Distortion

metric that assesses the performance of K-means (smaller values better)



$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \|x_n - \mu_k\|^2 \qquad \text{Goal: choose } r_{nk} \text{ and } \mu_k \text{ that minimizes } J$$

hard assignments!
$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_{j} ||x_n - \mu_j||^2 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{dJ}{d\mu_k} = 2\sum_{n=1}^{N} r_{nk}(x_n - \mu_k) = 0 \implies \mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}} = \frac{1}{N_k} \sum_n r_{nk} x_n$$

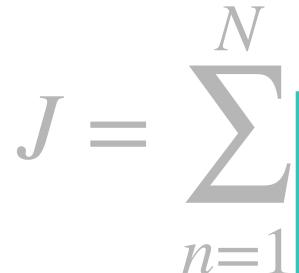
optimal value for μ_k minimizing our loss is the mean of all data points in that cluster

Distortion

metric that assesses the performance of K-means (smaller values better)



minimizes



/ actual data point n

- 1. choose **k** random points as cluster centers
- 2. for each data point, assign it the cluster whose centroid is the closest
- 3. using these assignments, recalculate the centers
- 4. reiterate from step (2) until convergence:
- cluster membership does not change
- center only changes very very little

$$\frac{dJ}{d\mu_{1}} = 2$$

hard assignmen

$$\sum_{n} r_{nk} = N_k = \sum_{n} r_{nk} x_n$$

optimal value for μ_k minimizing our loss is the mean of all data points in that cluster