

: تعریف کلی Clustering خوشه‌بندی:

Clustering یکی از روش‌های یادگیری بدون نظارت (**Unsupervised Learning**) در علم داده و یادگیری ماشین است که هدف آن تقسیم داده‌ها به دسته‌ها (خوشه‌ها) است به طوری که:

- داده‌های داخل هر خوشه بیشتر شبیه هم باشند
- داده‌های بین خوشه‌های مختلف کمتر به هم شبیه باشند

به زبان ساده‌تر:

Clustering یعنی اینکه داده‌هایی که شبیه هم هستند رو کنار هم بگذاریم، بدون اینکه از قبل بدانیم کدام داده متعلق به کدام دسته است.

K-Means Clustering

K-Means یک الگوریتم خوشه‌بندی (**Clustering**) در یادگیری ماشین است که:

داده‌ها را به **k خوشه‌ی مجزا** تقسیم می‌کند، به طوری که نقاط در هر خوشه به مرکز آن خوشه نزدیک‌تر از مراکز دیگر باشند. این الگوریتم یکی از ساده‌ترین روش‌های دسته‌بندی محسوب می‌شود. یعنی نقطه‌ای را پیدا می‌کنیم که در مرکز خوشه قرار گرفته باشد و به تمام داده‌ها در آن خوشه نزدیکترین حالت را داشته باشد.

در کتابخانه **sklearn** یک تابع برای انجام clustering وجود دارد.

```
from sklearn.datasets import make_blobs
```

make_blobs یک تابع برای تولید داده‌های خوشه‌ای شکل (**blobs**) در فضای n-بعدی است.

این تابع برای ساخت داده‌های تستی برای الگوریتم‌های **Clustering** طراحی شده است.

```
X, _ = make_blobs(random_state=42)
X.shape
```

این دستور یک دیتاست مصنوعی (ساختگی) می‌سازد که به صورت خوشه‌ای (**blobs**) است. و به طور پیش‌فرض، ۱۰۰ نقطه داده‌ای دوبعدی ایجاد می‌کند که به ۳ خوشه تقسیم شده‌اند.

پارامتر `random_state=42` باعث می‌شود که تولید داده‌ها تصادفی ولی تکرارپذیر باشد، یعنی هر بار که کد را اجرا کنید، دقیقاً همان داده‌ها ساخته شوند.

`X, _ = make_blobs(...)`

تابع `make_blobs` دو خروجی دارد:

- `X`: داده‌ها (مختصات نقاط در فضای دوبعدی)
- `Y`: برچسب هر نقطه (که خوشه‌اش را مشخص می‌کند)

در این خط فقط `X` را نگه می‌داریم و `Y` را با `_` دور می‌اندازیم، چون به آن احتیاجی نداریم.

نکته: مقدار 42 یک عدد رایج بین برنامه نویسان است و میتوان آن را تغییر داد. این عدد برگرفته از رمان "راهنمای مسافران کهکشان به کهکشان" است که در آن 42 پاسخ نهایی به زندگی، جهان و همه چیز است. بنابراین می‌توان این عدد را تغییر داد اما اگر این مقدار خالی باشد `None` برگردانده می‌شود و هر بار کد اجرا شود مقدار متفاوتی بازگردانده می‌شود.

```
plt.scatter(X[:, 0], X[:, 1]);
```

در ادامه یک نمودار دو بعدی رسم می‌کنیم.

`X[:, 0]`

یعنی: تمام ردیف‌ها، ستون اول از آرایه‌ی $X \rightarrow$ این می‌شود محور افقی (X-axis)

`X[:, 1]`

یعنی: تمام ردیف‌ها، ستون دوم از آرایه‌ی $X \rightarrow$ این می‌شود محور عمودی (Y-axis)

به عبارتی:

`X[:, 0]`

فرض کن `X` یک جدول (یا آرایه) از داده‌هاست. این دستور یعنی از همه‌ی ردیف‌ها، فقط ستون اول رو بردار ستون اول معمولاً مقدار محور `X` هست.

`X[:, 1]`

از همه‌ی ردیف‌ها، فقط ستون دوم رو بردار. ستون دوم معمولاً مقدار محور `Y` است.

```
kmeans = KMeans(n_clusters=3, random_state=42)
```

یک شی جدید از الگوریتم kmeans بساز که قرار است داده‌ها را به سه خوشه تقسیم کند و برای اینکه نتایج تصادفی قابل تکرار باشند مقدار random state را برابر با عدد چهل و دو قرار بدهیم.

عبارت n clusters برابر با سه یعنی ما می‌خواهیم سه گروه یا خوشه پیدا کنیم.

عبارت random state برابر با چهل و دو یعنی در اجرای دوباره کد نتایج یکسان بمانند چون الگوریتم kmeans به طور پیش‌فرض رفتار تصادفی دارد و با این عدد آن را کنترل می‌کنیم.

```
kmeans = KMeans(n_clusters=3, random_state=42)
```

این خط کد الگوریتم kmeans را روی داده‌های X اجرا می‌کند و به صورت هم‌زمان دو کار انجام می‌دهد. اول الگوریتم را روی داده‌ها آموزش می‌دهد یعنی خوشه‌ها را پیدا می‌کند و مراکز آن‌ها را مشخص می‌کند. دوم برای هر داده در X مشخص می‌کند که به کدام خوشه تعلق دارد و آن را برچسب‌گذاری می‌کند. خروجی این عملیات لیستی از اعداد است که نشان می‌دهد هر داده به کدام خوشه اختصاص یافته است. این لیست در متغیری به نام labels ذخیره می‌شود تا بعداً بتوانیم از آن برای تحلیل یا رسم نمودار استفاده کنیم.

```
plt.scatter(X[:, 0], X[:, 1], c=labels);
```

این خط کد یک نمودار نقطه‌ای از داده‌های X رسم می‌کند و برای هر نقطه رنگ متفاوتی بر اساس خوشه‌ای که به آن تعلق دارد در نظر می‌گیرد. محور افقی شامل مقادیر ستون اول داده‌ها و محور عمودی شامل مقادیر ستون دوم است. با استفاده از آرگومان c برابر با labels، هر داده با رنگ خوشه‌ای که به آن تعلق دارد نمایش داده می‌شود. این باعث می‌شود که خوشه‌بندی الگوریتم KMeans به صورت بصری قابل مشاهده باشد و بتوان به راحتی تشخیص داد که الگوریتم چگونه داده‌ها را گروه‌بندی کرده است.