1. As for the 2 type of PCA

① PCA based on projected variance maximization

let $v_1, v_2, \ldots, v_d$ denote the $d$ principal components,

$v_i \cdot v_j = 0$, $i \neq j$ and $v_i \cdot v_i = 1$, $i = j$.

And Assume that data is centered.

Let $X = [x_1, x_2, \ldots, x_n]$, we have $\frac{1}{n} \sum_{i=1}^{n} (v^T x_i)^2 = v^T X X^T v$

we need to find vector that maximizes sample variance of projected data

which is $\max_{v} (v^T X \bar{X}^T v)$ s.t. $v^T v = 1$

$$= Tr(v v^T X^T X v v^T) = Tr(X^T X v v^T)$$

② PCA based on projected error minimization,

$\frac{1}{n} \sum_{i=1}^{n} \| x_i - (v^T x_i) v \|^2$, we have that

$$\min_{v} \| X - X v v^T \|^2 = Tr((X - X v v^T)^T (X - X v v^T))$$
$$= Tr((I - v v^T)^T X^T X (I - v v^T))$$
$$= Tr(X^T X (I - v v^T))$$
$$= Tr(X^T X - X^T X v v^T)$$
$$= Tr(X^T X) - Tr(X^T X v v^T)$$

as the $Tr(X^T X)$ is a constant, the

$\max Tr(X^T X v v^T)$ is same as $\min -Tr(X^T X v v^T)$

so we get that PCA based on projected variance maximization is equivalent to PCA based on projected error minization.

2.

(a) Cost function:

$$\sum_{i=1}^{n}\sum_{j=1}^{k}\|X_i - \mu_j\|_2^2 \;\; Y_i^{\,j}, \quad \text{where } Y_i^{\,j}=1 \text{ when } y_i=j, \text{ and } 0 \text{ for otherwise}$$

(b) as defined in (a), we need to minimize the cost function, after each calculation, the mean will be updated, we denote the new mean is $\mu'$

$$f = \sum_{i=1}^{n}\sum_{j=1}^{k}\|X_i - \mu'_j\|_2^2 \, Y_i^{\,j}, \quad \text{where } Y_i^{\,j}=1, \text{ when } y_i=j, \; 0 \text{ for others}$$

$$\frac{\partial f}{\partial \mu'_j} = -2\sum_{i\in C_j}(X_i - \mu'_j)=0, \quad \text{then we get that}$$

$$|C_j|\,\mu'_j = \sum_{i\in C_j} X_i, \quad \text{where } |C_j| \text{ is the number of points in the corresponding}$$

$j$th cluster. So we get that $\mu'_j = \dfrac{\sum_{i\in C_j} X_i}{|C_j|}$, as the definition of centroid, we find that the each $\mu_j$ need to be the mean of the sample points assigned to the cluster i.

(c) The It do not follow the advice, there will be mainly 2 problems:

① When we met the case that a sample point has serval clusters means that are equally close, when we move to the next iteration, it not stay at the last cluster, we may fall into a case that we find a local minimum solution, which result the wrong solution.

② If we do not stay at the last cluster, the sample point may continuously move among the serval cluster, which will need cost more iterations. so will make the time longer or even not converge.

(d) CD