

K-Means 聚类算法

简介: ① unsupervised Learning Algo. 将所有近似的 example 聚集到一个 cluster 中。

算法步骤: (1) 随机初始化 K 个 centroids

往往是在原 dataset 中随机选 K 个 example, 而不是人为指定

(2) 对每个 $x^{(i)}$, 将其匹配到离其是近的那个 centroid, 使用 $idx(i)$ 来保存 centroid 的编号

(3) 更新 centroid 的值。即用该 centroid 内所有 $x^{(i)}$ 的平均值 (mean) 来更新该 centroid。

代价函数: $J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu^{(1)}, \mu^{(2)}, \dots, \mu^{(K)}) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu^{(idx(i))}\|^2$

就是保证每个 $x^{(i)}$ 与其对应的 centroid 的距离总和最小

如何选择 K : 一般使用 elbow method: 逐渐增大 K , 在

$J(x)$ 变化不大的那个 K 即取为最终的 K 。



K-Means 作图像 compression: 原本图像每个 pixel 24bit.
(RGB 表示, 每个 8bit, 表 0-255 之间 integer). 图应一个 picture
有 hundred of thousands colors. 现在, 将 ~~所有~~ picture 用
K 个 color 表示。即将所有 3D 中的颜色点进行聚类。最
后用 centroid 的 color 值替换每个 $x^{(i)}$ 的 color 值。只不过现
在每个 pixel 只需要存储 centroid 的索引即可。

