

第一回ベイズ統計学・機械学習研究会

～ベイズってなに？～

目次

- ベイズ統計学の歴史
- ベイズの定理
- ベイズ統計学の用途
- 最尤法とベイズ統計学
 - コインを使った例

ベイズ統計学の歴史

- 1740年ごろトーマス・ベイズがベイズの定理を発見
- フィッシャーとネイマン、ピアソンがベイズ統計学を否定
 - 事前分布の決め方に主観が入ることがあり数学的に厳密ではないから
- 1950年ドイツ軍の暗号をイギリスがベイズ統計学を用いて解読
- コンピューターの計算能力が向上し今まで手計算では太刀打ちできなかった問題が解けるように
 - マルコフ・チェーン・モンテカルロ (MCMC) の登場

ベイズの定理

- 使う道具はこれだけ

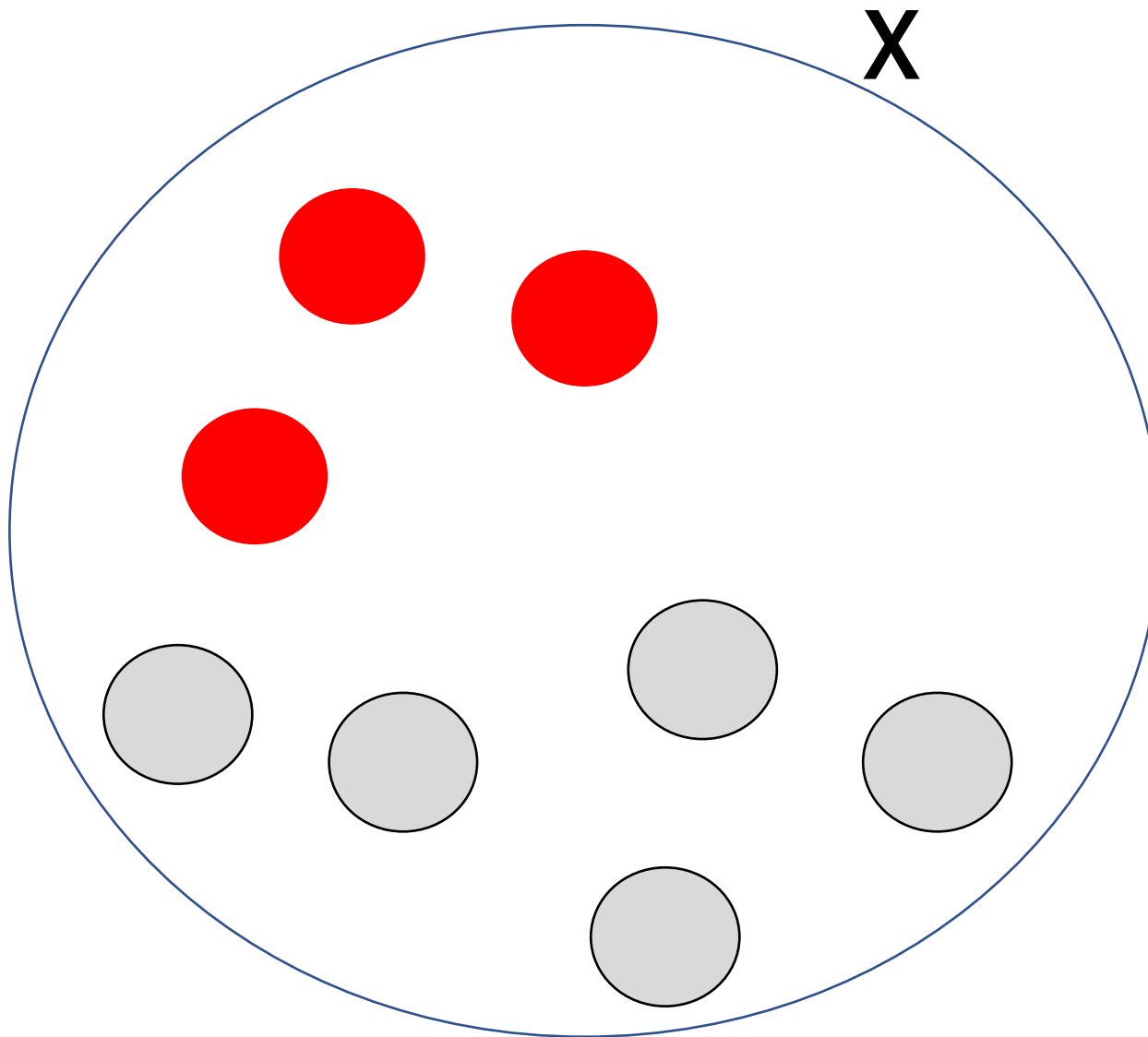
$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

例題

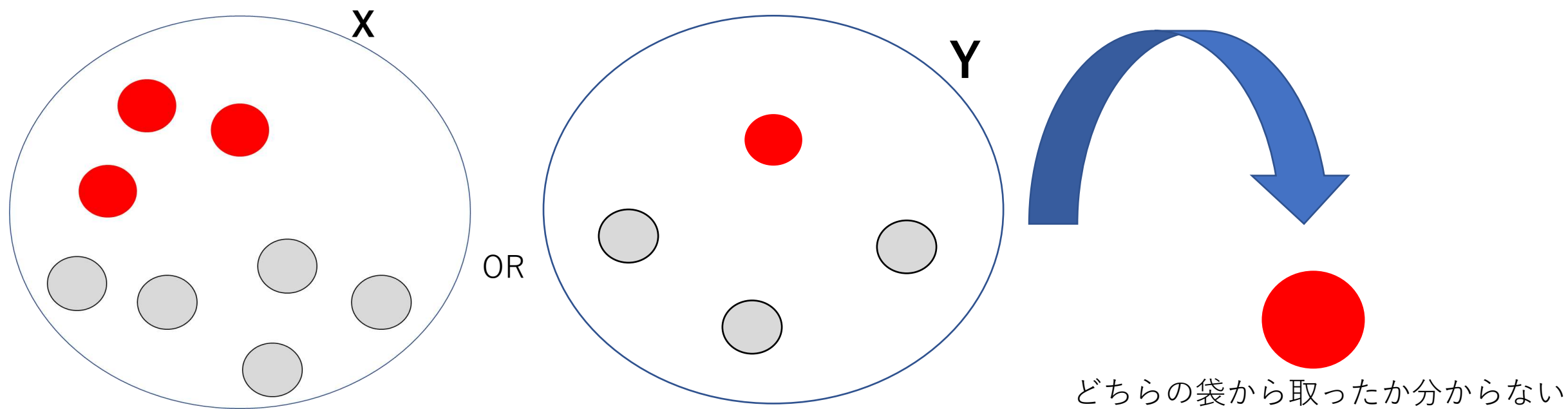
- XとYという二つの袋を用意する。Xには赤玉が3つ白玉が5つ、Yには赤玉が1つ白玉が3つ入っている。二分の一の確率でXかYを選び、その中の玉を一つとる。
- Xを選んだ時赤玉を取る確率は？
- 赤玉を取った時、それがXの袋に入っていたものである確率は？

- Xを選んだ時赤玉を取る確率は？

条件付確率
 $P(\text{赤玉}|X)=3/8$



- 赤玉を取った時、それがXの袋に入っていたものである確率は？



$$P(X|\text{赤玉}) = \frac{P(X) * P(\text{赤玉}|X)}{P(\text{赤})}$$

$$P(X)=1/2 \text{ (Xの袋を選ぶ確率)}$$

$$P(\text{赤})=(1/2*3/8)+(1/2*1/4) \\ =5/16 \text{ (赤い球を引く確率)}$$

$$P(\text{赤}|X)=3/8 \text{ (Xを選んだ時赤玉を取る確率)}$$

$$P(X|\text{赤玉})= \frac{P(X)*P(\text{赤玉}|X)}{P(\text{赤})}$$

$$=3/5$$

時間の流れ通りの情報から

時間の流れに逆らう分析ができる

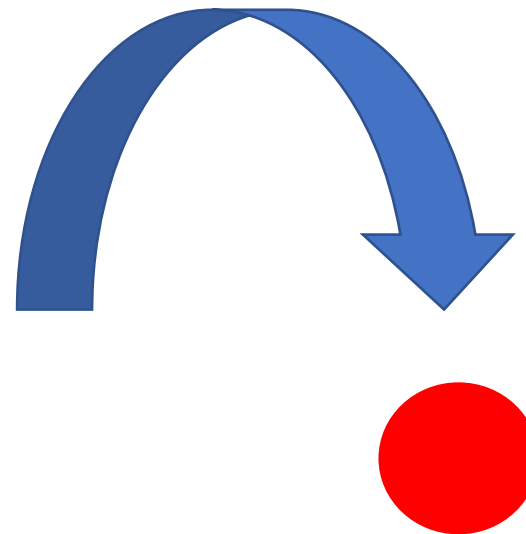
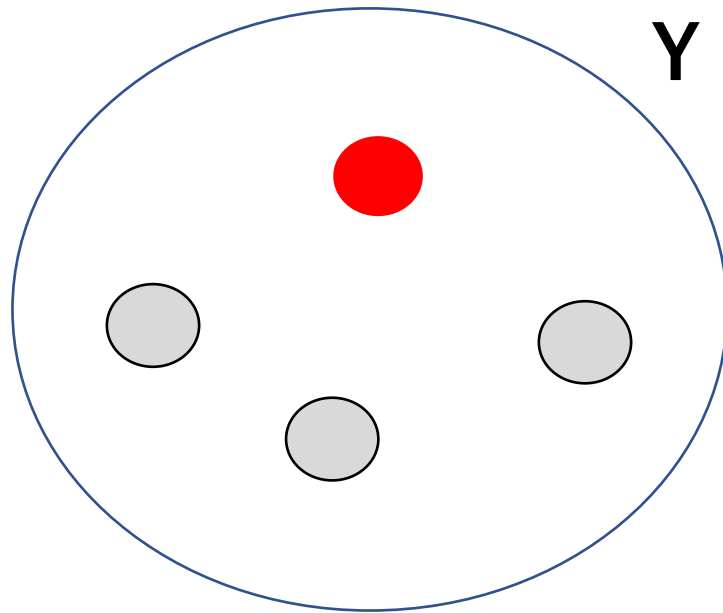
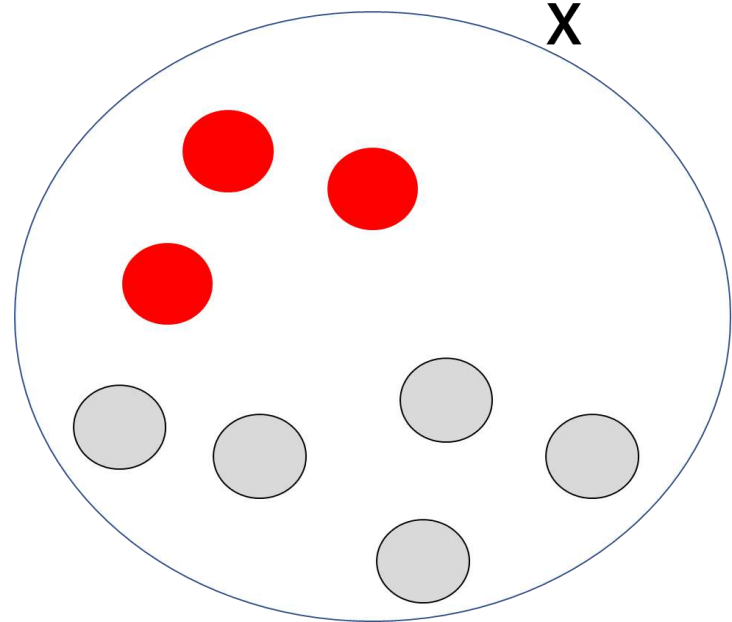
時間順行の確率

X

Y

OR

時間逆行の確率(ベイズの定理)



ベイズ統計の用途

- 結果から原因を推測
- データから母集団のパラメーターを予測
- 迷惑メールのフィルタリング
 - ユーザーに送られてきた迷惑メールからどんなメールが迷惑か推測
- 少ないデータからそれっぽい分析をする

最尤推定とベイズ統計学

- 最尤推定（さいゆうすいてい）
 - 与えられた観測値から、尤度を最大にするような母数や理論値を推定する方法
- 最もそれっぽい値を予測する方法（みなさんも使ってます）

例：n回投げてk回表が出たコインの表が出る確率 P_0 は

$P_0 = k/n$ と予測できる

Q. n 回投げて k 回表が出たコインの表が出る確率 P_0 を推測せよ。

コインの表が出る確率は、

$$\sum_k \binom{n}{k} p^k (1-p)^{n-k} \text{ と表せる。}$$

P_0 値によって確率は変わる

→これが最大になる P を求める

$$L(p) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$$

↑これを微分して最大になる P を求める

$L(p)$ が最大になる p も $\log L(p)$ が最大になる p も同じなので、対数をとって

$$\log L(p) = \log \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$$

$$= \log \binom{n}{k} + k \log p + (n-k) \log (1-p)$$

$$\frac{d}{dp} \log L(p) = \frac{k}{p} - \frac{(n-k)}{(1-p)}$$

これが 0 になるといいので

$$0 = \frac{k}{p} - \frac{(n-k)}{(1-p)}$$

$$p = \frac{k}{n}$$

おまけ

$(1-p) \log(1-p)$ の微分

$\log(1-p)$ を自然対数だと考えよう。

$$y = \log x \quad x = 1-p \text{ を } x \text{ として微分して、}$$

$$\frac{dy}{dx} = y' = \frac{1}{x} \quad x \text{ を代入して、}$$

$$= \frac{1}{1-p} \quad \text{... ①}$$

$$\frac{dx}{dp} = x' = -1 \quad \text{... ②}$$

①、②より

$$\frac{dy}{dx} \cdot \frac{dx}{dp} = \frac{dy}{dp} = \frac{1}{1-p} \times -1$$

$$= -\frac{1}{(1-p)}$$

最尤推定の限界

例題：AとBの二つのコインを用意する。Aは3回投げて2回表、Bは100回投げて60回表が出た。どちらのコインのほうが表が出やすいか答えよ。

最尤推定を用いて予測すると…

$$\begin{aligned} A: P(A) &= 2/3 \\ &= 66.7\% \end{aligned}$$

$$\begin{aligned} B: P(B) &= 60/100 \\ &= 60\% \end{aligned}$$

答え:Aのほうが表が出る確率が高い。

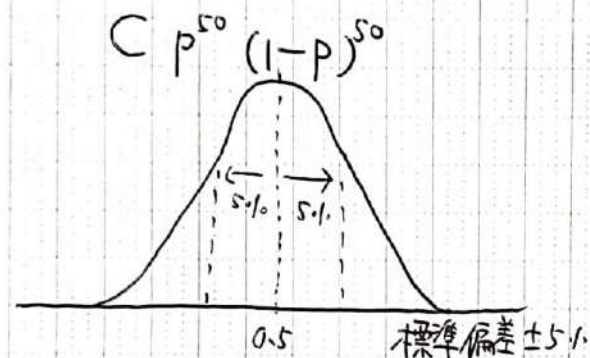
でも本当にそうだろうか…

ここでベイズ統計の力

- ベイズ統計では事前分布を使って確率を求める
- 事前分布とは、経験や一般常識、今まで集めたデータなどを確率分布の形であらわしたものだ。すでに経験的に知っていることをデータ分析に盛り込むことができる。
 - 事前分布の求め方は今後やります。

$$P(X|Y) = \frac{\overset{\text{尤度関数}}{P(Y|X)} \overset{\text{事前分布}}{P(X)}}{P(Y)}$$

・今回つかう事前分布



100回投げたとき、たいてい50%表が出るだろう。

コイ=A

事前分布 $C p^{50} (1-p)^{50}$ を仮定

$$p^{\text{Post}}(p) = C \times p^{50} \times (1-p)^{50} \times D \times p^2 \times (1-p)^1$$

$$= E \times p^{52} \times (1-p)^{51}$$



$$P_A = \frac{52}{52+51} = \frac{52}{103} = 50.49\%$$

コイ=B も同様 には

$$p^{\text{Post}}(p) = C p^{50} (1-p)^{50} \times D \times p^{60} \times (1-p)^{40}$$

$$= E \times p^{110} \times (1-p)^{90}$$



$$P_B = \frac{110}{110+90} = \frac{110}{200} = 55\%$$

ベイズ統計の力をつかうとい

コイ=A が表になる確率 50.49%

コイ=B 55 %

よりそれ、ほい分析ができる!!!

- 今日では例を示すため手計算でやりましたが、大体の計算はコンピューターがしてくれます
- コンピューターがどんな計算をしているのかを理解してベイズ統計学を使いこなしましょう
- 計算が理解できなくても、コンピューターに計算させるうちに理解できてくるようです
- **この講座が終わるころには皆さんを立派なベイジアンにします**