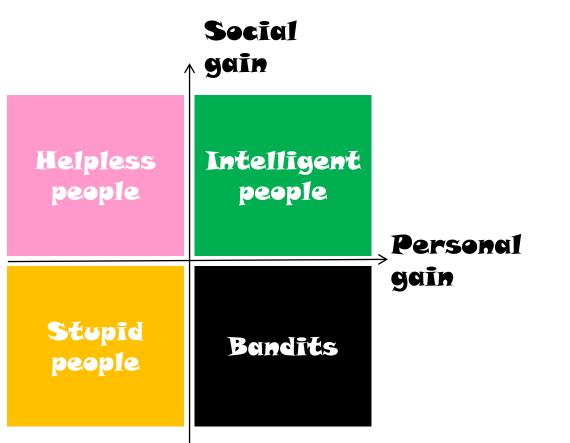
# Задача классификации данных



Классификация – нить <u>Ариадны</u> в лабиринте природы.

Жорж Санд

Cipolla C.M. The basic laws of human stupidity. Bologna: il Muli, 2011 (текст)

© М.Л. Цымблер 05.04.2025

# Содержание

- Основные понятия
- Деревья решений
- Байесовская классификация
- Классификация по ближайшим соседям
- Оценка качества классификации
- Ансамблевая классификация

# (Наивная) Байесовская классификация

- Идея
  - Предсказать вероятность принадлежности объекта классу
- Базис
  - Теорема Байеса
- Точность
  - Сопоставима с деревьями решений
- Инкрементальность
  - Возможность пополнения обучающей выборки
- Baseline (линия отсчета)
  - Минимум для сравнения других методов

#### Теорема Байеса

• P(A|B) — условная вероятность (вероятность наступления события A при условии, что событие B произошло):

$$P(A|B) = \frac{P(AB)}{P(B)}, P(B|A) = \frac{P(AB)}{P(A)}, AB$$
 — совместное событие



•  $AB \equiv BA$ :

$$P(AB) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

• Формула Байеса

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

• Формула полной вероятности для  $P(B) \neq 0$ :

$$P(B) = \sum_{i=1}^{n} P(A_i) \cdot P(B|A_i)$$

 $A_1, ..., A_n$  – несовместные события,  $\sum_{i=1}^n P(A_i) = 1$ 

# Теорема Байеса: пример

- A в бензобаке нет бензина В — авто не заводится
- $\bullet P(B|A) = 1$
- $P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)} = \frac{P(A)}{P(B)}$
- Если P(A) = 0.01, P(B) = 0.02 то для случайно выбранной машины вероятность того, что в бензобаке нет топлива, равна 0.5





# Теорема Байеса: пример

- Играют Трактор и АкБарс на поле АкБарс. Кто победит, если
  - Трактор выиграл 65% игр, АкБарс 35% игр
  - Трактор выиграл 30% игр на чужом поле, АкБарс выиграл 75% игр на своем поле
- В терминах вероятностей
  - $P(\Pi \circ E = T \circ E) = 0.65$ ,
  - $P(\Pi \circ E = A \kappa E \circ E) = 0.35$

  - $P(\Pi \text{оле} = \text{АкБарс} | \Pi \text{обеда} = \text{Трактор}) = 0.30$
  - Найти  $P(\Pi \circ E = A \kappa E \circ E)$





# Теорема Байеса: пример

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

- В терминах вероятностей
  - $P(\Pi \circ E = T \circ E) = 0.65, P(\Pi \circ E = E \circ E) = 0.35$
  - $P(\Pi \circ \pi) = A \kappa B \circ \pi \circ \pi = A \kappa B \circ \pi \circ \pi = 0.75$

  - Найти  $P(\Pi \circ \mathcal{E} = A \kappa \mathcal{E} \circ \mathcal{E})$
- Применение теоремы Байеса

$$P(\Pi \text{обеда} = \text{АкБарс} | \Pi \text{оле} = \text{АкБарс}) = \frac{P(\Pi \text{оле} = \text{АкБарс} | \Pi \text{обеда} = \text{АкБарс}) \cdot P(\Pi \text{обеда} = \text{АкБарс})}{P(\Pi \text{оле} = \text{АкБарс})}$$

$$= \frac{P(\Pi \text{оле} = \text{АкБарс} | \Pi \text{обеда} = \text{АкБарс}) \cdot P(\Pi \text{обеда} = \text{АкБарс})}{P(\Pi \text{оле} = \text{АкБарс, } \Pi \text{обеда} = \text{АкБарс}) + P(\Pi \text{оле} = \text{АкБарс, } \Pi \text{обеда} = \text{Трактор})} =$$

$$=\frac{P(\Pi \text{оле} = \text{АкБарс}|\Pi \text{обеда} = \text{АкБарс}) \cdot P(\Pi \text{обеда} = \text{АкБарс})}{P(\Pi \text{оле} = \text{АкБарс}|\Pi \text{обеда} = \text{АкБарс}) \cdot P(\Pi \text{обеда} = \text{АкБарс}) + P(\Pi \text{оле} = \text{АкБарс}|\Pi \text{обеда} = \text{Трактор}) \cdot P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{АкБарс}|\Pi \text{обеда} = \text{АкБарс}) \cdot P(\Pi \text{обеда} = \text{Трактор})}{P(\Pi \text{оле} = \text{КБарс}|\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{АкБарс}|\Pi \text{обеда} = \text{КБарс}) \cdot P(\Pi \text{обеда} = \text{Трактор})}{P(\Pi \text{оле} = \text{КБарс}|\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{КБарс}|\Pi \text{обеда} = \text{Трактор})}{P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{Трактор}) \cdot P(\Pi \text{обеда} = \text{Трактор})}{P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{Трактор}) \cdot P(\Pi \text{обеда} = \text{Трактор})}{P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{Трактор}) \cdot P(\Pi \text{обеда} = \text{Трактор})}{P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{Трактор}) \cdot P(\Pi \text{обеда} = \text{Трактор})}{P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{Трактор})}{P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{Tрактор})}{P(\Pi \text{обеда} = \text{Трактор})} = \frac{P(\Pi \text{оле} = \text{Tрактор})}{P(\Pi \text{обеда} = \text{Tрактор})} = \frac{P(\Pi \text{оле} = \text{Tрактор})}{P(\Pi \text{оле} = \text{Tрактор})} = \frac{P(\Pi \text{оле} = \text{Tрактор})}{P(\Pi \text{оле} = \text{Tрактор})} = \frac{P(\Pi \text{оле} = \text{Tрактор})}{P(\Pi \text{оле} = \text{Tрактор})} = \frac{P(\Pi \text{one} = \text{Tрактор})}{P(\Pi \text{one} = \text{Tpaken})} = \frac{P(\Pi \text{one} = \text{Tpaken})}{P(\Pi \text{one} = \text{Tpaken})} = \frac{P(\Pi \text{one} = \text{Tpaken})}{P(\Pi \text{one} = \text{Tpaken})} = \frac{P(\Pi \text{one} = \text{Tpaken})}{P(\Pi \text{one} = \text{Tpa$$

$$=\frac{0.75 \cdot 0.35}{0.75 \cdot 0.35 + 0.3 \cdot 0.65} = 0.5738$$

•  $P(\Pi \circ \mathcal{B} = A \kappa \mathcal{B} \circ \mathcal{B}) = \mathbf{57.4\%}$ 



# Применение теоремы Байеса для классификации

$$P(H|X) = \frac{P(X|H) \cdot P(H)}{P(X)}$$

X	Классифицируемый объект	Пациент
$A_1, \ldots, A_n$	Атрибуты объекта	(Возр, Пол, t <sup>0</sup> )
$x_1, \ldots, x_n$	Значения атрибутов объекта	(пожилой, м, высокая)
$\boldsymbol{c}_i$	Классы объектов	Здоров, Болен
Н	Гипотеза « $X \in C_i$ »	Пациент здоров
$P(\boldsymbol{H} \boldsymbol{X})$	Апостериорная вероятность <i>Н</i> при условии <i>X</i>	Вероятность, что пациент здоров, если это пожилой мужчина с высокой $t^0$
$P(\boldsymbol{H})$	Априорная вероятность Н	Вероятность, что пациент здоров
P(X H)	Апостериорная вероятность <i>X</i> при условии <i>H</i>	Вероятность, что здоровый пациент – это пожилой мужчина с высокой t <sup>0</sup>
$P(\boldsymbol{X})$	Априорная вероятность <i>X</i>	Вероятность, что пациент – это пожилой мужчина с высокой t <sup>0</sup>

# Байесовская классификация

• Принцип классификации

$$\bar{X} \in C_k \iff C_k = \arg\max_{C_i} P(C_i|X)$$

• Применим теорему Байеса

$$P(H|X) = \frac{P(X|H) \cdot P(H)}{P(X)}$$

P(X) одинаково для всех классов

• В итоге

$$X \in C_k \iff C_k = \arg\max_{C_i} (P(X|C_i) \cdot P(C_i))$$

#### Байесовская классификация

- $X \in C_k \iff C_k = \arg\max_{C_i} (P(X|C_i) \cdot P(C_i))$
- $P(C_i)$  доля кортежей класса  $C_i$  в обучающей выборке D
- Для вычисления  $P(X|C_i)$  предположим, что атрибуты  $A_1, \dots, A_n$  кортежа  $X=(x_1,\dots,x_n)$  независимы друг от друга. Тогда

$$P(X|C_i) = P(\forall j \ A_j = x_j | C_i) = \prod_{i=1}^{n} P(A_j = x_j | C_i)$$

#### Байесовская классификация

• 
$$P(X|C_i) = \prod_{j=1}^n P(A_j = x_j | C_i) = \prod_{j=1}^n \frac{|\{X \in C_i | X \cdot A_j = x_j\}|}{|C_i|}$$

• В итоге:

$$X \in C_k \iff C_k = \arg\max_{C_i} (P(C_i) \cdot P(X|C_i))$$

$$X \in C_k \iff C_k = \arg\max_{C_i} (\frac{|C_i|}{|D|} \cdot \prod_{j=1}^n \frac{|\{X \in C_i | X. A_j = x_j\}|}{|C_i|})$$

Доля объектов класса  $C_i$  в обучающей выборке

Доля объектов класса  $C_i$ , имеющих в атрибуте  $A_i$  значение  $x_i$ 

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

• *X*=(Возраст=Молодой, Пол=М, Гемоглобин=Норма, t<sup>0</sup>=Высокая) **Класс=?** 

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- X=(Возраст=Молодой, Пол=М, Гемоглобин=Норма,  $t^0$ =Высокая, **Класс=?**)
- $P(C_i)$ 
  - P(3доров) = 9/14 = 0.643
  - P(Болен) = 5/14 = 0.357

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- $X=(Bospacm=Moлoдoй, Пол=M, Гемоглобин=Норма, t^0=Высокая, Класс=?)$
- $P(X|C_i)$ 
  - P(Молодой|3доров) = 2/9 = 0.222

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- $X=(Bospacm=Moлoдoй, Пол=M, Гемоглобин=Hopma, t^0=Bысокая, Класс=?)$
- $P(X|C_i)$ 
  - P(Молодой|Здоров) = 0.222
  - *P*(Молодой|Болен) = 3/5 = 0.6

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- X=(Возраст=Молодой,  $\Pi$ ол=М, Гемоглобин=Норма,  $t^0$ =Высокая, Kласс=?)
- $P(X|C_i)$ 
  - Р(Молодой|Здоров) = 0.222
  - Р(Молодой|Болен) = 0.6
  - P(Мужчина|3доров) = 6/9 = 0.667

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- X=(Возраст=Молодой,  $\Pi$ ол=М, Гемоглобин=Норма,  $t^0$ =Высокая, Kласс=?)
- $P(X|C_i)$ 
  - Р(Молодой|Здоров) = 0.222
  - *P*(Молодой|Болен) =0.6
  - *P*(Мужчина|Здоров) = 0.667
  - P(Мужчина|Болен) = 1/5 = 0.2

Возраст	Пол	Гемо-	$t^{O}$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- X=(Возраст=Молодой, Пол=М,  $\Gamma$ емоглобин=Норма,  $t^0$ =Высокая, Kласс=?)
- $P(X|C_i)$ 
  - Р(Молодой|Здоров) = 0.222
  - *P*(Молодой|Болен) =0.6
  - Р(Мужчина|Здоров) = 0.667
  - Р(Мужчина|Болен) = 0.2
  - P(Норма|3доров) = 4/9 = 0.444

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- X=(Возраст=Молодой, Пол=М,  $\Gamma$ емоглобин=Норма,  $t^0$ =Высокая, Kласс=?)
- $P(X|C_i)$ 
  - P(Молодой|Здоров) = 0.222
  - *P*(Молодой|Болен) =0.6
  - P(Норма|3доров) = 0.444
  - Р(Мужчина|Здоров) = 0.667
  - Р(Мужчина|Болен) = 0.2
  - P(Норма|Болен) = 2/5 = 0.4

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- X=(Возраст=Молодой, Пол=М, Гемоглобин=Норма,  $t^0$ =Высокая, Класс=?)
- $P(X|C_i)$ 
  - P(Молодой|Здоров) = 0.222
  - Р(Молодой|Болен) = 0.6
  - Р(Мужчина|Здоров) = 0.667
  - Р(Мужчина|Болен) = 0.2
  - P(Норма|3доров) = 0.444
  - Р(Норма|Болен) = 0.4
  - P(Высокая|3доров) = 6/9 = 0.667

Возраст	Пол	Гемо-	$t^0$	Класс
		глобин		
Молодой	Ж	Высокий	Высокая	Болен
Молодой	Ж	Высокий	Норма	Болен
Средний	Ж	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Высокая	Здоров
Пожилой	M	Низкий	Высокая	Здоров
Пожилой	M	Низкий	Норма	Болен
Средний	M	Низкий	Норма	Здоров
Молодой	Ж	Норма	Высокая	Болен
Молодой	M	Низкий	Высокая	Здоров
Пожилой	M	Норма	Высокая	Здоров
Молодой	M	Норма	Норма	Здоров
Средний	Ж	Норма	Норма	Здоров
Средний	M	Высокий	Высокая	Здоров
Пожилой	Ж	Норма	Норма	Болен

- X=(Возраст=Молодой, Пол=М, Гемоглобин=Норма,  $t^0$ =Высокая, Класс=?)
- $P(X|C_i)$ 
  - P(Молодой|3доров) = 0.222
  - Р(Молодой|Болен) = 0.6
  - Р(Мужчина|Здоров) = 0.667
  - P(Мужчина|Болен) = 0.2
  - P(Норма|3доров) = 0.444
  - P(Норма|Болен) = 0.4
  - P(Высокая|Здоров) = 0.667
  - P(Высокая|Болен) = 2/5 = 0.4

- *X*=(Возраст=Молодой, Пол=М, Гемоглобин=Норма, t<sup>0</sup>=Высокая, **Класс=?**)
- $P(C_i)$ 
  - P(3доров) = 0.643
  - P(Болен) = 0.357
- $P(X|C_i)$ 
  - P(Молодой|Здоров) = 0.222
  - P(Мужчина|3доров) = 0.667
  - P(Норма|3доров) = 0.444
  - P(Высокая|Здоров) = 0.667
  - P(Молодой|Болен) = 0.6
  - P(Мужчина|Болен) = 0.2
  - P(Норма|Болен) = 0.4
  - *P*(Высокая|Болен) = 0.4

- $X \in C_k \Leftrightarrow$   $C_k = \arg \max_{C_i} (P(X|C_i) \cdot P(C_i))$
- **Класс=Здоров**  $P(X|3доров) \cdot P(3доров) = (0.222 \cdot 0.667 \cdot 0.444 \cdot 0.667) \cdot 0.643$
- **Класс=Болен**  $P(X|\text{Болен}) \cdot P(\text{Болен}) = (0.6 \cdot 0.2 \cdot 0.4 \cdot 0.4) \cdot 0.357$

- *X*=(Возраст=Молодой, Пол=М, Гемоглобин=Норма, t<sup>0</sup>=Высокая, **Класс=?**)
- $P(C_i)$ 
  - P(3доров) = 0.643
  - P(Болен) = 0.357
- $P(X|C_i)$ 
  - Р(Молодой|Здоров) = 0.222
  - P(Мужчина|3доров) = 0.667
  - P(Норма|3доров) = 0.444
  - P(Высокая|Здоров) = 0.667
  - P(Молодой|Болен) = 0.6
  - P(Мужчина|Болен) = 0.2
  - P(Норма|Болен) = 0.4
  - P(Высокая|Болен) = 0.4

- $X \in C_k \Leftrightarrow$   $C_k = \arg \max_{C_i} (P(X|C_i) \cdot P(C_i))$
- Класс=Здоров

```
P(X|3доров) \cdot P(3доров)
= (0.222 \cdot 0.667 \cdot 0.444 \cdot 0.667) \cdot 0.643
= \mathbf{0.028}
```

• Класс=Болен

```
P(X|\text{Болен}) \cdot P(\text{Болен})
= (0.6 \cdot 0.2 \cdot 0.4 \cdot 0.4) \cdot 0.357 = \mathbf{0.007}
```

• Итог

*X*=(Возраст=Молодой, Пол=М, Гемоглобин=Норма, t<sup>0</sup>=Высокая, **Класс=Здоров**)

#### Байесовская классификация: не все хорошо

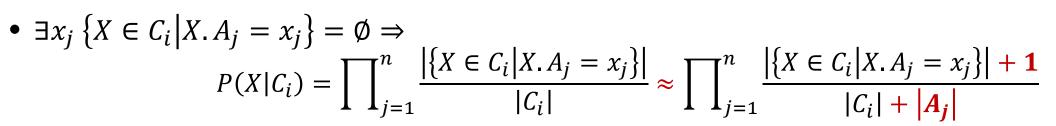
- Нулевая вероятность
  - Как быть, если в обучающей выборке нет объектов с таким же значением атрибута  $A_j$ , как у классифицируемого объекта (тогда  $P(X|C_i)=0$ )?
- Атрибуты с непрерывными значениями
  - Как вычислять апостериорную вероятность  $P(X|C_i)$ , если атрибут  $A_j$  имеет непрерывные (а не дискретные) значения?
- Независимость атрибутов
  - Насколько важно предположение о независимости атрибутов  $A_1$ , ...,  $A_n$ ?

#### Проблема нулевой вероятности

- $X \in C_k \iff C_k = \arg\max_{C_i} P(C_i) \cdot P(X|C_i) = \arg\max_{C_i} P(C_i) \cdot \prod_{j=1}^n \frac{|\{X \in C_i \mid X.A_j = x_j\}|}{|C_i|}$
- В обучающей выборке нет объектов с таким же значением атрибута  $A_j$ , как у классифицируемого объекта
  - $\bullet \ \exists x_i \ \{X \in C_i | X.A_i = x_i\} = \emptyset$
  - $P(X|C_i) = 0 \implies X \notin C_i$ , хотя это может быть не так

#### Нулевая вероятность: поправка Лапласа

- Предположение
  - Обучающая выборка настолько большая, что добавление нескольких объектов не повлияет на вычисление  $P(X|C_i)$





- Пример
  - Пусть в обучающей выборке 1000 здоровых пациентов:
     990 имеют нормальный гемоглобин и 10 высокий, 0 низкий
  - Тогда
    - $P(\text{Низкий}|3\text{доров}) = \frac{0+1}{1000+(1+1+1)} = 0.001$
    - $P(\text{Норма}|3\text{доров}) = \frac{990+1}{1000+(1+1+1)} = 0.988$
    - $P(\text{Высокий}|3\text{доров}) = \frac{10+1}{1000+(1+1+1)} = 0.011$

# Проблема некатегориальных атрибутов

• Если атрибут  $A_j$  — непрерывный, а не дискретный, то нельзя вычислить

$$P(X|C_i) = \prod_{j=1}^n \frac{|\{X \in C_i | X.A_j = x_j\}|}{|C_i|}$$

• Полагаем, что атрибут имеет нормальное (Гауссово) распределение

$$g(x,\mu,\sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

• 
$$P(A_j = x_j | C_i) = g(x_j, \mu_{C_i}, \sigma_{C_i})$$

# Гауссово (нормальное) распределение

• 
$$P(A_j = x_j | C_i) = g(x_j, \mu_{C_i}, \sigma_{C_i}) = \frac{1}{\sqrt{2\pi}\sigma_{C_i}} e^{-\frac{(x_j - \mu_{C_i})^2}{2\sigma_{C_i}^2}}$$
• Среднее значение по атрибуту  $A_i$  среди объекто

Среднее значение по атрибуту  $A_i$  среди объектов класса  $C_i$ 

$$\mu_{C_i} = \frac{1}{|C_i|} \sum_{k=1}^{|C_i|} A_{j_k}$$

• Стандартное отклонение по атрибуту  $A_i$  среди объектов класса  $C_i$ 

$$\sigma_{C_i} = \sqrt{\frac{1}{|C_i|} \sum_{k=1}^{|C_i|} A_{j_k}^2 - \mu_{C_i}^2}$$

- Пример
  - Пусть Bospacm непрерывный атрибут, здоровые пациенты имеют возраст 36..60 лет. Нужно найти P(Воspact = 45|3доров).
  - Тогда в формулу выше подставляем  $x_i = 45, \mu = 48, \sigma = 12.$

# Почему наивная Байесовская классификация?

- $P(X|C_i) = \prod_{j=1}^n P(A_j = x_j | C_i) \Leftrightarrow$  атрибуты  $A_1, \dots, A_n$  независимы друг от друга
- Но на практике это часто не так: атрибуты влияют друг на друга
  - Пример. У больных пациентов высокая температура часто наблюдается совместно с кашлем и низким гемоглобином: атрибуты  $t^0$ , Кашель, Гемоглобин не являются независимыми

#### Байесовская классификация: за и против

- Преимущества
  - Простая реализация
  - Точность классификации, сравнимая с деревьями решений и нейронными сетями в некоторых предметных областях
  - Используется как baseline (теоретическое обоснование для других классификаторов, которые явно не используют теорему Байеса)
- Недостатки
  - Потеря точности в случае зависимых друг от друга атрибутов

# Литература

- Han J., Kamber M., Pei J. Data Mining: Concepts and Techniques. 3rd Edition. Morgan Kaufmann, 2011. 740 p. ISBN 978-0123814791
  - 8.3. Bayes Classification Methods, pp. 350-355
- Tan P.-N., Steinbach M., Kumar V. Introduction to Data Mining. 1st Edition. Pearson, 2019. 732 p. ISBN: 978-1-292-02615-2
  - 5.3 Bayesian Classifiers, pp. 227-240