

# Clustering Report

Zijie Pan

## 1 Clustering

### 1.1 Experimental setup

**Notations** See [1].

**Dataset and models** Label column: “dec”. Feature column: “attr”, “sinc”, “intel”, “fun”, “amb”, “shar”. Models: k-means and GMM with 3 types of inference — expectation maximization (GMM-EM), variational inference (GMM-VI), and Gibbs sampling (GMM-GS).

**Visualization via PCA** Fig. 1 shows the distribution of the dataset, where pc1 denotes the first principal component and pc2 denotes the second.

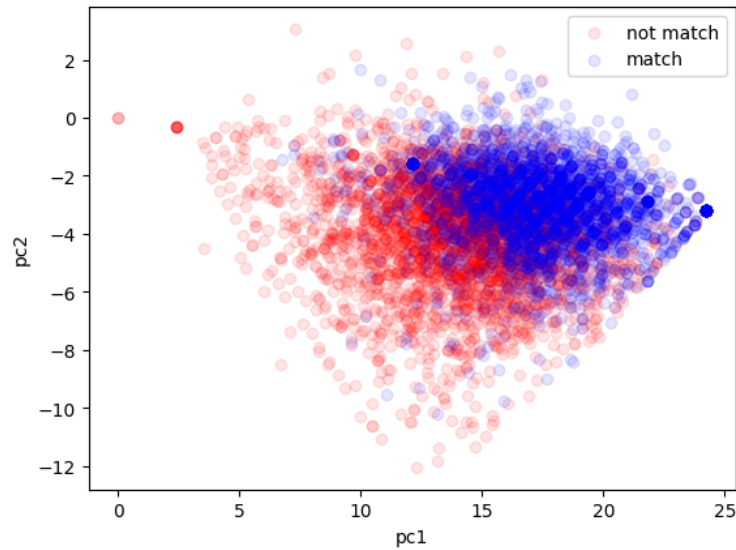


Figure 1: **Visualization via PCA.**

### 1.2 Results

**K-means**

**GMM-EM**

**GMM-VI**

**GMM-GS**

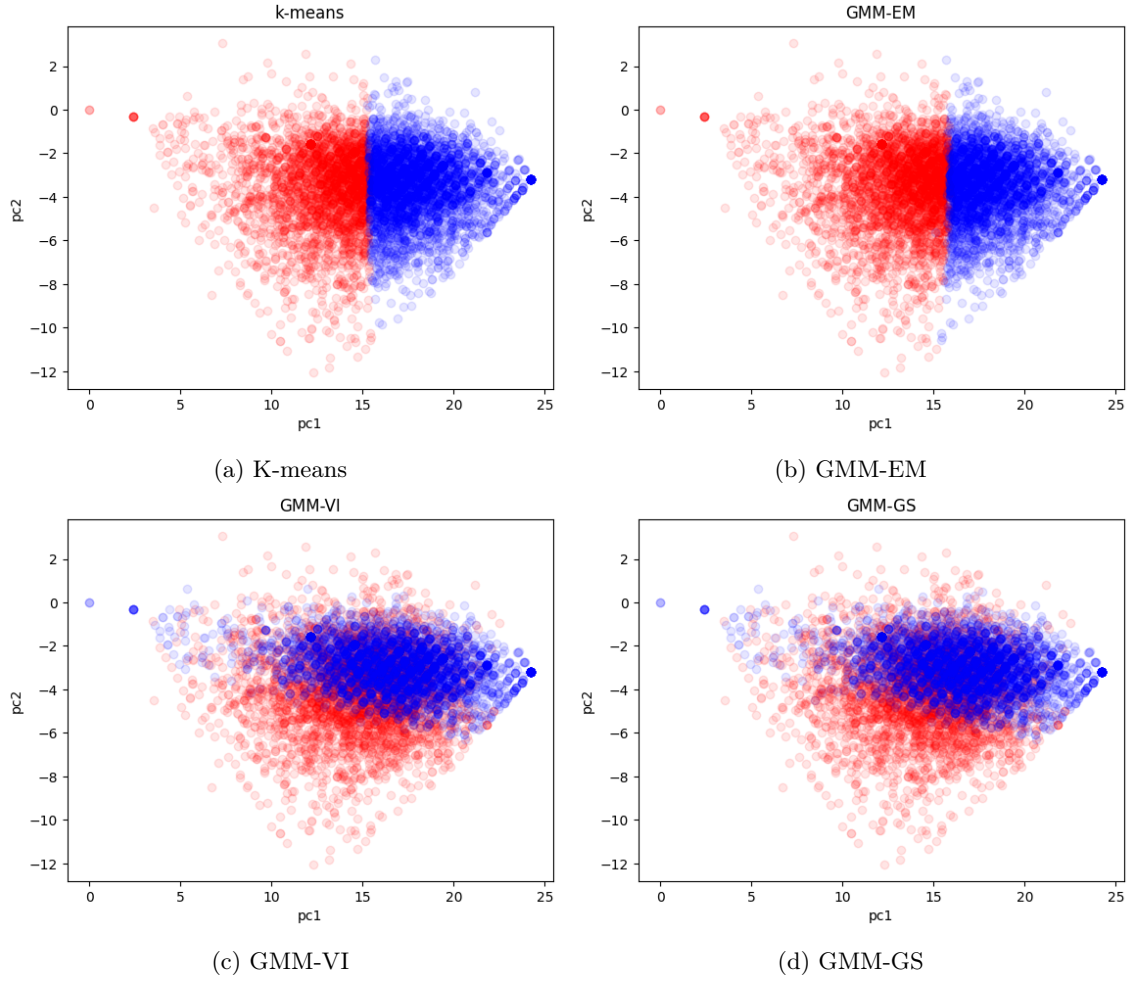


Figure 2: **Clustering results.**

## References

- [1] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.

---

**Algorithm 1** GMM-EM

---

```
1: input: initial random labels  $r_{ik}^0$ 
2: for  $t = 0, 1, \dots, T - 1$ 
3:   M-step
4:    $\pi_l^{t+1} = \frac{\sum_i r_{il}^t}{n}$ 
5:    $\mu_l^{t+1} = \frac{\sum_i r_{il}^t x_i}{\sum_i r_{il}^t}$ 
6:    $\Sigma_l^{t+1} = \frac{\sum_i r_{il}^t (x_i - \mu_l^{t+1})(x_i - \mu_l^{t+1})^\top}{\sum_i r_{il}^t}$ 
7:   E-step
8:    $r_{ik}^{t+1} = \frac{\pi_k^t N(x_i | \mu_k^t, \Sigma_k^t)}{\sum_l \pi_l^t N(x_i | \mu_l^t, \Sigma_l^t)}$ 
9: end for
```

---

---

**Algorithm 2** GMM-VI

---

```
1: input: initial random labels  $r_{ik}$  and prior parameters: Dirichlet parameter  $\alpha_0$ , Gaussian parameter  $m_0, \beta_0$ , Wishart parameter  $L_0, \nu_0$ 
2: for  $t = 0, 1, \dots, T - 1$ 
3:   M-like step
4:    $N_k = \sum_i r_{ik}$ 
5:    $\bar{x}_k = \frac{1}{N_k} \sum_i r_{ik} x_i$ 
6:    $S_k = \frac{1}{N_k} \sum_i r_{ik} (x_i - \bar{x}_k)(x_i - \bar{x}_k)^\top$ 
7:    $\beta_k = \beta_0 + N_k$ 
8:    $m_k = \frac{1}{\beta_k} (\beta_0 m_0 + N_k \bar{x}_k)$ 
9:    $L_k^{-1} = L_0^{-1} + N_k S_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{x}_k - m_0)(\bar{x}_k - m_0)^\top$ 
10:   $\nu_k = \nu_0 + N_k$ 
11:   $\alpha_k = \alpha_0 + N_k$ 
12:  E-like step
13:   $\tilde{\pi}_k = \exp \{ \psi(\alpha_k) - \psi(\sum_l \alpha_l) \}$ 
14:   $\tilde{\Lambda}_k = \exp \left\{ \sum_j \psi \left( \frac{\nu_k + 1 - j}{2} \right) + d \log 2 + \log \det L_k \right\}$ 
15:   $r_{ik} \propto \tilde{\pi}_k \tilde{\Lambda}_k^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} d \beta_k^{-1} - \frac{1}{2} \nu_k (x_i - m_k)^\top L_k (x_i - m_k) \right\}$ 
16: end for
```

---

---

**Algorithm 3** GMM-GS

---

```
1: input: initial  $\alpha, m_0, V_0, S_0, \nu_0$  and  $\pi, \mu_k, \Sigma_k$ 
2: for  $t = 0, 1, \dots, T - 1$ 
3:   Sampling  $z$ 
4:    $z_i \sim \text{Multinomial}(\pi_k N(x_i | \mu_k, \Sigma_k), k = 1, \dots, K)$ 
5:   Sampling  $\pi$ 
6:    $N_k = \sum_i 1(z_i = k)$ 
7:    $\pi \sim \text{Dir}(\alpha_1 + N_1, \dots, \alpha_K + N_K)$ 
8:   Sampling  $\mu_k$ 
9:    $V_k^{-1} = V_0^{-1} + N_k \Sigma_k$ 
10:   $m_k = V_k (\Sigma_k^{-1} \sum_{i: z_i = k} x_i + V_0^{-1} m_0)$ 
11:   $\mu_k \sim N(m_k, V_k)$ 
12:  Sampling  $\Sigma_k$ 
13:   $\bar{x}_k = \frac{1}{N_k} \sum_{i: z_i = k} x_i$ 
14:   $S_k = \sum_{i: z_i = k} (x_i - \bar{x}_k)(x_i - \bar{x}_k)^\top$ 
15:   $\Sigma_k \sim IW(S_0 + S_k, \nu_0 + N_k)$ 
16: end for
```

---