

Clustering algorithms

Konstantinos Koutroumbas

Unit 6

- Fuzzy CFO clustering algorithms
- Possibilistic CFO clust. Algorithms

Fuzzy CFO clustering algorithms

Fuzzy clustering algorithms:

Let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be a set of data points.

Each vector \mathbf{x}_i belongs to all clusters up to a certain degree, $u_{ij}, j = 1, \dots, m$,

Subject to the constraints

- $u_{ij} \in [0,1], i = 1, \dots, N, j = 1, \dots, m$
- $\sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$
- $0 < \sum_{i=1}^N u_{ij} < N, j = 1, \dots, m$

Each cluster is represented by a representative θ_j (point repr., hyperplane...).

Let $\Theta = \{\theta_1, \theta_2, \dots, \theta_m\}$

Define the cost function

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \theta_j), \quad (q > 1)$$

When $J_q(U, \Theta)$ is minimized?

When large u_{ij} 's are multiplied with small $d(\mathbf{x}_i, \theta_j)$'s.

Fuzzy CFO clustering algorithms

Minimizing the **cost function**

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \text{ s.t. } \sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$$

Since $\boldsymbol{\theta}_j$'s, u_{ij} 's are **continuous valued**, tools from analysis may be employed for **both** of them.

For **fixed $\boldsymbol{\theta}_j$'s**: Define the **Lagrangian function**

$$\mathcal{L}_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) - \sum_{i=1}^N \lambda_i \left(\sum_{j=1}^m u_{ij} - 1 \right)$$

Equating the **partial derivative** of $\mathcal{L}_q(U, \Theta)$ wrt u_{rs} to 0, it turns out that

$$\frac{\partial \mathcal{L}_q(U, \Theta)}{\partial u_{rs}} = 0 \Leftrightarrow u_{rs} = \frac{1}{\sum_{j=1}^m \left(\frac{d(\mathbf{x}_r, \boldsymbol{\theta}_s)}{d(\mathbf{x}_r, \boldsymbol{\theta}_j)} \right)^{\frac{1}{q-1}}}$$

For **fixed u_{ij} 's**: Solve the following **m** independent minimization problems

$$\boldsymbol{\theta}_j = \operatorname{argmin}_{\boldsymbol{\theta}_j} \sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j)$$

Fuzzy CFO clustering algorithms

Generalized Fuzzy Algorithmic Scheme (GFAS)

- Choose $\theta_j(0)$ as initial estimates for $\theta_j, j = 1, \dots, m$.

- $t = 0$

- Repeat

- For $i = 1$ to N % Determination of u'_{ij} s

- o For $j = 1$ to m

$$u_{ij}(t) = \frac{1}{\sum_{k=1}^m \left(\frac{d(\mathbf{x}_i, \theta_j(t))}{d(\mathbf{x}_i, \theta_k(t))} \right)^{\frac{1}{q-1}}}$$

- o End {For- j }

- End {For- i }

- $t = t + 1$

- For $j = 1$ to m % Parameter updating

- o Set

$$\theta_j(t) = \operatorname{argmin}_{\theta_j} \sum_{i=1}^N u_{ij}^q(t-1) d(\mathbf{x}_i, \theta_j), j = 1, \dots, m$$

- End {For- j }

- Until a termination criterion is met.

Fuzzy CFO clustering algorithms

Remarks:

- A candidate **termination condition** is

$$||\boldsymbol{\theta}(t) - \boldsymbol{\theta}(t - 1)|| < \varepsilon,$$

where $|| \cdot ||$ is any vector norm and ε a user-defined constant.

- GFAS may also be initialized from $U(0)$ instead of $\boldsymbol{\theta}_j(0)$, $j = 1, \dots, m$ and start iterations with computing $\boldsymbol{\theta}_j$ first.
- If a point \mathbf{x}_i **coincides** with one or more **representatives**, then it is shared arbitrarily among the clusters whose representatives coincide with \mathbf{x}_i , s.t. the constraint that the summation of all u_{ij} 's sum to 1.
- The degree of membership of \mathbf{x}_i in C_j cluster is related to the grade of membership of \mathbf{x}_i in rest $m - 1$ clusters.
- If $q = 1$, **no** fuzzy clustering is better than the best hard clustering in terms of $J_q(\boldsymbol{\theta}, U)$.
- If $q > 1$, **there are** fuzzy clusterings with lower values of $J_q(\boldsymbol{\theta}, U)$ than the best hard clustering.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The point representatives case

- **Point representatives** are used in the case of **compact clusters**.
- Each θ_j consists of l parameters.
- Any dissimilarity measure $d(\mathbf{x}_i, \theta_j)$ between two points can be used.
- Common choices for $d(\mathbf{x}_i, \theta_j)$ are

$$d(\mathbf{x}_i, \theta_j) = (\mathbf{x}_i - \theta_j)^T A (\mathbf{x}_i - \theta_j),$$

where A is symmetric and positive definite matrix.

It is:

$$\frac{\partial d(\mathbf{x}_i, \theta_j)}{\partial \theta_j} = 2A(\theta_j - \mathbf{x}_i)$$

In this case the problem

$$\theta_j = \operatorname{argmin}_{\theta_j} \sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \theta_j)$$

is solved as

$$\frac{\partial}{\partial \theta_j} \sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \theta_j) = 0 \Leftrightarrow 2A \sum_{i=1}^N u_{ij}^q (\theta_j - \mathbf{x}_i) = 0 \Leftrightarrow$$

$$\theta_j = \frac{\sum_{i=1}^N u_{ij}^q \mathbf{x}_i}{\sum_{i=1}^N u_{ij}^q}$$

Fuzzy CFO clustering algorithms

GFAS – The point representative with squared Mahalanobis distance

- Choose $\theta_j(0)$ as initial estimates for $\theta_j, j = 1, \dots, m$.
- $t = 0$
- Repeat

– For $i = 1$ to N % Determination of u'_{ij} s

o For $j = 1$ to m

$$u_{ij}(t) = \frac{1}{\sum_{k=1}^m \left(\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j(t))}{d(\mathbf{x}_i, \boldsymbol{\theta}_k(t))} \right)^{\frac{1}{q-1}}}$$

o End {For- j }

– End {For- i }

– $t = t + 1$

– For $j = 1$ to m % Parameter updating

o Set

$$\boldsymbol{\theta}_j(t) = \frac{\sum_{i=1}^N u_{ij}^q(t-1) \mathbf{x}_i}{\sum_{i=1}^N u_{ij}^q(t-1)}, j = 1, \dots, m$$

– End {For- j }

- Until a termination criterion is met.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The point representatives case

Remarks:

- GFAS with the Euclidean distance ($A = I$) is also known as **Fuzzy c-Means (FCM)** or **Fuzzy k-Means** algorithm.
- FCM **converges** to a **stationary point** of the cost function or it has at least one subsequence that converges to a stationary point. This point may be a local (or global) minimum or a saddle point.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The point representatives case

Example:

Generate and plot the data set X7, which consists of $N=216$ 2-dim. vectors. The first 100 stem from the normal distribution with mean $\mathbf{m}_1=[0, 0]^T$, the next 100 stem from the normal distribution with mean $\mathbf{m}_2=[13, 13]^T$. The other two groups of eight points each stem from the normal distribution with means $\mathbf{m}_3=[0, -40]^T$ and $\mathbf{m}_4=[-30, -30]^T$, respectively. The covariance matrices for all distributions are all equal to the 2x2 identity matrix. Obviously, the last two groups of points may be considered as outliers.

Apply the FCM on the data set X7 with $m=2$ clusters, plot the results and comment on the grade of memberships of the vectors to the two obtained clusters.

Apply also the k-means and the PAM on X7 and compare the results obtained from the three algorithms. ([SEE attached code](#))

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- Here the representatives are **quadric surfaces** (hyperellipsoids, hyperparaboloids, etc.)
- **First issue:** How to **represent** them?
- **General forms** of an equation describing a **quadric surface** Q :

1. $\mathbf{x}^T A \mathbf{x} + \mathbf{b}^T \mathbf{x} + c = 0,$

where A is an $l \times l$ symmetric matrix, \mathbf{b} is an $l \times 1$ vector, c is a scalar and $\mathbf{x} = [x_1, \dots, x_l]^T$.

For various choices of A , \mathbf{b} and c we obtain hyperellipses, hyperparabolas and so on.

2. $\mathbf{q}^T \mathbf{p} = 0,$

where

$$\mathbf{q} = [x_1^2, x_2^2, \dots, x_l^2, x_1 x_2, \dots, x_{l-1} x_l, x_1, x_2, \dots, x_l, 1]^T$$

and

$$\mathbf{p} = [p_1, p_2, \dots, p_l, p_{l+1}, \dots, p_r, p_{r+1}, \dots, p_s]^T$$

with $r = \frac{l(l+1)}{2}$ and $s = r + l + 1$.

NOTE: The above **representations** of Q are **equivalent**.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- **Second issue**: “Definition of the **distance** of a point **\mathbf{x}** to a quadric surface **Q** ”

Types of distances

- **Perpendicular distance**:

$$d_p^2(\mathbf{x}, Q) = \min_{\mathbf{z}} \|\mathbf{x} - \mathbf{z}\|^2,$$

subject to the constraint

$$\mathbf{z}^T \mathbf{A} \mathbf{z} + \mathbf{b}^T \mathbf{z} + c = 0$$

In words, $d_p^2(\mathbf{x}, Q)$ is the **distance** between **\mathbf{x}** and the **closest to \mathbf{x} point** that lies in **Q** .

Prove it for the $l = 2$ case.

- (Squared) **Algebraic distance**:

$$d_p^2(\mathbf{x}, Q) = (\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c)^2 \equiv \mathbf{p}^T \mathbf{M} \mathbf{p}$$

where $\mathbf{M} = \mathbf{q} \mathbf{q}^T$.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

– **Radial distance** (**only** when Q is a **hyperellipsoidal**):

For Q **hyperellipsoidal**, the representative **equation** can be written as

$$(\mathbf{x} - \mathbf{c})^T A (\mathbf{x} - \mathbf{c}) = 1$$

where \mathbf{c} is the **center** of the ellipse and A a **positive definite symmetric** matrix defining major axis, minor axis and **orientation**.

Then the **radial distance** is defined as

$$d_r^2(\mathbf{x}, Q) = \|\mathbf{x} - \mathbf{z}\|^2$$

subject to the **constraints**

$$(\mathbf{z} - \mathbf{c})^T A (\mathbf{z} - \mathbf{c}) = 1$$

and

$$(\mathbf{z} - \mathbf{c}) = a(\mathbf{x} - \mathbf{c}).$$

In words,

- the **intersection point** \mathbf{z} between the **line segment** $\mathbf{x} - \mathbf{c}$ and Q is determined
- the $d_r^2(\mathbf{x}, Q)$ is defined as the **squared Euclidean distance** between \mathbf{x} and \mathbf{z} .

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

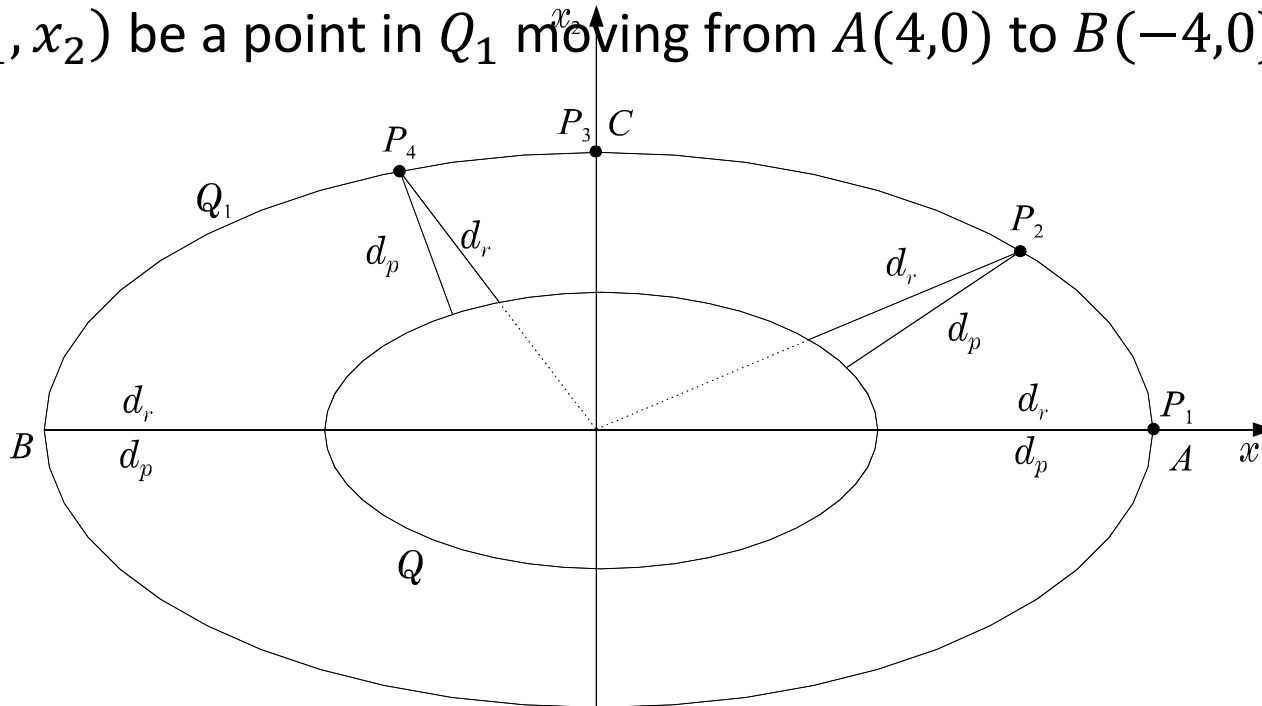
- (Squared) **Normalized radial distance** (only when Q is a **hyperellipsoidal**):

$$d_{nr}^2(\mathbf{x}, Q) = \left(((\mathbf{x} - \mathbf{c})^T A (\mathbf{x} - \mathbf{c}))^{1/2} - 1 \right)^2$$

$d_r^2(\mathbf{x}, Q) = d_{nr}^2(\mathbf{x}, Q) \|\mathbf{x} - \mathbf{z}\|^2$
 \mathbf{z} : **intersection** of $\mathbf{x} - \mathbf{c}$ and Q .

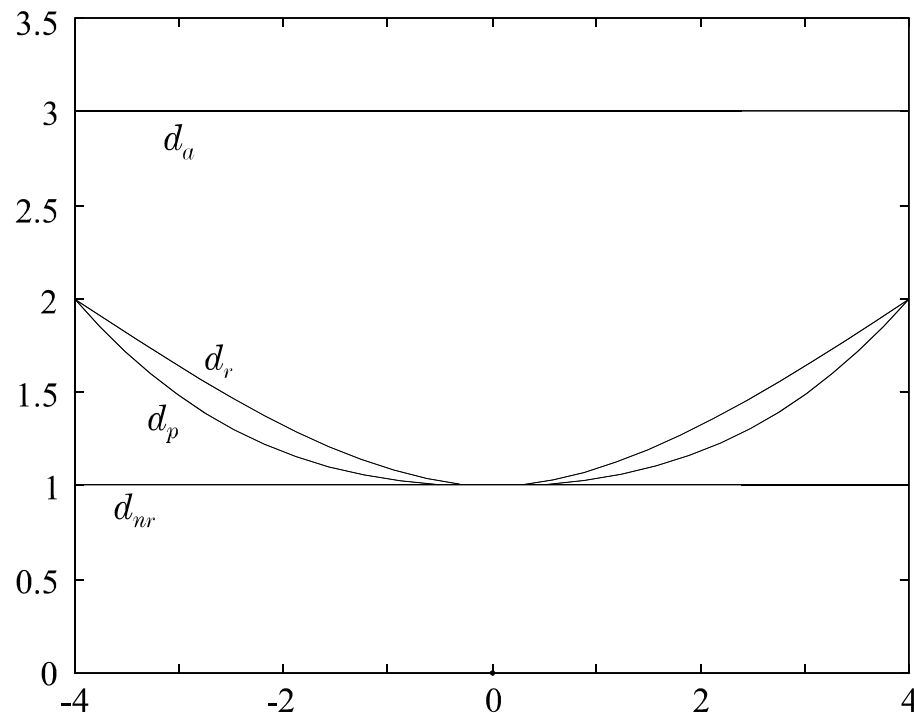
– Example 3:

- Consider two ellipses Q and Q_1 , centered at $\mathbf{c} = [0, 0]^T$, with $A = \text{diag}(\frac{1}{4}, 1)$ and $A_1 = \text{diag}(\frac{1}{16}, \frac{1}{4})$, respectively.
- Let $P(x_1, x_2)$ be a point in Q_1 moving from $A(4, 0)$ to $B(-4, 0)$, with $x_2 > 0$



Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case



Remarks:

- d_a and d_{nr} **do not vary** as P moves.
- d_r can be used as an **approximation** of d_p , when Q is a hyperellipsoid.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- **Third issue:** Choice of algorithm.

Recall that

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \text{ s.t. } \sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$$

- The **algorithms** in this case fall under the umbrella of **GFAS**.
- They all **share** the **same rule** for **updating** the matrix **U** .
- They **differ** on the choice of the **distance** between a **point** and the **representative** of a quadric surface.
 \Rightarrow they **differ** in the **representatives updating part**.
- At each iteration, the **updating** of the **representatives** is carried out by **setting** the **gradient** of J_q wrt them **equal** to **0** (for **fixed** u_{ij} 's) and **solving** (usually using iterative schemes) **for** the involved **parameters**.

Fuzzy CFO clustering algorithms

Generalized Fuzzy Algorithmic Scheme (GFAS)

- Choose $\theta_j(0)$ as initial estimates for $\theta_j, j = 1, \dots, m$.

- $t = 0$

- Repeat

- For $i = 1$ to N % Determination of u'_{ij} s

- o For $j = 1$ to m

$$u_{ij}(t) = \frac{1}{\sum_{k=1}^m \left(\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j(t))}{d(\mathbf{x}_i, \boldsymbol{\theta}_k(t))} \right)^{\frac{1}{q-1}}}$$

- o End {For- j }

- End {For- i }

- $t = t + 1$

- For $j = 1$ to m % Parameter updating

- o Set

$$\boldsymbol{\theta}_j = \operatorname{argmin}_{\boldsymbol{\theta}_j} \sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j), j = 1, \dots, m$$

- End {For- j }

- Until a termination criterion is met.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- **Third issue:** Choice of algorithm.

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \text{ s.t. } \sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$$

Algorithms:

- Fuzzy C Ellipsoidal Shells (FCES) Algorithm:
 - It **adopts** the **radial distance** between a vector and the surface representative
 - It **recovers** only **ellipsoidal clusters**.
- Fuzzy C Quadric Shells (FCQS) Algorithm:
 - It **adopts** the **algebraic distance** between a vector and the surf. repr. in the form $d_a^2(\mathbf{x}, Q) = \mathbf{p}^T M \mathbf{p}$, imposing **constraints** on vector **\mathbf{p}** .
 - It **recovers** **quadric clusters** of any kind (**ellipsoidal**, **hyperbolical**, **paraboloidal**, **pairs of lines**).

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- **Third issue:** Choice of algorithm.

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \text{ s.t. } \sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$$

Algorithms:

- Modified Fuzzy C Quadric Shells (MFCQS) Algorithm:
 - It **adopts** :
 - the **perpendicular distance** between a vector and the surface representative for the **updating** of matrix ***U***
 - The **algebraic distance** between a vector and the surface representative for the **updating** of the **cluster representatives**.
 - It **recovers** **quadric clusters** of any kind (**ellipsoidal**, **hyperbolical**, **paraboloidal**, **pairs of lines**).

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The hyperplane surfaces representatives case

- Here the representatives are **hyperplanes** (lines in the 2-D space, planes in the 3-D space etc.)
- **First issue:** How to **represent** them?

1. Via the **equation** of a **hyperplane** H :

$$H: \boldsymbol{\theta}^T \mathbf{x} + \theta_0 = 0,$$

where $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_l]^T$, $\mathbf{x} = [x_1, x_2, \dots, x_l]^T$.

2. Via a **center** \mathbf{c}_j and a **covariance matrix** Σ_j , that is, $\boldsymbol{\theta}_j = (\mathbf{c}_j, \Sigma_j)$.

NOTE: Another choice for **representing** such clusters is by using **line segments**. (only for the 2-D case).

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The hyperplane surfaces representatives case

- **Second issue**: “Definition of the *distance* of a point \mathbf{x} to a cluster”

Types of distances

– Distance of a point from a hyperplane:

$$d(\mathbf{x}, H) = \frac{|\boldsymbol{\theta}^T \mathbf{x} + \theta_0|}{\|\boldsymbol{\theta}\|}$$

– GK distance:

$$d_{GK}^2(\mathbf{x}, \boldsymbol{\theta}_j) = |\Sigma_j|^{1/l} (\mathbf{x} - \mathbf{c}_j)^T \Sigma_j^{-1} (\mathbf{x} - \mathbf{c}_j)$$

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The hyperplane surfaces representatives case

- **Third issue:** Choice of algorithm.

Recall that

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \text{ s.t. } \sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$$

- The **algorithms** in this case fall under the umbrella of **GFAS**.
- They all **share** the **same rule** for **updating** the matrix **U** .
- They **differ** on the choice of the **distance** between a **point** and the **representative** of a plane cluster.
 \Rightarrow they **differ** in the **representatives updating part**.
- At each iteration, the **updating** of the **representatives** is carried out by **setting** the **gradient** of J_q wrt them **equal** to **0** (for **fixed** u_{ij} 's) and **solving** (usually using iterative schemes) **for** the involved **parameters**.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- **Third issue:** Choice of algorithm.

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \text{ s.t. } \sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$$

Algorithms:

- Fuzzy C varieties (FCV) Algorithm:
 - It **adopts** the **classical distance** between a point and a hyperplane.
 - Disadvantages:
 - It tends to **recover very long clusters** and, thus, **collinear distinct clusters** may be **merged** to a single **one**.
 - If, at a certain iteration, a hyperplane representative **crosses two distinct clusters**, there is **no way to recover** from this situation.

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- **Third issue:** Choice of algorithm.

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \text{ s.t. } \sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$$

Algorithms:

- Gustafson-Kessel (GK) algorithm:
 - It **adopts** the **GK distance** between a point and a cluster.
 - The parameter updating takes place via the following two equations

$$\mathbf{c}_j(t) = \frac{\sum_{i=1}^N u_{ij}^q(t-1) \mathbf{x}_i}{\sum_{i=1}^N u_{ij}^q(t-1)}$$

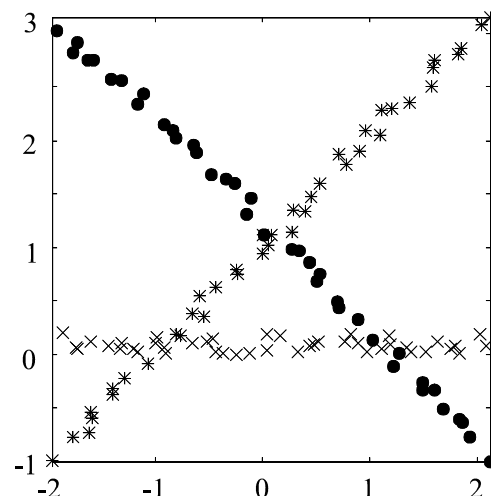
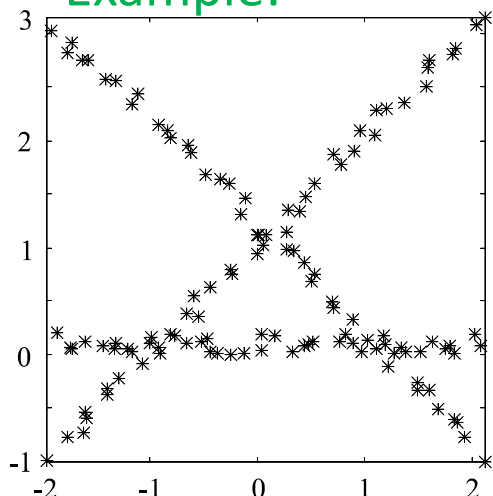
$$\Sigma_j(t) = \frac{\sum_{i=1}^N u_{ij}^q(t-1) (\mathbf{x}_i - \mathbf{c}_j(t)) (\mathbf{x}_i - \mathbf{c}_j(t))^T}{\sum_{i=1}^N u_{ij}^q(t-1)}$$

Fuzzy CFO clustering algorithms

Fuzzy Clustering – The quadric surfaces representatives case

- Gustafson-Kessel (GK) algorithm (cont.):

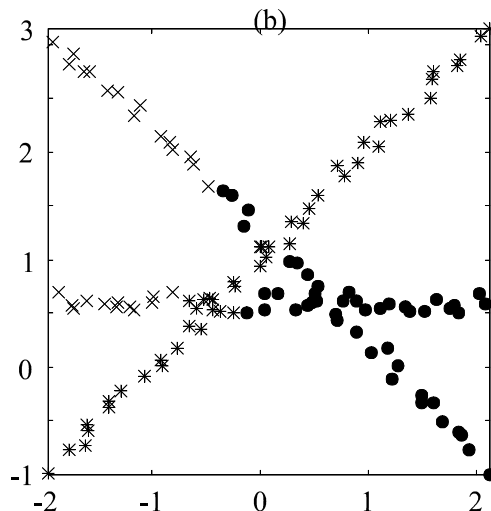
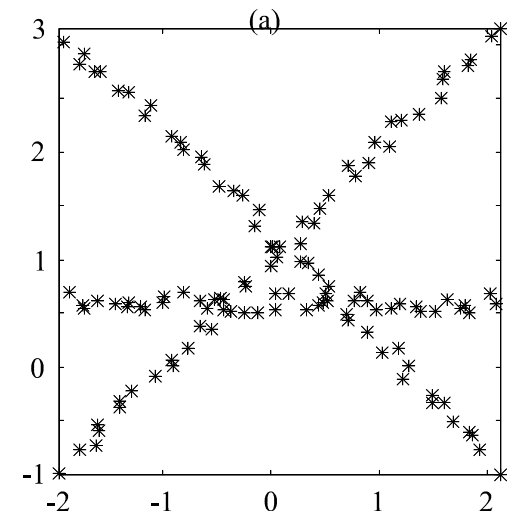
- Example:



Comments:

In the **first case**, the clusters are **well discriminated** and the **GK-algorithm recovers them correctly**.

In the **second case**, the clusters are **not well discriminated** and the **GK-algorithm fails to recover them correctly**.



(a)

(b)

Possibilistic CFO clustering algorithms

Possibilistic clustering algorithms:

Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be a set of data points.

For each vector \mathbf{x}_i its **degree of compatibility** with **all clusters**, $u_{ij}, j = 1, \dots, m$, is considered.

The **constraints** on u_{ij} 's are

- $u_{ij} \in [0,1], i = 1, \dots, N, j = 1, \dots, m$
- $0 < \sum_{i=1}^N u_{ij} < N, j = 1, \dots, m$

Each **cluster** is **represented** by a representative $\boldsymbol{\theta}_j$ (point repr., hyperplane...).

Let $\Theta = \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_m\}$

Define the **cost function**

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j)$$

When $J_q(U, \Theta)$ is **minimized**?

Possibilistic CFO clustering algorithms

Possibilistic clustering algorithms:

Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be a set of data points.

For each vector \mathbf{x}_i its **degree of compatibility** with **all clusters**, $u_{ij}, j = 1, \dots, m$, is considered.

The **constraints** on u_{ij} 's are

- $u_{ij} \in [0,1], i = 1, \dots, N, j = 1, \dots, m$
- $0 < \sum_{i=1}^N u_{ij} < N, j = 1, \dots, m$

Each **cluster** is **represented** by a representative $\boldsymbol{\theta}_j$ (point repr., hyperplane...).

Let $\Theta = \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_m\}$

Define the **cost function**

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j)$$

When $J_q(U, \Theta)$ is **minimized**?

When **all** u_{ij} 's are (very close to) **zero**.

Possibilistic CFO clustering algorithms

How to **avoid** the trivial **zero u_{ij} 's solution**?

Add a **suitable term** that discourages the zero solution.

A possible scenario:

Minimize the **cost function**

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) + \sum_{j=1}^m \eta_j \sum_{i=1}^N (1 - u_{ij})^q$$

where η_j 's are suitably defined **constants** (one for each cluster), **associated** with the **variance** of the **clusters**.

Since $\boldsymbol{\theta}_j$'s, u_{ij} 's are **continuous valued**, tools from analysis may be employed.

For **fixed $\boldsymbol{\theta}_j$'s**: Equating the **partial derivative** of **$J_q(U, \Theta)$** wrt **u_{ij}** to 0 we obtain

$$\frac{\partial J_q(U, \Theta)}{\partial u_{ij}} = 0 \Leftrightarrow u_{ij} = \frac{1}{1 + \left(\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j)}{\eta_j} \right)^{\frac{1}{q-1}}}$$

Notes: (a) u_{ij} **depends exclusively** on $\boldsymbol{\theta}_j$.

(b) It is $u_{ij} \in [0,1]$

Possibilistic CFO clustering algorithms

How to **avoid** the trivial **zero u_{ij} 's solution**?

Add a **suitable term** that discourages the zero solution.

A possible scenario:

Minimize the **cost function**

$$J_q(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) + \sum_{j=1}^m \eta_j \sum_{i=1}^N (1 - u_{ij})^q$$

where η_j 's are suitably defined **constants** (one for each cluster), **associated** with the **variance** of the **clusters**.

Since $\boldsymbol{\theta}_j$'s, u_{ij} 's are **continuous valued**, tools from analysis may be employed.

For **fixed u_{ij} 's**: Solve the following **m** independent minimization problems

$$\boldsymbol{\theta}_j = \operatorname{argmin}_{\boldsymbol{\theta}_j} \sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j)$$

Possibilistic CFO clustering algorithms

Generalized Possibilistic Algorithmic Scheme (*GPAS1*)

- Fix η_j 's, $j = 1, \dots, m$.
- Choose $\theta_j(0)$ as initial estimates for θ_j , $j = 1, \dots, m$.
- $t = 0$

- Repeat

- For $i = 1$ to N % Determination of u'_{ij} s

- o For $j = 1$ to m

$$u_{ij}(t) = \frac{1}{1 + \left(\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j(t))}{\eta_j} \right)^{\frac{1}{q-1}}}$$

- o End {For- j }

- End {For- i }

- $t = t + 1$

- For $j = 1$ to m % Parameter updating

- o Set

$$\boldsymbol{\theta}_j(t) = \operatorname{argmin}_{\boldsymbol{\theta}_j} \sum_{i=1}^N u_{ij}^q(t-1) d(\mathbf{x}_i, \boldsymbol{\theta}_j), j = 1, \dots, m$$

- End {For- j }

- Until a termination criterion is met.

Possibilistic CFO clustering algorithms

Remarks:

- A candidate **termination condition** is

$$||\boldsymbol{\theta}(t) - \boldsymbol{\theta}(t-1)|| < \varepsilon,$$

where $|| \cdot ||$ is any vector norm and ε a user-defined constant.

- GFAS may also be initialized from $U(0)$ instead of $\boldsymbol{\theta}_j(0)$, $j = 1, \dots, m$ and start iterations with computing $\boldsymbol{\theta}_j$ first.
- Based on GPAS, a possibilistic algorithm can be derived, for each fuzzy clustering algorithm derived previously.
- **High values of q :**
 - In **possibilistic clustering** **cause** almost **equal contributions** of all vectors to all clusters
 - In **fuzzy clustering** **cause increased sharing** of the vectors among all clusters.

Possibilistic CFO clustering algorithms

Three observations

- **Decomposition of $J(\theta, U)$:**

Since for each vector \mathbf{x}_i , u_{ij} 's, $j = 1, \dots, m$ are **independent** from each other, $J(\theta, U)$ can be written as

$$\begin{aligned} J(\theta, U) &= \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \theta_j) + \sum_{j=1}^m \eta_j \sum_{i=1}^N (1 - u_{ij})^q \\ &= \sum_{j=1}^m \left[\sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \theta_j) + \eta_j \sum_{i=1}^N (1 - u_{ij})^q \right] \equiv \sum_{j=1}^m J_j \end{aligned}$$

where

$$J_j = \sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \theta_j) + \eta_j \sum_{i=1}^N (1 - u_{ij})^q$$

Each J_j is **associated** with a different **cluster** and minimization of $J(\theta, U)$ with respect to u_{ij} 's can be carried out separately for each J_j .

Possibilistic CFO clustering algorithms

Three observations

- About η_j 's:
 - They **determine** the **relative significance** of the **two terms** in $J(\Theta, U)$.
 - They are **related** to the “**variance**” of the points of C_j 's, $j = 1, \dots, m$, around their centers.
 - Two scenarios for the estimation of η_j 's, for the **point representatives** case, are the following:
 - o **Run** the related FCM algorithm and after its convergence estimate η_j 's as

$$\eta_j = \frac{\sum_{i=1}^N u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j)}{\sum_{i=1}^N u_{ij}^q} \quad \text{or} \quad \eta_j = \frac{\sum_{u_{ij} > a} d(\mathbf{x}_i, \boldsymbol{\theta}_j)}{\sum_{u_{ij} > a} 1}$$

- o **Set** $\eta_j = \eta = \frac{\beta}{q\sqrt{m}}$, where $\beta = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2$ and $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$

Possibilistic CFO clustering algorithms

Three observations

- The mode-seeking property

Unlike Hard and fuzzy clustering algorithms which are **partition algorithms** (they terminate with the predetermined number of clusters no matter how many physical clusters are naturally formed in X), GPAS is a **mode-seeking algorithm** (it searches for dense regions of vectors in X).

Advantage: The number of clusters need not be a priori known.

If the number of clusters in GPAS, m , is greater than the true number of clusters k in X , some representatives will coincide with others. If $m < k$, **some** (and not all) of the clusters will be identified.

Disadvantage: Need for estimating η_j .

Possibilistic CFO clustering algorithms

How to **avoid** the trivial **zero u_{ij} 's solution**?

Add a **suitable term** that discourages the zero solution.

Another possible scenario:

Minimize the **cost function**

$$J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} d(\mathbf{x}_i, \boldsymbol{\theta}_j) + \sum_{j=1}^m \eta_j \sum_{i=1}^N (u_{ij} \ln u_{ij} - u_{ij})$$

where η_j 's are suitably defined **constants** (one for each cluster), **associated** with the **variance** of the **clusters**.

Since **$\boldsymbol{\theta}_j$'s, u_{ij} 's** are **continuous valued**, tools from analysis may be employed.

For **fixed $\boldsymbol{\theta}_j$'s**: Equating the **partial derivative** of **$J(U, \Theta)$** wrt **u_{ij}** to 0 we obtain

$$\frac{\partial J_q(U, \Theta)}{\partial u_{ij}} = 0 \Leftrightarrow u_{ij} = \exp\left(-\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j)}{\eta_j}\right)$$

Notes: (a) **u_{ij}** depends exclusively on **$\boldsymbol{\theta}_j$** .

(b) It is $u_{ij} \in [0,1]$

Possibilistic CFO clustering algorithms

How to **avoid** the trivial **zero u_{ij} 's solution**?

Add a **suitable term** that discourages the zero solution.

A possible scenario:

Minimize the **cost function**

$$J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} d(\mathbf{x}_i, \boldsymbol{\theta}_j) + \sum_{j=1}^m \eta_j \sum_{i=1}^N (u_{ij} \ln u_{ij} - u_{ij})$$

where η_j 's are suitably defined **constants** (one for each cluster), **associated** with the **variance** of the **clusters**.

Since $\boldsymbol{\theta}_j$'s, u_{ij} 's are **continuous valued**, tools from analysis may be employed.

For **fixed u_{ij} 's**: Solve the following **m** independent minimization problems

$$\boldsymbol{\theta}_j = \operatorname{argmin}_{\boldsymbol{\theta}_j} \sum_{i=1}^N u_{ij} d(\mathbf{x}_i, \boldsymbol{\theta}_j)$$

Possibilistic CFO clustering algorithms

Generalized Possibilistic Algorithmic Scheme (GPAS2)

- Fix η_j 's, $j = 1, \dots, m$.
- Choose $\theta_j(0)$ as initial estimates for θ_j , $j = 1, \dots, m$.
- $t = 0$

- Repeat

- For $i = 1$ to N % Determination of u'_{ij} s

- o For $j = 1$ to m

$$u_{ij}(t) = \exp\left(-\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j(t))}{\eta_j}\right)$$

- o End {For- j }

- End {For- i }

- $t = t + 1$

- For $j = 1$ to m % Parameter updating

- o Set

$$\boldsymbol{\theta}_j(t) = \operatorname{argmin}_{\boldsymbol{\theta}_j} \sum_{i=1}^N u_{ij}(t-1) d(\mathbf{x}_i, \boldsymbol{\theta}_j), j = 1, \dots, m$$

- End {For- j }

- Until a termination criterion is met.

Optimization theory – Basic concepts (supp. material)

Let $J(\mathbf{w})$ be a continuous function of \mathbf{w} .

Problem (P1): Determine the **position** \mathbf{w}^* where the function $J(\mathbf{w})$ achieves its **minimum** value.

A simple method for solving **(P1)** is that of **gradient descent**.

-Initialize $\mathbf{w} = \mathbf{w}(0)$

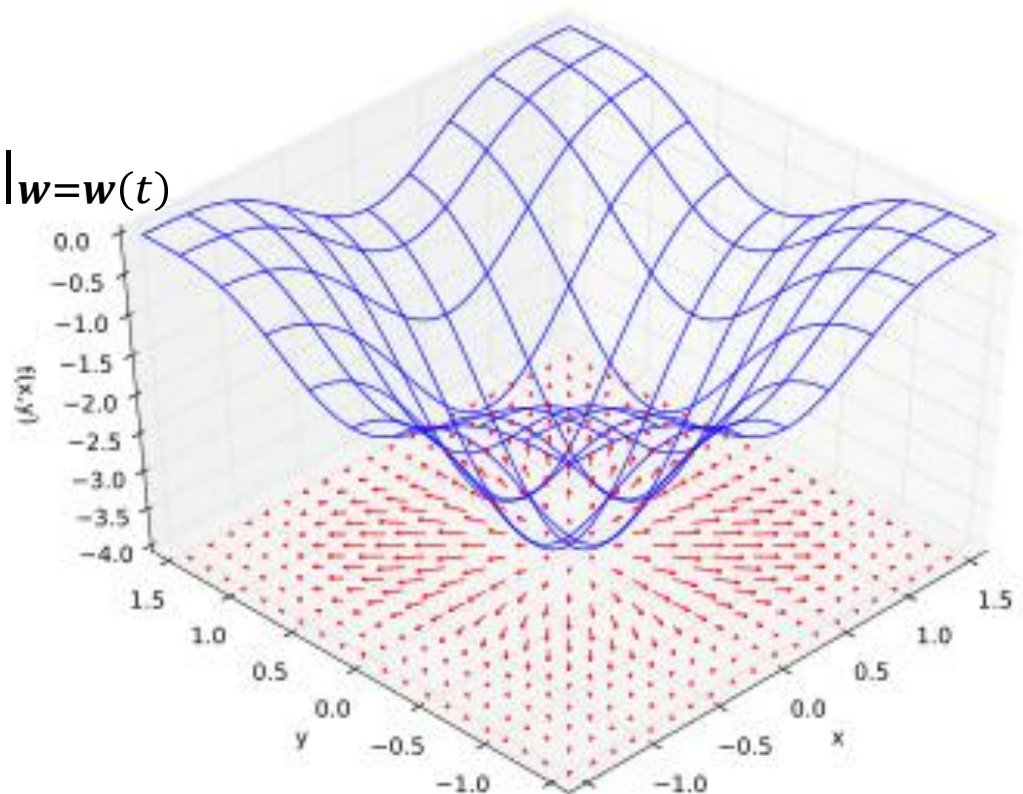
- $t = 0$

-Repeat

$$- \mathbf{w}(t + 1) = \mathbf{w}(t) - \mu \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}(t)}$$

$$- t = t + 1$$

-Until convergence



Optimization theory – Basic concepts (supp. material)

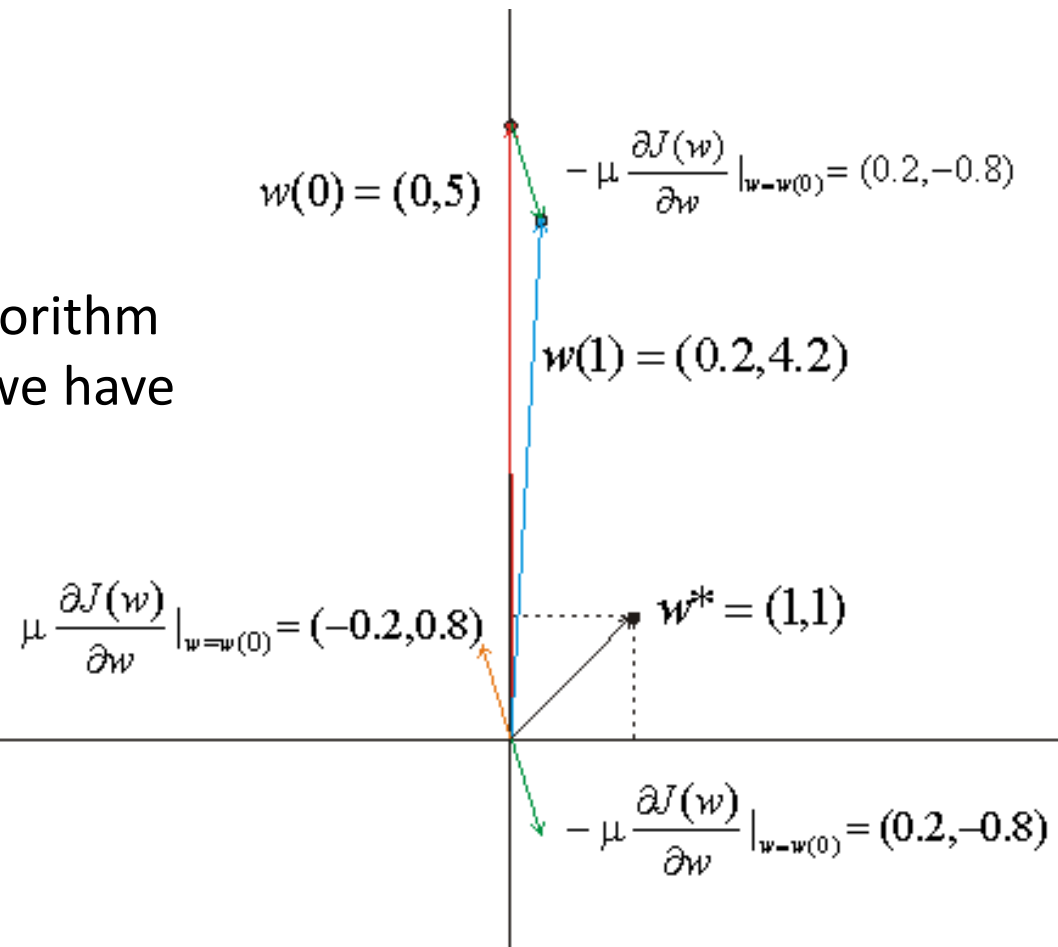
-An example: Let $\mathbf{w} = [w_1, w_2]^T$ and $J(\mathbf{w}) = (w_1 - 1)^2 + (w_2 - 1)^2$. Clearly, the minimum value of $J(\mathbf{w})$ is met at $\mathbf{w}^* = [1, 1]^T$.

-It is
$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \begin{bmatrix} 2w_1 - 2 \\ 2w_2 - 2 \end{bmatrix}$$

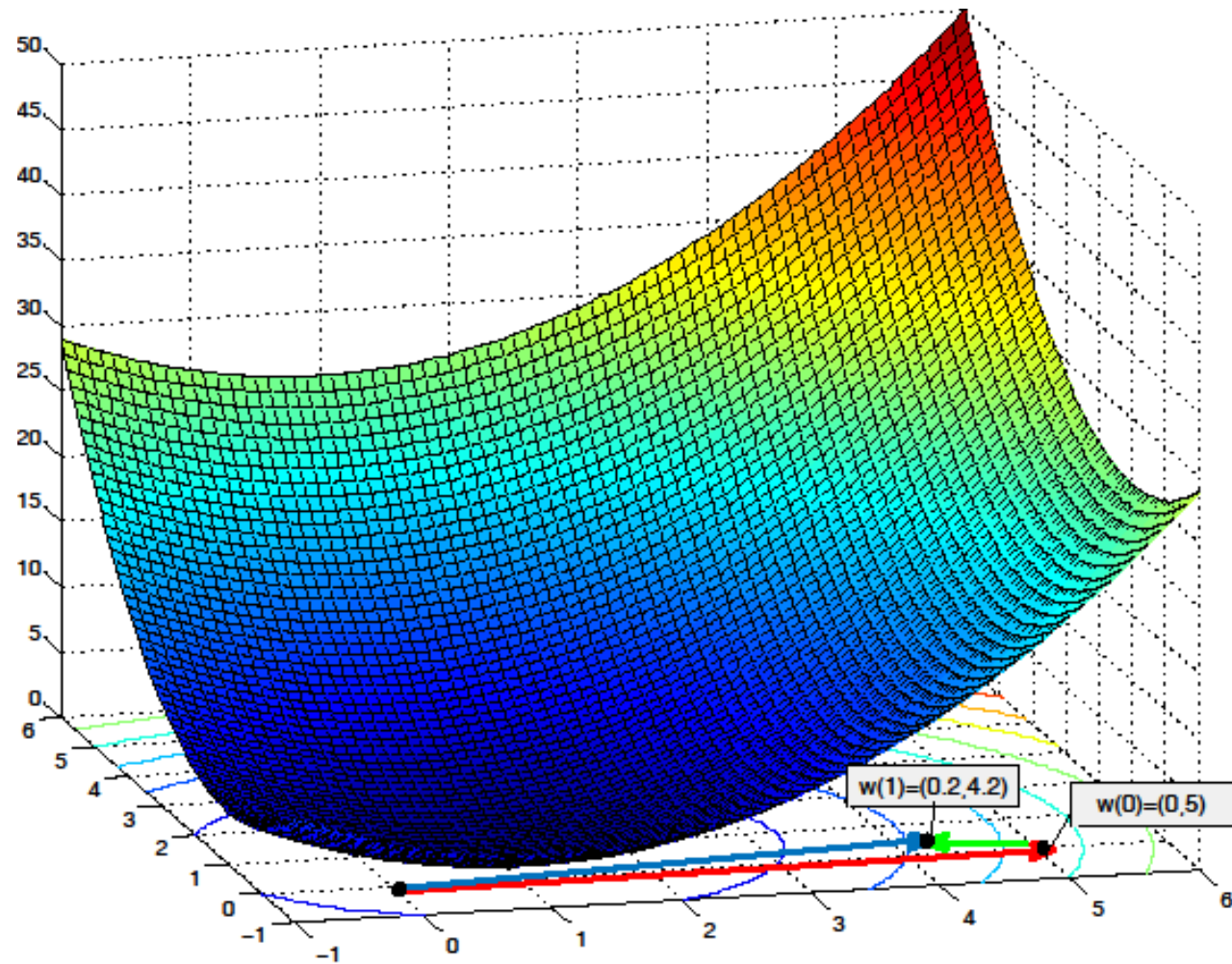
-Applying the gradient descent algorithm for $\mathbf{w}(0) = [0, 5]^T$, and $\mu = 0.1$, we have

$$\mathbf{w}(1) = \begin{bmatrix} 0 \\ 5 \end{bmatrix} - 0.1 \begin{bmatrix} -2 \\ 8 \end{bmatrix} = \begin{bmatrix} 0.2 \\ 4.2 \end{bmatrix}$$

-Thus, $\mathbf{w}(1)$ comes closer to \mathbf{w}^* .



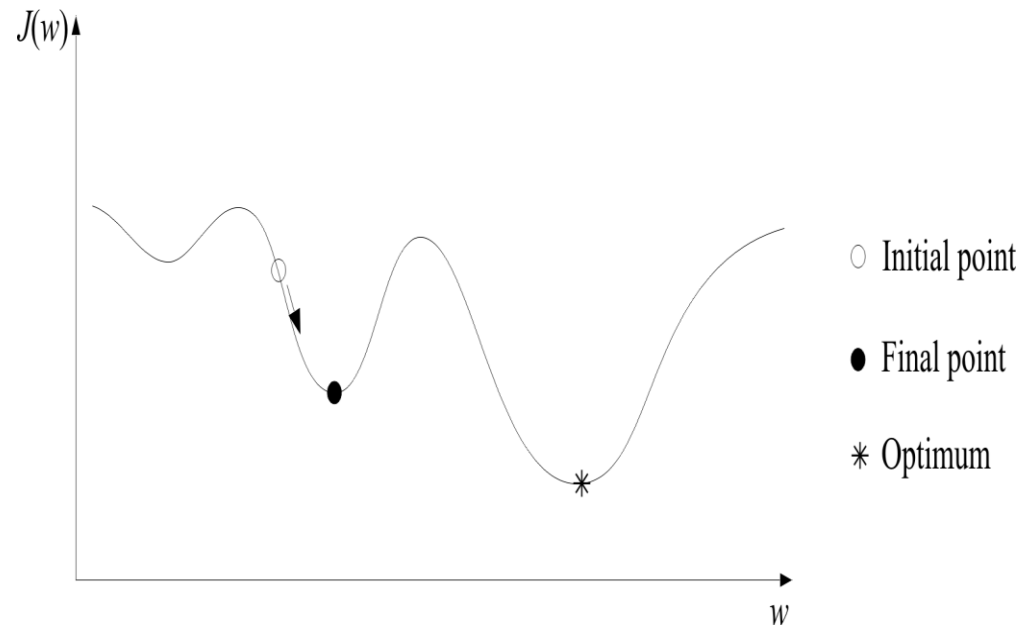
Optimization theory – Basic concepts (supp. material)



Optimization theory – Basic concepts (supp. material)

Remarks for gradient descent:

- The value of μ should be chosen **not too large**, in order to avoid oscillations around the minimum and **not too small** in order to avoid unnecessary delays in the convergence
- If $J(\mathbf{w})$ has **more than one local minima**, the gradient descent will converge (in general) to the one that is closest to $\mathbf{w}(0)$.
- If the algorithm is trapped to a **local minimum** that correspond to a poor solution, the only way to **escape** from it is to **re-initialize** the algorithm from another initial position.
- It can be proved that, under certain conditions, the algorithm **converges asymptotically** to a **local minimum** of $J(\mathbf{w})$.



Optimization theory – Basic concepts (supp. material)

Let $J(\mathbf{w})$ be a continuous function of \mathbf{w} .

Problem (P2): Determine the **position \mathbf{w}^*** where the function $J(\mathbf{w})$ achieves its **minimum** value, under the constraint that \mathbf{w} satisfies some **equality constraints**.

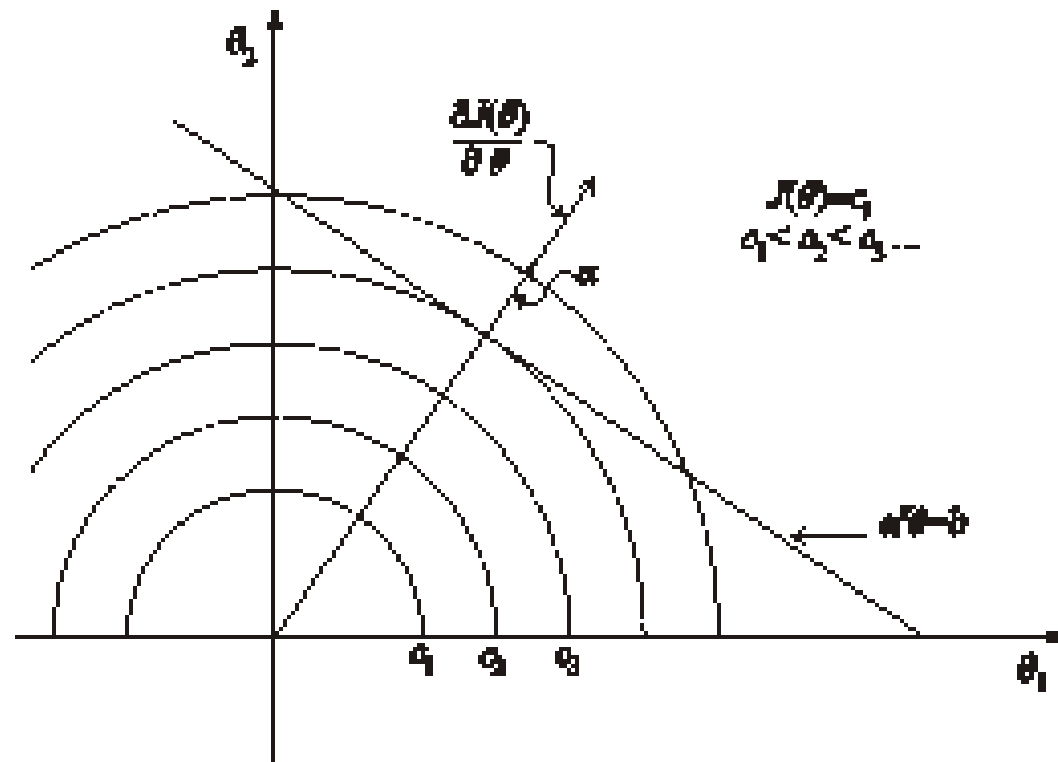
For **linear equality constraints**, the problem is stated as follows

- Minimize $J(\mathbf{w})$
- Subject to the constraints $A\mathbf{w} = \mathbf{b}$, where A an $m \times l$ matrix and \mathbf{b} an m -dim. Vector.

Solution: **Lagrange multipliers**

Minimize

- $L(\mathbf{w}) = J(\mathbf{w}) + \boldsymbol{\lambda}^T(A\mathbf{w} - \mathbf{b})$
- $\boldsymbol{\lambda}$ is an m -dim vector that is estimated through the constraints $A\mathbf{w} = \mathbf{b}$



Optimization theory – Basic concepts (supp. material)

Let $J(\mathbf{w})$ be a continuous function of \mathbf{w} .

Problem (P3): Determine the **position** \mathbf{w}^* where the function $J(\mathbf{w})$ achieves its **minimum** value, under the constraint that \mathbf{w} satisfies some **inequality constraints**.

For **linear inequality constraints**, the problem is stated as follows

- Minimize $J(\mathbf{w})$
- Subject to the constraints $A\mathbf{w} \geq \mathbf{b}$, where A an $m \times l$ matrix and \mathbf{b} an m -dim. Vector.

