

Търсене и извличане на информация. Приложение на дълбоко машинно обучение

Стоян Михов



Лекция 7: Клъстеризация във векторно пространство. Вероятно приблизително коректно обучение (РАС-обучение).

План на лекцията

- 1. Формалности за курса (5 мин)**
2. Клъстеризация (10 мин)
3. k-means (15 мин)
4. Варианти и подобрения на K-means (15 мин)
5. РАС-обучение (20 мин)
6. Пример за РАС-обучение (25 мин)

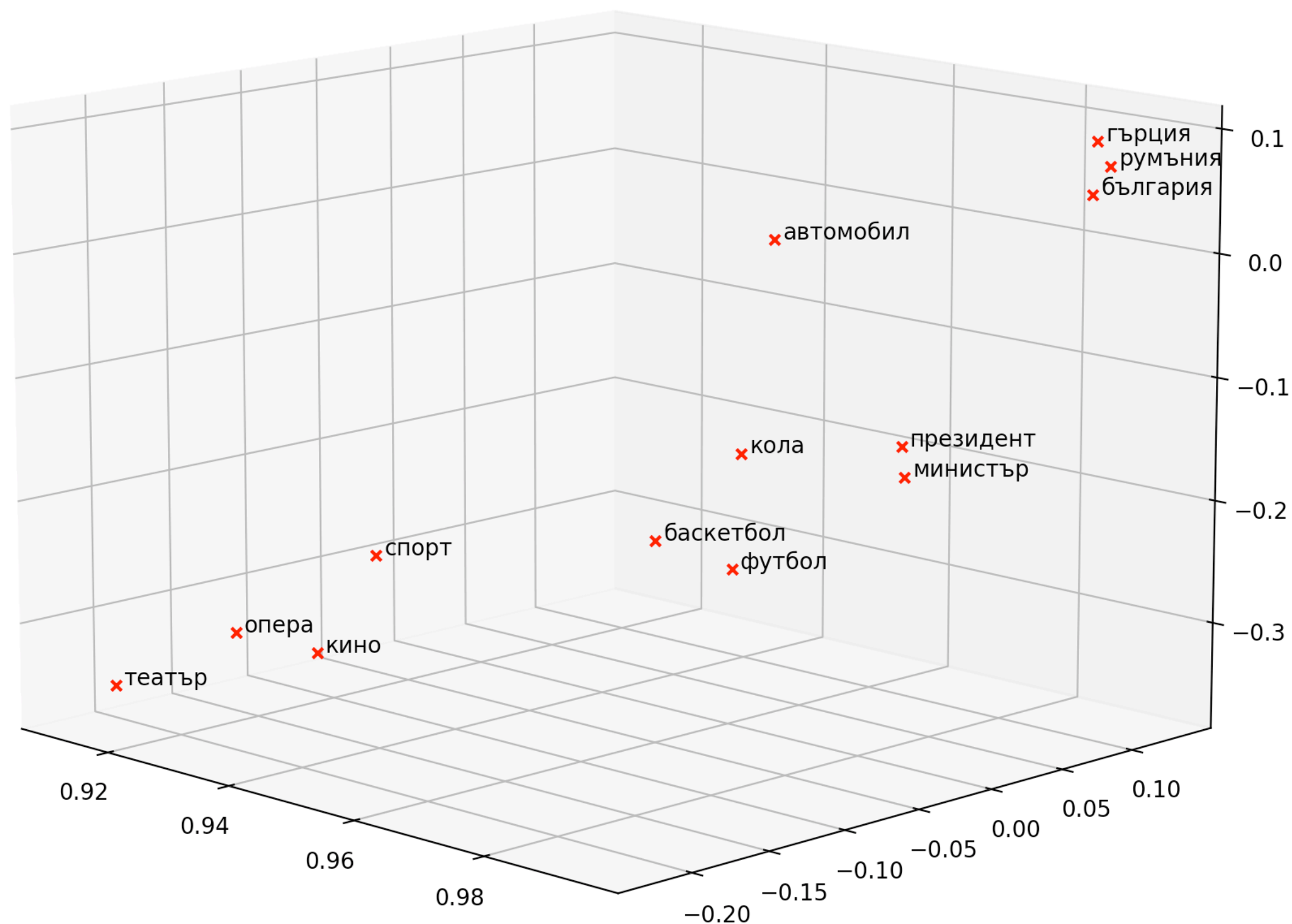
Формалности

- Засега ще провеждаме занятията онлайн всяка сряда от 8:15 до 12:00 часа.
- Засега ще използваме платформата Google meet: meet.google.com/hue-frfx-axb
- Днес ще използваме едновременно слайдове и бяла дъска. Моля следете съответния екран.
- В Moodle на 19.11.2021 г. ще бъде публикувано домашно задание, което следва да бъде предадено до края на деня на 28.11.2021 г.
- Седмата лекция се базира на глава 16 от първия учебник и секция 10.4 от втория учебник.

План на лекцията

1. Формалности за курса (5 мин)
- 2. Клъстеризация (10 мин)**
3. k-means (15 мин)
4. Варианти и подобрения на K-means (15 мин)
5. РАС-обучение (20 мин)
6. Пример за РАС-обучение (25 мин)

Семантично пространствени релации



Клъстеризация

- Задачата на клъстеризацията е да намерим “естествено” групиране на обектите в “клъстери”
- Целта е:
 - В рамките на един клъстер обектите да са близки
 - Обектите от различни клъстери да са далеч един от друг

Приложение на клъстеризацията в търсенето на информация

- При търсене в Интернет често се връщат много хиляди резултати, като потребителя може да разгледа едва няколко от тях.
- Поради многозначността на езика резултатите могат да бъдат от различни области.
- Чрез клъстеризиране на резултатите от заявката на първата страница се връщат по няколко резултата от всеки от клъстерите.

Google

видове корона

AI Images Maps More Settings Tools

Any time Verbatim Clear

botanic.cc › coronavirus-coid-2019-i... Translate this page
Коронавирус | COVID-19 и други видове коронавирус ...
Mar 2, 2020 — Начало / За Здравето / Коронавирус | COVID-2019 и други видове ...
„Корона“ на латински както и на български означава „ореол“ или ...

d2ouvy59p0dg6k.cloudfront.net › o... PDF Translate this page
Определител на растителните видове за оценка на гори с ...
кийския дъб се запазват, а дърветата от други видове в съседство се запазват и ...
костилка, светлочервени, на върха с корона от прави ча- шелистчета.

www.bgfermer.bg › Article Translate this page
Свободно растящата корона еподходяща за овощните ...
Dec 20, 2019 — Формирането на дърветата като свободно растяща корона се препоръчва за почти всички овощни видове, присадени на умерено или ...

www.bgfermer.bg › Article Translate this page
Видове корони за дръвчетата | Овошки | Български Фермер
Feb 27, 2019 — Този тип корона може да се прилага при почти всички овощни видове. Принципното различие при формирането на тази корона в ...

www.booking.com › city › corona-d... Translate this page
Най-добрите хотели в Корона дел Мар, САЩ (на цена от ...
Located in Corona del Mar in the California region, 310 1/2 Iris Ave. 3 Bedroom Condo has a balcony. This property offers access to a patio, free private parking ...

www.booking.com › corona.bg.html Translate this page
Най-добрите 10 за хотела с басейни в Корона, САЩ ...
Корона. Featuring a heated outdoor pool and a fitness centre, Ayres Lodge & Suites Corona West is just 9 minutes' drive from Fender Center. Free WiFi access ...

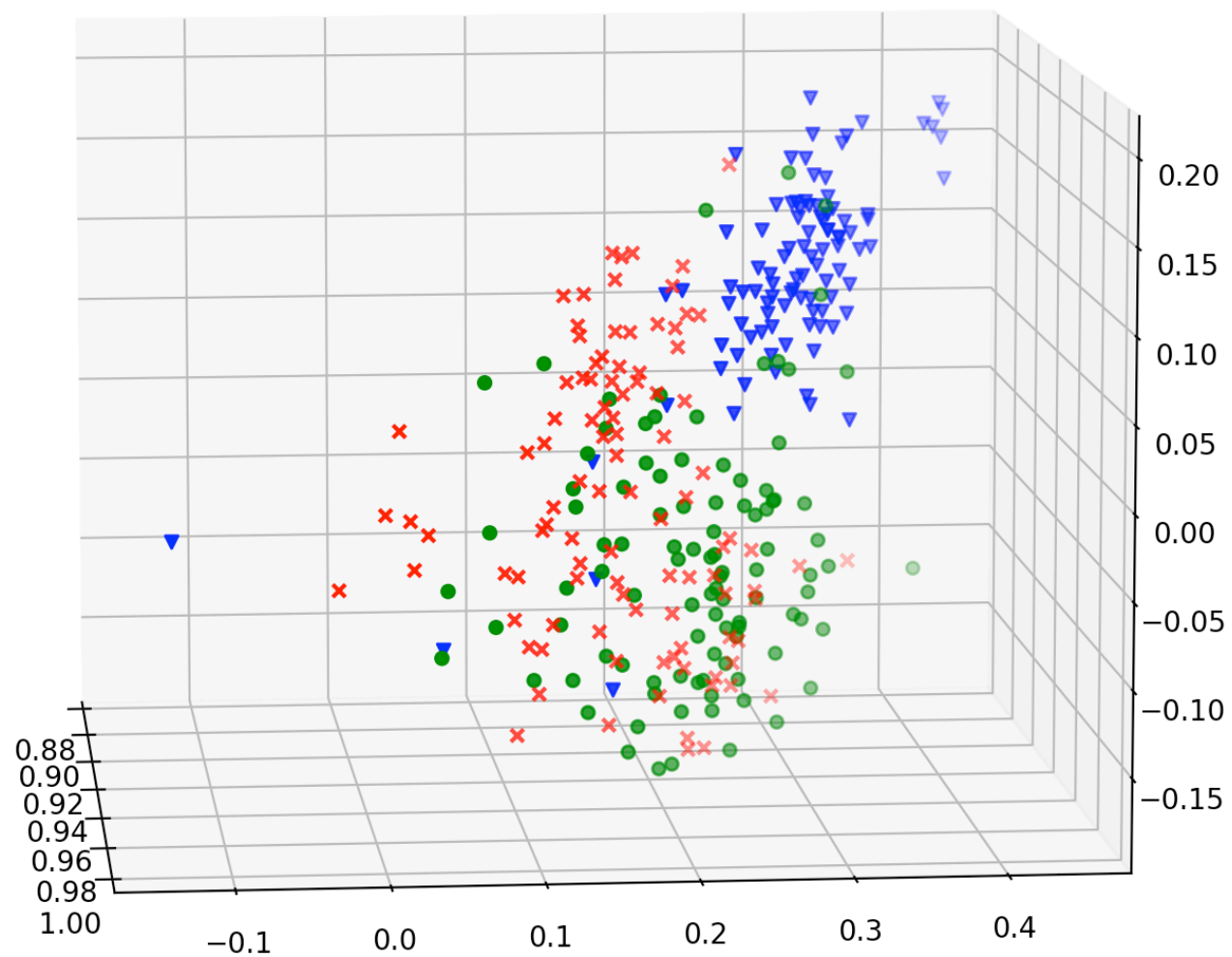
www.medicaldent.bg › Блог › Други Translate this page
Кога идва нуждата от коронка и какви видове има?
... корони превъзхождат метало-керамичните по естетични показатели, като циркониевата корона превъзхожда останалите и по здравина. Изцяло ...

Класификация \Leftrightarrow клъстеризация

- При класификацията имаме предварително зададени класове и база от класифицирани документи \Leftrightarrow при клъстеризацията нямаме нито зададени класове, нито техния брой, нито класифицирани документи.
- При класификацията се стремим да намерим функция (класификатор), която да определя класа на даден документ по подобие на класифицираните документи \Leftrightarrow при клъстеризацията се търси разбиване на базата на имплицитните закономерности в базата от документи.
- Задачата за класификация е пример за обучение “с учител” (supervised learning) \Leftrightarrow задачата за клъстеризация е пример за обучение “без учител” (unsupervised learning)

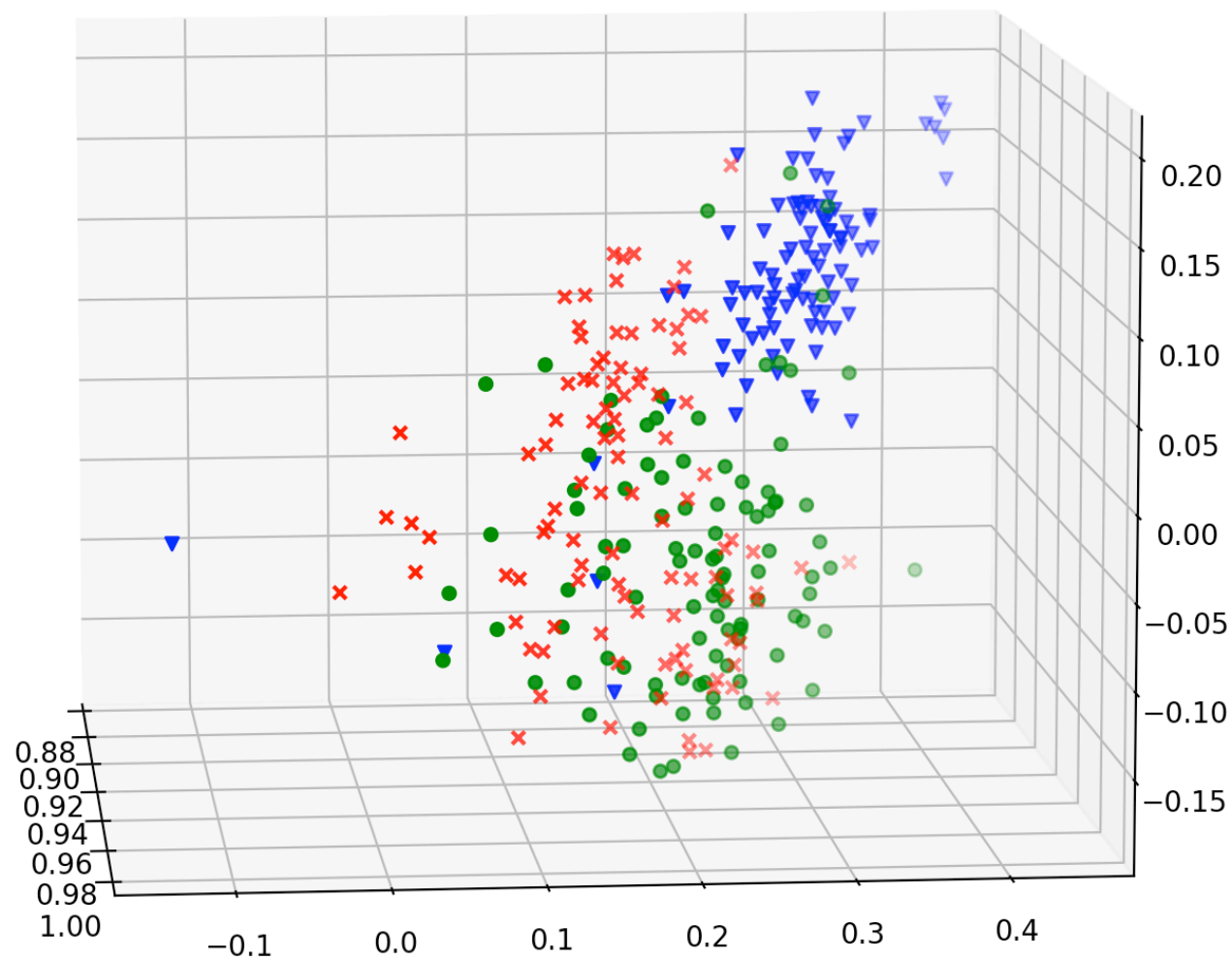
Класификация \Leftrightarrow клъстеризация

Класове от документи

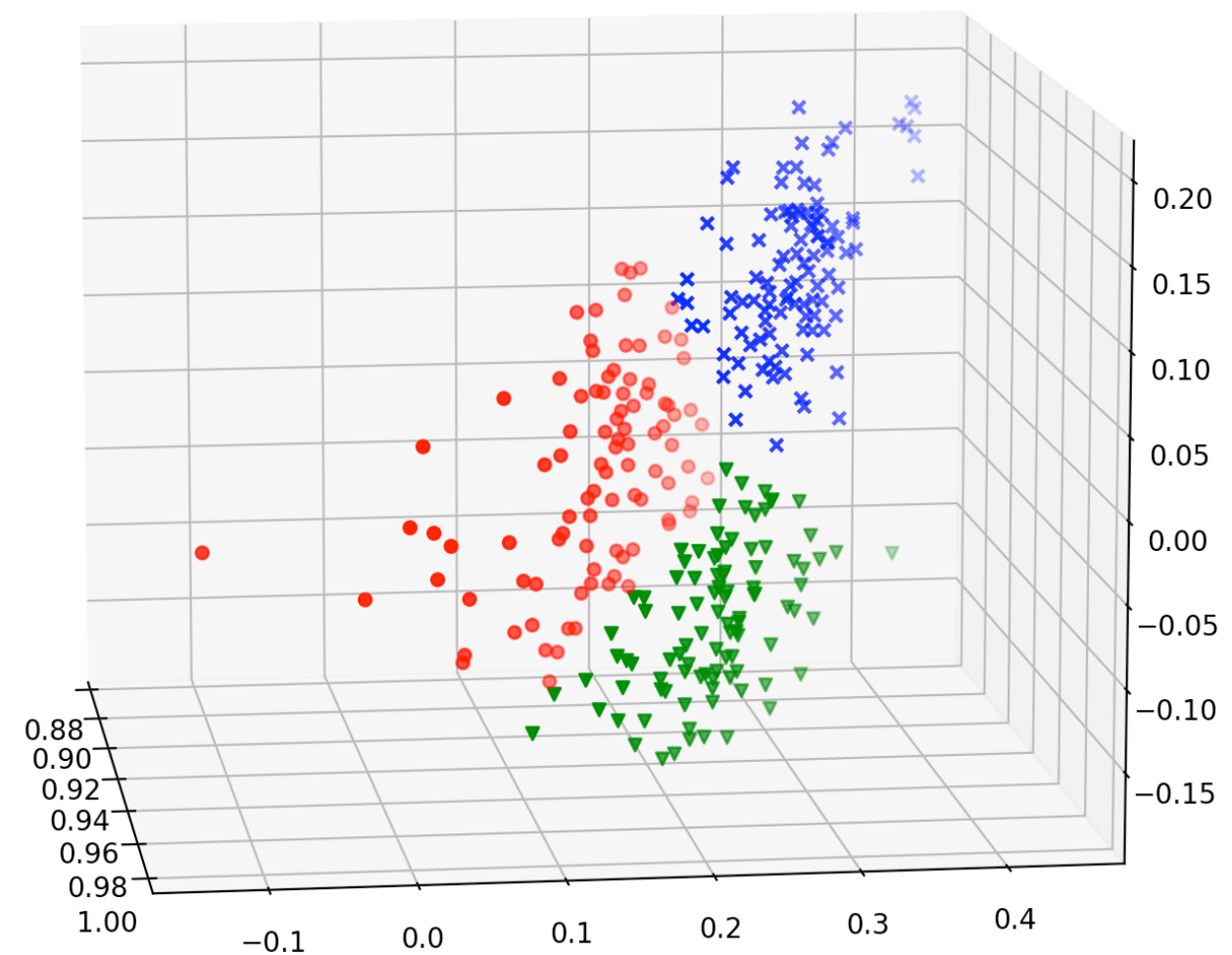


Класификация \Leftrightarrow клъстеризация

Класове от документи



Клъстери от документи



Други приложения на клъстеризацията

- Групиране на резултатите от търсене
- Групиране на поток от документи — новини, мейлове, съобщения, ...
- Ускоряване на търсене по подобие
- Търсене чрез разбиване-събиране (gather-scatter)

План на лекцията

1. Формалности за курса (5 мин)
2. Клъстеризация (10 мин)
- 3. k-means (15 мин)**
4. Варианти и подобрения на K-means (15 мин)
5. РАС-обучение (20 мин)
6. Пример за РАС-обучение (25 мин)

K-means

- Дадено е множество от вектори $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_S\} \subset \mathbb{R}^N$, и число $K \in \mathbb{N}$

- Търсим разбиване на \mathbf{X} в K клъстера $W_1, W_2, \dots, W_K \subset \mathbf{X}$, $W_1 \cup W_2 \cup \dots \cup W_K = \mathbf{X}$, така че:

$$\text{RSS}(W_1, W_2, \dots, W_K) = \sum_{k=1}^K \sum_{\mathbf{x} \in W_k} \|\mathbf{x} - \mu_k\|^2 \text{ е минимално,}$$

където $\mu_k = \frac{1}{|W_k|} \sum_{\mathbf{x} \in W_k} \mathbf{x}$ за $k = 1, 2, \dots, K$.

- μ_k е центроида (центъра на тежестта) на W_k

Алгоритъм K-means

1. Започваме с първоначални центроиди $\mu_1, \mu_2, \dots, \mu_K$.

2. За всеки от центроидите μ_k намираме клъстера W_k
$$W_k = \{\mathbf{x} \in \mathbf{X} \mid \arg \min_i \|\mathbf{x} - \mu_i\|^2 = k\}.$$

3. Намираме новите стойности на центроидите:

$$\mu_k = \frac{1}{|W_k|} \sum_{\mathbf{x} \in W_k} \mathbf{x}$$

4. Докато не се изпълни условие за край повтаряме стъпките 2-4

Алгоритъм K-means

```
K-means ( {x[1],...,x[S] } , K )
1  (μ[1],...,μ[K] ) <- SelectSeeds({x[1],...,x[S]}, K )
2  while stopping criterion has not been met do
3      for k <- 1 to K do
4          ω[k] <- {}
5          for i <- 1 to S do
6              k <- argminj |μ[j] - x[i]|
7              ω[k] <- ω[k] ∪ {x[i]}
8          for k <- 1 to K do
9              μ[k] <- 1/|ω[k]| ∑x∈ω[k] x
10  return{μ[1],...,μ[K]}
```

Коректност

Твърдение: Нека са дадени вектори $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^N$ и $W_1, W_2, \dots, W_K \subset \mathbf{X}$ са дефинирани като $W_k = \{\mathbf{x} \in \mathbf{X} \mid \arg \min_i \|\mathbf{x} - \mu_i\|^2 = k\}$, за $k = 1, 2, \dots, K$. Нека $W'_1, W'_2, \dots, W'_K \subset \mathbf{X}$ е произволно разбиване на \mathbf{X} . Тогава: $\text{RSS}(W_1, W_2, \dots, W_K) \leq \text{RSS}(W'_1, W'_2, \dots, W'_K)$.

Доказателство:

Нека $k(\mathbf{x}) = \arg \min_i \|\mathbf{x} - \mu_i\|^2$ и $l(\mathbf{x}) = l \leftrightarrow \mathbf{x} \in W'_l$. Тогава:

$$\begin{aligned} \text{RSS}(W'_1, W'_2, \dots, W'_K) &= \sum_{k=1}^K \sum_{\mathbf{x} \in W'_k} \|\mathbf{x} - \mu_k\|^2 = \sum_{\mathbf{x} \in \mathbf{X}} \|\mathbf{x} - \mu_{l(\mathbf{x})}\|^2 \geq \sum_{\mathbf{x} \in \mathbf{X}} \|\mathbf{x} - \mu_{k(\mathbf{x})}\|^2 \\ &= \sum_{k=1}^K \sum_{\mathbf{x} \in W_k} \|\mathbf{x} - \mu_k\|^2 = \text{RSS}(W_1, W_2, \dots, W_K) \end{aligned}$$

Твърдение: Нека е дадено множество $W \subset \mathbb{R}^N$. Тогава:

$$\arg \min_{\mathbf{y} \in \mathbb{R}^N} \sum_{\mathbf{x} \in W} \|\mathbf{x} - \mathbf{y}\|^2 = \frac{1}{|W|} \sum_{\mathbf{x} \in W} \mathbf{x}$$

Доказателство:

$$\frac{\partial}{\partial \mathbf{y}} \sum_{\mathbf{x} \in W} \|\mathbf{x} - \mathbf{y}\|^2 = 2|W|\mathbf{y} - 2 \sum_{\mathbf{x} \in W} \mathbf{x} = 0$$

Следствие: На всяка стъпка от алгоритъма **RSS** намалява.

План на лекцията

1. Формалности за курса (5 мин)
2. Клъстеризация (10 мин)
3. k-means (15 мин)
- 4. Варианти и подобрения на K-means (15 мин)**
5. РАС-обучение (20 мин)
6. Пример за РАС-обучение (25 мин)

Условия за край

- Функцията RSS не е изпъкнала. Намирането на глобален екстремум е NP пълен проблем.
- Евристично решение:
 - Когато RSS спре да се подобрява
 - Когато подобрението на RSS е под определен праг
 - След извършване на предварително фиксиран брой итерации

Начални центроиди

Оказва се, че резултатът от клъстеризирането с k-means силно зависи от началните центроиди.

Варианти:

1. Избираме първите K вектора от \mathbf{X} и изпълняваме алгоритъма k-means — най-просто, но наивно.
2. Избираме с равномерно случайно разпределение K вектора от \mathbf{X} и изпълняваме алгоритъма k-means.
3. Повтаряме няколко пъти точка 2 и избираме резултата с най-добър RSS.
4. Избираме началните центроиди, така че да ги раздалечим — виж k-means++

K-means++

1. Избираме първия центроид μ_1 с равномерно случайно разпределение от \mathbf{X} .
2. Нека сме избрали центроиди $\mu_1, \mu_2, \dots, \mu_l$. Нека $D(\mathbf{x}) = \min_{i=1}^l \|\mathbf{x} - \mu_i\|$. Дефинираме случайно разпределение върху \mathbf{X} като за всеки вектор $\mathbf{x} \in \mathbf{X}$ дефинираме
$$\Pr_l[\mathbf{x}] = \frac{D(\mathbf{x})^2}{\sum_{\mathbf{x}' \in \mathbf{X}} D(\mathbf{x}')^2}.$$
Избираме центроида μ_{l+1} от \mathbf{X} със случайно разпределение $\Pr_l[\mathbf{x}]$.
3. Повтаряме стъпки 2-3, докато изберем K центроида.
4. С избраните центроиди изпълняваме алгоритъма k-means

Доказва се, че очакваното за RSS при k-means++ е по малко от $O(\log K)$ по глобалния минимум за RSS:

Arthur, D.; Vassilvitskii, S. (2007). "k-means++: the advantages of careful seeding". Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms. Society for Industrial and Applied Mathematics Philadelphia, PA, USA. pp. 1027–1035.

План на лекцията

1. Формалности за курса (5 мин)
2. Клъстеризация (10 мин)
3. k-means (15 мин)
4. Варианти и подобрения на K-means (15 мин)
- 5. РАС-обучение (20 мин)**
6. Пример за РАС-обучение (25 мин)

Какво ще разбираме под “машинно обучение”

- Съществуват различни подходи (алгоритми), за машинно обучение: вероятностни, градиентни, алгебрични, комбинаторни, ...
- Машинното обучение може да се разглежда като “с учител” и “без учител”
- За формалното изследване на различните алгоритми е целесъобразно въвеждането на формална рамка, която да ни позволява да сравняваме и оценяваме свойствата на различните алгоритми за машинно обучение.
- Ще разгледаме един от разпространените подходи за формализиране на понятието машинно обучение — рамката “Вероятно приблизително коректно” обучение (Probably Approximately Correct PAC-обучение)

Формализиране на задачата за машинно обучение

- Нека с \mathcal{X} означим множеството от всички възможни **наблюдения** (примери) и ще наричаме входно пространство или домейн.
- Нека с \mathcal{Y} означим множеството от възможни **класове**. Ще разглеждаме само крайни множества от класове.
- Функция $c : \mathcal{X} \rightarrow \mathcal{Y}$ ще наричаме **класификатор (концепт)**. Когато $\mathcal{Y} = \{0,1\}$ ни е даден бинарен класификатор, който може да разглеждаме като подмножеството \mathcal{X} , за което c , връща стойност 1.
- Под **клас от класификатори** ще разбираме конкретно множество от класификатори, които ще искаме да научим и ще означаваме с \mathcal{C} . Например, ако \mathcal{X} е двумерното пространство и разглеждаме бинарни класификатори, то класът от класификатори би могъл да бъде множеството от всички триъгълници в равнината.
- Предполагаме, че ни е дадено фиксирано, но неизвестно **вероятностно разпределение** \mathcal{D} върху \mathcal{X} , и че всички наблюдения от \mathcal{X} , които правим са независими и еднакво разпределени с разпределение \mathcal{D} .

Задачата на машинното обучение

- Обучаемият получава извадка $S = (x_1, x_2, \dots, x_m)$ от независими и идентично разпределени с разпределение \mathcal{D} наблюдения, заедно със съответни етикети $(c(x_1), c(x_2), \dots, c(x_m))$ относно конкретен класификатор $c \in \mathcal{C}$, който следва да се научи.
- Задачата на обучаемият е въз основа на наблюденията S класификатор да избере h_S , с минимална грешка при обобщение спрямо целевия класификатор $c \in \mathcal{C}$.
- **Грешка при обобщение** $R(h)$ или истинска грешка между класификатор $h \in \mathcal{H}$ и целевия класификатор $c \in \mathcal{C}$ дефинираме: $R(h) = \Pr_{x \sim \mathcal{D}}[\{x \mid h(x) \neq c(x)\}]$.

Дефиниция на PAC-обучение

Казваме, че класът от класификатори \mathcal{C} е **PAC-обучаем** ако:

- съществува алгоритъм \mathcal{A} , връщащ на извадка от наблюдения S класификатор $h_S = \mathcal{A}(S)$, и
- полиномиална функция $poly : \mathbb{R}^2 \rightarrow \mathbb{R}$,

така че за всяко $\varepsilon > 0$ и всяко $\delta > 0$ ако $m > poly(1/\varepsilon, 1/\delta)$ и S е извадка с поне m наблюдения, то:

$$\Pr_{S \sim \mathcal{D}^m}[\{S \mid R(h_S) \leq \varepsilon\}] \geq 1 - \delta.$$

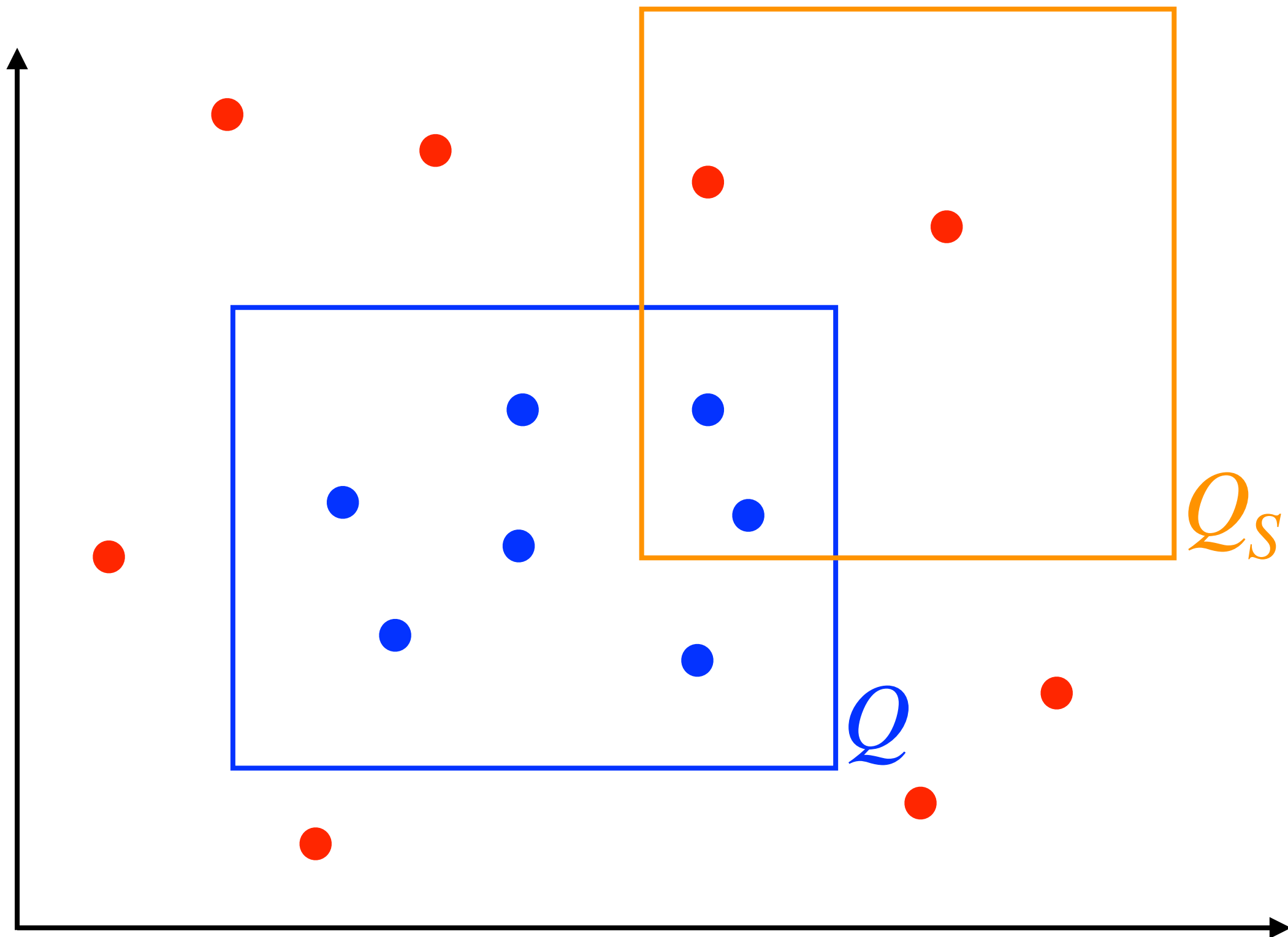
Т.е. Клас от класификатори \mathcal{C} е PAC-обучаем, ако хипотезата, върната от алгоритъма след извадка от наблюдения, чийто брой е полиномиален спрямо $1/\varepsilon$ и $1/\delta$, е приблизително правилна (грешка най-много ε) с голяма вероятност (поне $1 - \delta$), което оправдава терминологията “Вероятно приблизително коректно”.

План на лекцията

1. Формалности за курса (5 мин)
2. Клъстеризация (10 мин)
3. k-means (15 мин)
4. Варианти и подобрения на K-means (15 мин)
5. РАС-обучение (20 мин)
- 6. Пример за РАС-обучение (25 мин)**

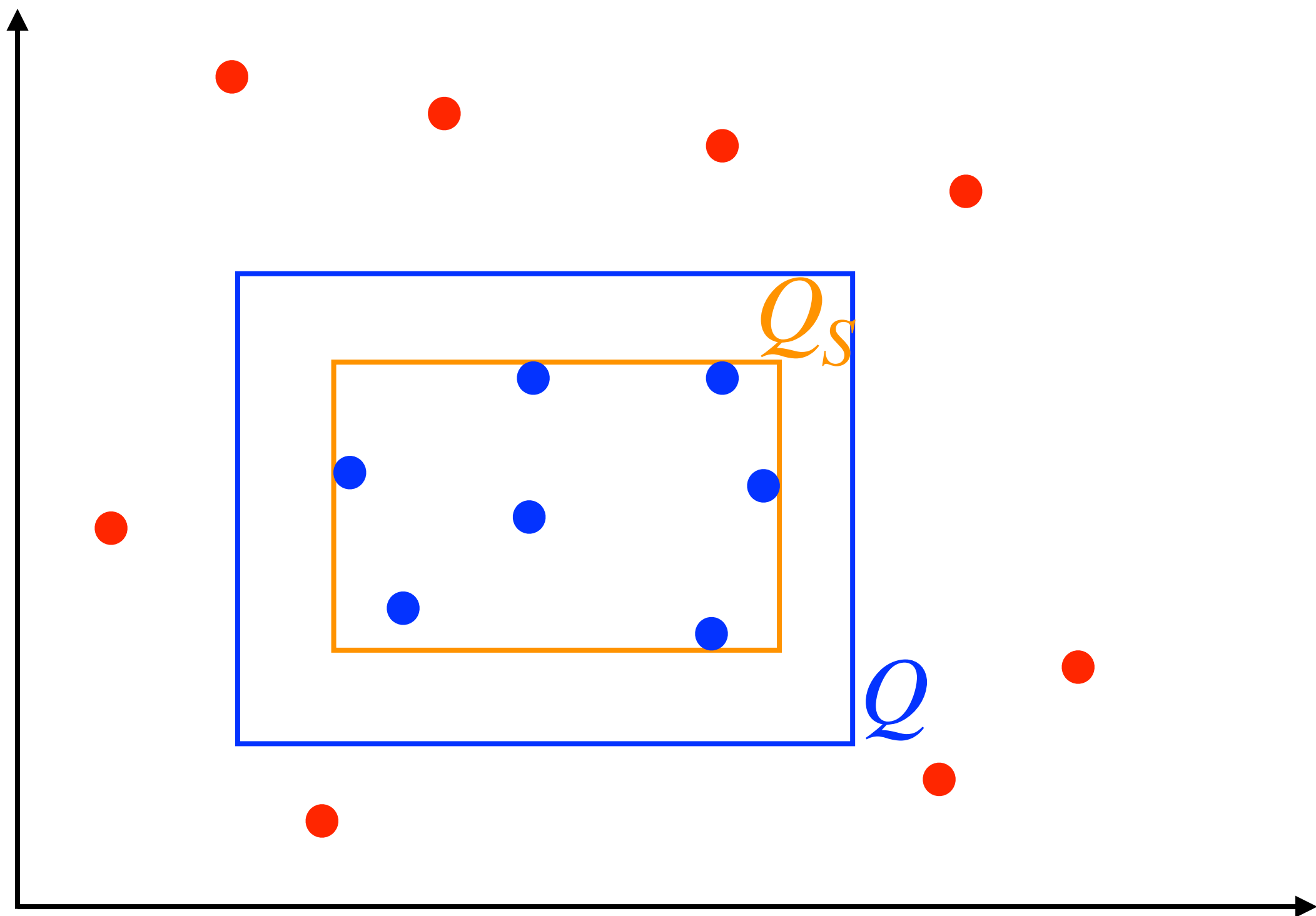
Пример за RAS-обучаем класификатор

- Нека $\mathcal{X} = \mathbb{R}^2$ — пространството от всички точки в равнината.
- Нека \mathcal{C} множеството на всички правоъгълници, чиито страни са успоредни на координатните оси. Т.е. всеки класификатор $c \in \mathcal{C}$ е множеството от влизащи в правоъгълник със страни успоредни на осите.
- Задачата на обучаемия е по извадка от наблюдения с техните етикети да избере правоъгълник, който е максимално близък до целевия.



Решение

- Нека за дадена наблюдения $S = (x_1, x_2, \dots, x_m)$ връщаме най-малкия правоъгълник от \mathcal{C} , който съдържа всички положителни наблюдения. Т.е. ако $P = \{x \in S \mid c(x) = 1\}$ то $Q_S = [\min \text{Proj}_1(P), \max \text{Proj}_1(P)] \times [\min \text{Proj}_2(P), \max \text{Proj}_2(P)]$ т.е. $h_S(x) = \delta_{x \in Q_S}$
- Ясно е, че ако $Q = \{x \in \mathbb{R}^2 \mid c(x) = 1\}$ то $Q \supset Q_S$.
- Следователно: $c(x) \neq h_S(x) \implies x \in Q \setminus Q_S$. Т.е. всички грешки ще са вътре в Q .



- За да докажем, че нашия алгоритъм удовлетворява условието за РАС-обучаемост, нека е дадено $\varepsilon > 0$. Нека за търсения класификатор Q , вероятността дадена точка да влезе в Q спрямо разпрелението \mathcal{D} бележим с $\Pr[Q] = \Pr[\{x \mid x \in Q\}]$. Ако $\varepsilon \geq \Pr[Q]$ то вероятността за грешка

$$R(h_S) = \Pr[\{x \mid h_S(x) \neq c(x)\}] = \Pr[Q \setminus Q_S] \leq \Pr[Q] \leq \varepsilon. \text{ Т.е.}$$

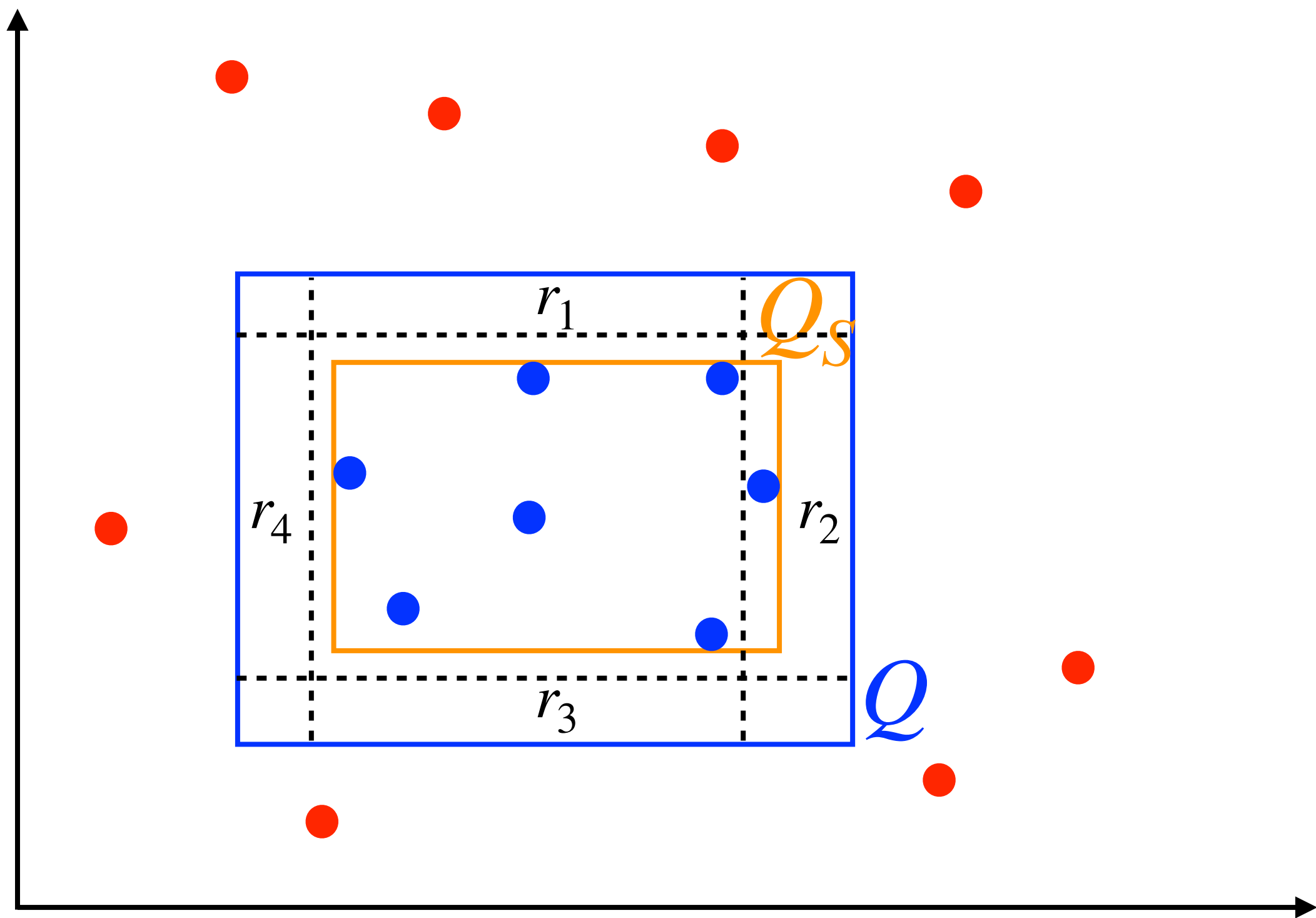
винаги (с вероятност 1) грешката ще е по-малка от ε и твърдението е изпълнено.
- Нека сега $\varepsilon < \Pr[Q]$. Нека $Q = [l, r] \times [b, t]$. Дефинираме правоъгълниците:

$$r_1 = [l, r] \times [b', t], b' = \sup\{s \mid \Pr[[l, r] \times [s, t]] \geq \varepsilon/4\}$$

$$r_2 = [l', r] \times [b, t], l' = \sup\{s \mid \Pr[[s, r] \times [b, t]] \geq \varepsilon/4\}$$

$$r_3 = [l, r] \times [b, t'], t' = \inf\{s \mid \Pr[[l, r] \times [b, s]] \geq \varepsilon/4\}$$

$$r_4 = [l, r'] \times [b, t], r' = \inf\{s \mid \Pr[[l, s] \times [b, t]] \geq \varepsilon/4\}$$



- Ако допуснем, че Q_S има непразно сечение с всеки от правоъгълниците r_1, r_2, r_3, r_4 , то за грешката получаваме:

$$R(h_S) = \Pr[Q \setminus Q_S] \leq \Pr\left[\bigcup_{i=1}^4 r_i\right] \leq \sum_{i=1}^4 \Pr[r_i] = \varepsilon.$$

- Ако допуснем, че Q_S има празно сечение с поне един от правоъгълниците r_1, r_2, r_3, r_4 , то за грешката получаваме:
- $$\Pr_{S \sim \mathcal{D}^m}[\{S \mid R(h_S) > \varepsilon\}] \leq \Pr_{S \sim \mathcal{D}^m}\left[\bigcup_{i=1}^4 \{S \mid Q_S \cap r_i = \emptyset\}\right]$$

$$\begin{aligned} &\leq \sum_{i=1}^4 \Pr_{S \sim \mathcal{D}^m}[\{S \mid Q_S \cap r_i = \emptyset\}] \\ &\leq 4(1 - \varepsilon/4)^m \end{aligned}$$

- Помощно неравенство: $1 - x \leq e^{-x}$. Разглеждаме $f(x) = e^{-x} + x - 1$, $f'(x) = 1 - e^{-x}$, $f''(x) = e^{-x}$.
Следователно функцията е изпъкнала и има единствен минимум при $x = 0$. Но $f(0) = 0$ откъдето следва неравенството.
- $\Pr_{S \sim \mathcal{D}^m}[\{S \mid R(h_S) > \varepsilon\}] \leq 4(1 - \varepsilon/4)^m \leq 4e^{-m\varepsilon/4}$.
- Следователно, за дадено $\delta > 0$, ако изберем $m > \frac{4}{\varepsilon} \log \frac{4}{\delta}$ получаваме: $\delta > 4e^{-m\varepsilon/4} \geq \Pr_{S \sim \mathcal{D}^m}[\{S \mid R(h_S) > \varepsilon\}]$,
откъдето следва $\Pr_{S \sim \mathcal{D}^m}[\{S \mid R(h_S) \leq \varepsilon\}] \geq 1 - \delta$, което трябваше да се покаже. Освен това m е в порядък от $O\left(\frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$.

Заклучение

- Понятието “Вероятно приблизително коректно” РАС-обучение е рамка за математически анализ на алгоритмите за машинното обучение.
- Чрез това и свързаните с него понятия става възможно теоретичното изследване на свойствата и ограниченията на методите за машинно обучение.
- По-задълбочено този подход се разглежда в курса “Теория на машинното обучение и някои нейни приложения в невронните мрежи“
- Няма много практически приложения на този подход, поради което няма да го разглеждаме по-нататък в курса.