

سوال ۱ — تابع هدف K-Means و فرمول مرکز خوش

چرا مرکز خوش = میانگین نقاط؟

میخواهیم:

- این داده‌ها را به K گروه (خوش) تقسیم کنیم
- طوری که نقاط هر خوش تا حد ممکن به هم نزدیک باشند (نزدیک بودن = فاصله کم)

تعريف تابع هدف (Objective Function)

حالا کل ایده K-Means در یک فرمول خلاصه می‌شود:

$$J = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|^2$$

یه نقطه رو به عنوان مرکز هر گروه بگیر
بعد بین مجموع فاصله‌ی همهی نقاطها تا این مرکز چقدر میشه
حالا مرکز رو جوری جابه‌جا کن که این مجموع فاصله کمترین مقدار ممکن باشه»

ابن «کمترین کردن» همون چیزیه که اسمش هست:
تابع هدف:

- برای هر خوش
 - فاصله تمام نقاطش تا مرکز خوش را حساب کن
 - مربع کن
- ❖ چرا مربع؟
❖ فاصله منفی نمی‌شه
❖ نقاط خیلی دور «خیلی بدتر» حساب می‌شن
❖ محاسبه‌اش راحت‌تره

- جمع بزن
- کل این عدد را کمینه کن

این همان چیزی است که در کد `sklearn` می‌بینی به اسم:

`inertia = WCSS`

در K-Means :

عدد بزرگ → خوشبندی بد

عدد کوچک → خوشبندی خوب

تمرکز روی یک خوشه (تفکیک مسئله)

فرض کن:

- خوشه‌ی k ام را داریم
- نقاطش:

$$C_k = \{x_1, x_2, \dots, x_m\}$$

تابع هدف فقط برای این خوشه:

$$J_k(\mu_k) = \sum_{i=1}^m \|x_i - \mu_k\|^2$$

اینجا فقط μ_k متغیر است. نقاط ثابت‌اند

$$(x_i - \mu)^2 = x_i^2 - 2x_i\mu + \mu^2$$

$$J(\mu) = \sum_{i=1}^m (x_i^2 - 2x_i\mu + \mu^2)$$

$$J(\mu) = \sum_{i=1}^m x_i^2 - 2\mu \sum_{i=1}^m x_i + \sum_{i=1}^m \mu^2$$

μ^2 برای همه‌ی آها یکیه

$$J(\mu) = \underbrace{\sum_{i=1}^m x_i^2}_{\text{ثابت نسبت به } \mu} - 2\mu \sum_{i=1}^m x_i + m\mu^2$$

قسمت‌های مهم همون‌هایی‌اند که μ توشون هست
برای کمینه کردن یک تابع، معمولاً مشتق می‌گیریم و برابر صفر می‌ذاریم:

$$\frac{dJ}{d\mu} = 0 - 2 \sum_{i=1}^m x_i + 2m\mu$$

$$m\mu = \sum_{i=1}^m x_i$$

$$\boxed{\mu = \frac{1}{m} \sum_{i=1}^m x_i}$$

مثلاً اگر داده 2 بعدی باشد:

$$x_i = (x_{i1}, x_{i2})$$

آن‌وقت:

$$\mu = \left(\frac{1}{m} \sum x_{i1}, \frac{1}{m} \sum x_{i2} \right)$$

ارتباط با K-Means در عمل

دو کار را تکرار می‌کند: K-Means

Assign .1: هر نقطه میره نزدیکترین مرکز

Update .2: برای هر خوش، مرکز جدید = میانگین نقاط خوش (همین چیزی که اثبات کردیم)

و چون در مرحله 2 واقعاً بهترین مرکز را برای خوش‌هی فعلی انتخاب می‌کنیم، مقدار هزینه کم می‌شود.

(2) روش Elbow در K-Means

«چطور K مناسب رو انتخاب کنیم؟»

چرا Elbow لازم داریم؟

در K-Means یک مشکل اساسی داریم:

K (تعداد خوشه‌ها) را از کجا بفهمیم؟

- K-Means خودش K رو پیدا نمی‌کنه
- ما باید قبل از اجرا بگیم K چند تاست

در روش Elbow دقیقاً چه چیزی رسم می‌شود؟

محور افقی (x-axis):

...K=1,2,3

یعنی:

تعداد خوشه‌ها

محور عمودی (y-axis):

Inertia یا WCSS

یعنی:

- مجموع مربع فاصله نقاط تا مرکز خوشه‌ها
- همان تابع هدف K-Means

در کد: kmeans.inertia_

WCSS همیشه با افزایش K کم می‌شود

وقتی WCSS را بر حسب K رسم می‌کنیم، معمولاً انتظار داریم:

- اول: افت شدید
- بعد: افت ملایم
- نمودار شبیه یک بازوی خم شده (Elbow)

اسم روش Elbow دقیقاً از همین شکل آمده.

چرا این روش دقیق ریاضی نیست؟

چون:

- «آرنج» تعریف ریاضی دقیق ندارد
- بیشتر یک قضاوت بصری است
- دو نفر ممکن است دو K متفاوت انتخاب کند

به همین دلیل معمولاً:

- Elbow + Silhouette
- دانش مسئله + Elbow
- یا

چه زمانی نمودار Elbow آرنج واضح ندارد؟

اگر Elbow واضح نباشد، این‌ها را می‌توان نتیجه گرفت:

- یا داده ذاتاً خوش ندارد
- یا K-Means فرض مناسبی برای داده نیست
- یا باید از معیارهای دیگر کمک گرفت
- یا تعداد خوش‌ها «ذاتی» نیست

داده ساختار خوش‌ای واضح ندارد (داده تقریباً یکنواخت باشدا خوش‌ها خیلی در هم باشند)

خوش‌ها همان‌دازه و خیلی مشابه‌اند (همه خوش‌ها تراکم مشابه دارند)

داده پیچیده یا غیر کروی است (اگر داده: هلالی و کشیده و یا با نویز زیاد

(WCSS رفتار روشی نشان نمی‌دهد)

ابعاد داده زیاد است (فاصله‌ها شبیه هم می‌شوند)

سوال ۳ — چرا به Hungarian نیاز داریم؟

ما دو چیز متقابل داریم:

(1) خروجی K-Means

- برچسب‌ها: ۰, ۱, ۲, ...
- این اعداد فقط اسماند
- هیچ معنایی ندارند

مثلًا:

- خوشی ۰ ← می‌تواند « عدد ۷ » باشد

(2) برچسب‌های واقعی (Ground Truth)

- مثلًا در دیتاست digits :

- کلاس ۰ یعنی عدد « ۰ »
- کلاس ۱ یعنی عدد « ۱ »

- این‌ها معنا دارند

خطاهای؟

(1) خطای مفهومی

چون:

- فرض می‌کند شماره‌ی خوشی معنا دارد
- در حالی که ندارد

(2) خطای عددی / ارزیابی

چون:

- عدد Precision یا Accuracy کاملاً گمراهنده می‌شود
- مدل خوب، بد گزارش می‌شود

Hungarian Algorithm چی کار می‌کند؟

بهترین نگاشت یکبیک بین خوشه‌ها و کلاس‌ها را پیدا می‌کند طوری که مجموع خط‌کمینه (یا تطبیق بیشینه) شود

یعنی:

- خوشه‌ها را بهترین‌طور ممکن به کلاس‌ها نگاشت کنیم

نقش Hungarian در ارزیابی دقیقاً چیست؟

(1) حذف اثر جایه‌جایی بر چسب‌ها

- دیگر مهم نیست خوش 0 باشد یا 7
- فقط تطبیق واقعی مهم است

(2) ایجاد یک معیار عدالت‌نه

بعد از نگاشت:

- Accuracy معنی‌دار می‌شود
- Precision/Recall قابل تفسیر می‌شوند

نگاشت (mapping) یعنی:

این تصمیم که بگیم:

- خوشی 0 ↔ کلاس 7
- خوشی 1 ↔ کلاس 2
- خوشی 2 ↔ کلاس 9

یعنی: هر خوش نماینده‌ی کدام کلاس واقعی فرض شود؟
این نگاشت از قبل معلوم نیست و باید پیدا شود.

اگر Hungarian نباشد چه می‌شود؟

- ممکن است نگاشت اشتباه انتخاب شود
- خیلی پایین گزارش شود
- در حالی که خوشبندی در واقع خوب بوده

یعنی:

خطای ارزیابی، نه خطای مدل

اگر Hungarian کمک می‌کند بفهمیم اگر هر خوشه را به «بهترین» کلاس ممکن وصل کنیم، عملکرد واقعی خوشبندی چقدر خوب بوده است.

اگر اشتباه مقایسه کنیم چه می‌شود؟

اگر بگیم:

- خوشه 0 ↔ گربه
- خوشه 1 ↔ سگ

آن وقت:

- همه داده‌ها غلط حساب می‌شوند
- $\text{Accuracy} = \text{صفر}$

در حالی که مدل بی‌نقص کار کرده!

پس راه درست چیه؟

باید اول این سوال رو جواب بدیم:

هر خوشه بیشتر شبیه کدام کلاس واقعیه؟

یعنی:

- خوشه‌ای که بیشتر گربه داره → گربه
- خوشه‌ای که بیشتر سگ داره → سگ

سوال ۴ — ARI و NMI خیلی ساده و مفهومی

ما داریم خوشبندی بدون ناظر انجام می‌دهیم، ولی:

- برچسب واقعی داریم (برای ارزیابی)
- Accuracy به تنهایی قابل اعتماد نیست
- شماره‌ی خوشها معنا ندارن

ARI دقیقاً چی رو اندازه می‌گیره؟

آیا الگوریتم، جفتداده‌ها را مثل واقعیت کنار هم گذاشته یا جدا کرده؟ «نتیجه رو نسبت به شанс صحیح می‌کنه

: بازه‌ی ARI

- 1 → تطابق کامل
- 0 → تقریباً تصادفی
- منفی → بدتر از تصادفی

یعنی چی؟

به هر دو نقطه نگاه می‌کنه و می‌گه:

- آیا در واقعیت با هم بودن؟
- آیا در خوشبندی هم با هم افتادن؟

NMI دقیقاً چی رو اندازه می‌گیره؟

: NMI می‌پرسه

«چقدر اطلاعات خوشبندی، درباره‌ی کلاس‌های واقعی به من می‌ده؟»

دقت کن:

- نمی‌پرسه لیبل خوشه چند است
- می‌پرسه رابطه‌ی بین دو دسته‌بندی چطوره

: NMI بازه‌ی

- 0 → هیچ اطلاعات مشترکی
- 1 → تطابق کامل

یعنی:

- اگر برچسب خوش رو بدونم
- چقدر می‌تونم حدس بزنم کلاس واقعی چی بوده؟

حساسیت به جایه‌جایی برچسب‌ها

سؤال: اگر شماره‌ی خوش‌ها عوض شود چه می‌شود؟

مثالاً:

- خوش 0 ↔ خوش 3
- خوش 1 ↔ خوش 7

پاسخ:

- ARI: کاملاً پایدار
- NMI: کاملاً پایدار

هیچ‌کدام به «اسم خوش» نگاه نمی‌کنند

فقط به روابط بین داده‌ها نگاه می‌کنند

چرا ممکن است ARI و NMI برای یک خوشبندی فرق زیادی داشته باشند؟

تعداد خوش‌ها با تعداد کلاس‌ها فرق دارد

مثالاً:

- کلاس واقعی: 3 تا
- خوشبندی: 6 تا خوش (ریزتر)

: NMI ◆

- ممکنه هنوز بالا باشه
- چون اطلاعات کلاس‌ها تا حدی حفظ شده

:ARI ◆

- ممکنه پایین بیاد
- چون خیلی از جفت‌ها برخلاف واقعیت جدا شده‌اند

تفسیر:

خوشبندی اطلاعات خوبی دارد، ولی بیش‌از‌حد خوش‌سازی کرده

خوش‌ها همپوشانی دارند

- بعضی داده‌ها مرزی‌اند
- خوشبندی کاملاً تمیز نیست

:ARI ◆

- حساس‌تر به این اشتباهات
- سریع‌تر افت می‌کند

:NMI ◆

- ممکن است هنوز متوسط یا بالا بماند

تفسیر: *

ساختار کلی حفظ شده، ولی مرزها دقیق نیستند

کلاس‌ها نامتوازن هستند

مثالاً:

- یک کلاس خیلی بزرگ
- چند کلاس خیلی کوچک

:ARI ◆

- به روابط جفتی حساس است
- ممکن است شدیداً تحت‌تأثیر کلاس بزرگ قرار گیرد

:NMI ◆

- متوازن‌تر رفتار می‌کند

اگر NMI بالا و ARI پایین بود:

- معمولاً یعنی Over-clustering یا حفظ "اطلاعات کلی" ولی خراب بودن روابط دقیق
- خوشها خالص‌اند، اما یک کلاس را تکه‌تکه کرده‌ای

اگر ARI بالا و NMI پایین (کمتر رایج، ولی ممکن):

- می‌تواند یعنی برخی روابط جقتی خوب حفظ شده ولی توزیع اطلاعات بین خوشها/کلاس خوب نیست
- یا تعداد خوشها خیلی کم/زیاد و توزیع‌ها نامناسب است
(در عمل این حالت کمتر از "NMI بالا، ARI پایین" دیده می‌شود)

اگر هر دو پایین:

- خوشبندی واقعاً با کلاس‌ها همخوان نیست
- یا داده خوشپذیر نیست
- یا الگوریتم/فاصله/اسکیلینگ مناسب نیست

۵(تفاوت Accuracy و ARI و NMI)

«چرا Accuracy برای خوشه‌بندی مناسب نیست؟»

ما داریم درباره خوشه‌بندی بدون ناظر (Clustering) صحبت می‌کنیم، ولی برچسب واقعی هم داریم (برای ارزیابی).

حالا سه معیار داریم:

- Accuracy
- ARI
- NMI

سؤال اینه:

چرا Accuracy معیار خوبی نیست،
ولی NMI و ARI مناسب‌ترند؟

Accuracy می‌پرسد:

«چند درصد نمونه‌ها دقیقاً همان برچسب درست را گرفته‌اند؟»

$$Accuracy = \frac{\text{تعداد پیش‌بینی درست}}{\text{کل نمونه‌ها}}$$

نکته‌ی خیلی مهم:

Accuracy فرض می‌کند:

- عدد برچسب مهم است
- کلاس ۰ یعنی ۰
- کلاس ۱ یعنی ۱

این فرض در Classification درست است 
اما در Clustering کاملاً غلط است 

چرا Accuracy برای خوشبندی بدون ناظر مناسب نیست؟

در خوشبندی:

شماره‌ی خوشها دلخواه و بدون معنا هستند

یعنی:

- خوش 0 می‌توانست خوش 5 باشد
- خوش 1 می‌توانست خوش 2 باشد

حتی با Accuracy هم Hungarian محدود است

فرض کنیم Hungarian استفاده کردیم و:

- بهترین نگاشت خوش \leftrightarrow کلاس را پیدا کردیم

حالا Accuracy عدد معقولی می‌دهد،
اما باز هم یک مشکل بزرگ دارد

فقط می‌پرسد: Accuracy

«در نهایت چند تا درست شد؟»

اما نمی‌پرسد:

- ساختار خوشبندی چقدر خوب بوده؟
- آیا خوشها بیش از حد ریز شده‌اند؟
- آیا داده‌ها تصادفی کنار هم افتاده‌اند؟

ARI چه ویژگی خاصی دارد؟

«آیا روابط بین جفتداده‌ها درست حفظ شده یا نه؟»

ARI مهم ویژگی

خوشبندی تصادفی $\rightarrow ARI \approx 0$

بیش از حد خوشسازی \rightarrow جرمیه می‌شود

NMI چه ویژگی خاصی دارد؟

«دانستن خوشها، چقدر درباره کلاس‌های واقعی به ما اطلاعات میدهد؟»

اگر هر خوش تقریباً فقط یک کلاس باشد → NMI بالا

اگر خوشها قاطی باشند → NMI پایین

ویژگی مهم NMI

- نسبت به جابه‌جایی برچسب‌ها پایدار است
- به ساختار کلی داده حساس است
- به عدد برچسب‌ها کاری ندارد

اگر فقط به Accuracy تکیه کنیم چه سوءبرداشت‌هایی ایجاد می‌شود؟

سوءبرداشت 1

خوشبندی خوب است
در حالی که فقط برچسب‌ها اتفاقی همنام شده‌اند

سوءبرداشت 2

خوشبندی بد است
در حالی که فقط شماره خوشها جابه‌جا شده‌اند

سوءبرداشت 3

مدل ساده بهتر از مدل پیچیده است
چون Accuracy بیشتری دارد
(در حالی که ساختار داده را نابود کرده)

۶) تفاوت مفهومی K-Means، DBSCAN و Hierarchical Clustering

→ « تقسیم‌بندی بر اساس فاصله تا مرکز »

فرض اصلی:

خوش‌ها کروی (گرد) هستند

همهی خوش‌ها اندازه و تراکم تقریباً مشابه دارند

چرا؟

چون از فاصله اقلیدسی تا مرکز (میانگین) استفاده می‌کند

نتیجه:

برای خوش‌های کشیده، هلالی، پیچ‌خورده → بد عمل می‌کند

حتماً باید K را از قبل بدانیم

اگر K را اشتباه انتخاب کنیم:

- نتیجه کاملاً عوض می‌شود

- بسیار حساس به مقیاس

- حتماً نیاز به:

- normalization ○

- standardization ○ یا

پارامتر مهم:

K •
initialization •

نویز را تشخیص نمی‌دهد

- هر نقطه‌ای باید به یک خوش برود
- می‌تواند مرکز خوش را جابه‌جا کند outlier •

→ «پیدا کردن نواحی پرتراکم» → **DBSCAN**

فرض اصلی:

خوشها نواحی با چگالی بالا هستند

شکل خوشمه مهم نیست (کروی، هلالی، مارپیچ...)

نتیجه:

برای خوشها غیر کروی عالی است

میتواند شکل واقعی داده را کشف کند

• بهشدت حساس به:

(*epsilon* (ε) ○
min_samples ○

اگر مقیاس درست نباشد:

• ع معنی خود را از دست می‌دهد

• کل خوشبندی خراب می‌شود

DBSCAN در scaling  است

نویز را بهصورت طبیعی تشخیص می‌دهد

• نقاط کمچگالی → *noise*

1- • برچسب

 یکی از بزرگترین مزیت‌های DBSCAN

نیازی به *K* ندارد

• تعداد خوشها بهصورت خودکار از داده در می‌آید

• حتی ممکن است:

○ خوش θ

○ یا خوشها کم/زیاد

«ساختن ساختار درختی از داده‌ها» → **Hierarchical**

فرض اصلی:

داده‌ها دارای ساختار سلسله‌مراتبی هستند

خوشه‌بندی را می‌شود در سطوح مختلف دید

شكل خوش:

به نوع **linkage** بستگی دارد

می‌تواند کروی یا غیرکروی باشد

حساس به **linkage**:

→ زنجیره‌ای •

→ فشرده •

→ متعادل •

→ کروی‌تر •

به مقیاس هم حساس است (بسته به **linkage**)

نویز را صریح تشخیص نمی‌دهد

• ولی در **dendrogram**

◦ نقاط دورافتاده دیرتر وصل می‌شوند

• می‌توان با برش مناسب اثر نویز را کم کرد

سؤال ۷ — K-Means بودن تابع هدف non-convex

تابع هدف K-Means چه اثری روی بهینه‌سازی دارد؟

تابع هدف K-Means چی بود؟

مجموع مربع فاصله‌ی نقاط تا مراکز خوش‌ها (WCSS)

وقتی:

- همه‌ی مراکز خوش‌ه را همزمان در نظر می‌گیریم

این تابع:

- چندین مینیمم محلی دارد
- فقط یک مینیمم سراسری دارد (بهترین حالت ممکن)

چرا initialization های مختلف جواب‌های مختلف می‌دهند؟

زیرا نقطه‌ی شروع مسیر را تعیین می‌کند

مشکل این است که:

- K-Means فقط پایین رفتن را بلد است
- بالا رفتن و دوباره پایین آمدن بلد نیست

یعنی:

- اگر به یک مینیمم محلی رسید
- دیگر راه فرار ندارد

در عمل چه کار می‌کنیم که این مشکل کمتر شود؟

راهکار ۱: اجرای چندباره الگوریتم (`n_init`)

رایج‌ترین و ساده‌ترین راه:

- K-Means را چند بار اجرا می‌کنیم
- هر بار با مراکز اولیه‌ی متفاوت
- بهترین نتیجه (کمترین WCSS) را انتخاب می‌کنیم

راهکار 2: استفاده از ++K-Means

بهجای انتخاب کاملاً تصادفی مراکز اولیه:

- مرکز اول تصادفی
- بقیه مراکز:
 - دور از هم
 - با احتمال بیشتر برای نقاط دورتر

نتیجه: 

- شروع بهتر
- افتادن کمتر در مینیمم‌های بد

در :sklearn

`init="k-means++"`

راهکار 3: نرمال‌سازی / داده Scaling

چون:

- K-Means به فاصله حساس است
- اگر یک ویژگی خیلی بزرگ باشد، همه چیز را خراب می‌کند

:پس

- StandardScaler
- MinMaxScaler

این باعث می‌شود: 

سطح تابع هدف «منطقی‌تر» شود

راهکار 4: انتخاب K مناسب

- K خیلی کم → خوشبندی بد
- K خیلی زیاد → مینیمم‌های محلی زیاد

Elbow و Silhouette کمک می‌کنند:

- در ناحیه‌ی معقول K کار کنیم

راهکار 5: چند الگوریتم یا چند Seed مختلف

در پروژه‌های جدی:

- چند بار اجرا
- با seed های مختلف
- نتایج مقایسه می‌شود

