

курс «Прикладные задачи анализа данных»

**Функции ошибки /
функционалы качества**
**Часть 3: многоклассовые задачи,
ранжирование, кластеризация**

Александр Дьяконов

23 октября 2020 года

План на несколько лекций

задача регрессии

задача бинарной классификации

- **чёткая классификация**
- **скоринговые функции**

задача классификации с несколькими классами

задачи ранжирования

задачи кластеризации

Weighted kappa

Если есть разумные веса ошибок за конкретные несогласованности

Когда это бывает?

$$\kappa = 1 - \frac{\sum_{i=1}^l \sum_{j=1}^l w_{ij} m_{ij}}{\sum_{i=1}^l \sum_{j=1}^l w_{ij} s_{ij}} \in [-1, +1]$$

**матрица случайных
ответов**

$$s_{ij} = m_{i:} m_{:j} = \sum_j m_{ij} \sum_i m_{ij}$$

$$s_{ij} \leftarrow \frac{s_{ij}}{m^2} m = \frac{s_{ij}}{\sum_{tr} s_{tr}} \sum_{tr} m_{tr}$$

квадратичные веса

$$w_{ij} = \frac{(i - j)^2}{(l - 1)^2}$$

м.б. любая весовая схема

Вычисление Quadratic Weighted Кэрра

ответы

	у	а
0	1	1
1	1	1
2	1	2
3	2	1
4	2	3
5	3	2
6	3	3
7	3	3
8	1	2
9	2	2

матрица ошибок

у	а	1	2	3
1	2	2	0	
2	1	1	1	
3	0	1	2	

матрица случайных ответов

	0	1	2
0	12	16	12
1	9	12	9
2	9	12	9

матрица весов

	0	1	2
0	0.00	0.25	1.00
1	0.25	0.00	0.25
2	1.00	0.25	0.00

после нормировки

	0	1	2
0	1.2	1.6	1.2
1	0.9	1.2	0.9
2	0.9	1.2	0.9

WK = 0.615

Вычисление Quadratic Weighted Кappa

```
from sklearn.metrics import cohen_kappa_score  
cohen_kappa_score(df.y, df.a, weights='quadratic')
```

```
n = pd.crosstab(df.y, df.a)  
n = n.values  
m = np.outer(n.sum(axis=1) , n.sum(axis=0))  
m = m / m.sum() * n.sum()  
w = (np.arange(1, 4)[:,np.newaxis] -  
      np.arange(1, 4)) ** 2 / ((3-1)*(3-1))  
1 - (np.sum(n*w) / np.sum(m*w))
```

Потом **докажем**:

Каппа Коэна частный случай Weighted Кappa для 2х классов

Quadratic Weighted Kappa

Применяется в задачах, где классы упорядочены
«ранжирование»

	y	1.0	0.83	0.83	0.33	0.8	0.0	-1.0
0	0	0	0	0	0	0	0	2
1	0	0	0	0	0	0	1	2
2	0	0	1	0	2	0	2	2
3	1	1	1	1	1	0	0	1
4	1	1	1	1	1	0	1	1
5	1	1	0	2	1	0	2	1
6	2	2	2	2	2	2	0	0
7	2	2	2	2	2	2	1	0
8	2	2	2	1	0	2	2	0

Многоклассовая задача «Multi-label»

матрица классификаций

$$\| y_{ij} \|_{m \times l}$$

	class 1	class 2	class 3
0	1	0	0
1	0	1	0
2	0	0	1
3	1	1	0

матрица ответов

$$\| a_{ij} \|_{m \times l}$$

	class 1	class 2	class 3
0	0.75	0.00	0.25
1	0.00	0.50	0.25
2	0.25	1.00	0.25
3	0.00	0.25	0.75

По сути, надо сравнить матрицы на похожесть

микро-подход	можно сравнивать матрицы как векторы
макро-подход	можно сравнивать столбцы матриц
по объектам	можно сравнивать строки матриц и усреднять

Многоклассовая задача: Hamming Loss

Число (процент) ошибок

$$a_{ij} \in \{0,1\}$$

$$\text{HL} = \frac{1}{ml} \sum_{i=1}^m \sum_{j=1}^l I[y_{ij} \neq a_{ij}]$$

(1 – точность)

Многоклассовая задача: Log Loss (cross-entropy)

**Естественное обобщение
логистической ошибки**

$$a_{ij} \in [0,1]$$

$$\text{logloss} = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^l y_{ij} \log a_{ij}$$

(тонкость: лучше для непересекающихся классов)

Многоклассовая задача: Mean Probability Rate

(это функционал качества, $a_{ij} \in [0,1]$ ~ распределение)

Есть макро-версия, см. дальше

$$\text{MPR} = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^l y_{ij} a_{ij}$$

$$\text{MAPR} = \frac{1}{l} \sum_{j=1}^l \frac{\sum_{i=1}^m y_{ij} a_{ij}}{\sum_{i=1}^m y_{ij}}$$

Многоклассовая задача: MSE, MAE

$$\text{MSE} = \frac{1}{ml} \sum_{i=1}^m \sum_{j=1}^l (y_{ij} - a_{ij})^2$$

$$\text{MAE} = \frac{1}{ml} \sum_{i=1}^m \sum_{j=1}^l |y_{ij} - a_{ij}|$$

это всё вариации на тему схожести / различия бинарного и вещественного вектора

Многоклассовый AUCROC: Макро-усреднение

$$\text{AUC} = \frac{1}{l} \sum_{j=1}^l \text{AUC}_j$$

AUC_j – значение функционала в задаче бинарной классификации

«j-й класс / не j-й класс»

$$\{(x_i, I[y(x_i)_{[j]} = 1])\}_{i=1}^m$$

Многоклассовый AUCROC: Весовое макро-усреднение

$$\text{AUC} = \frac{\sum_{j=1}^l P_j \text{AUC}_j}{\sum_{j=1}^l P_j}$$

P_j – вероятность j-го класса

(процент «1» в столбце матрицы классификации)

Многоклассовый AUCROC: Микро-усреднение

значение функционала в задаче

$$\{((x_i, j), I[y(x_i)_{[j]} = 1])\}_{i=1, j=1}^{m, l}$$

«вытягиваем матрицу ответов в вектор»

Многоклассовый AUCROC: Усреднение по объектам

$$AUC = \frac{1}{m} \sum_{i=1}^m AUC'_i$$

AUC'_i – **значение функционала в задаче**

$$\{((x_i, j), I[y(x_i)_{[j]} = 1])\}_{j=1}^l$$

«решение задачи по строкам»

Многомерный AUC: минутка кода

```
from sklearn.metrics import roc_auc_score

roc_auc_score(y, a, average='macro')
# эквивалентно:
auc_pclass = [roc_auc_score(y[:,i], a[:,i]) for i in range(1)]
auc_pclass, mean(auc_pclass)

roc_auc_score(y, a, average='micro')
# эквивалентно:
roc_auc_score(y.ravel(), a.ravel())

roc_auc_score(y, a, average='weighted')
# эквивалентно:
w = y.sum(axis=0)
sum(np.array(auc_pclass) * w) / sum(w)

roc_auc_score(y, a, average='samples')
# эквивалентно:
auc_pinstance = [roc_auc_score(y[i,:], a[i,:]) for i in
range(m)]
auc_pinstance, mean(auc_pinstance)
```

Многомерный AUC ROC

матрица классификаций

	class 1	class 2	class 3
0	1	0	0
1	0	1	0
2	0	0	1
3	1	1	0

macro	micro	weighted	samples
0.49	0.53	0.52	0.56

матрица ответов

	class 1	class 2	class 3
0	0.75	0.00	0.25
1	0.00	0.50	0.25
2	0.25	1.00	0.25
3	0.00	0.25	0.75

	class 0	class 1	class 2
AUC_per_class	0.62	0.5	0.33
P_per_class	0.50	0.5	0.25

	class 0	class 1	class 2	class 3
AUC_per_instance	1.0	1.0	0.25	0.0

Точность: сравнение макро- и микро- усреднения

$$P_j = \frac{TP_j}{TP_j + FP_j}$$

$$P_{\text{macro-mean}} = \frac{1}{l} \sum_{j=1}^l P_j$$

$$P_{\text{micro-mean}} = \frac{\sum_{j=1}^l TP_j}{\sum_{j=1}^l TP_j + \sum_{j=1}^l FP_j}$$

**Кстати, макро-усреднение делают по-разному.
Часто: среднее геометрическое.**

Точность: сравнение макро- и микро- усреднения
При вычислении каких-то функционалов, например точности

МАКРО	<ul style="list-style-type: none">• вычислить точность для каждого класса• усреднить полученные точности $(H.TP / (H.TP + H.FP)) .mean ()$	0.344
МИКРО	<ul style="list-style-type: none">• вычислить точность сразу для всех классов $\begin{aligned} TP &= H.TP.sum () + H.TP.sum () \\ FP &= H.FP.sum () + H.FP.sum () \\ P &= TP / (TP + FP) \end{aligned}$	0.246

	TP	FP		P
class 1	2	2	class 1	0.50
class 2	5	10	class 2	0.33
class 3	10	40	class 3	0.20

Точность: сравнение макро- и микро- усреднения

	TP	FP
class 1	2	2
class 2	5	10
class 3	100	400

не изменились точности по классам \Rightarrow не изменилась макро-точность 0.344

изменились TP, FP по классам \Rightarrow микро-точность смещается в сторону «большого» класса: 0.206 (вместо 0.246)

где какое усреднение лучше использовать?

Совет: смотреть дисперсию показателей по классам

F-мера – ещё больше вариантов усреднения

Макро F-мера

$$F_{\text{macro-mean}} = \frac{1}{l} \sum_{j=1}^l F_j$$

на основе других макро-параметров

$$F = \frac{2}{\frac{1}{P_{\text{macro-mean}}} + \frac{1}{R_{\text{macro-mean}}}}$$

ДЗ провести сравнение!

Сбалансированная точность «Balanced accuracy» – макро-усреднение полноты

Сбалансированная точность (accuracy)
не есть усреднение точностей (precision)

$$BA = \frac{1}{l} \sum_{j=1}^l R_j = \frac{1}{l} \sum_{j=1}^l \frac{\sum_{t=1}^m I[y(x_t)_{[j]} = 1] I[a(x_t)_{[j]} = 1]}{\sum_{t=1}^m I[y(x_t)_{[j]} = 1]}$$

```
from sklearn.metrics import balanced_accuracy_score
```

Другие (неэквивалентные) определения:

$$BA = \frac{1}{l} \sum_{j=1}^l \min[P_j, R_j]$$

$$BA = \frac{1}{l} \sum_{j=1}^l \min[\text{sens}_j, \text{spec}_j]$$

Д3 Сравните разные подходы

Пример: соревнование LSHTC

Id, Predicted			
1,	12	35	200
2,	54	55	
3,	11		
4,	1	7	101
...			

$$\tilde{F} = \frac{2\tilde{P}\tilde{R}}{\tilde{P} + \tilde{R}}$$

$$\tilde{P} = \frac{1}{l} \sum_{j=1}^l \frac{TP_j}{TP_j + FP_j}$$

$$\tilde{R} = \frac{1}{l} \sum_{j=1}^l \frac{TP_j}{TP_j + FN_j}$$

Решающее правило с отсечкой:

$$\alpha_{ij} = \begin{cases} 1, & \gamma_{ij} \geq \min(c, \max\{\gamma_{ij}\}_{j=1}^l), \\ 0, & \text{иначе.} \end{cases}$$

Решать задачу по вертикали / по горизонтали

Оценка результатов поиска/рекомендаций



Задача с бинарной релевантностью

$$x_1 \prec x_2 \prec \dots \prec x_m$$

$y_i = 1$ – релевантный объект

$y_i = 0$ – нерелевантный объект

Задача ранжирования

Целевой признак может быть бинарным, но это не задача классификации

Precision at n

Точность на первых n элементах

$$p @ n = \frac{y_1 + \dots + y_n}{n}$$

Average Precision at n

Средняя точность на первых n элементах

$$ap @ n = \sum_{k=1}^n \frac{P(k)}{\min(n, r)}$$

**r – мощность множества релевантных объектов
(товаров, документов)**

n – сколько рекомендаций будет учитываться

$$P(k) = \begin{cases} p @ k, & y_k = 1, \\ 0, & y_k = 0, \end{cases}$$

y_i – бинарное значение релевантности

Mean Average Precision

– усреднение $ap @ n$ по всем пользователям

Average Precision at n

Примеры (три релевантных объекта, $r = 3$):

$0 \prec 0 \prec 0$

$$ap @ 3 = \frac{1}{3} [0 + 0 + 0]$$

$0 \prec 0 \prec 1$

$$ap @ 3 = \frac{1}{3} \left[0 + 0 + \frac{1}{3} \right]$$

$0 \prec 1 \prec 1$

$$ap @ 3 = \frac{1}{3} \left[0 + \frac{1}{2} + \frac{2}{3} \right]$$

$1 \prec 0 \prec 0$

$$ap @ 3 = \frac{1}{3} \left[\frac{1}{1} + 0 + 0 \right]$$

$0 \prec 0 \prec 1 \prec 1 \prec 1$

$$ap @ 5 = \frac{1}{3} \left[0 + 0 + \frac{1}{3} + \frac{2}{4} + \frac{3}{5} \right]$$

$1 \prec 1 \prec 1 \prec 0 \prec 0$

$$ap @ 5 = \frac{1}{3} \left[\frac{1}{1} + \frac{2}{2} + \frac{3}{3} + 0 + 0 \right]$$

Concordant – Discordant ratio

$$\frac{|\{(i, j) \mid y_i > y_j, 1 \leq i < j \leq m\}|}{|\{i \mid y_i = 1\}| \cdot |\{j \mid y_j = 0\}|}$$

Упорядочили: E, D, C, B, A (по убыванию релевантности)

На самом деле: B, E – релевантные

Пары «релевантный» – «нерелевантный»:

BA

EA

BC

EC

BD

ED

Качество упорядочивания: 4 / (2 + 4)

~ AUC ROC

**Что ещё может встретиться...
в задачах рекомендации**

$$\frac{1}{|Z|} \sum_{z \in Z} \frac{|\{x_1, \dots, x_z\} \cap \{x'_1, \dots, x'_z\}|}{z}$$

x_1, \dots, x_n – **упорядоченный** список ответов

x'_1, \dots, x'_m – **все релевантные**

$$Z \subseteq \{1, 2, \dots, n\}$$

$$Z = \{5, 10, 15, 20, 25, 30\}$$

когда логично применить?

Рекомендации



Новинка

Ritmix RH-126M, Black Green наушники

[К сравнению](#) [В избранное](#) [Поделиться](#)

Цвет: черный, зеленый

Тип: Наушники

Модель: 15119162

Тип соединения: Проводные


Вид наушников: Вставные (внутриканальные), Спортивные

Конструкция наушников: Динамические, С микрофоном


[Перейти к описанию](#)




Рекомендуем также




от 115 ₽
Ritmix RH-120M наушники
[Все варианты](#)




от 282 ₽
Ritmix RH-180M наушники
[Все варианты](#)




289 ₽
Ritmix RH-158, Dark Venge наушники
[В корзину](#)




от 145 ₽
Ritmix RH-125 наушники
[Все варианты](#)




183 ₽
Ritmix RH-115M Luminous наушники
[В корзину](#)




от 220 ₽
Ritmix RH-150M наушники
[Все варианты](#)




709 ₽
Бестхарий и чувства
[В корзину](#)




1 032 ₽
Ritmix RH-567M Gaming игровые наушники
[В корзину](#)




934 ₽
Ritmix RH-565M Gaming игровые наушники
[В корзину](#)




65 ₽
Ritmix RH-011 наушники
[Все варианты](#)




47 ₽
Ritmix RH-004 наушники
[Все варианты](#)




65 ₽
Ritmix RH-011, Black наушники
[В корзину](#)



45 ₽
Ritmix RH-003 наушники
[Все варианты](#)



4 790 ₽
Ritmix RSK-615 электронная книга
[В корзину](#)



1 199 ₽
Ritmix RCH-106 WH автомобильный держатель
[В корзину](#)

Mean Reciprocal Rank (MRR)

- это усреднение Reciprocal rank (RR) по всем ранжированиям, который сделал алгоритм.

$$RR = \frac{1}{\min\{i : y_i = 1\}}$$

Часто оптимизируют именно его!

Классические функционалы в поиске

Случай небинарной релевантности

Выдали id документов/товаров/..., а их ценность (релевантность):

$$y_1, \dots, y_m$$

Cumulative Gain at n

$$CG@n = y_1 + \dots + y_n$$

Discounted Cumulative Gain at n

$$DCG@n = \sum_{i=1}^n \frac{2^{y_i} - 1}{\log_2(i + 1)}$$

Ещё вариант:

$$DCG@n = y_1 + \sum_{i=2}^n \frac{y_i}{\log_2(i)} = y_1 + y_2 + \frac{y_3}{\log_2 3} + \dots + \frac{y_n}{\log_2 n}$$

Цена ошибок за неправильное ранжирование

$$\frac{1}{\log_2(1+1)} - \frac{1}{\log_2(1+2)} \approx 0.37$$

$$\frac{1}{\log_2(1+10)} - \frac{1}{\log_2(1+11)} \approx 0.01$$

$$\frac{1}{\log_2(1+10)} - \frac{1}{\log_2(1+20)} \approx 0.06$$

Normalized DCG

$$nDCG = \frac{DCG}{IDCG}$$

IDCG = ideal DCG**для того, чтобы не было зависимости от длины выдачи**

Ещё подход к сравнению порядков:

Пусть алгоритм выдал

$$x_1 \prec x_2 \prec \dots \prec x_m$$

Правильный порядок

$$x_{i_1} \prec x_{i_2} \prec \dots \prec x_{i_m}$$

Надо сравнить:

$$(1, 2, \dots, m)$$

$$(i_1, i_2, \dots, i_m)$$

Ранговые корреляции...

Ещё подход к оценке ранжирования

Известны вероятности того, что объект является релевантным

$$p_i = p(x_i)$$

~ пользователь выберет ссылку

Expected reciprocal rank (ERR)

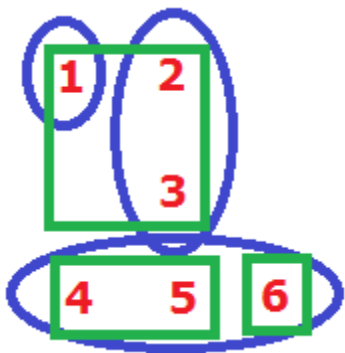
$$ERR @ n = \frac{1}{n} \sum_{k=1}^n \frac{1}{k} p_k \prod_{i < k} (1 - p_i)$$

Как интерпретировать?

Редакторское расстояние

Операции

- добавление к кластеру
- создание кластера с одним объектом
- удаление из кластера
- удаление кластера с одним объектом



```
1 2 3;4 5;6
1 2 3; 4 5 [delC]
2 3; 4 5 [del]
2 3; 4 5; 1 [insC]
2 3; 4 5 6; 1 [ins]
```

	2 3	4 5 6	1
1 2 3	1	6	2
4 5	4	1	3
6	3	2	2

Редакторское расстояние

- Плохо заносить не в тот кластер (целых две операции на перенос)
 - Плохо создавать неправильный кластер

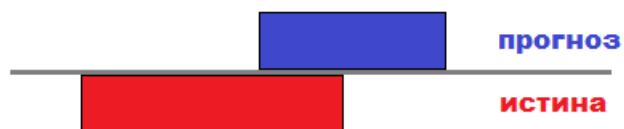
⇒ осторожный алгоритм



- Много зависит от операций...

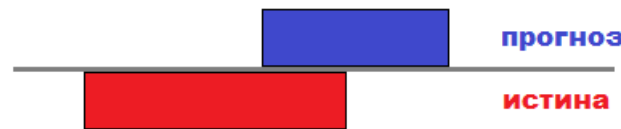
Задача с «неклассическим целевым вектором»

Надо предсказывать не значение,
а интервал $[a, b]$



Как измерить качество?

Задача с интервальным целевым вектором



Интервал – это множество!

Коэффициент Жаккара (Jaccard)

$$\frac{|A \cap B|}{|A \cup B|}$$

**коэффициент Шимкевича-Симпсона
(Szymkiewicz, Simpson)**

$$\frac{|A \cap B|}{\min(|A|, |B|)}$$

коэффициент Браун-Бланке (Braun-Blanquet)

$$\frac{|A \cap B|}{\max(|A|, |B|)}$$

См. википедию «Коэффициент сходства» для переноса идеи Колмогорова об обобщённом среднем...

Вариации на тему усреднения...

коэффициент Сёренсена (Sørensen)

$$\frac{2|A \cap B|}{|A| + |B|}$$

коэффициент Кульчинского (Kulczynsky)

$$\frac{|A \cap B|}{2} (1/|A| + 1/|B|)$$

коэффициент Отиаи (Ochiai)

$$\frac{|A \cap B|}{\sqrt{|A| \cdot |B|}}$$

Меры включения

$$\frac{\frac{|A \cap B|}{|A|}}{\frac{|A \cap B|}{2|A| - |A \cap B|}}$$

$$\frac{\frac{|A \cap B|}{|B|}}{\frac{|A \cap B|}{2|B| - |A \cap B|}}$$

Как решать задачи с интервалами? Потом вернёмся...

Оценка результатов кластеризации

Если знаем верную кластеризацию... внешняя оценка (External evaluation)

Вопрос: когда?

**ничего не знаем \Rightarrow согласованность с данными
внутренняя оценка (Internal evaluation)**

Оценка результатов кластеризации: «Internal evaluation»

Пусть чёткая (нет пересечений) кластеризация $U = u_1 \cup \dots \cup u_{|U|}$
множества $X = \{x_1, \dots, x_m\}$

Davies–Bouldin index

Использует центроиды и дисперсии

$$DB = \frac{1}{|U|} \sum_{i=1}^{|U|} \max_{j \neq i} \left(\frac{\sigma_i + \sigma_j}{d(c_i, c_j)} \right)$$

Dunn index =

min между кластерами / max внутри
кластерами

$$D = \frac{\min_{1 \leq i < j \leq |U|} d(u_i, u_j)}{\max d_{\text{in}}(u_i)}$$

Silhouette

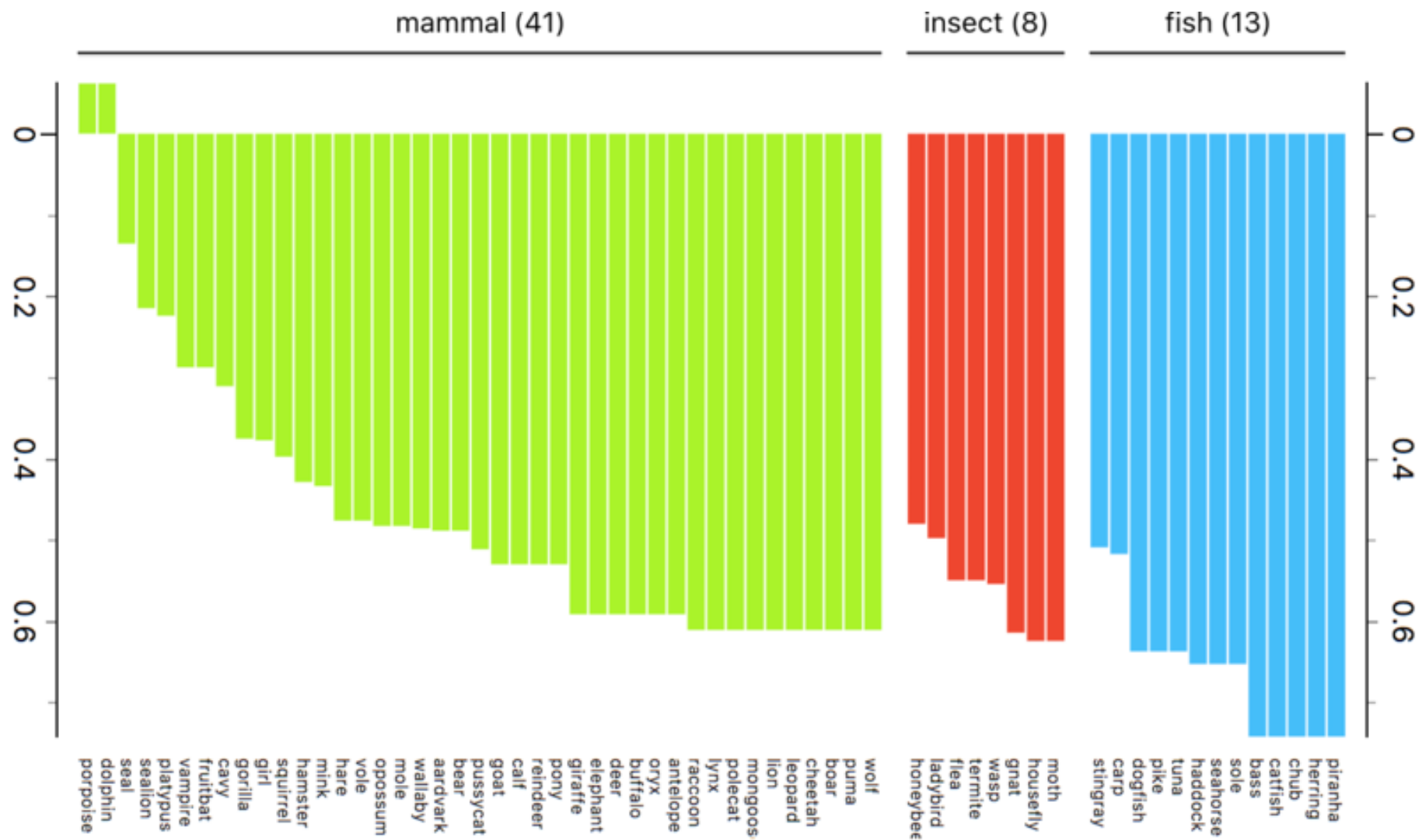
$$x_i \in u_1, d(x_i, u_2) \leq d(x_i, u_3) \leq \dots$$

Расстояние считается как среднее до
всех точек кластера

$$\text{silhouette}(x_i) = \frac{d(x_i, u_2) - d(x_i, u_1)}{\max(d(x_i, u_2), d(x_i, u_1))}$$

Можно усреднять по точкам

Оценка результатов кластеризации: «Internal evaluation»



[https://en.wikipedia.org/wiki/Silhouette_\(clustering\)](https://en.wikipedia.org/wiki/Silhouette_(clustering))

Calinski-Harabasz Index (Variance Ratio Criterion) ■

$$\frac{\text{trace}\left(\frac{1}{|U|-1} \sum_{i=1}^{|U|} |U_i| (x - c_i)(x - c_i)^T\right)}{\text{trace}\left(\frac{1}{m - |U|} \sum_{i=1}^{|U|} \sum_{x \in U_i} (x - c_i)(x - c_i)^T\right)}$$

след матрицы межклассовой ковариации / след матрицы внутриклассовой ковариации

лучше подходит для выпуклых кластеров и евклидовой метрики

External evaluation: взаимная информация

Пусть чёткие (нет пересечений) кластеризации

$$U = u_1 \cup \dots \cup u_{|U|}$$

$$V = v_1 \cup \dots \cup v_{|V|}$$

множества $X = \{x_1, \dots, x_m\}$

$$p_i = \frac{|u_i|}{m}$$

$$H(U) = - \sum_{i=1}^{|U|} p_i \log p_i$$

Аналогично $H(V)$

$$p_{ij} = \frac{|u_i \cap v_j|}{m}$$

$$MI = \sum_{i=1}^{|U|} \sum_{j=1}^{|V|} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

**потом MI ~ насколько более чётко определена U при знании V
уже её можно использовать...**

External evaluation: скорректированная взаимная информация Adjusted mutual information

$$AMI(U, V) = \frac{MI(U, V) - \mathbf{E}(MI(U, V))}{\max(H(U), H(V)) - \mathbf{E}(MI(U, V))}$$

1 – если кластеризации равны

~0 – если кластеризации случайны

матожидание можно вычислить аналитически

**нужно калибровать, т.к. чем больше кластеров в кластеризациях,
тем больше значение MI**

```
from sklearn.metrics import mutual_info_score # MI
from sklearn.metrics import normalized_mutual_info_score # [0, 1]
from sklearn.metrics.cluster import adjusted_mutual_info_score
adjusted_mutual_info_score([0, 0, 1, 1], [0, 0, 1, 1])
```

https://en.wikipedia.org/wiki/Adjusted_mutual_information

Д3 параметрический пример, в котором AMI меняется в своих пределах

External evaluation: V-мера

V – среднее гармоническое homogeneity и completeness

homogeneity ~ каждый кластер содержит только объекты отдельного класса

completeness ~ все объекты конкретного класса отнесены в один кластер

```
from sklearn.metrics.cluster import homogeneity_score
from sklearn.metrics.cluster import completeness_score
from sklearn.metrics.cluster import v_measure_score

v_measure_score([0, 0, 1, 1], [0, 0, 1, 1])
```

External evaluation: Adjusted Rand index

Аналогичная «Adjusted» идея, но проще...
поскольку кластеризация задаёт отношение эквивалентности

Rand index

$$R = \frac{|\{i, j : (i \sim_U j) \& (i \sim_V j)\}| + |\{i, j : (i \not\sim_U j) \& (i \not\sim_V j)\}|}{C_m^2}$$

теперь калибровка под случайную кластеризацию:

$X \backslash Y$	Y_1	Y_2	\dots	Y_s	Sums
X_1	n_{11}	n_{12}	\dots	n_{1s}	a_1
X_2	n_{21}	n_{22}	\dots	n_{2s}	a_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
X_r	n_{r1}	n_{r2}	\dots	n_{rs}	a_r
Sums	b_1	b_2	\dots	b_s	

Adjusted Index
 \widehat{ARI}

=

$\overbrace{\sum_{ij} \binom{n_{ij}}{2}}$
Index

$-\underbrace{[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}]/\binom{n}{2}}_{\text{Expected Index}}$

$\frac{1}{2}[\underbrace{\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}}_{\text{Max Index}}] - \underbrace{[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}]/\binom{n}{2}}_{\text{Expected Index}}$

```
from sklearn.metrics.cluster import adjusted_rand_score
adjusted_rand_score([0, 0, 1, 1], [0, 0, 1, 1])
```

https://en.wikipedia.org/wiki/Rand_index#Adjusted_Rand_index

External evaluation: общий подход

Кластеризация ~ классификация пар

$$\{x_1, \dots, x_m\} \rightarrow \{(1,1), \dots, (i,j), \dots, (m,m)\}$$

$$a_U(i, j) = 1 \Leftrightarrow i \sim_U j$$

Можно сравнивать классификации a_U и a_V

Пример, Rand index:
$$RI = \frac{TP + TN}{TP + FP + FN + TN}$$

Fowlkes-Mallows index (FMI)

– среднее геометрическое точности и полноты

$$FMI = \frac{TP}{\sqrt{(TP + FP)(TP + FN)}}$$

```
from sklearn.metrics.cluster import fowlkes_mallows_score
fowlkes_mallows_score([0, 0, 1, 1], [0, 0, 1, 1])
```

есть и много других...

Литература

**К.Д. Маннинг, П. Рагхаван, Х. Шютце «Введение в информационный поиск». —
Вильямс, 2011.**