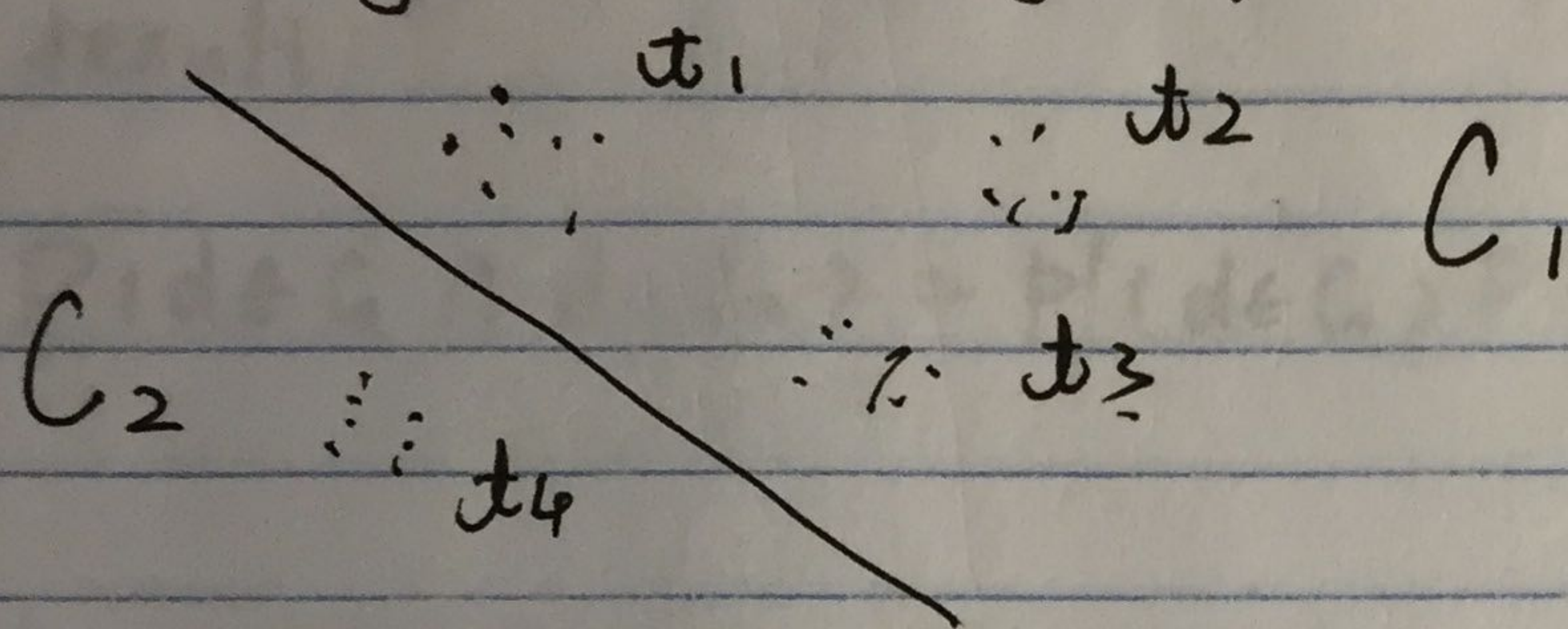# Model:

The data are clustered in nature, the observed categories are groups of the underlying clusters.

$t_1, t_2, t_3, t_4$ are the actual clusters,

$C_1$ and $C_2$ are the observed categories.

Taking naive bayes as an example, the probability for a data point $d$ to be in $C_1$ is

$$P(d \in C_1) = \sum_{i=1}^{3} p(d \in t_i) = \sum_{i=1}^{3} \prod_{w \in d} P(w | t_i) \cdot P(t_i)$$

The directly trained model is:

$$P'(d \in C_1) = \prod_{w \in d} P(w | C_1) \cdot P(C_1) \neq$$

$$= \prod_{w \in d} \sum_{i=1}^{3} P(w | t_i) P(t_i) \cdot \left( \sum_{i=1}^{3} P(t_i) \right)$$

These two can be different.

The previously proposed cross domain classification.

The previously proprosed model gives the following result:

$$P(d \in G \cap d \in C_2) = p'(d \in G) p'(d \in C_2) \frac{|G \cap C_2|}{|G \cup C_2|}$$

it is hard to show that this is less biased.

Thus, it is better to discover the latent clusters directly.

But there is an issue that $P(t_i \in C_2) = 0$, this will cause problems in EM, I haven't figure out how to treat it properly.

But if we want some preliminary results, we may try the following two algorithms first.

(Detailed direvation will not be shown.)

# Algorithm 1:

Input. $D = \{d_1, d_2, \cdots\}$     $C = \{c_1, c_2 \cdots\}$

~~data~~ $\{C_d = c_i, d \in c_i\}$

Initialized     $t = C$,         flag $=$ true.

while ( flag ):

$$P(t) = \frac{\mathbb{1}\{t_d = t\}}{M} \;;\quad P(w|t) = \frac{\sum \mathbb{1}\{t_d = t_i, w \in d\}}{\sum \mathbb{1}\{t_d = t\}}$$

E-M iterations :

$$E:\quad Q(c, t, d) = \frac{P(t)\, P(c|t)\, \prod_w P(w|t)}{\sum_{\hat{c},\hat{t}} P(\hat{t})\, P(\hat{c}|\hat{t})\, \prod_w P(w|\hat{t})}$$

$$M:\quad p(t_i) = \frac{\sum Q(c, t_i | d)}{\sum Q(c, t | d)}$$

$$P(w|t_i) = \frac{\sum Q(c, t_i | d)\, \mathbb{1}\{w \in d\}}{\sum Q(c, t_i | d)}$$

flag $=$ false :

For $c_i$ in $C$ :
    $IG = H(c_i) - H(c_i | t_c)$

    If $IG > IG$ pre:
        $IG$ pre $= IG$
        flag $=$ true
        $T_{c_i}$. add $t_{new}$
        $t_{new} = \{d_i ; t_d \notin T_{c_i}, d \in c_i\}$

# Algorithm 2:

$t \leftarrow C.$

while $|\{t\}| < maxNum$:

$$P(t_i) = \frac{\Sigma 1\{t_d = t_i\}}{|\{t\}|}.$$

$$P(w|t_i) = \frac{\Sigma 1\{t_d = t_i, w \in d\}}{\Sigma 1\{t_d = t_i\}}$$

$$P(c|t_i) = \frac{\Sigma 1\{t_d = t_i, d \in C\} + 1}{\Sigma 1\{t_d = t_i\} + |\{t\}|}$$

E-M iteration.
  $Q(t, c|d)$. $P(t_i)$, $P(w|t_i)$, $P(c|t_i)$.

For $t_i$ in $\{t\}$:

$t_i \in C$   it $P(c|t_i) \ge P(c'|t_i)$ $\forall c'$.

Record $\{t\}$.

$t_{new} = \{d; t_d \notin T_{c_i}, d \in c_i\}.$

Choose the best $\{t\}$. manually through validation.