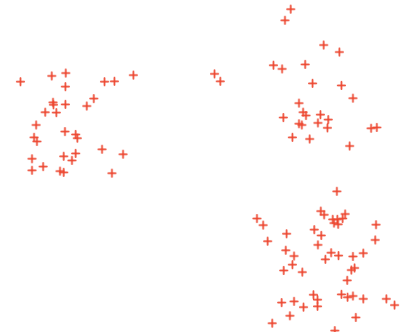


#### 1. *K-means Clustering*

- Dessinez trois ou quatre groupes de points en dimension 2, par exemple selon le modèle ci-contre.
- Trouvez le composant permettant d'exécuter l'algorithme des *k-means* pas à pas, c'est-à-dire de manière interactive et connectez-le.
- Indiquez le bon nombre de clusters, puis cliquez plusieurs fois sur **Randomize Positions**.
- Pour exécuter une itération de l'algorithme, il suffit de cliquer sur l'algorithme **Reassign Membership**. Comptez le nombre d'itérations nécessaire à l'obtention d'une partition stable.
- On peut aussi initialiser en positionnant les centres initiaux. Trouvez une configuration qui vous semble difficile et exécutez de nouveau pas à pas.
- Que va-t-il se passer si on demande un nombre de groupes différent de ce que suggère le visuel (un de plus ou un de moins) ? Testez.
- Trouvez le composant **k-Means** et connectez-le aux données. Vous offre-t-il une solution pour le problème du nombre de groupes ?



#### 2. *Clustering Hiérarchique*

- Chargez les données du fichier **TD4-xyzt.xlsx**, trouvez le moyen de calculer le tableau de distances (euclidienne usuelle) entre tous les points, puis de le visualiser.
- Calculez les distances de *Manhattan* et vérifiez visuellement (nuage dessiné en TD) ainsi qu'à l'aide du composant.
- Trouvez le moyen de calculer la distance cosinus entre les deux variables.
- Trouvez le composant permettant de réaliser une classification hiérarchique, puis retrouvez la hiérarchie de Ward vue en TD.
- Quelle semble être la meilleure partition ? Trouvez le moyen de récupérer une variable indicatrice du résultat.
- Utilisez cette indicatrice pour visualiser les points.
- Changez de distance et constatez le résultat.
- Exécutez les *K-means* avec  $K = 2$ , puis trouvez le moyen de comparer la partition obtenue avec celle déduite de la hiérarchie de Ward. Si c'est la même partition, changez un paramètre de sorte d'en comparer deux différentes.

#### 3. Données Réelles

- Chargez les données **wine.tab**, et réalisez leur classification hiérarchique.
  - Obtenez une partition en un nombre de groupes égal au nombre de qualités des 178 vins.
  - Comparez le résultat (indicatrice de *cluster*) et la vérité-terrain (donnée d'expertise).
  - Calculez (manuellement) le taux de désaccord.
  - Les méthodes d'Analyse de Données ne sont pas mutuellement exclusives. Reprenez l'ACP de ces données, puis la réalisez le partitionnement par *k-means* du nuage réduit à un nombre judicieux de composantes principales.
  - Recommencez la comparaison avec la vérité terrain (taux de désaccord).
  - Le taux de désaccord est-il meilleur ? Peut-on affirmer que c'est toujours ainsi ?
- (HW) Reprenez la même analyse (a) à (g) en inversant l'ordre : *k-means* dans l'espace initial et *hierarchical clustering* dans l'espace des composantes principales.
- Cherchez, par *k-means* et/ou par *hierarchical clustering*,

#### 4. Digits Images

- Décompressez le fichier **digits0137.zip** et importez les images qu'il contient à l'aide du composant approprié.
- Connectez un **Data Table** afin de savoir ce qu'il en est, puis trouvez le moyen de visualiser les images.
- Un réseau de neurones profond a été codé dans le composant **Image Embedding** afin d'extraire 1 000 caractéristiques numériques de chaque image quelle que soit sa taille. Connectez ce composant et regardez le tableau de données résultant.
- Réalisez la *meilleure* classification hiérarchique au sens de la vérité-terrain en jouant sur la distance entre points et celle entre groupes. Vous pourrez visualiser les images par groupe.
- Comment réduire intelligemment le nombre de variables en amont ? Testez.

#### 5. Ponctuation

Chargez le fichier **auteurs.xlsx**, et... débrouillez-vous ;)