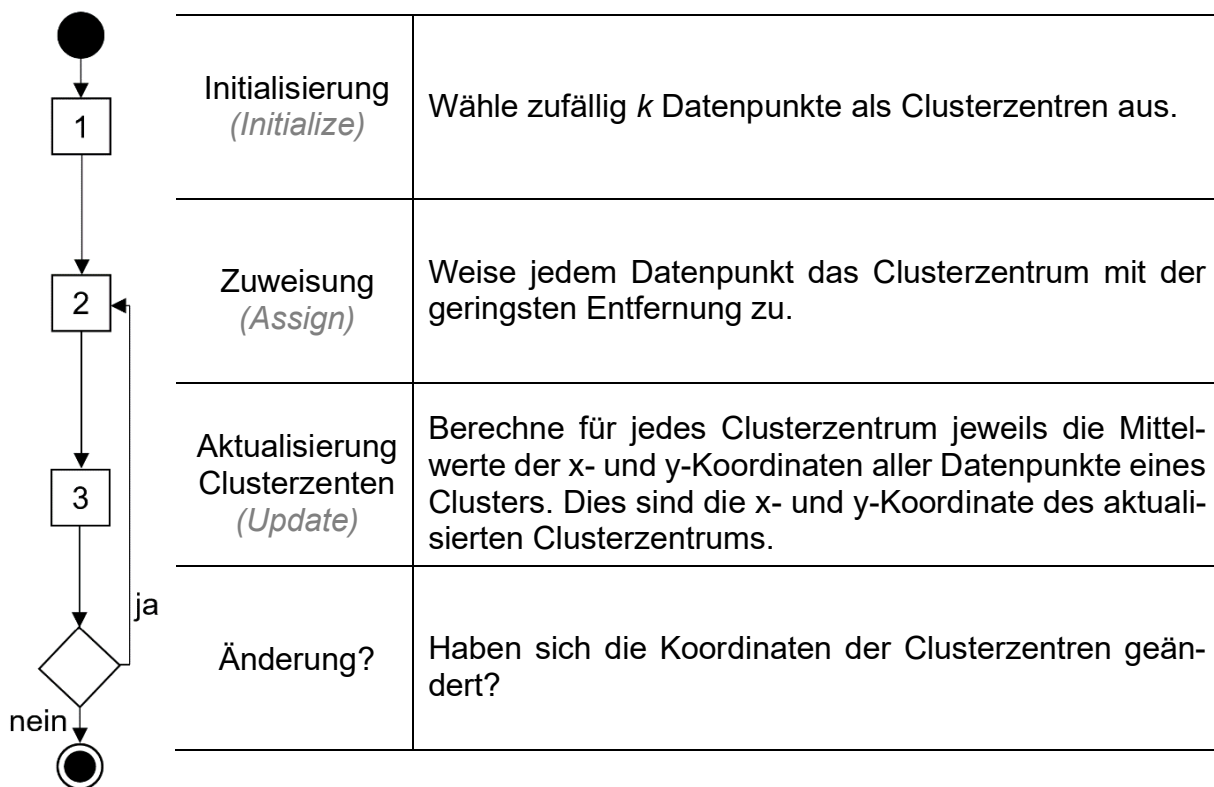


***k*-Means-Algorithmus**

Einführung:

Bei KI-Systemen unterscheidet man zwischen dem wissensbasierten Ansatz und dem datenbasierten Ansatz (auch *maschinelles Lernen*). Teilbereiche des *maschinellen Lernens* sind u.a. das *Reinforcement Learning*, das *Supervised Learning* und das *Unsupervised Learning*. Beim *Unsupervised Learning* sind die Daten *ungelabelt* und der Algorithmus teilt diese in Gruppen, die sog. *Cluster*, ein. Ein typischer Algorithmus aus dem Bereich des *Unsupervised Learning* ist der *k*-Means-Algorithmus, um den es in dieser Unterrichtseinheit geht.

Bei der Anwendung des *k*-Means-Algorithmus gehst du wie folgt vor:



Arbeitsauftrag 1 (Clustering):

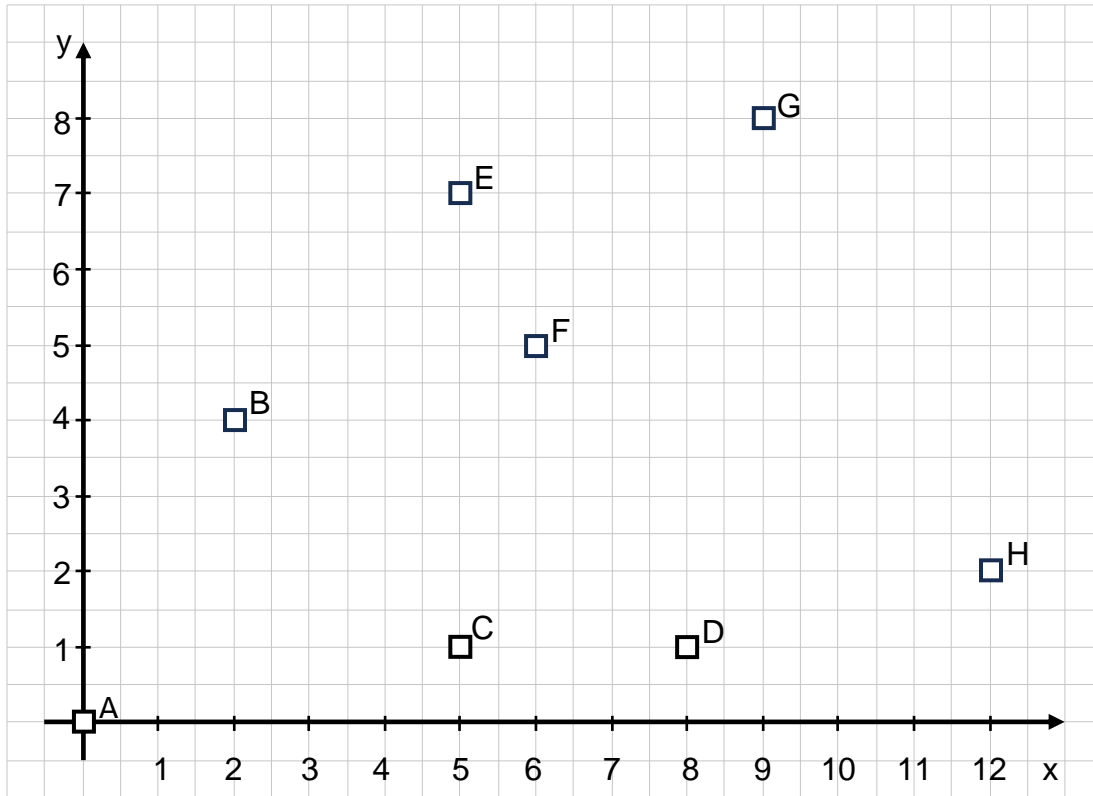
Teile mit Hilfe des *k*-Means-Algorithmus die Datenpunkte **A** bis **H** in **zwei** Cluster! Fülle dazu die folgenden Lücken aus.

1 Initialisierung

$$k = 2$$

Wir gehen davon aus, dass per Zufall als Clusterzentrum $Z_{\text{grün}}$ der Punkt **C** und als Clusterzentrum Z_{gelb} der Punkt **F** ausgewählt wurde. Trage diese als grünes bzw. gelbes „X“ im folgenden Koordinatensystem ein.

1. Durchlauf



2 Zuweisung

Messe alle Entfernungen zwischen den Datenpunkten **A** bis **H** und den Clusterzentren $Z_{\text{grün}}$ und Z_{gelb} und trage diese in die Tabelle ein.

	A (0 0)	B (2 4)	C (5 1)	D (8 1)	E (5 7)	F (6 5)	G (9 8)	H (12 2)
$Z_{\text{grün}}$ ()								
Z_{gelb} ()								
$Z_{\text{grün}}? Z_{\text{gelb}}?$								

Das Cluster zum Clusterzentrum $Z_{\text{grün}}$ enthält die Datenpunkte _____

Das Cluster zum Clusterzentrum Z_{gelb} enthält die Datenpunkte _____

Markiere im vorherigen Koordinatensystem die Datenpunkte in der Farbe des zugehörigen Clusterzentrums.

3 Aktualisierung der Clusterzentren

Berechne die neuen Koordinaten von $Z_{\text{grün}}$!

x = _____

y = _____

Berechne die neuen Koordinaten von Z_{gelb} !

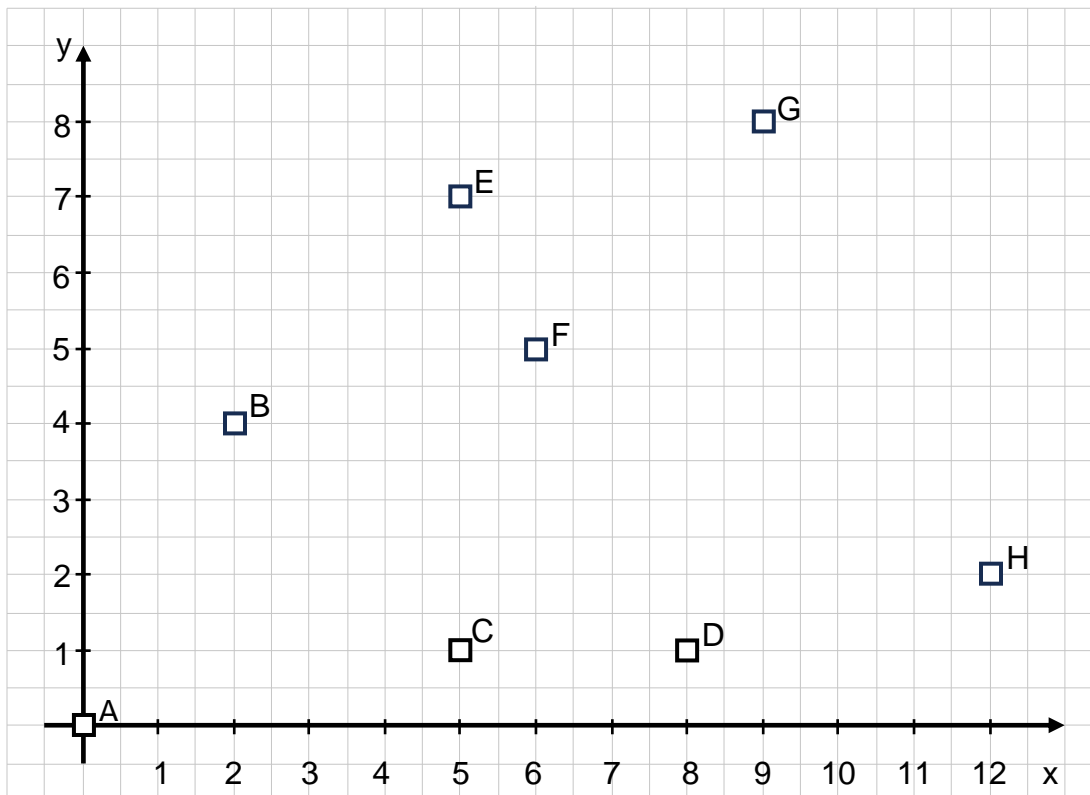
x = _____

y = _____

Haben sich die Koordinaten der Clusterzentren geändert? ja nein

2. Durchlauf

Trage die aktualisierten Clusterzentren als grünes bzw. gelbes „X“ im folgenden Koordinatensystem ein.



2 Zuweisung zu den aktualisierten Clusterzentren

	A (0 0)	B (2 4)	C (5 1)	D (8 1)	E (5 7)	F (6 5)	G (9 8)	H (12 2)
Z_{grün} ()								
Z_{gelb} ()								
Z_{grün}? Z_{gelb}?								

Das Cluster zum Clusterzentrum **Z_{grün}** enthält die Datenpunkte _____

Das Cluster zum Clusterzentrum **Z_{gelb}** enthält die Datenpunkte _____

Markiere im vorherigen Koordinatensystem die Datenpunkte in der Farbe des zugehörigen Clusterzentrums.

3 Aktualisierung der Clusterzentren

Berechne die neuen Koordinaten von **Z_{grün}**!

x = _____

y = _____

Berechne die neuen Koordinaten von **Z_{gelb}**!

x = _____

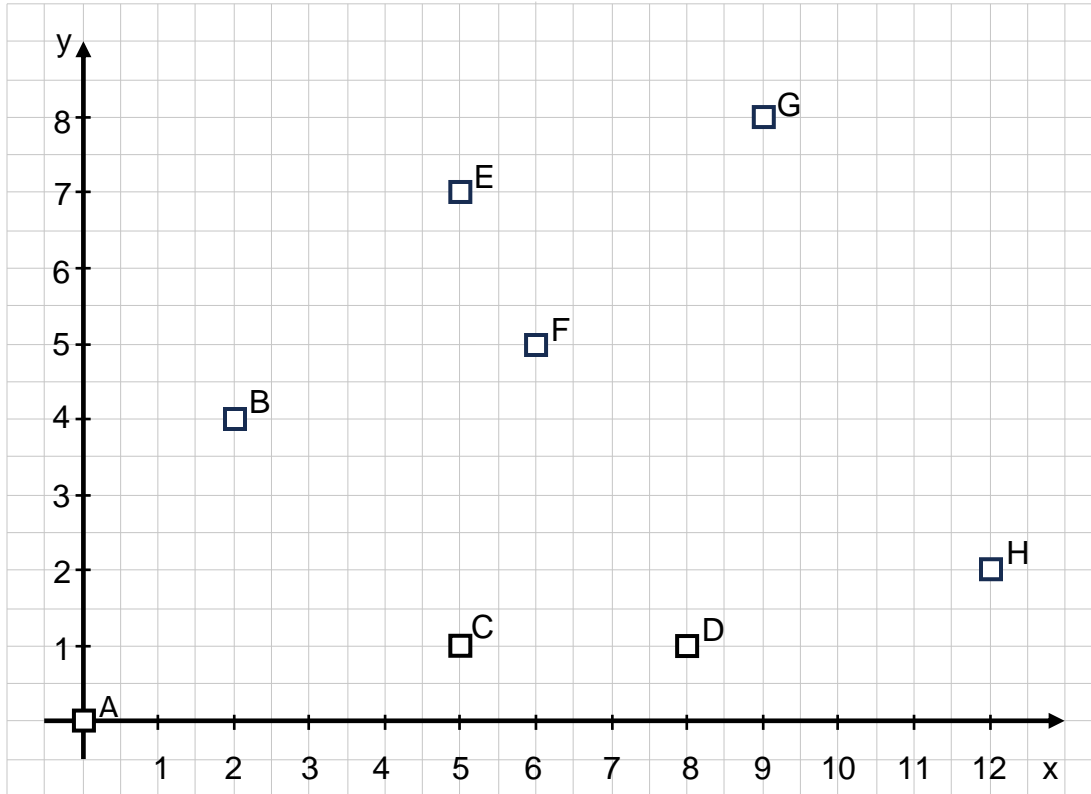
y = _____

Haben sich die Koordinaten der Clusterzentren geändert? ja nein



3. Durchlauf

Trage die aktualisierten Clusterzentren als grünes bzw. gelbes „X“ im folgenden Koordinatensystem ein.



2 Zuweisung zu den aktualisierten Clusterzentren

	A (0 0)	B (2 4)	C (5 1)	D (8 1)	E (5 7)	F (6 5)	G (9 8)	H (12 2)
$Z_{\text{grün}}$ ()								
Z_{gelb} ()								
$Z_{\text{grün}}?$ $Z_{\text{gelb}}?$								

Das Cluster zum Clusterzentrum $Z_{\text{grün}}$ enthält die Datenpunkte _____

Das Cluster zum Clusterzentrum Z_{gelb} enthält die Datenpunkte _____

Markiere im vorherigen Koordinatensystem die Datenpunkte in der Farbe des zugehörigen Clusterzentrums.

3

Aktualisierung der ClusterzentrenBerechne die neuen Koordinaten von $Z_{\text{grün}}$!

x = _____

y = _____

Berechne die neuen Koordinaten von Z_{gelb} !

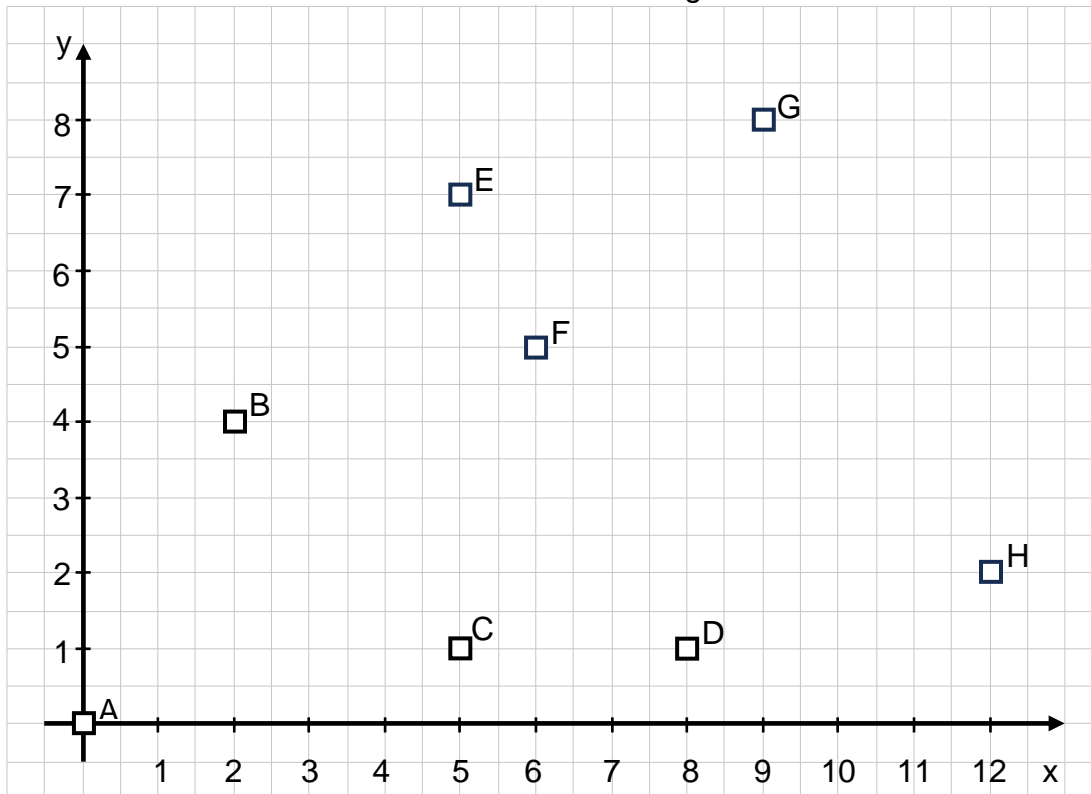
x = _____

y = _____

Haben sich die Koordinaten der Clusterzentren geändert? ja nein**Ergebnis**

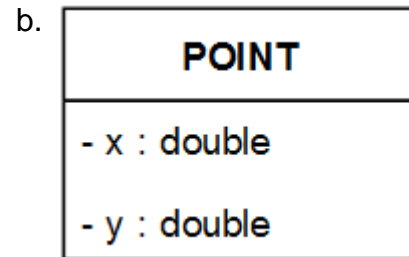
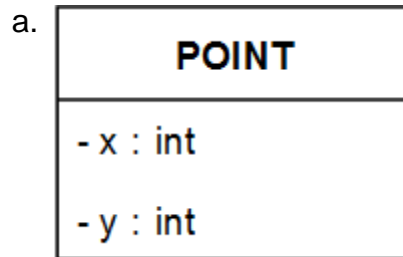
Arbeitsauftrag 2 (Cluster-Hülle):

Zeichne die Cluster-Hüllen ein! Tipp: Verbinde zunächst alle Datenpunkte eines Clusters und zeichne dann nur die äußersten Linien farbig nach.



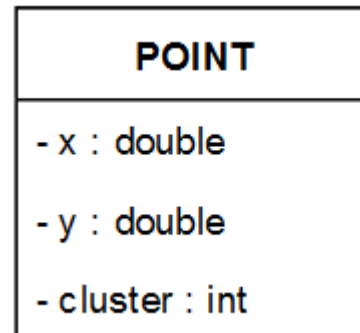
Arbeitsauftrag 3 (Modellierung):

- a) Für jeden Datenpunkt müssen die x- und y-Koordinaten gespeichert werden. Diskutiere: Welches Klassendiagramm würdest du dafür bevorzugen?

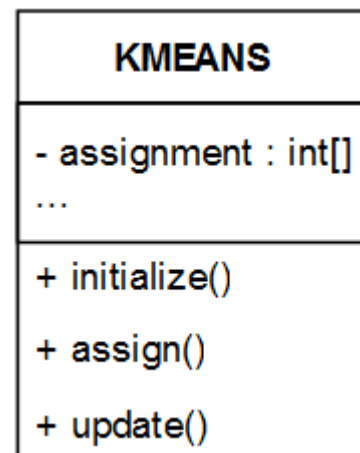


- b) Für jeden Datenpunkt muss es möglich sein, dessen Zuweisung zu einem Cluster zu speichern. Welche der drei untenstehenden Möglichkeiten ermöglicht eine effektive Zuweisung, die keine Prinzipien der Objektorientierung verletzt?

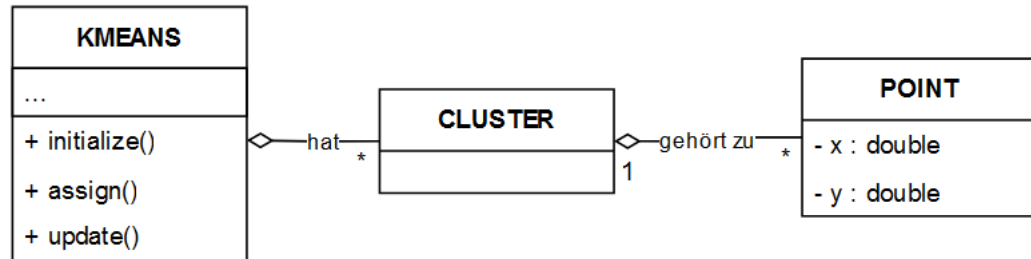
- a. Speicherung in der Klasse POINT



- b. Speicherung in der Klasse KMEANS als Array der Länge #Eingabedaten.



- c. Speicherung in separater Klasse CLUSTER, die die zum Cluster zugehörigen Punkte referenziert.



- c) Die Datenpunkte der Eingabedaten müssen in einer Datenstruktur gehalten werden. Überlege, welche Datenstrukturen du bereits kennst. Für welche würdest du dich entscheiden?

Arbeitsauftrag 4 (Pseudocode):

Ergänze den Pseudocode für folgende Methoden

```
def initialize()
```

Wiederhole k mal:

Wähle einen _____ Punkt aus _____

aus.

Erstelle an dessen Koordinaten ein _____.

```
def assign()
```

Für _____ in der Liste _____:

Speichere die Zugehörigkeit zu dem Zentrum mit

_____.

```
def update()
```

Für alle _____ in der Liste _____:

Berechne jeweils den _____ der _____ und

_____ Koordinaten aller Punkte, die zu diesem

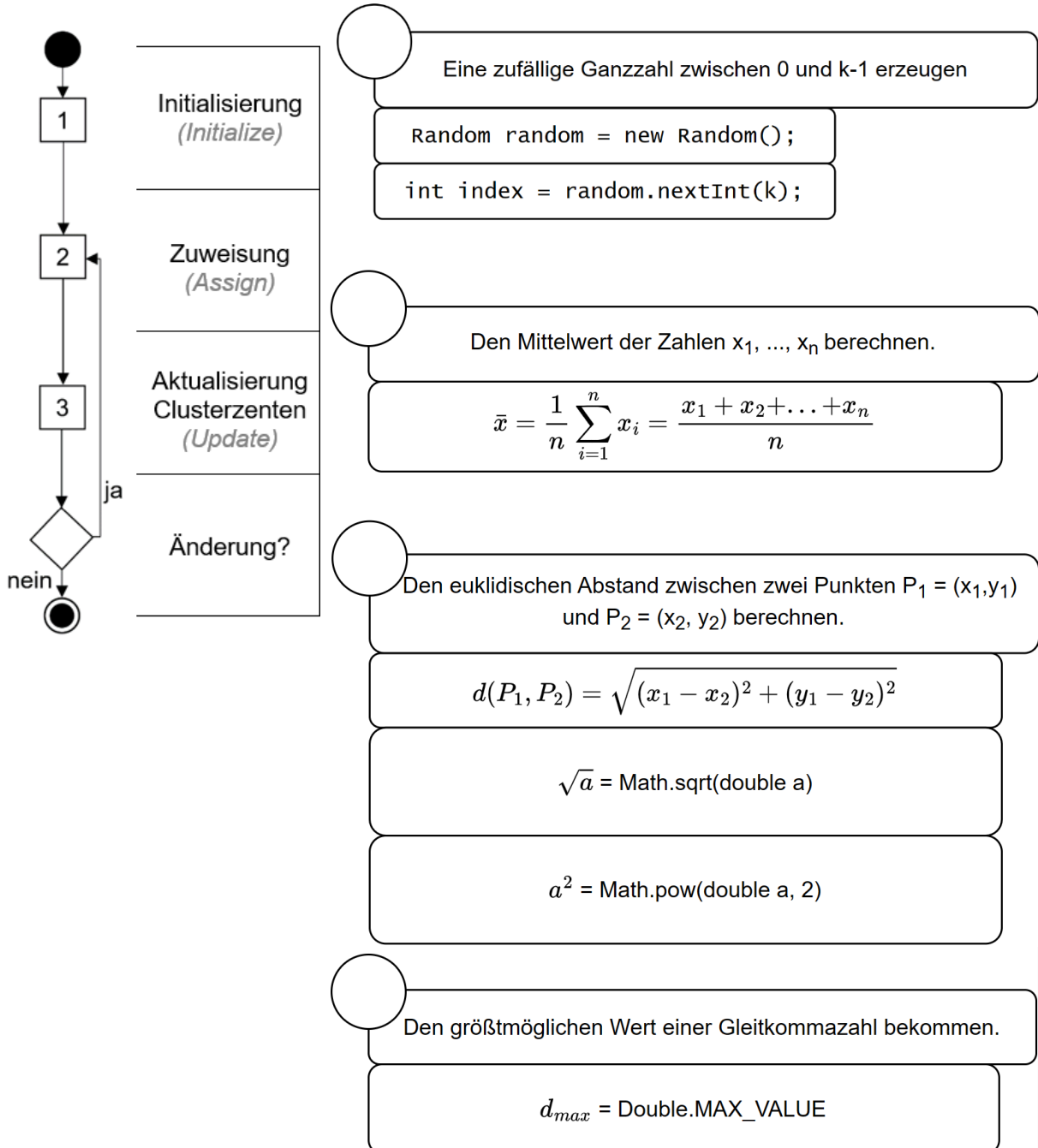
_____ gehören.

Setze die x- und y- Koordinaten des _____

Zentrums auf die berechneten _____.

Arbeitsauftrag 5 (Hilfestellung fürs Programmieren):

Ordne die folgenden Ablaufschritte den einzelnen Phasen des Algorithmus zu. Schreibe die Schrittnummer in die Kreise.



Arbeitsauftrag 6 (Implementierung in BlueJ):

Öffne das *BlueJ*-Projekt „kMeans_BlueJ_Vorlage“. Dort findest du sieben Klassen. Deine Aufgabe ist es die fehlenden Methoden der Klasse KMEANS zu ergänzen. Sie entsprechen den einzelnen Schritten des *k-Means-Algorithmus*.

- a) Ergänze die Methode `initialize()` der Klasse KMEANS. Achte darauf, dass du das Zentrum als einen **neuen** Punkt anlegst, da dessen x- und y-Koordinaten im Verlauf des Algorithmus geändert werden.
- b) Ergänze die Methode `assign()`. Das Array `assignment` ist zu Beginn komplett mit „-1“ gefüllt und muss von dir auf die Clusterzahl (= Index des entsprechenden Zentrums in `centers`) des nächsten Zentrums gesetzt werden.

Hinweis: Die Methode `getDistance(POINT p)` der Klasse `POINT` berechnet den euklidischen Abstand zwischen dem Punkt selbst und einem anderen Punkt `p`.

- c) Ergänze die Methode `update()` und bestimme die Laufzeit der von dir implementierten Lösung. Zusatzaufgabe: Kannst du die Laufzeit noch verbessern?

Testen: Zum Testen deiner Methoden kannst du entweder den Konstruktor der Klasse GUI ohne Parameter aufrufen oder ein Objekt der Klasse DATAHANDLER erstellen und dieses zusammen mit deinem gewünschten *k* dem Konstruktor der Klasse GUI übergeben. Mittels der aufgehenden Visualisierung siehst du die Auswirkungen deines Codes.