

# بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

درس : یادگیری ماشین

استاد : خانم دکتر زرین بال ماسوله

موضوع : *clustering validity index*

دانشجو : فرزاد محسنی

مقطع : کارشناسی ارشد

رشته تحصیلی : مهندسی آیتی

کد دانشجویی :

دانشگاه : صنعتی امیرکبیر

<b><i>Method</i></b> -----	<b><i>Page</i></b>
<b><i>Silhouette Score</i></b> -----	<b><i>3</i></b>
<b><i>Calinski-Harabasz Index</i></b> -----	<b><i>5</i></b>
<b><i>Davies-Bouldin Index</i></b> -----	<b><i>7</i></b>
<b><i>Adjusted Rand Index</i></b> -----	<b><i>9</i></b>
<b><i>Normalized Mutual Information – NMI</i></b> -----	<b><i>11</i></b>
<b><i>Comparison Table of Clustering Validity Indices</i></b> -----	<b><i>13</i></b>

## Silhouette Score

### امتیاز سیلوئت (Silhouette Score) چیست؟

- Silhouette Score یک معیار برای ارزیابی کیفیت خوشه بندی است. این معیار بررسی می کند که هر نقطه داده (Data Point) تا چه حد به خوشه ی خودش تعلق دارد و چقدر از خوشه های دیگر جداست.

### چه زمانی از Silhouette Score استفاده می شود؟

- برای ارزیابی کیفیت خوشه بندی بدون نیاز به برچسب های واقعی (Unsupervised).
- برای مقایسه ی تعداد خوشه ها: مثلاً اگر نمی دانی چند خوشه مناسب تر است، می توانی سیلوئت را برای مقادیر مختلف K محاسبه کنی و بهترین را انتخاب کنی.
- در الگوریتم های خوشه بندی مثل K-Means، DBSCAN، Agglomerative Clustering و غیره کاربرد دارد.

### مزایا

- بدون نیاز به برچسب (Label-Free)
- ساده و قابل فهم
- امکان بررسی خوشه بندی به صورت نقطه به نقطه

### معایب

- محاسبه اش می تواند برای دیتاست های بزرگ کند باشد.
- به معیار فاصله ی مورد استفاده حساس است (مثلاً فاصله اقلیدسی، کسینوسی، و غیره).

### فرمول امتیاز سیلوئت

برای هر نقطه داده  $i$ ، امتیاز سیلوئت به صورت زیر محاسبه می‌شود:

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

که در آن:

- $a(i)$ : میانگین فاصله‌ی نقطه  $i$  تا تمام نقاط در همان خوشه (درونی بودن خوشه).
- $b(i)$ : کمترین میانگین فاصله‌ی نقطه  $i$  تا نقاط خوشه‌ی نزدیک دیگر (بیرونی بودن خوشه).

### مقدار امتیاز سیلوئت

- $s(i)$  همیشه بین  $-1$  تا  $+1$  است:
  - اگر نزدیک به  $1$  باشد: یعنی نقطه به خوشه‌ی خودش بسیار خوب تعلق دارد.
  - اگر نزدیک به  $0$  باشد: یعنی نقطه در مرز بین دو خوشه است.
  - اگر منفی باشد: یعنی احتمالاً نقطه به خوشه‌ی اشتباهی تخصیص داده شده است.

## Calinski-Harabasz Index

### شاخص کالینسکی-هاراباز چیست؟

- شاخص Calinski-Harabasz که به آن Variance Ratio Criterion نیز می‌گویند، یک معیار برای ارزیابی کیفیت خوشه بندی است. این شاخص نسبت پراکندگی بین خوشه ها (بین خوشه‌ای) به پراکندگی درون خوشه ها (درون خوشه‌ای) را اندازه‌گیری می‌کند.

### چه زمانی استفاده می‌شود؟

- زمانی که می‌خواهی خوشه بندی بدون برچسب واقعی را ارزیابی کنی (unsupervised evaluation).
- برای تعیین تعداد بهینه خوشه ها (با مقایسه‌ی مقدار شاخص برای مقادیر مختلف  $k$ ).
- کاربرد در الگوریتم‌هایی مثل K-Means, K-Medoids و الگوریتم‌های سلسله مراتبی (hierarchical clustering).

### مزایا

- محاسبه سریع و مؤثر
- بدون نیاز به داده های برچسب خورده
- قابل تفسیر ساده: مقدار بیشتر = خوشه بندی بهتر

### معایب

- به شکل خوشه ها و اندازه‌ی داده ها حساس است
- فرض می‌کند خوشه ها شکل کروی و اندازه مشابهی دارند

فرمول شاخص کالینسکی-هاریابازش

$$CH = \frac{\text{Tr}(B_k)}{\text{Tr}(W_k)} \cdot \frac{N - k}{k - 1}$$

که در آن:

- $\text{Tr}(B_k)$  مجموع مربعات بین خوشه ها (Between-cluster dispersion)
- $\text{Tr}(W_k)$  مجموع مربعات درون خوشه ها (Within-cluster dispersion)
- $N$ : تعداد کل نمونه ها
- $k$ : تعداد خوشه ها

تفسیر شاخص

- مقدار بالاتر بهتر است.
- یعنی خوشه ها به خوبی از هم جدا هستند و داده ها درون هر خوشه فشرده و نزدیک به هم اند.
- شاخص CH مقدار منفی ندارد و می تواند از صفر تا بی نهایت باشد.
- اگر خوشه ها با هم ادغام یا پراکنده باشند، مقدار شاخص کاهش می یابد.

## Davies-Bouldin Index

### شاخص دیویس-بولدین چیست؟

- Davies-Bouldin Index یا DBI یکی دیگر از معیارهای ارزیابی کیفیت خوشه بندی است. این شاخص بررسی می کند که خوشه ها چقدر از هم جدا هستند و چقدر درون خود فشرده اند.

### چه زمانی استفاده می شود؟

- برای ارزیابی کیفیت خوشه بندی در حالت بدون نظارت (بدون نیاز به برچسب ها).
- برای مقایسه خوشه بندی ها با تعداد مختلف خوشه و انتخاب بهترین K.
- کاربرد در الگوریتم هایی مثل K-Means، DBSCAN، Agglomerative Clustering و غیره.

### مزایا

- مستقل از برچسب
- محاسبه ی نسبتاً ساده
- ترکیب خوبی از فشردگی داخلی و جدایی بین خوشه ها را می سنجد

### معایب

- به شکل و توزیع خوشه ها حساس است
- اگر یک خوشه بسیار بزرگ یا پراکنده باشد، می تواند مقدار کل شاخص را خراب کند

فرمول شاخص دیویس-بولدین

$$DBI = \frac{1}{k} \sum_{i=1}^k \max_{j \neq i} \left( \frac{s_i + s_j}{d_{ij}} \right)$$

که در آن:

- $k$ : تعداد خوشه ها
- $s_i$  میانگین فاصله نقاط در خوشه  $i$  تا مرکز خوشه (درون خوشه ای یا فشرده‌گی خوشه)
- $d_{ij}$  فاصله بین مراکز خوشه های  $i$  و  $j$

تفسیر شاخص DBI

- مقدار کمتر بهتر است
- چون نشان می دهد خوشه ها از یکدیگر جدا هستند و داده های هر خوشه نزدیک به هم هستند.
- مقدار DBI همیشه بزرگ تر یا مساوی صفر است.
- در حالت ایده آل، هر خوشه باید کمترین هم پوشانی را با خوشه های دیگر داشته باشد.



## Adjusted Rand Index

### شاخص Adjusted Rand Index (ARI) چیست؟

- شاخص رند تعدیل شده معیاری برای ارزیابی کیفیت خوشه بندی است در صورتی که برچسب های واقعی (Ground Truth Labels) در دسترس باشند.
- این شاخص بررسی می کند که چقدر خوشه بندی به دست آمده با برچسب های واقعی هماهنگ و منطبق است.

### تفاوت با Rand Index معمولی؟

- شاخص رند ساده فقط میزان تطابق بین دو دسته بندی را بررسی می کند، ولی ARI تعدیل شده برای تطابق تصادفی است.
- یعنی اگر دو دسته بندی کاملاً تصادفی باشند، ARI نزدیک به صفر خواهد بود، نه مقدار بالا.

### چه زمانی استفاده می شود؟

- فقط زمانی که برچسب های واقعی داده ها موجود باشند (برای ارزیابی دقیق).
- مقایسه ی عملکرد الگوریتم های خوشه بندی در داده های labeled.
- بسیار پرکاربرد در ارزیابی K-Means، Agglomerative Clustering و سایر الگوریتم ها با ground truth.

### مزایا

- دقت بالا در مقایسه با معیار های ساده
- تعدیل شده برای جلوگیری از نتایج گمراه کننده در حالت تصادفی
- قابل اعتماد برای مقایسه با ground truth

### معایب

- فقط در صورت وجود برچسب های واقعی قابل استفاده است
- ممکن است محاسبات برای داده های بسیار بزرگ کمی کند شود

فرمول کلی

$$ARI = \frac{\text{Index} - \text{Expected Index}}{\text{Max Index} - \text{Expected Index}}$$

جزئیات کامل فرمول با استفاده از ماتریس احتمال هم قابل نمایش است، اما ایده‌ی اصلی این است:

•  $ARI = 1 \rightarrow$  خوشه بندی کاملاً مطابق با برجسب های واقعی

•  $ARI \approx 0 \rightarrow$  تطابقی بیش تر از حالت تصادفی ندارد

•  $ARI < 0 \rightarrow$  بدتر از تصادف (نادر)

ویژگی ها

• دامنه: از -1 تا 1

• ۱: تطابق کامل

• ۰: تصادفی

• کمتر از صفر: بدتر از حالت تصادفی

## Normalized Mutual Information – NMI

### NMI چیست؟

- (NMI) یک معیار برای مقایسه دو خوشه بندی است، معمولاً بین خوشه بندی حاصل از الگوریتم و برچسب های واقعی (Ground Truth).
- این معیار می سنجد که چه مقدار اطلاعات بین دو دسته بندی مشترک است.

### ایده ی اصلی:

- اطلاعات متقابل (Mutual Information - MI) میزان اشتراک اطلاعات بین دو متغیر تصادفی (در اینجا: دو خوشه بندی) را اندازه می گیرد.
- اما MI به تعداد نمونه ها و اندازه ی خوشه ها حساس است. برای رفع این مشکل، آن را نرمال سازی (Normalization) می کنند.

### چه زمانی استفاده می شود؟

- زمانی که برچسب های واقعی در دسترس هستند
- برای مقایسه ی خوشه بندی ها با ground truth یا حتی بین دو الگوریتم خوشه بندی
- کاربرد در ارزیابی K-Means, Spectral Clustering و سایر الگوریتم ها

### مزایا

- مقیاس بندی شده و قابل مقایسه بین داده های مختلف
- برابر است در صورت تعویض برچسب خوشه ها (invariant to label permutations)
- مناسب برای خوشه بندی های با تعداد متفاوت خوشه

## معایب

- فقط برای داده های دارای برچسب قابل استفاده است
- نیاز به محاسبه آنترופی و اطلاعات متقابل (کمی پیچیده تر از برخی معیارها)

## فرمول NMI

$$NMI(U, V) = \frac{2 \cdot I(U; V)}{H(U) + H(V)}$$

## که در آن:

- **U**: خوشه بندی پیش بینی شده
- **V**: خوشه بندی واقعی
- **I(U;V)**: اطلاعات متقابل بین
- **H(U)** و **H(V)**: آنترופی خوشه بندی ها (میزان عدم قطعیت)

## دامنه و تفسیر

- $NMI \in [0, 1]$
- ۱ = تطابق کامل بین خوشه بندی و برچسب های واقعی
- ۰ = هیچ ارتباط اطلاعاتی بین دو دسته بندی وجود ندارد

## جدول مقایسه شاخص های اعتبار خوشه بندی

شاخص	نوع ارزیابی	نیاز به برچسب واقعی؟	دامنه مقدار	معیار بهتر بودن	کاربرد اصلی	مزایا	معایب
<b>Silhouette Score</b>	درونی (Internal)	ندارد	-1 تا +1	مقدار بیشتر بهتر	ارزیابی فشردگی درون خوشه‌ای و جدایی بین خوشه‌ای	شهودی، بدون نیاز به برچسب، قابل تفسیر	در دیتاست های بزرگ کند می‌شود
<b>Calinski-Harabasz Index</b>	درونی (Internal)	ندارد	0 تا $\infty$	مقدار بیشتر بهتر	انتخاب تعداد بهینه خوشه ها	سریع، ساده، بدون نیاز به برچسب	فرض کروی بودن خوشه ها
<b>Davies-Bouldin Index</b>	درونی (Internal)	ندارد	0 تا $\infty$	مقدار کمتر بهتر	بررسی جدایی و فشردگی خوشه ها	ترکیب فشردگی و جدایی، بدون نیاز به برچسب	به خوشه های پراکنده حساس است
<b>Adjusted Rand Index</b>	بیرونی (External)	دارد	-1 تا +1	مقدار بیشتر بهتر	ارزیابی دقت نسبت به برچسب واقعی	تعدیل شده برای تصادف، دقیق	فقط در صورت داشتن برچسب کاربرد دارد
<b>Normalized Mutual Information - NMI</b>	بیرونی (External)	دارد	0 تا 1	مقدار بیشتر بهتر	مقایسه خوشه بندی با برچسب یا خوشه دیگر	نرمال شده، مستقل از ترتیب برچسب ها	نیاز به محاسبه آنترופی و MI