

exercise 4

fun with Frank-Wolfe

solutions due

until **December 9, 2025** at **23:30** via **ecampus**

general remarks

The following sheets contain a lot of text with a lot of “unconventional” information not found in the books. All this text is to provide you with the theoretical background of the tasks you are supposed to solve.

These tasks will all be about applications of the Frank-Wolfe algorithm for convex optimization over convex sets. In fact, they are (almost) all about computing global or local means of data sets. While you may be familiar with 1-mean or k -means computation, it is likely that you will not yet have seen the methods we employ here.

If you work carefully through this exercise, you will gain incredibly deep insights into the nature of k -means clustering. This should be fun and will come in handy later in the course . . .

task 4.1 [5 points]**sample means (part 1)**

Consider a sample $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ of n data points $\mathbf{x}_j \in \mathbb{R}^m$ which we may gather in a data matrix

$$\mathbf{X} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_n \\ | & | & \cdots & | \end{bmatrix} \in \mathbb{R}^{m \times n}$$

In this task, we are interested in estimating the mean $\hat{\boldsymbol{\mu}} \in \mathbb{R}^m$ of such a sample. To begin with, we recall a lesser known fact, namely that the familiar maximum likelihood estimator

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

is the minimizer of the following convex function

$$f(\boldsymbol{\mu}) = \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j - \boldsymbol{\mu}\|^2 = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^\top \mathbf{x}_j - 2 \mathbf{x}_j^\top \boldsymbol{\mu} + \boldsymbol{\mu}^\top \boldsymbol{\mu}$$

In other words, **the sample mean of $\mathcal{X} \subset \mathbb{R}^m$ is the particular point $\hat{\boldsymbol{\mu}} \in \mathbb{R}^m$ that has the smallest average distance to all the points contained in \mathcal{X} .**

background info

Let's quickly verify that

$$\hat{\boldsymbol{\mu}} = \underset{\boldsymbol{\mu}}{\operatorname{argmin}} f(\boldsymbol{\mu})$$

We have the gradient

$$\nabla f = \frac{df}{d\boldsymbol{\mu}} = \frac{1}{n} \sum_{j=1}^n 2\boldsymbol{\mu} - 2\mathbf{x}_j$$

which, when equated to 0 (and then divide by 2), yields the equivalencies

$$\frac{1}{n} \sum_{j=1}^n \boldsymbol{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \Leftrightarrow \frac{1}{n} n \boldsymbol{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \Leftrightarrow \boldsymbol{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

Now, given what we saw in earlier exercises, we realize that we can also write the sample mean like this

$$\hat{\mu} = \frac{1}{n} \mathbf{X} \mathbf{1} \equiv \mathbf{X} \hat{\mathbf{w}}$$

where the entries of vector $\hat{\mathbf{w}} = \frac{1}{n} \mathbf{1} \in \mathbb{R}^n$ are non-negative and sum to one

$$\begin{aligned} \hat{\mathbf{w}} &\geq \mathbf{0} \\ \mathbf{1}^\top \hat{\mathbf{w}} &= 1 \end{aligned}$$

But this is then to say that we may also compute the sample mean using a rather elaborate two-step procedure, namely

Step 1:

$$\begin{aligned} \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \quad & \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j - \mathbf{X} \mathbf{w}\|^2 \\ \text{s.t.} \quad & \mathbf{1}^\top \mathbf{w} = 1 \\ & \mathbf{w} \geq \mathbf{0} \end{aligned} \tag{1}$$

Step 2:

$$\hat{\mu} = \mathbf{X} \hat{\mathbf{w}} \tag{2}$$

Below, your task will be to implement this procedure. Indeed, based on what we did in earlier exercises, everybody should be able to solve the problem in (1) using `scipy.optimize` methods. However, there also is another, arguably better way ...

To see this, we first of all note that the objective function of the constrained minimization problem in (1) is

$$g(\mathbf{w}) = \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j - \mathbf{X} \mathbf{w}\|^2 = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j^\top \mathbf{x}_j - 2 \mathbf{x}_j^\top \mathbf{X} \mathbf{w} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w})$$

for which we have the following gradient

$$\begin{aligned} \nabla g &= \frac{dg}{d\mathbf{w}} = \frac{1}{n} \sum_{j=1}^n (2 \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 \mathbf{X}^\top \mathbf{x}_j) \\ &= \frac{2}{n} \mathbf{X}^\top \mathbf{X} \mathbf{w} - \frac{2}{n} \mathbf{X}^\top \sum_{j=1}^n \mathbf{x}_j \\ &= 2 \mathbf{X}^\top \mathbf{X} \mathbf{w} - \frac{2}{n} \mathbf{X}^\top \mathbf{X} \mathbf{1} \\ &= 2 \mathbf{X}^\top \mathbf{X} \left[\mathbf{w} - \frac{1}{n} \mathbf{1} \right] \end{aligned} \tag{3}$$

Second of all, we note that any $\mathbf{w} \in \mathbb{R}^n$ which obeys the constraints $\mathbf{w} \geq \mathbf{0}$ and $\mathbf{1}^\top \mathbf{w} = 1$ is an element of the **standard simplex**

$$\Delta^{n-1} = \left\{ \sum_{i=1}^n w_i \mathbf{e}_i \mid \mathbf{w} \in \mathbb{R}^n, \mathbf{w} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{w} = 1 \right\}$$

where \mathbf{e}_i denotes the i -th standard basis vector of \mathbb{R}^n .

Third of all, since Δ^{n-1} is a *compact convex set* and $g(\mathbf{w})$ is a convex objective function, we realize that (1) constitutes a convex minimization problem over a compact convex set and can thus be solved by running the following (problem specific) version of the **Frank-Wolfe algorithm**

procedure FW_SAMPLE_MEAN_WEIGHTS($\mathbf{X} \in \mathbb{R}^{m \times n}, t_{\max} \in \mathbb{N}$)

guess a feasible $\mathbf{w}_0 \in \Delta^{n-1}$, say

$$\mathbf{w}_0 = \mathbf{e}_1$$

for $t = 0, \dots, t_{\max} - 1$

determine step direction

$$\mathbf{s}_t = \operatorname{argmin}_{\mathbf{s} \in \Delta^{n-1}} \mathbf{s}^\top \nabla g(\mathbf{w}_t)$$

determine step size

$$\eta_t = \frac{2}{t+2}$$

update current guess

$$\mathbf{w}_{t+1} = (1 - \eta_t) \mathbf{w}_t + \eta_t \mathbf{s}_t$$

return \mathbf{w}_{t+1}

Note: $\mathbf{s}^\top \nabla g(\mathbf{w}_t)$ is a linear function over Δ^{n-1} . Since Jensen's inequality tells us that the minimum of a convex function over a convex set is attained at a of said vertex, we have

$$\mathbf{s}_t = \operatorname{argmin}_{\mathbf{s} \in \Delta^{n-1}} \mathbf{s}^\top \nabla g(\mathbf{w}_t) \quad \Leftrightarrow \quad \mathbf{s}_t = \operatorname{argmin}_{\mathbf{e}_i \in \mathbb{R}^n} \mathbf{e}_i^\top \nabla g(\mathbf{w}_t)$$

an you should use the very convenient fact in your solutions for this task.

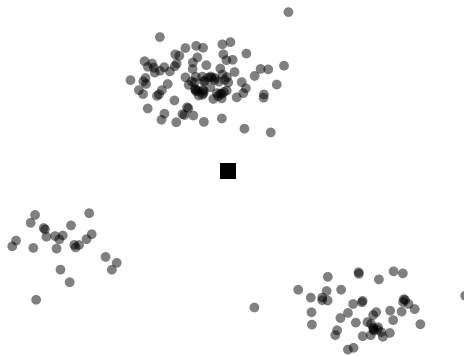
here is your task

The file

```
threeBlobs.csv
```

contains a data matrix $\mathbf{X} \in \mathbb{R}^{2 \times 175}$ of 175 data points $\mathbf{x}_j \in \mathbb{R}^2$. Read this matrix into memory and—to get a ground truth result— (properly) use the [numpy](#) function [mean](#) to compute the sample mean $\bar{\mathbf{x}}$ of the \mathbf{x}_j . Have a look at the result.

Next, using the gradient in (3), implement the Frank-Wolfe algorithm for solving the problem in (1). Run the algorithm for $T = 10000$ iterations to get an estimate of $\hat{\mathbf{w}} \in \mathbb{R}^{175}$. Given $\hat{\mathbf{w}}$, compute $\hat{\boldsymbol{\mu}} = \mathbf{X}\hat{\mathbf{w}}$ and plot the data (say as round dots) together with your estimated mean (say as a square). Your plot should look something like this



Finally, here is a statement w/o proof: “After T iterations, the Frank-Wolfe algorithm will produce a solution that is $O(\frac{1}{T})$ away from an optimum”.

Practically investigate what this means! Proceed as follows: Run your Frank-Wolfe implementation for $T \in \{10, 100, 1000, 10000\}$ iterations. This should result in different $\hat{\mathbf{w}}$ and therefore in different mean estimates $\hat{\boldsymbol{\mu}} = \mathbf{X}\hat{\mathbf{w}}$. Compare each of these mean estimates to the ground truth $\bar{\mathbf{x}}$ which you computed above. Discuss what you observe.

task 4.2 [5 points]**sample means (part 2)**

Let x_j , X , μ , $\hat{\mu}$, w , \hat{w} , and $\mathbf{1}$ be as in the previous task and (try to) convince yourself that the following equivalencies hold true

$$\sum_{j=1}^n \|x_j - \mu\|^2 \Leftrightarrow \|X - \mu \mathbf{1}^\top\|^2 \Leftrightarrow \|X - Xw\mathbf{1}^\top\|^2$$

Working with the rightmost expression, we can thus compute the mean $\hat{\mu}$ of the columns x_j of X using this elaborate two-step procedure

Step 1:

$$\begin{aligned} \hat{w} = \underset{w}{\operatorname{argmin}} \quad & \|X - Xw\mathbf{1}^\top\|^2 \\ \text{s.t.} \quad & \mathbf{1}^\top w = 1 \\ & w \geq \mathbf{0} \end{aligned} \tag{4}$$

Step 2:

$$\hat{\mu} = X\hat{w} \tag{5}$$

Indeed, the objective of the minimization problem in (4) can be written as

$$\begin{aligned} h(w) &= \|X - Xw\mathbf{1}^\top\|^2 \\ &= \operatorname{tr} \left[(X - Xw\mathbf{1}^\top)^\top (X - Xw\mathbf{1}^\top) \right] \\ &= \operatorname{tr} \left[X^\top X - 2X^\top Xw\mathbf{1}^\top + \mathbf{1}w^\top X^\top Xw\mathbf{1}^\top \right] \\ &= \text{const} - 2 \operatorname{tr} [X^\top Xw\mathbf{1}^\top] + \operatorname{tr} [\mathbf{1}w^\top X^\top Xw\mathbf{1}^\top] \end{aligned}$$

At this point, we can resort to two results which you are supposed to prove below. The two traces can be rewritten such that our objective function reads

$$\begin{aligned} h(w) &= \text{const} - 2 \mathbf{1}^\top X^\top Xw + \mathbf{1}^\top \mathbf{1} w^\top X^\top Xw \\ &= n w^\top X^\top Xw - 2 \mathbf{1}^\top X^\top Xw + \text{const} \end{aligned}$$

This, in turn, gives us the following gradient

$$\nabla h = \frac{dh}{dw} = 2n X^\top Xw - 2 X^\top X \mathbf{1}$$

Convince yourself that equating ∇h to 0 and then solving for w leads to the familiar result

$$\hat{w} = \frac{1}{n} [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{1} = \frac{1}{n} \mathbf{1}$$

If you are wondering why we derived this result yet again, then note that it was to familiarize ourselves with the kind of linear algebra needed in the subsequent tasks ...

here is your task

Let X and w be as above and $z \in \mathbb{R}^n$ be any n -dimensional vector. Prove the following two identities¹

$$\text{tr}[\mathbf{X}^\top \mathbf{X} w z^\top] = z^\top \mathbf{X}^\top \mathbf{X} w$$

$$\text{tr}[z w^\top \mathbf{X}^\top \mathbf{X} w z^\top] = z^\top z \cdot w^\top \mathbf{X}^\top \mathbf{X} w$$

¹In case you are tackling this task with an AIs please tell us about your experiences! Did it work? Did the AI reason correctly?

task 4.3 [10 points]**finding maximally different data points**

In the previous tasks, you merely had to do what the problem specifications asked you to do. In this task, you need to get creative!

Let $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ be a sample of n data points $\mathbf{x}_j \in \mathbb{R}^m$ which we may gather in a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$. Now, *invent* an algorithm that selects $2 \leq k < n$ of the n columns of \mathbf{X} such that the selected data vectors are as far apart as possible.

Mathematically, the problem considered in this task is a subset selection problem and we can formalize it as follows: given $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, solve

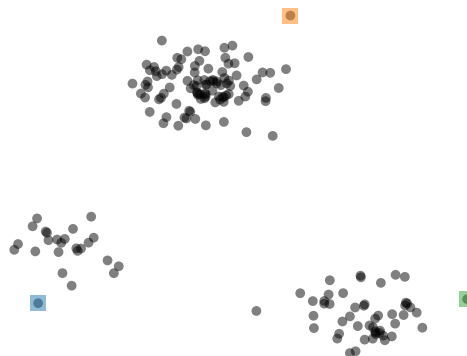
$$\hat{\mathcal{S}} = \underset{\mathcal{S} \subset \mathcal{X}}{\operatorname{argmax}} \sum_{\mathbf{x}_i \in \mathcal{S}} \sum_{\mathbf{x}_j \in \mathcal{S}} \|\mathbf{x}_i - \mathbf{x}_j\|^2$$

s.t. $|\mathcal{S}| = k$

Note: this problem is actually way more difficult than it may appear at first sight; in fact, it is NP-hard and an *efficient* algorithm that is *guaranteed* to find the optimal $\hat{\mathcal{S}}$ for very large sets \mathcal{X} and arbitrary choices of k remains elusive to this date. Having said this, one can definitely come up with efficient and well working heuristics and your instructors are looking forward to your ideas ...

task 4.3.1 [5 points]

Implement your algorithm and run it on the data in file `threeBlobs.csv`. Test your algorithm for several choices of k , say $k \in \{3, 4, 5\}$, and plot the data (say as black dots) as well as the selected columns (say as colored squares). For instance, for $k = 3$, your result could look like this:



task 4.3.2 [5 points]

Next, run your algorithm on some more interesting data. The file

`faceMatrix.npy`

contains a data matrix $\mathbf{X} \in \mathbb{R}^{361 \times 2429}$ whose columns x_j represent tiny face images of size 19×19 pixels. To read this matrix into memory and check its size, you may use

```
matX = np.load('faceMatrix.npy').astype('float')
m, n = matX.shape
print(m, n)
```

To have a look at one of the images it contains, say x_{15} , you may use this

```
vecX = matX[:,14].reshape(19,19)
plt.imshow(vecX, cmap='gray')
plt.xticks([])
plt.yticks([])
plt.show()
```

Run your algorithm on this new data matrix \mathbf{X} and consider several choices for k , say $k \in \{4, 9, 16, 25, 49, 100\}$. Visualize the extracted columns in terms of tiny images. For instance, for $k = 25$, your result could look like this:



Note: We ordered the extracted maximally diverse data points / images according to their overall brightness. You do not have to this. But if your results look completely different from ours, your algorithm / implementation may have issues.

task 4.4 [10 points]**unconventional k -means clustering (part 1)**

Note: the following background material is for those who are not familiar with k -means clustering. If you know about it, you may safely skip ahead.

background

k -means clustering is without a doubt the most popular idea for clustering a data set $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^m$ into k clusters $\mathcal{C}_1, \dots, \mathcal{C}_k \subset \mathcal{X}$ which are disjoint, i.e. $\mathcal{C}_i \cap \mathcal{C}_j = \emptyset$ for $i \neq j$, and cover the data, i.e. $\mathcal{C}_1 \cup \mathcal{C}_2 \cup \dots \cup \mathcal{C}_k = \mathcal{X}$. The two defining properties of k -means clustering are

- a) it represents each cluster \mathcal{C}_i in terms of its *centroid* or mean $\boldsymbol{\mu}_i \in \mathbb{R}^m$
- b) it clusters according to Euclidean distances between centroids and data points; that is, a point \mathbf{x}_j will be assigned to cluster \mathcal{C}_i if its distance to $\boldsymbol{\mu}_i$ is less than that to any other centroid.

A common formalization of the k -means clustering problem is as follows

$$\hat{\boldsymbol{\mu}}_1, \dots, \hat{\boldsymbol{\mu}}_k = \underset{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k}{\operatorname{argmin}} \sum_{i=1}^k \sum_{\mathbf{x}_j \in \mathcal{C}_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \quad (6)$$

Observe that this poses a chicken-and-egg problem: in order to compute cluster means, we need clusters and in order to compute clusters, we need cluster means.

Further observe that k -means clustering is NP-hard so that any presently known algorithm for solving (6) is but a heuristic. A well known k -means heuristic is Lloyd's algorithm

"randomly" initialize $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$

repeat

for $i = 1, \dots, k$

$$\mathcal{C}_i = \left\{ \mathbf{x} \in \mathcal{X} \mid \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 \leq \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\}$$

// compute clusters

for $i = 1, \dots, k$

$$\boldsymbol{\mu}_i = \frac{1}{|\mathcal{C}_i|} \sum_{\mathbf{x} \in \mathcal{C}_i} \mathbf{x}$$

// update means

until convergence

which is used by the k -means functions available in [scipy.cluster.vq](#).

an alternative k -means formalization

The objective function in (6) can be expressed differently. Indeed, (try to) convince yourself that we have the following equivalence

$$\sum_{i=1}^k \sum_{\mathbf{x}_j \in \mathcal{C}_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \quad \Leftrightarrow \quad \sum_{i=1}^k \sum_{j=1}^n z_{ij} \cdot \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

where the additional variables $z_{ij} \in \{0, 1\}$ are **indicator variables** such that

$$z_{ij} = \begin{cases} 1, & \text{if } \mathbf{x}_j \in \mathcal{C}_i \\ 0, & \text{otherwise} \end{cases}$$

The k -means clustering problem can therefore also be formalized as

$$\begin{aligned} \hat{\boldsymbol{\mu}}_1, \dots, \hat{\boldsymbol{\mu}}_k = \operatorname{argmin}_{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k} & \sum_{i=1}^k \sum_{j=1}^n z_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \\ & z_{ij} \in \{0, 1\} \\ \text{s.t.} & \sum_{i=1}^k z_{ij} = 1 \end{aligned} \quad (7)$$

where the sum-to-one constraints enforce the goal of assigning each data point \mathbf{x}_j to only one cluster \mathcal{C}_i .

an unconventional k -means formalization

Next, things will get more interesting. One can prove² the following equality

$$\sum_{i=1}^k \sum_{j=1}^n z_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 = \|\mathbf{X} - \mathbf{M}\mathbf{Z}\|^2$$

where

$\mathbf{X} \in \mathbb{R}^{m \times n}$ is a matrix of data vectors $\mathbf{x}_j \in \mathbb{R}^m$

$\mathbf{M} \in \mathbb{R}^{m \times k}$ is a matrix of cluster centroids $\boldsymbol{\mu}_i \in \mathbb{R}^m$

$\mathbf{Z} \in \mathbb{R}^{k \times n}$ is a matrix of binary indicator variables

$$Z_{ij} = \begin{cases} 1, & \text{if } \mathbf{x}_j \in \mathcal{C}_i \\ 0, & \text{otherwise.} \end{cases}$$

²C. Bauckhage, “ k -Means Clustering Is Matrix Factorization”, arXiv:1512.07548, 2015

Note that we have two crucial restrictions on matrix Z . First of all, it must be binary, i.e. $Z \in \{0, 1\}^{k \times n}$. Second of all, since each data point can only belong to one cluster, each of its columns must sum to one. If we introduce the vectors $\mathbf{1}_k$ and $\mathbf{1}_n$ of all ones in \mathbb{R}^k and \mathbb{R}^n , respectively, then this second constrained can be formalized as $\mathbf{1}_k^\top Z = \mathbf{1}_n^\top$.

With all this in place, we can also formalize the k -means problem as

$$\begin{aligned} \hat{M}, \hat{Z} = \underset{M, Z}{\operatorname{argmin}} \quad & \|X - MZ\|^2 \\ \text{s.t.} \quad & Z \in \{0, 1\}^{k \times n} \\ & Z \geq \mathbf{0}_{k \times n} \\ & \mathbf{1}_k^\top Z = \mathbf{1}_n^\top \end{aligned} \tag{8}$$

where we write $\mathbf{0}_{k \times n}$ to denote the $k \times n$ matrix of all zeros.

Speaking of this matrix, we are well aware that the constraints $Z \in \{0, 1\}^{k \times n}$ and $Z \geq \mathbf{0}_{k \times n}$ are redundant. Our reason for introducing this redundancy will become apparent later . . .

For the objective function in (8), we have the following gradients w.r.t. Z and M

$$\begin{aligned} \nabla_Z \|X - MZ\|^2 &= G_Z = 2 [M^\top MZ - M^\top X] \\ \nabla_M \|X - MZ\|^2 &= G_M = 2 [MZZ^\top - XZ^\top] \end{aligned}$$

Equating these to zero and solving for Z and M , respectively, we find

$$\begin{aligned} Z &= [M^\top M]^{-1} M^\top X \\ M &= XZ^\top [ZZ^\top]^{-1} \end{aligned}$$

Naïvely thinking, this would thus suggest the following k -means algorithm

“randomly” initialize M

for $T = 1, \dots, T_{\max}$

$$Z = [M^\top M]^{-1} M^\top X \quad \text{// compute clusters}$$

$$M = XZ^\top [ZZ^\top]^{-1} \quad \text{// update means}$$

Alas, this naïve algorithm will produce unreasonable results (try it) since it does not pay attention to the constraints on Z . So how can we incorporate those?

Assuming we had reasonable guesses for M and Z , we can generalize what we learned above and observe that

$$t_{\max} = 1$$

iterations of the following Frank-Wolfe procedure

procedure FW_UPDATE_Z($X, M, Z = [z_1, \dots, z_n], t_{\max}$)

for $t = 0, \dots, t_{\max} - 1$

$$G_Z = 2 [M^\top M Z - M^\top X]$$

for $j = 1, \dots, n$

$$o = \operatorname{argmin}_i [G_Z]_{ij}$$

$$z_j = z_j + \frac{2}{t+2} [e_o - z_j]$$

return Z

will “naturally” produce a matrix Z that meets the constraints in problem 8. Based on this insight, we can then use the following procedure for k -means clustering

procedure FW_KMEANS_VERSION1($X \in \mathbb{R}^{m \times n}, k, T_{\max}$)

“randomly” initialize matrix $M \in \mathbb{R}^{m \times k}$

for $T = 0, \dots, T_{\max} - 1$

$$Z = \frac{1}{k} \mathbf{1}_{k \times n}$$

$$Z = \text{FW_UPDATE_Z}(X, M, Z, t_{\max} = 1)$$

$$M = X Z^\top [Z Z^\top]^{-1}$$

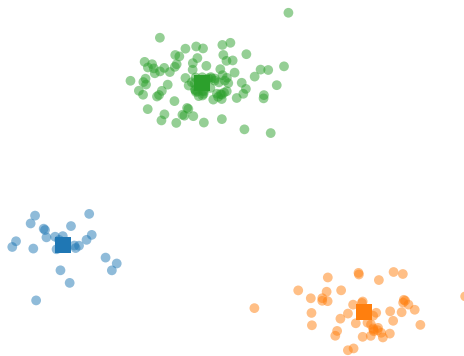
return M, Z

where we write $\mathbf{1}_{k \times n}$ to denote the $k \times n$ matrix of all ones.

Finally ...

task 4.4.1 [8 points]

Implement functions `FW_update_Z` and `FW_kMeans_Version1` and test this approach to k -means clustering on the data in file `threeBlobs.csv`. Visualize your clustering results. For example, for $k = 3$, your results could look like this



Note: if you are clever, you can avoid the `for` loop over the j in function `FW_update_Z`. Spend some time on thinking of how this could be achieved; your code will be much faster . . .



Note: k -means clustering is often very sensitive w.r.t. the initialization of the cluster centers $M = [\mu_1, \dots, \mu_k]$. A reasonable idea is to randomly select k columns of the data matrix X .

However . . . don't be surprised if your program crashes. What we are doing in procedure `FW_KMEANS_VERSION1` involves matrix inversion which can sometimes lead to numerical instabilities. A more elaborate but generally safer idea for initializing M would be to use your code from task 4.3 to extract k very diverse columns from X .

task 4.4.2 [2 points]

Try to run your code on the data in `faceMatrix.npy` and visualize the resulting mean faces. The numerical stability issues we mentioned above may be a headache here. Nevertheless, using initialization with maximally diverse data points, your instructors managed to get the following result for $k = 16$



task 4.5 [10 points]**unconventional k -means clustering (part 2)**

Let's dial things up a notch! We just considered the objective $\|X - MZ\|^2$ and found that M can be written as $M = XZ^\top [ZZ^\top]^{-1}$. But this is to say that we may also formalize the k -means clustering problem as

$$\begin{aligned} \hat{Z} = \underset{Z}{\operatorname{argmin}} \quad & \|X - XZ^\top [ZZ^\top]^{-1} Z\|^2 \\ \text{s.t.} \quad & Z \in \{0, 1\}^{k \times n} \\ & Z \geq \mathbf{0}_{k \times n} \\ & \mathbf{1}_k^\top Z = \mathbf{1}_n^\top \end{aligned} \quad (9)$$

If we could solve this problem, we could finally obtain $\hat{M} = X\hat{Z}^\top [\hat{Z}\hat{Z}^\top]^{-1}$. Alas, the problem in (10) is a truly formidable problem (to see why, just try to compute the gradient w.r.t. Z).

There are complicated ways for addressing this issue, but there also is a relatively simple one. We may introduce the $n \times k$ matrix

$$Y = Z^\top [ZZ^\top]^{-1}$$

for which one can show the following: If Z is stochastic ($Z \geq \mathbf{0}_{k \times n}$ and $\mathbf{1}_k^\top Z = \mathbf{1}_n^\top$), then Y is stochastic ($Y \geq \mathbf{0}_{n \times k}$ and $\mathbf{1}_n^\top Y = \mathbf{1}_k^\top$).

But this is to say that there is yet another formalization of the k -means problem, namely

$$\begin{aligned} \hat{Y}, \hat{Z} = \underset{Y, Z}{\operatorname{argmin}} \quad & \|X - XYZ\|^2 \\ \text{s.t.} \quad & Z \in \{0, 1\}^{k \times n} \\ & Z \geq \mathbf{0}_{k \times n} \\ & \mathbf{1}_k^\top Z = \mathbf{1}_n^\top \\ & Y \geq \mathbf{0}_{n \times k} \\ & \mathbf{1}_n^\top Y = \mathbf{1}_k^\top \end{aligned} \quad (10)$$

For the objective function in (10), we have the following gradient w.r.t. Y

$$\nabla_Y \|X - XYZ\|^2 = G_Y = 2 [X^\top XYZZ^\top - X^\top XZ^\top]$$

This allows us to design the following Frank-Wolfe procedure for updating Y based on reasonable guesses for Y and Z ...

```

procedure FW_UPDATE_Y( $\mathbf{X}$ ,  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_k]$ ,  $\mathbf{Z}$ ,  $t_{\max}$ )

  for  $t = 0, \dots, t_{\max} - 1$ 

     $\mathbf{G}_Y = 2 [\mathbf{X}^\top \mathbf{X} \mathbf{Y} \mathbf{Z} \mathbf{Z}^\top - \mathbf{X}^\top \mathbf{X} \mathbf{Z} \mathbf{Z}^\top]$ 

    for  $i = 1, \dots, k$ 

       $o = \operatorname{argmin}_j [\mathbf{G}_Y]_{ji}$ 

       $\mathbf{y}_i = \mathbf{y}_i + \frac{2}{t+2} [\mathbf{e}_o - \mathbf{y}_i]$ 

  return  $\mathbf{Y}$ 

```

This, in turn, can be used in the following procedure for k -means clustering

```

procedure FW_KMEANS_VERSION2( $\mathbf{X} \in \mathbb{R}^{m \times n}$ ,  $k$ ,  $T_{\max}$ )

  “randomly” initialize matrix  $\mathbf{M} \in \mathbb{R}^{m \times k}$ 

  for  $T = 0, \dots, T_{\max} - 1$ 

     $\mathbf{Z} = \frac{1}{k} \mathbf{1}_{k \times n}$ 

     $\mathbf{Z} = \text{FW\_UPDATE\_Z}(\mathbf{X}, \mathbf{M}, \mathbf{Z}, t_{\max} = 1)$ 

     $\mathbf{Y} = \frac{1}{n} \mathbf{1}_{n \times k}$ 

     $\mathbf{Y} = \text{FW\_UPDATE\_Y}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}, t_{\max} = 100)$ 

     $\mathbf{M} = \mathbf{X} \mathbf{Y}$ 

  return  $\mathbf{M}, \mathbf{Y}, \mathbf{Z}$ 

```

where we observe that the t_{\max} in the call of FW_UPDATE_Z differs from the t_{\max} in the call of FW_UPDATE_Y.

We also observe that procedure FW_kMeans_Version2 does not involve matrix inversion steps anymore!!!!

here is your task

Reuse your above code for function `FW_update_Z` and implement two new functions `FW_update_Y` and `FW_kMeans_Version2`. Test your combined new code on the data in file `threeBlobs.csv`. Visualize your clustering results. They should look like those in task 4.4.1.

Also test your code on the data in `faceMatrix.npy` and visualize the resulting mean faces. These should look like those in task 4.4.2.

Overall, we expect that you will find that this approach towards k -means works more robustly than the one in task 4.4.

task 4.6 [10 points]**archetypal analysis**

Above, we expressed the k -means clustering as the following constrained optimization problem

$$\begin{aligned} \hat{\mathbf{Y}}, \hat{\mathbf{Z}} = \operatorname{argmin}_{\mathbf{Y}, \mathbf{Z}} \quad & \|\mathbf{X} - \mathbf{X}\mathbf{Y}\mathbf{Z}\|^2 \\ & \mathbf{Z} \in \{0, 1\}^{k \times n} \\ & \mathbf{Z} \geq \mathbf{0}_{k \times n} \\ \text{s.t.} \quad & \mathbf{1}_k^\top \mathbf{Z} = \mathbf{1}_n^\top \\ & \mathbf{Y} \geq \mathbf{0}_{n \times k} \\ & \mathbf{1}_n^\top \mathbf{Y} = \mathbf{1}_k^\top \end{aligned}$$

Now, let us ever so slightly but very profoundly change the problem to

$$\begin{aligned} \hat{\mathbf{Y}}, \hat{\mathbf{Z}} = \operatorname{argmin}_{\mathbf{Y}, \mathbf{Z}} \quad & \|\mathbf{X} - \mathbf{X}\mathbf{Y}\mathbf{Z}\|^2 \\ & \mathbf{Z} \geq \mathbf{0}_{k \times n} \\ \text{s.t.} \quad & \mathbf{1}_k^\top \mathbf{Z} = \mathbf{1}_n^\top \\ & \mathbf{Y} \geq \mathbf{0}_{n \times k} \\ & \mathbf{1}_n^\top \mathbf{Y} = \mathbf{1}_k^\top \end{aligned} \tag{11}$$

Try to spot the difference! The problem in (11) is called **archetypal analysis** and is a powerful tool in data exploration. To solve it, we suggest

procedure FW_ARCHETYPAL_ANALYSIS($\mathbf{X} \in \mathbb{R}^{m \times n}, k, T_{\max} = 100$)

“randomly” initialize matrix $\mathbf{A} \in \mathbb{R}^{m \times k}$

for $T = 0, \dots, T_{\max} - 1$

$\mathbf{Z} = \text{FW_UPDATE_Z}(\mathbf{X}, \mathbf{A}, \frac{1}{k} \mathbf{1}_{k \times n}, t_{\max} = 100)$

$\mathbf{Y} = \text{FW_UPDATE_Y}(\mathbf{X}, \frac{1}{n} \mathbf{1}_{n \times k}, \mathbf{Z}, t_{\max} = 100)$

$\mathbf{A} = \mathbf{X}\mathbf{Y}$

return $\mathbf{A}, \mathbf{Y}, \mathbf{Z}$

Observe that parameter t_{\max} in the call of FW_UPDATE_Z has changed!

here is your task

Reuse your above code for *FW_update_Z* and *FW_update_Y* and implement *FW_Archetypal_Analysis*. Test your combined code on the data in file `threeBlobs.csv`. Plot the data in X (say as black dots) and the **archetypes** in A (say as colored squares). Also test your code on the data in `faceMatrix.npy` and visualize the resulting archetypal faces.

We won't spoil the fun and show how your results should look like. But we suggest you (visually) compare them against your results from tasks 4.3 and 4.5 ... What do you observe?

task 4.7

submission of presentation and code

As always, prepare a presentation / set of slides on your solutions for all the mandatory tasks and submit them to eCampus. Also as always, put your code into a ZIP file and submit it to eCampus.