

## Learning Feature Representations with K-means

*Lecturer: Krishnan Srinivasan**Scribe: Jonathon Cai*

## 1 Today

– Paper Discussion

## 2 Derivation of arg max term

We would like to minimize

$$\min_{s^{(i)}} \|Ds^{(i)} - x^{(i)}\|_2^2$$

given  $D \in \mathbb{R}^{n \times k}$ ,  $x^{(i)} \in \mathbb{R}^n$ , and  $s^{(i)} \in \mathbb{R}^k$ . The  $j$ -th column vector of  $D$  is denoted by  $D^{(j)}$ , where  $\forall j$ , the  $l_2$  norms of  $D^{(j)}$  are 1 and  $s^{(i)}$  has at most one nonzero coordinate. It turns out that both these assumptions – the restriction on  $D^{(j)}$  and the fact that  $\|s^{(i)}\|_0 \leq 1$  – are important for the analysis.

For the moment assume that  $s^{(i)}$  has a nonzero coordinate at  $r = s_j^{(i)}$ , where  $r$  is a constant number. Then  $Ds^{(i)} = rD^{(j)}$ . So we would like to solve the equation:

$$\frac{\partial}{\partial r} \|rD^{(j)} - x^{(i)}\|_2^2 = 0$$

Expanding yields

$$\frac{\partial}{\partial r} \sum_k (rd_k - x_k)^2 = 0$$

$$\sum_k 2(rd_k - x_k)d_k = 0$$

$$r \sum_k d_k^2 = \sum_k x_k d_k$$

$$r = \sum_k x_k d_k$$

$$r = D^{(j)T} x^{(i)}$$

normalized  $D^{(j)}$  assumption

Now we will consider the minimum of the norm over all  $j$ . Let  $r_j$  denote the  $r$  corresponding to  $j$ . Note that:

$$\begin{aligned}
\sum_k (r_j d_k - x_k)^2 &= \sum_k r_j^2 d_k^2 + x_k^2 - 2r_j d_k x_k \\
&= r_j^2 + \sum_k x_k^2 - 2r_j^2 \\
&= C - r_j^2
\end{aligned}$$

$C = l_2$  norm squared of  $x^{(i)}$

Thus, for  $\|s^{(i)}\|_0 \leq 1$ , the norm is minimized when we maximize  $|r_j|$ , or when  $j = \arg \max_j |r_j|$ .