

Hand on session 2: Clustering

Exercise 2: On the consistency of K-Means

$$\text{crit}(G) = \sum_{c=1}^K \sum_{a \in G_c} \|x_a - \bar{x}_{G_c}\|^2 \quad \text{with} \quad \bar{x}_{G_c} = \frac{1}{|G_c|} \sum_{b \in G_c} x_b$$

$$1. \text{crit}(G) = \sum_{c=1}^K \sum_{a \in G_c} \left\langle x_a - \frac{1}{|G_c|} \sum_{b \in G_c} x_b, x_a - \bar{x}_{G_c} \right\rangle$$

$$= \sum_{c=1}^K \frac{1}{|G_c|} \sum_{a, b \in G_c} \langle x_a - x_b, x_a - \bar{x}_{G_c} \rangle$$

$$= \sum_{c=1}^K \frac{1}{|G_c|} \sum_{a, b \in G_c} \left(\langle x_a - x_b, x_a \rangle - \langle x_a - x_b, \bar{x}_{G_c} \rangle \right)$$

$$\text{and} \quad \sum_{a, b \in G_c} \langle x_a - x_b, \bar{x}_{G_c} \rangle = \sum_{c \in G_c} \langle |G_c| x_a - |G_c| \bar{x}_{G_c}, \bar{x}_{G_c} \rangle$$
$$= 0$$

$$\text{So, } \boxed{\text{crit}(G) = \sum_{c=1}^K \frac{1}{|G_c|} \sum_{a, b \in G_c} \langle x_a, x_a - x_b \rangle}$$

$$\forall a, b \in G_c, \|x_b - x_a\|^2 = \langle x_a, x_a - x_b \rangle + \langle x_b, x_b - x_a \rangle$$

$$\text{So, } 2 \text{crit}(G) = \sum_{c=1}^K \frac{1}{|G_c|} \sum_{a, b \in G_c} \|x_b - x_a\|^2$$

$$\text{So } \boxed{\text{crit}(G) = \frac{1}{2} \sum_{c=1}^K \frac{1}{|G_c|} \sum_{a, b \in G_c} \|x_a - x_b\|^2}$$

$$2. \quad x_a = \mu_a + \varepsilon_a, \quad \varepsilon_1, \dots, \varepsilon_n \text{ centered and independent}$$

$$\forall a, b, \|x_a - x_b\|^2 = \|\mu_a - \mu_b\|^2 + 2\langle \varepsilon_a, \mu_a - \mu_b \rangle + 2\langle \varepsilon_b, \mu_a - \mu_b \rangle$$
$$+ \|\varepsilon_a\|^2 + \|\varepsilon_b\|^2 - 2\langle \varepsilon_a, \varepsilon_b \rangle$$

$$S_0, \forall a, b, E(\|X_b - X_a\|^2) = \|\mu_a - \mu_b\|^2 + 2\langle E(\varepsilon_a), \mu_a - \mu_b \rangle + 2\langle E(\varepsilon_b), \mu_a - \mu_b \rangle \\ + E(\|\varepsilon_a\|^2) + E(\|\varepsilon_b\|^2) - 2\langle E(\varepsilon_a), E(\varepsilon_b) \rangle \\ \text{since } \varepsilon_a \text{ and } \varepsilon_b \text{ are indept if } a \neq b \text{ (0 otherwise)}$$

$$S_0, \forall a, b, E(\|X_b - X_a\|^2) = (\|\mu_a - \mu_b\|^2 + v_a + v_b) \mathbb{1}_{a \neq b} \\ \text{where } v_a = \text{tr}(\text{cov}(X_a)) = E(\|\varepsilon_a\|^2) \\ \text{since } E(\varepsilon_a) = E(\varepsilon_b) = 0$$

$$S_0, \left| E(\text{cnt}(G)) = \frac{1}{2} \sum_{c=1}^K \frac{1}{|G_c|} \sum_{a, b \in G_c} (\|\mu_a - \mu_b\|^2 + v_a + v_b) \mathbb{1}_{a \neq b} \right|$$

3 $\forall \mu: \exists m_1, \dots, m_K \in \mathbb{M}^p, \exists \sigma_1, \dots, \sigma_K > 0$

a) $\text{st } \forall c \in [1, K], \forall a \in G_c^*, \mu_a = m_c \text{ and } v_a = \sigma_c^2$

$$\text{then } E(\text{cnt}(G^*)) = \frac{1}{2} \sum_{c=1}^K \frac{1}{|G_c^*|} \sum_{a, b \in G_c^*} (\|m_c - m_c\|^2 + \sigma_c + \sigma_c) \mathbb{1}_{a \neq b} \\ = \frac{1}{2} \sum_{c=1}^K \frac{1}{|G_c^*|} 2 \sigma_c (|G_c^*| - 1) |G_c^*|$$

$$S_0, \left| E(\text{cnt}(G^*)) = \sum_{c=1}^K \sigma_c (|G_c^*| - 1) \right|$$

b) $\text{If } \sigma_1 = \dots = \sigma_K, G = G^* \text{ minimizes } E(\text{cnt}(G))$

$$c) \|m_2 - m_3\|^2 = \|(0, 1, 0) - (0, 1 - \tau, \sqrt{1 - (1 - \tau)^2})\|^2 \\ = (1 - (1 - \tau))^2 + 1 - (1 - \tau)^2 = 1 + \tau^2 - (1 - \tau)^2 = 2\tau$$

$$S_0, \underline{\|m_2 - m_3\|^2 = 2\tau}$$

$$d) \underline{E(\text{cnt}(G^*)) = (\sigma_+ + 2\sigma_-)(S - 1)}$$

$$e) \sum_{a, b \in G_1} (\|\mu_a - \mu_b\|^2 + \nu_a + \nu_b) \mathbb{1}_{a \neq b} = \sum_{a, b \in G_1} (\nu_a + \nu_b) \mathbb{1}_{a \neq b} \\ = \binom{s}{2} \left(\frac{s}{2} - 1 \right) (2\sigma_+)$$

$$\sum_{a, b \in G_2} (\|\mu_a - \mu_b\|^2 + \nu_a + \nu_b) \mathbb{1}_{a \neq b} = \binom{s}{2} \left(\frac{s}{2} - 1 \right) (2\sigma_+)$$

$$\sum_{a, b \in G_2} (\|\mu_a - \mu_b\|^2 + \nu_a + \nu_b) \mathbb{1}_{a \neq b} =$$

$$s(s-1)2\sigma_- + s(s-1)2\sigma_- + s^2(2\tau + 2\sigma_-) + s^2(2\tau + 2\sigma_-)$$

$$So, E(\text{crit}(G')) = \frac{1}{2} [2(s-2)\sigma_+ + (s-1)2\sigma_- + (s-1)2\sigma_- + 2s(2\tau + 2\sigma_-)]$$

$$E(\text{crit}(G')) = s(\sigma_+ + 2\sigma_- + \tau) - (2\sigma_+ + \sigma_-)$$

$$f) E(\text{crit}(G^*)) < E(\text{crit}(G')) \iff (\sigma_+ + 2\sigma_-)(s-1) < s(\sigma_+ + 2\sigma_- + \tau) - (2\sigma_+ + \sigma_-)$$

$$\iff \tau > \frac{1}{s} [(\sigma_+ + 2\sigma_-)(s-1) + 2\sigma_+ + \sigma_-] - \sigma_+ - 2\sigma_-$$

$$\iff \tau > (\sigma_+ - \sigma_-) / s$$

g) The optimal solution of k-News for the MM is the generative model only if the model is well-separable.