

* تعريف خوارزمية k_means :

هي خوارزمية تقسيم البيانات إلى مجموعات k-means خوارزمية متجانسة بحيث يكون كل مجموعة تحتوي على نقاط متشابهة. (clusters) الهدف من هذه الخوارزمية هو تقليل المسافة بين نقاط البيانات في نفس المجموعة وزيادة المسافة بين نقاط المجموعات المختلفة. يتم تحديد عدد المجموعات المطلوبة مسبقاً وتقوم الخوارزمية بتحديد مراكز هذه المجموعات بناءً على البيانات المدخلة

* طرق حساب k في k_means

1. Elbow Method: هذه الطريقة تتضمن تشغيل خوارزمية k-means وحساب متوسط المربعات البيضاء k عدة مرات مع قيم مختلفة لـ (Sum of Squared Errors) يتم رسم هذه القيم على الرسم البياني k لكل قيمة من (Squared Errors) على الرسم البياني وهي (elbow) التي تظهر "كوع" k ويتم اختيار قيمة القيمة التي يعتبر بها انخفاض معدل الخطأ أبطأ

* مثال Elbow Method :

في k يتم تحديد قيمة مناسبة لـ (Elbow Method) طبقاً للطريقة الكوعية . (Sum of Squared Errors) عن طريق رسم خطأ التجميع k-means خوارزمية يتم اختيار (clusters) مقابل عدد مختلف من المجموعات (Errors - SSE) مما يشبه شكل كوع ، حيث يحدث تغيير حاد في الخطأ التجميعي k قيمة : لنفترض أن لدينا مجموعة بيانات تحتوي على نقاط مثل التالي

x	y	نقطة
1	2	3

2	3 4
3	2 5
4	6 7
5	7 6
6	8 9

بإستخدام الطريقة الكوعية. يمكننا تطبيق k نريد تحديد قيمة مناسبة لـ SSE وحساب k (مثلاً من 1 إلى 6) لعدة قيم مختلفة لـ k -means خوارزمية (k) مقابل عدد المجموعات (SSE) لكل قيمة. ثم يتم رسم خطأ التجميع

نبحث عن النقطة التي تشبه كوعاً في الرسم البياني. هذه النقطة تشير ،بعد ذلك تقع في هذا الموضع k إلى أنه يجب اختيار قيمة

$K = 1$: Inertia = [2560](#)

$K = 2$: Inertia = [1800](#)

$K = 3$: Inertia = [1200](#)

$K = 4$: Inertia = [900](#)

$K = 5$: Inertia = [800](#)

$K = 6$: Inertia = [750](#)

$k = 4$ أفضل قيمة عندما تكون

2. **Silhouette Method:** تستخدم هذه الطريقة مقياس الـ k لعدة قيم مختلفة من k -means التقييم جودة التجزئة. يتم تشغيل خوارزمية التي تعطي k يتم اختيار قيمة k . لكل قيمة من k Silhouette ويتم حساب قيمة Silhouette. أعلى قيمة

3. **Gap Statistics:** هذه الطريقة تقارن متوسط الـ SS (Sum of Squared Errors) للبيانات الفعلية بالقيم المتوقعة إذا كانت البيانات توزعت للبيانات SS التي تعطي أكبر فجوة بين متوسط الـ k عشوائياً. يتم اختيار قيمة الفعلية والبيانات المتوقعة.

4. **AIC (Akaike Information Criterion) أو BIC (Bayesian Information Criterion):** تستخدم هذه الطرق معاقبة إضافية لعدد التي تحقق أفضل توازن بين دقة k وتحاول اختيار قيمة k المجموعات النموذج وبساطته.

-