DMML, 20 Feb 2020

5H5T: $P(5H5T \mid \theta = 0.60) = P_1$
$P(5H5T \mid \theta = 0.50) = P_2$

$$\frac{P_1}{P_1 + P_2} \quad \frac{P_2}{P_1 + P_2}$$

0.45   0.55

**E-step** ②

HTTTHHTHTH
HHHHTHHHHH
HTHHHHHTHH
HTHTTTHHTT
THHHTHHHTH

0.45 x Ⓐ    0.55 x Ⓑ
0.80 x Ⓐ    0.20 x Ⓑ
0.73 x Ⓐ    0.27 x Ⓑ
0.35 x Ⓐ    0.65 x Ⓑ
0.65 x Ⓐ    0.35x Ⓑ

| Coin A | Coin B |
|---|---|
| ≈ 2.2 H, 2.2 T | ≈ 2.8 H, 2.8 T |
| ≈ 7.2 H, 0.8 T | ≈ 1.8 H, 0.2 T |
| ≈ 5.9 H, 1.5 T | ≈ 2.1 H, 0.5 T |
| ≈ 1.4 H, 2.1 T | ≈ 2.6 H, 3.9 T |
| ≈ 4.5 H, 1.9 T | ≈ 2.5 H, 1.1 T |
| ≈ 21.3 H, 8.6 T | ≈ 11.7 H, 8.4 T |

$\hat{\theta}_A^{(0)} = 0.60$

$\hat{\theta}_B^{(0)} = 0.50$

①

$$\hat{\theta}_A^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{\theta}_B^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$

③ **M-step**

④

$\hat{\theta}_A^{(10)} \approx 0.80$

$\hat{\theta}_B^{(10)} \approx 0.52$

Bottleneck of supervised learning — need for labelled
training data

Volume + Timeliness

Bootstrapping — Semi Supervised Learning

1. EM — Expectation Maximization

   Text Classification — Naïve Bayes

$$P(w_i | c_j)$$

$$P(d_i | c_j)$$

$$P(w_i|c_j) = \frac{\text{Occ. of } w_i \text{ in } c_j}{\text{All words in } c_j}$$

$$= \frac{\sum_{d_k \in c_j} n_{ik}}{\sum_{w_\ell \in V} \sum_{d_k \in c_j} n_{\ell k}}$$

$n_{ik}$ = # of occurrences of $w_i$ in $d_k$

$\begin{cases} 0 \text{ if } d_k \notin c_j \\ 1 \text{ if } d_k \in c_j \end{cases}$

Instead $\dfrac{\sum_{d_k} n_{ik} P(d_k|c_j)}{\sum_{w_\ell \in V} \sum_{d_k} n_{\ell k} P(d_k|c_j)}$

# Semi Supervised Learning

Supervised case — compute $P(w_i | c_j)$, $P(d_k | c_j)$

Given a new $d$

Compute $P(c_i | d)$ for all $c_i$

$$\underset{i}{\arg\max} \ P(c_i | d)$$

Fractional assignments of categories that we discard after arg max

Semi Supervised

Label, say, 10% of documents "randomly selected"
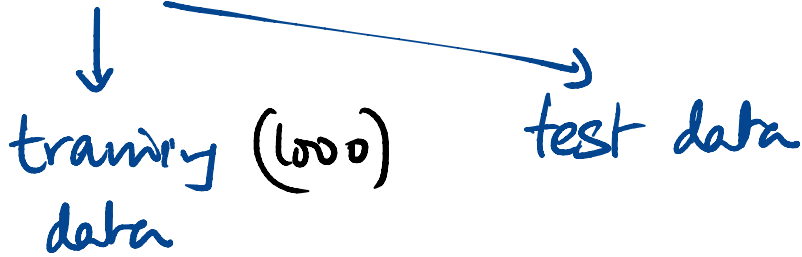
$\downarrow$

Naive Bayes model

$\downarrow$

Fractional <u>topic</u> allocation per document
class

$\downarrow$

Now $P(d_\kappa | c_i)$ makes sense!

New estimates of $P(w_i | c_j)$, $P(d_n | c_i)$

# Another example

MNIST digit images

→ training (1000) data

→ test data

label only 50

Cluster remaing 950 using these 50 as centroids
- labels are inherited from centroid

Use distance from centroid for logistic regression

So labelled point $\xrightarrow[1000]{\text{extend to}}$ 1000 labelled point

↓

Model

↓

Test data accuracy

$A_1$

Extend labels to "nearby" point

Nearest 20% pts in each cluster

↳ Model $A_3$

↓

Model

↓

Accuracy

$A_2$

$A_1 < A_2 < A_3$