を条件付確率でモデル化:  $P(x/t)$
- 目的は  $x$  が観測されたときの  $t$  の確率すなわち事後確率  $P(t/x)$  を最大化する情報源の確率。

$$\begin{aligned}\hat{t} &= \arg \max_t P(t | x) && \text{ここでベイズの定理により} \\ &= \arg \max_t P(x | t) P(t)\end{aligned}$$


- $P(t)$  は情報源記号列の既知の統計的性質が利用できる
- $P(x/t)$  は情報源記号列  $t$  が情報変換およびnoisy channelの雑音によって  $x$  毎に変化する確率。
- この確率は多数の  $\langle t, x \rangle$  対の観測データにより計算する

# 情報変換過程モデルの適用例

## ➤ 例：機械翻訳

元言語

$x$ : 私がリンゴを食べる

機械  
翻訳

翻訳先言語

$t$ : I eat an apple

- $P(t/x)$  は元言語のテキスト  $x$  (既知) が翻訳先言語のテキスト  $t$  に翻訳される確率
- $P(x/t)$  は  $t$  という翻訳結果に対する元言語のテキストが  $x$  である確率
- $P(t)$  翻訳先言語におけるテキスト  $t$  の自然さ。例えば、 $N$  単語列のコーパスにおける 単語3-gram 確率
- 以上の設定で下の式 は機械翻訳の出力

$$\hat{t} = \arg \max_t P(t | x) = \arg \max_t P(x | t) P(t)$$

- この考え方を元にしたのが現在主流となってきた統計的機械翻訳 (IBM で 1993 年に開発された)

# 日英機械翻訳の例

- $P(\text{リンゴを食べる} | \text{eats an apple}) = 0.3$
- $P(\text{リンゴを食べる} | \text{eats apples}) = 0.2$
- $P(\text{彼は} | \text{He}) = 1.0$
- $P(\text{He eats apples}) = 0.2$ 、 $P(\text{He eats an apple}) = 0.5$
- $P(\text{He eats an apple} | \text{彼はリンゴを食べる})$ 
  - $= 1.0 \times 0.3 \times 0.5 = 0.15$
- $P(\text{He eats apples} | \text{彼はリンゴを食べる})$ 
  - $= 1.0 \times 0.2 \times 0.2 = 0.04$
- “He eats an apple” のほうが良い英訳
- 事前知識として  $P(\text{He eats apple}) = 0.0$  があれば
- 非文  $P(\text{He eats apple}) = 0$  にできるのがベイズの強み

## ➤ 例：文書分類

- $P(t/x)$  において  $x$  が与えられた文書、 $t$  がカテゴリ

$$\text{推定されたカテゴリ} : \hat{t} = \arg \max_t P(x | t) P(t)$$

- $P(t)$  はカテゴリ  $t$  の文書の出現確率
- $P(x/t)$  はカテゴリ  $t$  において文書  $x$  が出現する確率
- このモデル化にはいろいろな方法があるが、簡単なのは、出現する単語  $w_1, \dots, w_N$
- $P(x/t) = P(w_1, \dots, w_N/t)$  だが、このままでは計算しにくいので  $w_1, \dots, w_N$  が独立だとすると

$$P(w_1, \dots, w_N | t) = \prod_{n=1}^N P(w_n | t)$$

Why?

- これを naïve Bayse 分類とよぶ。



## 文書分類の例: 長澤まさみ vs 上野樹里

- 「長澤まさみ」関連の文書に高い確率で出現する単語
  - 主演、映画、東宝、吉田礼、薬師丸ひろ子、サッカー、
- 「上野樹里」関連の文書に高い確率で出現する単語
  - 主演、のため、カンタービレ、ドラマ、ラスト、フジテレビ、
- 分類したい文書: Dの含む単語は
  - 主演、ラスト、フレンズ
  - $P(\text{主演} | \text{長澤}) = 0.1$ 、 $P(\text{主演} | \text{上野}) = 0.1$
  - $P(\text{ラスト} | \text{長澤}) = 0.2$ 、 $P(\text{ラスト} | \text{上野}) = 0.2$
  - $P(\text{フレンズ} | \text{長澤}) = 0.2$ 、 $P(\text{フレンズ} | \text{上野}) = 0.2$

Googleのヒット数から推定したところ、  
 $P(\text{長澤}) = 0.6$ 、 $P(\text{上野}) = 0.4$

$$\begin{aligned} \text{➤ } P(\text{長澤} \mid D) &= P(D \mid \text{長澤}) P(\text{長澤}) \\ &= P(\text{主演} \mid \text{長澤}) P(\text{ラスト} \mid \text{長澤}) P(\text{フレンズ} \mid \text{長澤}) P(\text{長澤}) \\ &= 0.1 \times 0.2 \times 0.2 \times 0.6 = 0.0024 \end{aligned}$$

$$\begin{aligned} \text{➤ } P(\text{上野} \mid D) &= P(D \mid \text{上野}) P(\text{上野}) \\ &= P(\text{主演} \mid \text{上野}) P(\text{ラスト} \mid \text{上野}) P(\text{フレンズ} \mid \text{上野}) P(\text{上野}) \\ &= 0.1 \times 0.2 \times 0.2 \times 0.4 = 0.0016 \end{aligned}$$

よって、文書Dは長澤に分類

しかし、Dに「カンタービレ」という単語も含まれ、  
 $P(\text{カン..} \mid \text{長澤}) = 0.1$ 、 $P(\text{カン..} \mid \text{上野}) = 0.8$ だと  
 $P(\text{長澤} \mid D) = 0.00024$ 、 $P(\text{上野} \mid D) = 0.00128$   
で文書Dは上野に分類。 **直感にあっているようだ！**

# 教師あり学習

- 上記の例では、情報源のモデルである $P(t)$ や $P(x|t)$ は単に出現確率だったが、ここで適切な確率分布を考えることが可能
- すると、その分布を決めるパラメータを推定する必要が出てくる。
- そのために $\langle t, x \rangle$ という情報源の状態と出力データの対データが多数入手できれば利用する。
- この $\langle t, x \rangle$ を教師データ(あるいは観測データ)と呼ぶ。
- すると、機械学習の中心となる教師あり学習は、

確率分布 $P(t)$ 、 $P(x/t)$ のパラメータを  
教師データ $\langle t, x \rangle$ を利用して求める

という問題になる。

# 教師なし学習

- 教師あり学習では教師データ $\langle t, x \rangle$ の集合が与えられた状態で、 $P(t)$ や $P(x|t)$ のパラメータを求めた。
- しかし、データ $\langle x \rangle$ の集合だけが与えられていて( $t$ は与えられていない)ときはどうする？

➤ データ $\langle x \rangle$ の集合から、 $P(x)$ のパラメータだけを求めることになる。

- 直観的には、データ $\langle x \rangle$ を類似したものにグループ化する
  - → クラスタリングと言い、グループのことをクラスタと呼ぶ。
- これを教師なし学習と呼ぶ。

# 識別モデルと生成モデル

- 入力データ $x$ に対応する予測値 $t$ を求める
- 識別モデル(discriminative model):  $p(t|x)$ を直接モデル化する。この $p(t|x)$ によって、未知の $x$ に対する $t$ を予測(あるいは推定)する方法
  - $t=f(x)$ となる関数を直接求めるものもあり。
- 生成モデル(generative model): ベイズの定理で  $p(t|x)$ を  $p(x|t)p(t)/p(x)$ に変換。 $p(x|t)$ を学習。 $p(t)$ を事前データから求める。これと既知の $\langle x, t \rangle$ のペアのデータから $p(x|t)$ のパラメータを更新。これによって、未知の $x$ に対する $t$ を求める $p(t|x)$ の確率分布をモデル化する。

事前  
分布

観測データが知られて後の $p(x|t)$ の事後分布

# 最尤推定とMAP推定

## ➤ 最尤推定

- 分布  $P(X | \theta)$  のパラメータ  $\theta$  の推定値  $\hat{\theta}$  を以下の式で求める

$$\hat{\theta} = \arg \max_{\theta} P(X_1, \dots, X_N | \theta)$$

- あるいは対数を取り推定: 対数尤度の最大化

$$\hat{\theta} = \arg \max_{\theta} \log P(X_1, \dots, X_N | \theta)$$

## ➤ MAP推定(事後確率の最大化)

- 事前確率  $P(\theta)$  が与えられていたときには、次式のように事後分布の確率を最大化するパラメータを求める

$$\hat{\theta} = \arg \max_{\theta} \log P(X_1, \dots, X_N | \theta) P(\theta)$$

ただし、 $X_1, \dots, X_N$  は  $N$  個の観測データ

◆問題1  $P(X) = \theta^X (1-\theta)^{1-X}$  ( $X_i$ は、0か1)で定義されるベルヌーイ試行を独立にN回繰り返したとき、0が $m$ 回、1が $N-m$ 回観測されたとする。最尤推定して  $\theta$ を求めよ

また、事前分布として、 $P(\theta) = b\theta$ ただし、 $0 \leq \theta \leq 1$ のときのMAP推定した  $\theta$  を求めよ。  
この場合の結果の意味を考察せよ。

◆問題2 次式の多項分布  $P(X) = \frac{N!}{X_1! \cdots X_K!} \theta_1^{X_1} \cdots \theta_K^{X_K} \quad \sum_{k=1}^K \theta_k = 1$

において最尤推定して  $\theta_i$  を求めよ。

事前分布が、 $P(\theta) = \frac{\Gamma(\alpha_1 + \cdots + \alpha_K)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \theta_1^{\alpha_1-1} \cdots \theta_K^{\alpha_K-1}$

の場合の、MAP推定した  $\theta_i$  を求めよ

# データの性質

- 今までは、情報源の記号 $t$ と出力記号列(= 直接に観測されたデータ) $x$ は、機械学習において**直接**に計算の対象としていた。
- この仮定が成立する場合も多い。
  - 身長、体重、薬の濃度、価格、などの(連続)数値データ
  - 人数、個数、などの整数をとる数値データ
  - 割合、%など
  - 男女、国籍など属性が記号の場合(整数に変換すれば数値として処理可能)
- しかし、必ずしも**直接**に観測されたデータだけを使える場合ばかりではない。



# 観測データを表す情報の次元

## ➤ 観測データ点が人間の場合の例

- $x = (\text{身長、体重、血圧、収入金額})^T \rightarrow$  数値だけなので簡単。単位は外部知識とする。

Ex (170, 50, 120, 10,000,000)

- 確率分布としては正規分布など。

- $x = (\text{職業、発熱})^T \rightarrow$  記号。2つの方法

- 記号に番号を与える。Ex 無職=0, 学生=1,...、発熱無=0、有=1

- 確率分布としては離散数値をとる分布など。数値の意味付けが難しい。

- 記号の種別ごとに1次元を与える(次のページ参照)

# 記号の種類ごとに次元を割り当てる方法

➤  $x = (\text{訪問国1}, \dots, \text{訪問国N})^\top$      $\text{ex}(\text{USA}, \text{UK}, \text{Italy})^\top$

➤ 対策: 国を番号つける。

(USA=1, UK=2, Japan=3, China=4, Italy=5, ...)

➤ この番号がベクトルの何番目の要素かを示すとして、数値のベクトルとして表現: **Bernoulli分布**:  $\text{Bern}(x | \mu) = \mu^x (1 - \mu)^{1-x}$

➤ 上のexは  $(1, 1, 0, 0, 1, \dots)^\top$

➤ このベクトルの次元は世界中の国の数だけあるため、かなり大きい。しかし、観測データには0が多く、スパースなデータ

➤ 記号の出現回数のある場合

$x = ((\text{訪問国1}, \text{滞在日数1}), \dots, (\text{訪問国N}, \text{滞在日数}))^\top$

➤  $\text{ex}((\text{USA}, 15), (\text{UK}, 5), (\text{Italy}, 3))^\top \rightarrow (15, 5, 0, 0, 3, \dots)^\top$

➤ **多項分布**:  $\text{Mult}(15, 5, 0, 0, 3, \dots | \mu_{\text{USA}}, \mu_{\text{UK}}, \mu_{\text{Japan}}, \mu_{\text{China}}, \mu_{\text{Italy}}, \dots)$

$$\propto \mu_{\text{USA}}^{15} \mu_{\text{UK}}^5 \mu_{\text{Japan}}^0 \mu_{\text{China}}^0 \mu_{\text{Italy}}^3 \dots$$

# 次元の大きさ

- 国と滞在日数の例と同じタイプの問題を、テキストデータで考えてみよう。
- あるテキストを表現するには、そのテキストに出現した各単語の個数で表現する。
  - 次元は語彙数 日本語の新聞では約40万語。固有名詞や複合語まで入れると、100万以上。→ **100万次元のベクトルを扱う必要あり！**
  - 個々の単語だけを対象にすれば済むのか？
    - ABC証券、ABC証券株式会社、...、総理が失言、総理が訂正、...、というような単語の連鎖で見ないと分からない場合は？
    - N単語の連鎖(=N-gram)の種類数は、100万のN乗！！
    - **しかし、このような多次元がすべて重要な情報だとも思えない**
    - 次元圧縮の技術が有望 i.e. Singular Value Decomposition (SVD)とかLatent Semantic Indexing(LSI)

# 特殊性を表すデータ 1

- これまでに示したデータ点の数値は、観測された数値（出現回数など）を直接使っていた。
- 観測データ全体の構造を利用した**tf\*idf**と呼ばれる数値も有力
- データ点頻度 Data point Frequency : DF
- ただし、 $DF(j)$ はj番目の次元のデータが0でないデータ点の数
- また、観測データ点の総数を $N$ とする。

## 特殊性を表すデータ 2

- データ点頻度 Data point frequency:DF
- ただし、DF(j)はj番目の次元のデータが0でないデータ点の数
- また、観測データ点の総数をNとする。
- $IDF(j)=1/DF(j)$
- $TF(i,j)$ =観測データ点iで第j次元のデータの出現回数
- TF\*IDFの定義:

$$w_{i,j} = TF * IDF(i, j) = TF(i, j) \cdot \log \frac{N}{DF(j)}$$

# 例

## ➤ データ例

旅行者a: (USA=10, UK=2, Japan=3, China=0, Italy=0)

旅行者b: (USA=0, UK=2, Japan=0, China=4, Italy=0)

旅行者c: (USA=5, UK=0, Japan=2, China=0, Italy=0)

旅行者d: (USA=2, UK=0, Japan=1, China=2, Italy=1)

➤  $DF(USA)=3$ ,  $DF(UK)=2$ ,  $DF(JP)=3$ .  $DF(CH)=2$ ,  $DF(IT)=1$

➤  $N/DF(..)$ は  $USA=4/3$ ,  $UK=4/2$ ,  $JP=4/3$ ,  $CH=4/2$ ,  $IT=4/1$

➤  $TF*IDF(USA,a)=10*\log(4/3)=4.114$ ,  $TF*IDF(USA,b)=0$

$TF*IDF(UK,a)=2*\log(4/2)=2$

$TF*IDF(IT,d)=1*\log(4/1)=2$

# 特殊性を表すデータ 3

- TF\*IDFの定義:  $w_{i,j} = \text{TF*IDF}(i,j) = \text{TF}(i,j) \cdot \log \frac{N}{\text{DF}(j)}$
- TF\*IDF(i,j)は、データ点: iだけで特別に多く現れる次元: jの数値を表す。
  - 例えば、新聞の1記事を観測データ点とし、次元を単語とすると、TF\*IDF(i,j)の大きな単語iは、偏りのある特殊ないし専門の単語、小さな単語は一般的な単語といえる。
- TF\*IDFを用いて観測データ点を表現しなおすと、いろいろなことが見えてくることがある。

# 距離の定義

- 観測データ点を多次元空間中の点と定義
  - そこで2つの問題
  - 各次元は観測データ点からどのように定義するか
    - 次元のことを**feature**あるいは**素性**(そせい)と呼ぶ
    - この問題をfeature design : 素性設計と呼ぶ。例えば、
      - 2つの素性の比を新たな素性とする ex 身長/体重
      - 2つの素性の連続したもの ex 日本・銀行、日本・沈没、
    - しかし、これは個別適用分野に応じて工夫すべし。
  - 多次元空間における2点間の距離の定義
    - ユークリッド距離ばかりではないのだ！