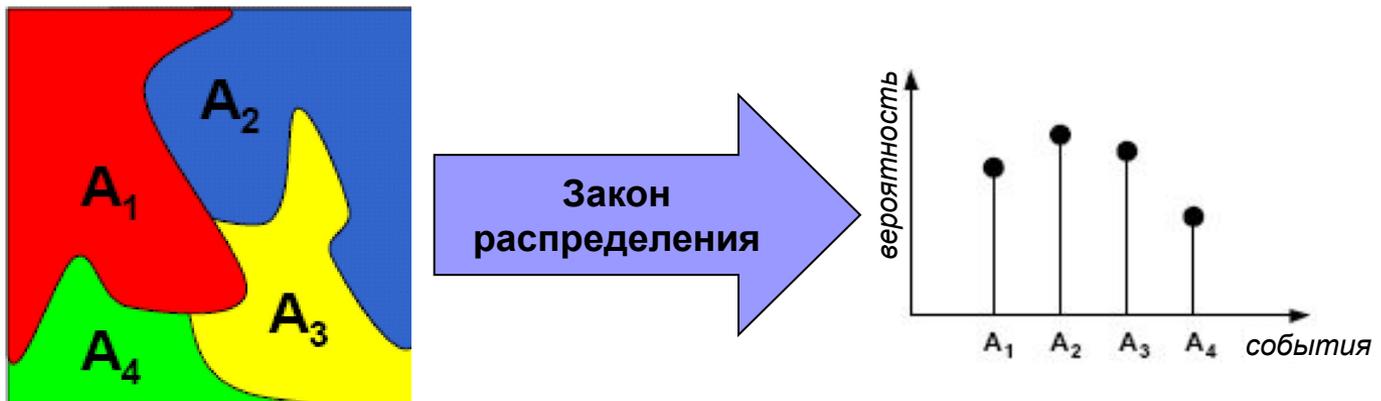




# Байесовская классификация

# Вероятность: основные понятия

- Определения (неформальные)
  - **Вероятность** – число, сопоставляемое событию и показывающее «насколько часто» будет происходить это событие при проведении случайного эксперимента
  - **Закон распределения вероятностей** в случайном эксперименте – правило, сопоставляющее вероятности событиям в эксперименте
  - **Пространство исходов  $S$**  случайного эксперимента – множество всех возможных итогов эксперимента



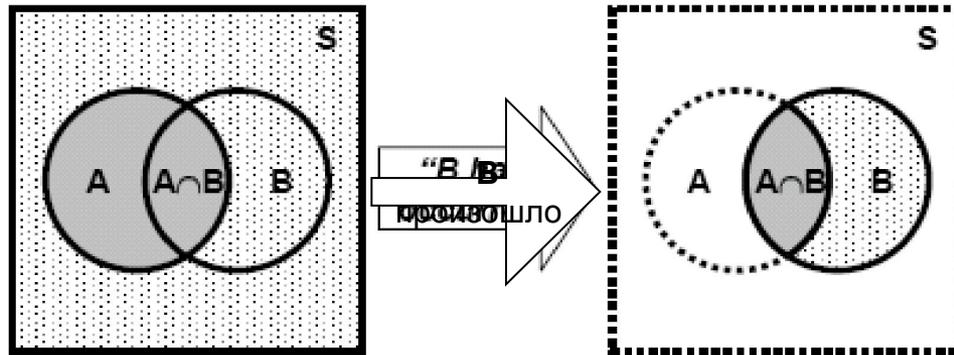
# Вероятность: свойства

- Свойство 1:  $P[A^C] = 1 - P[A]$
- Свойство 2:  $P[A] \leq 1$
- Свойство 3:  $P[\emptyset] = 0$
- Свойство 4: Для  $\{A_1, A_2, \dots, A_N\}$ , если  $\{A_i \cap A_j = \emptyset \ \forall i, j\}$ , то  $P[\bigcup_{k=1}^N A_k] = \sum_{k=1}^N P[A_k]$
- Свойство 5:  $P[A_1 \cup A_2] = P[A_1] + P[A_2] - P[A_1 \cap A_2]$
- Свойство 6:  $P[\bigcup_{k=1}^N A_k] = \sum_{k=1}^N P[A_k] - \sum_{j < k} P[A_j \cap A_k] + \dots + (-1)^{N+1} P[A_1 \cap A_2 \cap \dots \cap A_N]$
- Свойство 7: Если  $A_1 \subset A_2$ , то  $P[A_1] \leq P[A_2]$

# Условная вероятность

- Если  $A$  и  $B$  – события, то вероятность события  $A$  при условии, что  $B$  уже произошло, определяется соотношением

$$P[A | B] = \frac{P[A \cap B]}{P[B]} \quad \text{для} \quad P[B] > 0$$

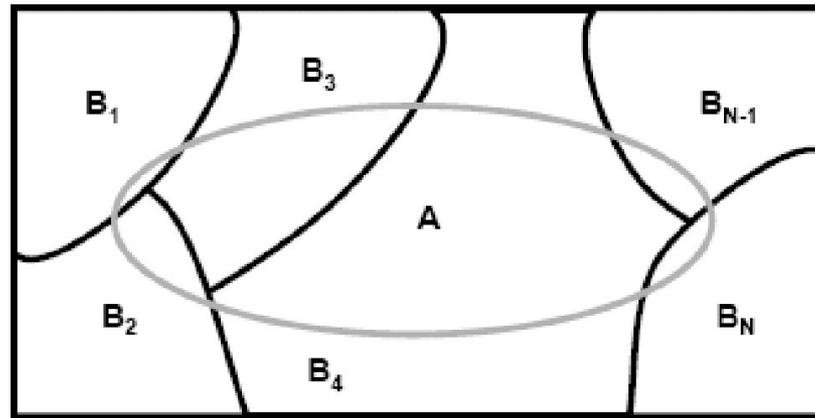


- Интерпретация
  - Новое обстоятельство « $B$  произошло» имеет следующий эффект
    - Исходное вероятностное пространство  $S$  (весь квадрат) сужается до  $B$  (правый круг)
    - Событие  $A$  становится  $A \cap B$
  - Таким образом,  $P[B]$  просто нормирует вероятность событий, происходящих совместно с  $B$

# Формула полной вероятности

- Пусть  $B_1, B_2, \dots, B_N$ , – взаимоисключающие события, которые в объединении дают пространство исходов  $S$  (т.е разбиение  $S$ )
- Событие  $A$  может быть представлено как

$$A = A \cap S = A \cap (B_1 \cup B_2 \cup \dots \cup B_N) = (A \cap B_1) \cup (A \cap B_2) \cup \dots \cup (A \cap B_N)$$



- Поскольку  $B_1, B_2, \dots, B_N$  – взаимоисключающие, то

$$P[A] = P[A \cap B_1] + P[A \cap B_2] + \dots + P[A \cap B_N]$$

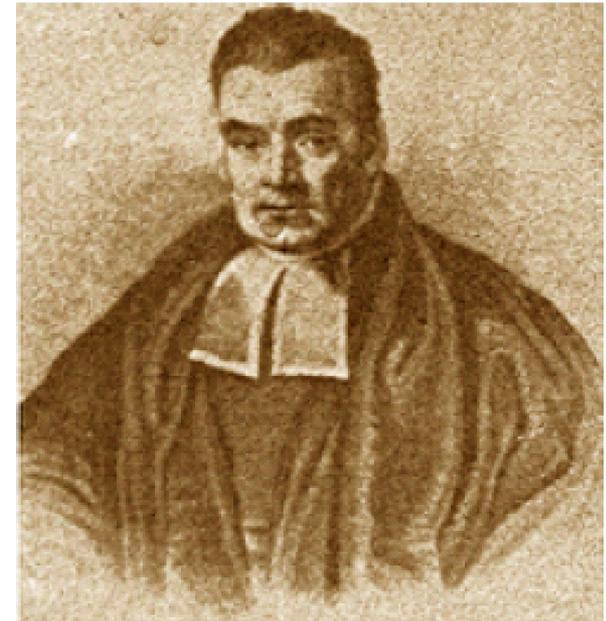
- и, следовательно,

$$P[A] = P[A | B_1]P[B_1] + \dots + P[A | B_N]P[B_N] = \sum_{k=1}^N P[A | B_k]P[B_k]$$

# Формула Байеса

- Пусть  $B_1, B_2, \dots, B_N$  дают разбиение  $S$ . Предположим, что произошло  $A$ . Какова вероятность  $B_j$ ?
  - Из определения условной вероятности и формулы полной вероятности следует

$$P[B_j | A] = \frac{P[A \cap B_j]}{P[A]} = \frac{P[A | B_j] \cdot P[B_j]}{\sum_{k=1}^N P[A | B_k] \cdot P[B_k]}$$



Томас Байес (Thomas Bayes)  
1702-1761

- Это соотношение известно как формула Байеса (правило Байеса, теорема Байеса) и является одним из самых полезных соотношений в теории вероятностей и статистике
  - В распознавании образов формула Байеса – один из фундаментальных результатов

# Формула Байеса и статистическое распознавание образов

- Для решения задачи классификации формула Байеса может быть переписана следующим образом:

$$P[\omega_j | \mathbf{x}] = \frac{P[\mathbf{x} | \omega_j] \cdot P[\omega_j]}{\sum_{k=1}^N P[\mathbf{x} | \omega_k] \cdot P[\omega_k]} = \frac{P[\mathbf{x} | \omega_j] \cdot P[\omega_j]}{P[\mathbf{x}]}$$

- где  $\omega_j$  –  $j$ -й класс,  $\mathbf{x}$  – вектор характеристик образа
- Типичное решающее правило: выбрать класс  $\omega_j$ , для которого  $P[\omega_j | \mathbf{x}]$  – наибольшая
  - Интуитивно: выбираем класс, который наиболее ожидаемо даст  $\mathbf{x}$
- Каждый член в формуле Байеса имеет специальное название
  - $P[\omega_j]$  **априорная вероятность** (класса  $\omega_j$ )
  - $P[\omega_j | \mathbf{x}]$  **апостериорная вероятность** (класса  $\omega_j$  при условии, что в результате наблюдения получен  $\mathbf{x}$ )
  - $P[\mathbf{x} | \omega_j]$  **правдоподобие** (условная вероятность появления наблюдения  $\mathbf{x}$ , при условии, что наблюдается класс  $\omega_j$ )
  - $P[\mathbf{x}]$  нормализующая константа, не влияющая на выбор решения

# Простой пример

- Диагностическая задача: необходимо решить, болен ли пациент, основываясь на *несовершенном* тесте заболевания
  - Некоторые больные могут быть не распознаны (ложно-отрицательный результат теста)
  - Некоторые здоровые могут быть отнесены к больным (ложно-положительный результат теста)
- Терминология
  - Вероятность  $P[\text{ОТР}|\text{ЗДОРОВ}]$  истинно-отрицательного результата теста – **избирательность**
  - Вероятность  $P[\text{ПОЛ}|\text{БОЛЕН}]$  истинно-положительного результата теста – **чувствительность**

# Простой пример

## ■ Задача

- Имеется выборка из 10000 человек, в которой больными являются по 1 человеку из каждых 100
- Используется тест заболевания с избирательностью – 98% и чувствительностью – 90%
- Допустим для Вас тест – положительный.
- **Какова вероятность, что Вы больны?**
  - Решение 1: Заполнить таблицу сопряженности
  - Решение 2: Воспользоваться формулой Байеса

	Тест положителен	Тест отрицателен	Итого
Болен	Истинно-полож. $P[\text{ПОЛ} \text{БОЛЕН}]$	Ложно-отриц. $P[\text{ОТР} \text{БОЛЕН}]$	
Здоров	Ложно-полож. $P[\text{ПОЛ} \text{ЗДОРОВ}]$	Ложно-отриц. $P[\text{ОТР} \text{ЗДОРОВ}]$	
Итого			

# Простой пример

## ■ Задача

- Имеется выборка из **10000** человек, в которой больными являются по **1** человеку из каждых **100**
- Используется тест заболевания с избирательностью – **98%** и чувствительностью – **90%**
- Допустим для Вас тест – положительный.
- **Какова вероятность, что Вы больны?**
  - Решение 1: Заполнить таблицу сопряженности
  - Решение 2: Воспользоваться формулой Байеса

	Тест положителен	Тест отрицателен	Итого
Болен	Истинно-полож. $P[\text{ПОЛ} \text{БОЛЕН}]$ <b><math>100 \times 0.90</math></b>	Ложно-отриц. $P[\text{ОТР} \text{БОЛЕН}]$ <b><math>100 \times (1 - 0.90)</math></b>	<b>100</b>
Здоров	Ложно-полож. $P[\text{ПОЛ} \text{ЗДОРОВ}]$ <b><math>9900 \times (1 - 0.98)</math></b>	Ложно-отриц. $P[\text{ОТР} \text{ЗДОРОВ}]$ <b><math>9900 \times 0.98</math></b>	<b>9900</b>
Итого	<b>288</b>	<b>9712</b>	<b>10000</b>

# Простой пример

## ■ Задача

- Имеется выборка из 10000 человек, в которой больными являются по 1 человеку из каждых 100
- Используется тест заболевания с избирательностью – 98% и чувствительностью – 90%
- Допустим для Вас тест – положительный.
- **Какова вероятность, что Вы больны?**
  - Решение 1: Заполнить таблицу сопряженности
  - **Решение 2: Воспользоваться формулой Байеса**

$$P[\text{БОЛЕН}|\text{ПОЛ}] =$$

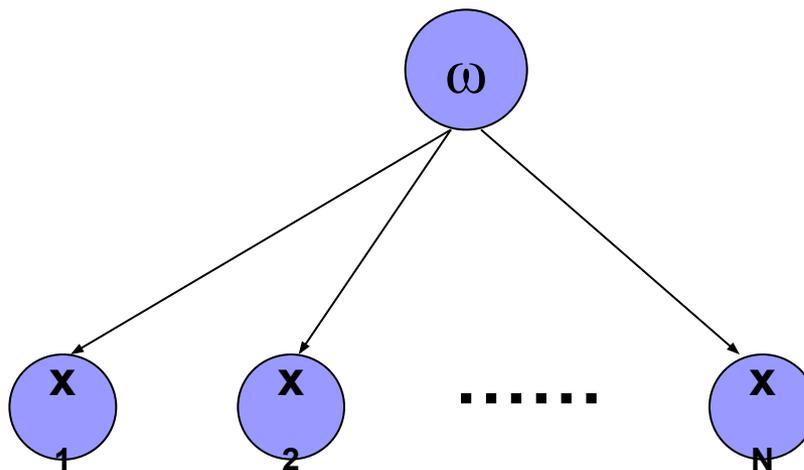
$$\frac{P[\text{ПОЛ}|\text{БОЛЕН}] \cdot P[\text{БОЛЕН}]}{P[\text{ПОЛ}]} =$$

$$\frac{P[\text{ПОЛ}|\text{БОЛЕН}] \cdot P[\text{БОЛЕН}]}{P[\text{ПОЛ}|\text{БОЛЕН}] \cdot P[\text{БОЛЕН}] + P[\text{ПОЛ}|\text{ЗДОРОВ}] \cdot P[\text{ЗДОРОВ}]} =$$

$$\frac{0.90 \cdot 0.01}{0.90 \cdot 0.01 + (1 - 0.98) \cdot 0.99} = 0.3125$$

# Наивный байесовский классификатор

- Наивный байесовский классификатор – простой вероятностный классификатор, основанный
  - на формуле Байеса
  - на предположении **независимости** наблюдаемых признаков (что наивно)
- Графически
  - $\omega$  – класс (метка, номер класса), требуется оценить  $P[\omega | x_1, x_2, \dots, x_N]$
  - $x_i$  – наблюдаемые признаки
  - все признаки независимы



# Наивный байесовский классификатор

- Вероятностная модель наивного байесовского классификатора

$$P[\omega | x_1, \dots, x_N] = \frac{P[\omega] \cdot P[x_1, \dots, x_N | \omega]}{P[x_1, \dots, x_N]}$$

- Строго говоря,

$$\begin{aligned} P[\omega | x_1, \dots, x_N] &= \\ &= P[\omega] \cdot P[x_1 | \omega] \cdot P[x_2 | \omega, x_1] \cdot P[x_3 | \omega, x_1, x_2] \cdot P[x_4, \dots, x_N | \omega, x_1, x_2, x_3] \end{aligned}$$

- Но в виду условной независимости признаков  $P[x_i | \omega, x_j] = P[x_i | \omega]$  и

$$P[\omega | x_1, \dots, x_N] = P[\omega] \cdot P[x_1 | \omega] \cdot P[x_2 | \omega] \cdot \dots \cdot P[x_N | \omega]$$

- Тогда распределение условной вероятности класса  $\omega$  есть

$$P[\omega | x_1, \dots, x_N] = \frac{1}{Z} P[\omega] \prod_{i=1}^N P[x_i | \omega]$$

где  $Z$  – константа, не зависящая от  $\omega$

# Наивный байесовский классификатор

- Наивный байесовский классификатор – комбинация вероятностной модели и решающего правила «максимум апостериорной вероятности»

$$\text{classify}(x_1, \dots, x_N) = \arg \max_j P[\omega = \omega_j] \prod_{i=1}^N P[x_i = f_i | \omega = \omega_j]$$

- Решающее правило
  - **Содержательно:** «Выбрать тот класс, который наиболее ‘вероятен’ для измерения  $\mathbf{x}$ »
  - **Формально:** Вычислить апостериорную вероятность для каждого класса  $P[\omega_j | \mathbf{x}]$  и выбрать класс с ее наибольшим значением
- Для построения наивного байесовского классификатора
  - Оценить  $P[\omega = \omega_j]$ , как долю объектов в обучающей выборке, для которых  $\omega = \omega_j$
  - Оценить  $P[x_i = f_i | \omega = \omega_j]$ , как долю объектов, для которых  $\omega = \omega_j$  и при этом  $x_i = f_i$
  - Для предсказания значения  $\omega$ , найти то  $j$ , которое дает максимальное значение апостериорной вероятности

$$j^* = \arg \max_j P[\omega = \omega_j] \prod_{i=1}^N P[x_i = f_i | \omega = \omega_j]_{14}$$

# Наивный байесовский классификатор

- **Задача.** Отнести текстовые документы к одному из predetermined классов (спорт, политика, экономика,...)
- **Дано.** Имеется выборка из  $N$  текстов, разбитая на  $K$  групп
- **Решение**
  - Вычислить априорные вероятности классов
    - Посчитать количество текстов в каждом классе (группе) –  $N_j$
    - Оценить априорную вероятность  $P[\omega_j] = N_j / N, j=1, \dots, K$
  - Вычислить условные вероятности появления слов в текстах разных классов
    - Пусть есть словарь из  $M$  слов:  $x_1, x_2, \dots, x_M$
    - Вычислить  $c_{ij}$  – сколько раз слово  $x_i$  встретилось в текстах класса  $\omega_j$
    - Вычислить  $n_j$ , сколько слов из словаря встречается в текстах класса  $\omega_j$
    - Условные вероятности слов есть  $P[x_i | \omega_j] = c_{ij} / n_j, j=1, \dots, K$

# Наивный байесовский классификатор

- **Задача.** Отнести текстовые документы к одному из predetermined классов (спорт, политика, экономика,...)
- **Дано.** Имеется выборка из  $N$  текстов, разбитая на  $K$  групп
- **Решение**
  - Классифицировать новый текст  $t$ :
    - Вычислить признаки  $t$  (посчитать сколько раз входит слово  $x_i$  в  $t$ )
    - $P[\omega_j | t] = P[\omega_j | x_1, x_2, \dots, x_M] = P[\omega_j] P[x_1 | \omega_j] P[x_2 | \omega_j] \dots P[x_M | \omega_j]$
    - Отнести  $t$  к классу  $\omega_j$  с максимальной апостериорной вероятностью  $P[\omega_j | t]$

# Наивный байесовский классификатор

- **Задача.** Отнести текстовые документы к одному из predetermined классов (спорт, политика, экономика,...)
- **Дано.** Имеется выборка из  $N$  текстов, разбитая на  $K$  групп
- **Решение**
  - Предобработка
    - Устранение пунктуации
    - Устранение чисел
    - Приведение регистра (все буквы - строчные)
    - Устранение слов короче 4 символов
  - Построение словаря и оценка встречаемости слов производится с помощью хэш-таблицы
  - Как выбрать ключевые слова (признаки)?
    - Взять  $k$  наиболее популярных слов в каждом классе
    - Взять  $k$  наиболее популярных слов в выборке
    - Объединить все эти слова. Вектор признаков – количества этих слов.

# Наивный байесовский классификатор

- **Задача.** Отнести текстовые документы к одному из predetermined классов (спорт, политика, экономика,...)
- **Дано.** Имеется выборка из  $N$  текстов, разбитая на  $K$  групп
- **Решение**
  - Прочие нюансы
    - Нулевые вероятности «убивают» байесовский классификатор
      - Если слово встречается только в одном классе, то условные вероятности такого слова «обнулятся»
      - Для предотвращения этого условные вероятности следует оценивать как  $P[x_i | \omega_j] = \varepsilon / n_j$ , где  $\varepsilon$  – малая настраиваемая константа
    - Все вероятности имеет смысл логарифмировать, чтобы избежать переполнения (особенно при вычислении длинного произведения)
    - Непрерывные признаки имеет смысл квантовать (дискретизировать)

# Критерий отношения правдоподобия

- Пусть объект классифицируется на основании измерения (вектора характеристик)  $\mathbf{x}$
- Разумное решающее правило:
  - «Выбрать тот класс, который наиболее ‘вероятен’ для измерения  $\mathbf{x}$ »
    - Более формально: вычислить апостериорную вероятность для каждого класса  $P(\omega_j | \mathbf{x})$  и выбрать класс с ее наибольшим значением

# Критерий отношения правдоподобия

- В случае 2-классовой задачи классификации:

- Если  $P(\omega_1 | \mathbf{x}) > P(\omega_2 | \mathbf{x})$ , выбрать  $\omega_1$ , иначе выбрать  $\omega_2$

- Или, короче: 
$$P(\omega_1 | \mathbf{x}) \underset{\omega_2}{\overset{\omega_1}{\gtrless}} P(\omega_2 | \mathbf{x})$$

- Используя формулу Байеса, получим

$$\frac{P(\mathbf{x} | \omega_1)P(\omega_1)}{P(\mathbf{x})} \underset{\omega_2}{\overset{\omega_1}{\gtrless}} \frac{P(\mathbf{x} | \omega_2)P(\omega_2)}{P(\mathbf{x})}$$

- $P(\mathbf{x})$  не влияет на решающее правило и может быть исключено

$$\Lambda(\mathbf{x}) = \frac{P(\mathbf{x} | \omega_1)}{P(\mathbf{x} | \omega_2)} \underset{\omega_2}{\overset{\omega_1}{\gtrless}} \frac{P(\omega_2)}{P(\omega_1)}$$

- $\Lambda(\mathbf{x})$  называется **отношением правдоподобия**, а решающее правило известно как **критерий отношения правдоподобия (КОП)**

# Критерий отношения правдоподобия: пример

- Используя критерий отношения правдоподобия, построить решающее правило для следующих условных плотностей классов (полагая априорные вероятности равными):

$$P(x|\omega_1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-4)^2}$$

$$P(x|\omega_2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-10)^2}$$

## Решение

- КОП для заданных плотностей:

$$\Lambda(x) = \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-4)^2}}{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-10)^2}} \begin{matrix} \omega_1 > 1 \\ \omega_2 < 1 \end{matrix}$$

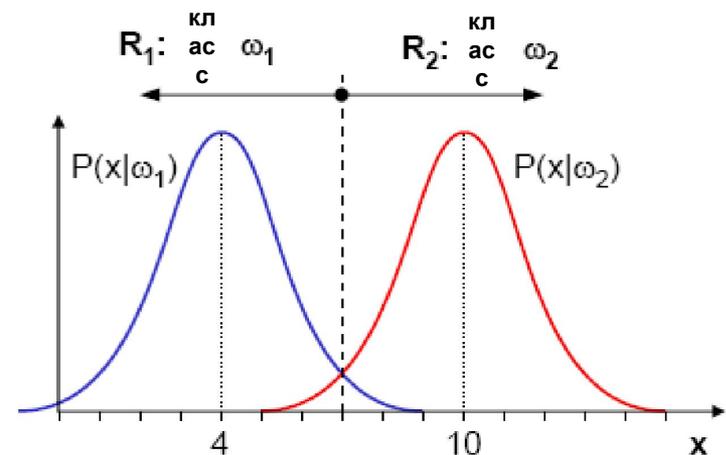
- Упрощаем:  $\Lambda(x) = \frac{e^{-\frac{1}{2}(x-4)^2}}{e^{-\frac{1}{2}(x-10)^2}} \begin{matrix} \omega_1 > 1 \\ \omega_2 < 1 \end{matrix}$

- Меняем знаки и логарифмируем:

$$(x-4)^2 - (x-10)^2 \begin{matrix} \omega_1 < 0 \\ \omega_2 > 0 \end{matrix}$$

- В итоге:

$$x \begin{matrix} \omega_1 < 7 \\ \omega_2 > 7 \end{matrix}$$



## Критерий отношения правдоподобия: пример

- Предыдущий пример понятен с интуитивной точки зрения, т.к. правдоподобия идентичны и отличаются только средними значениями
- Как изменится критерий отношения правдоподобия, если, например, априорные вероятности такие, что  $P(\omega_1) = 2P(\omega_2)$ ?

# Вероятность ошибки

- Мерой качества любого решающего правила может выступать с **вероятность ошибки**  $P[\text{error}]$
- В соответствии с формулой полной вероятности может быть представлена как

$$P[\text{error}] = \sum_{i=1}^C P[\text{error} | \omega_i] P[\omega_i]$$

- Условная вероятность ошибки класса  $P[\text{error} | \omega_i]$ :

$$P[\text{error} | \omega_i] = P[\text{выбр} \neq \omega_j | \omega_i] = \int_{R_j} P(x | \omega_i) dx$$

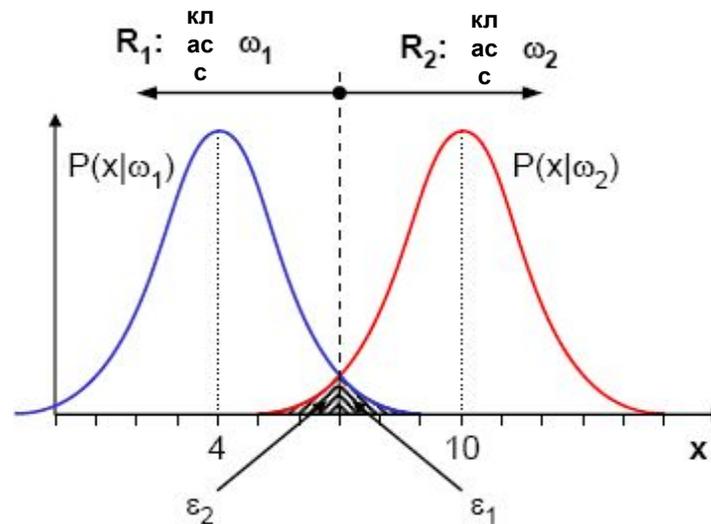
- Для 2-классовой задачи вероятность ошибки примет вид

$$P[\text{error}] = P[\omega_1] \underbrace{\int_{R_2} P(x | \omega_1) dx}_{\varepsilon_1} + P[\omega_2] \underbrace{\int_{R_1} P(x | \omega_2) dx}_{\varepsilon_2}$$

- где  $\varepsilon_i$  – интеграл правдоподобия  $P(x | \omega_i)$  по области  $R_j$ , соответствующей выбору  $\omega_j$

# Вероятность ошибки

- Для решающего правила из рассмотренного примера интегралы  $\varepsilon_1$  и  $\varepsilon_2$  показаны на рисунке
  - Поскольку априорные вероятности полагались равными, то  $P[\text{error}] = (\varepsilon_1 + \varepsilon_2)/2$



# Вероятность ошибки

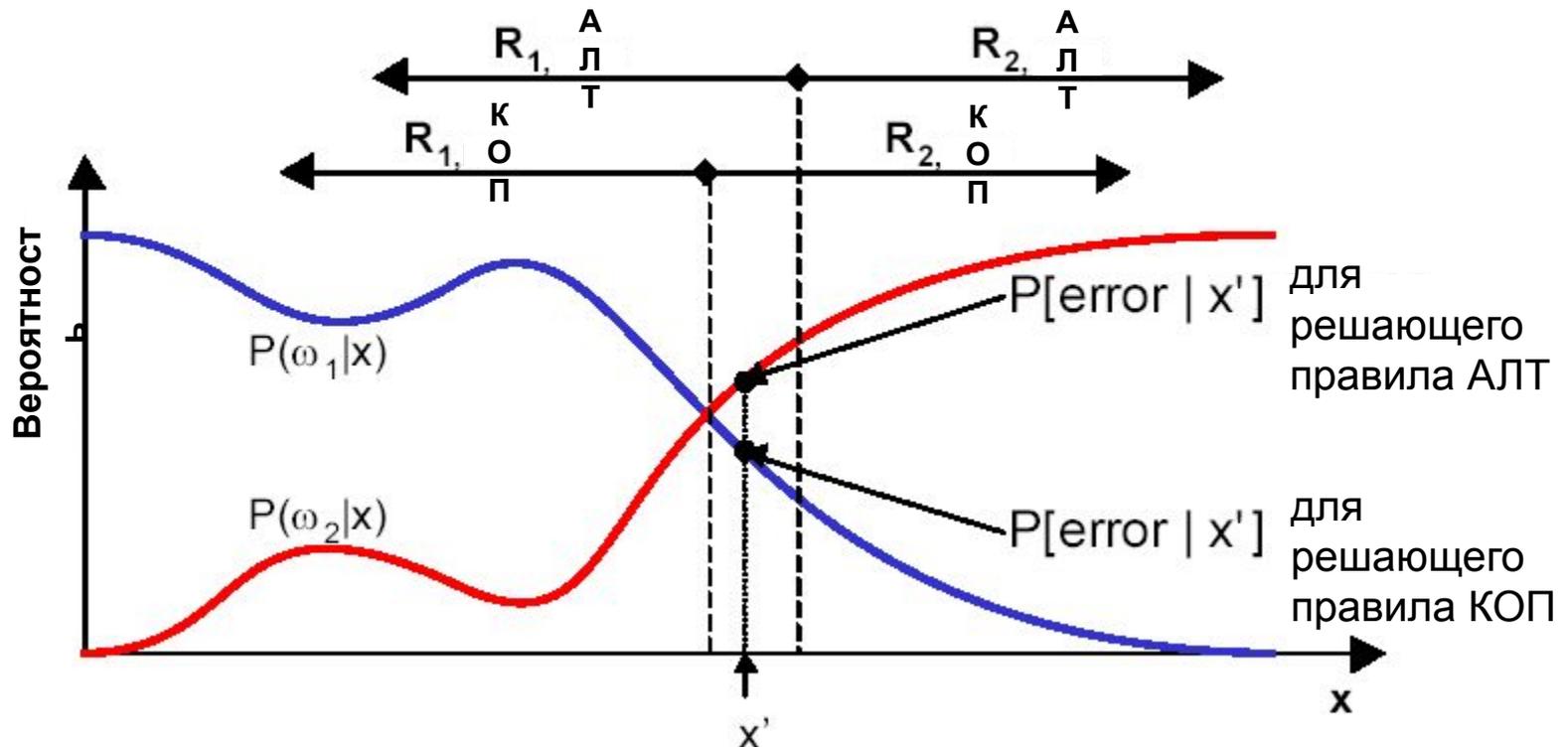
- Выясним, насколько хорош критерий отношения правдоподобия в смысле вероятности ошибки
  - Для этого удобно выразить  $P[\text{error}]$  в терминах апостериорной вероятности  $P[\text{error} | x]$

$$P[\text{error}] = \int_{-\infty}^{+\infty} P[\text{error} | x] P(x) dx$$

- Оптимальное решающее правило – правило, минимизирующее  $P[\text{error} | x]$  для любого  $x$ , чтобы минимизировать весь интеграл
- В каждой точке  $x'$  вероятность  $P[\text{error} | x'] = P(\omega_i | x')$ , если выбран другой класс  $\omega_j$

# Вероятность ошибки

- В каждой точке  $x'$  вероятность  $P[\text{error} | x'] = P(\omega_j | x')$ , если выбран другой класс  $\omega_j$



# Вероятность ошибки

Для любой задачи минимальная вероятность ошибки достигается, если в качестве решающего правила используется критерий отношения правдоподобия.

Эта вероятность ошибки называется **байесовским уровнем ошибки** и является **наилучшим** значением среди всех возможных классификаторов.

# Байесовский риск

- До сих пор полагалось, что цена ошибочного отнесения к классу  $\omega_1$  образца, принадлежащего классу  $\omega_2$ , совпадает с ценой противоположной ошибки. В общем случае это не верно.
  - Ошибочное отнесение больного раком к здоровым имеет несравненно более печальные последствия, чем обратная ошибка
- Это формализуется введением функции стоимости ошибок  $C_{ij}$ 
  - $C_{ij}$  – цена отнесения образца к классу  $\omega_i$  если на самом деле он принадлежит классу  $\omega_j$
- Байесовский риск определяется как мат. ожидание стоимости

$$\mathfrak{R} = E[C] = \sum_{i=1}^2 \sum_{j=1}^2 C_{ij} \cdot P[\text{выбран } \omega_i \text{ при } x \in \omega_j] = \sum_{i=1}^2 \sum_{j=1}^2 C_{ij} \cdot P[x \in R_i | \omega_j] \cdot P[\omega_j]$$

# Байесовский риск

- Какое решающее правило минимизирует байесовский риск?

- Заметим, что

$$P[x \in R_i | \omega_j] = \int_{R_i} P(x | \omega_j) dx$$

- Байесовский риск выражается как

$$\mathfrak{R} = \int_{R_1} [C_{11} \cdot P[\omega_1] \cdot P(x | \omega_1) + C_{12} \cdot P[\omega_2] \cdot P(x | \omega_2)] dx + \\ \int_{R_2} [C_{21} \cdot P[\omega_1] \cdot P(x | \omega_1) + C_{22} \cdot P[\omega_2] \cdot P(x | \omega_2)] dx$$

- Для любого правдоподобия

$$\int_{R_1} P(x | \omega_1) dx + \int_{R_2} P(x | \omega_1) dx = \int_{R_1 \cup R_2} P(x | \omega_1) dx = 1$$

# Байесовский риск

- Подставим последнее уравнение в выражение байесовского риска

$$\mathfrak{R} = \begin{array}{l} \boxed{C_{11}P[\omega_1] \int_{R_1} P(x | \omega_1) dx} + \boxed{C_{12}P[\omega_2] \int_{R_1} P(x | \omega_2) dx} + \\ \boxed{+ C_{21}P[\omega_1] \int_{R_2} P(x | \omega_1) dx} + \boxed{+ C_{22}P[\omega_2] \int_{R_2} P(x | \omega_2) dx} + \\ \boxed{+ C_{21}P[\omega_1] \int_{R_1} P(x | \omega_1) dx} + \boxed{+ C_{22}P[\omega_2] \int_{R_1} P(x | \omega_2) dx} + \\ \boxed{- C_{21}P[\omega_1] \int_{R_1} P(x | \omega_1) dx} - \boxed{C_{22}P[\omega_2] \int_{R_1} P(x | \omega_2) dx} \end{array}$$

- Избавимся от интегралов по  $R_2$

$$\mathfrak{R} = \boxed{C_{21}P[\omega_1]} + \boxed{C_{22}P[\omega_2]} + \boxed{+ (C_{12} - C_{22})P[\omega_2] \int_{R_1} P(x | \omega_2) dx} - \boxed{(C_{21} - C_{11})P[\omega_1] \int_{R_1} P(x | \omega_1) dx}$$

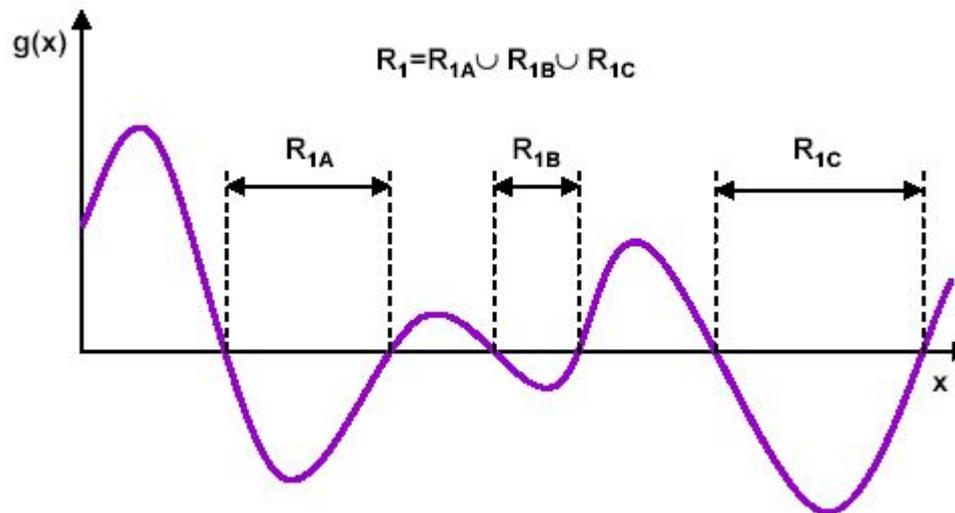
- Т.к. первые два слагаемых – константы, ищем  $R_1$ , минимизирующее остальное:

$$\begin{aligned} R_1 &= \operatorname{argmin} \left\{ \int_{R_1} [(C_{12} - C_{22})P[\omega_2]P(x | \omega_2) - (C_{21} - C_{11})P[\omega_1]P(x | \omega_1)] dx \right\} \\ &= \operatorname{argmin} \left\{ \int_{R_1} g(x) dx \right\} \end{aligned}$$

# Байесовский риск

- На уровне интуиции: к каким областям  $R_1$  приводит минимизация байесовского риска?

- Отыскивается  $R_1$ , минимизирующая интеграл  $\int_{R_1} g(x) dx$
- Иными словами, ищутся области, в которых  $g(x) < 0$



- Таким образом,  $R_1$  выбирается так, чтобы

$$(C_{21} - C_{11})P[\omega_1]P(x | \omega_1) > (C_{12} - C_{22})P[\omega_2]P(x | \omega_2)$$

# Байесовский риск

- На уровне интуиции: к каким областям  $R_1$  приводит минимизация байесовского риска?

- Отыскивается  $R_1$ , минимизирующая интеграл  $\int_{R_1} g(x) dx$

- Иными словами, ищутся области, в которых  $g(x) < 0$

- Таким образом,  $R_1$  выбирается так, чтобы

$$(C_{21} - C_{11})P[\omega_1]P(x | \omega_1) > (C_{12} - C_{22})P[\omega_2]P(x | \omega_2)$$

- Или, после нехитрых преобразований:

$$\frac{P(x | \omega_1)}{P(x | \omega_2)} > \frac{(C_{12} - C_{22}) P[\omega_2]}{(C_{21} - C_{11}) P[\omega_1]}$$

- Итак, минимизация байесовского риска снова приводит к **критерию отношения правдоподобий**

# Байесовский риск: пример

- Классификация на два класса с функциями правдоподобия

$$P(x|\omega_1) = \frac{1}{\sqrt{2\pi}\sqrt{3}} e^{-\frac{1x^2}{2 \cdot 3}}$$

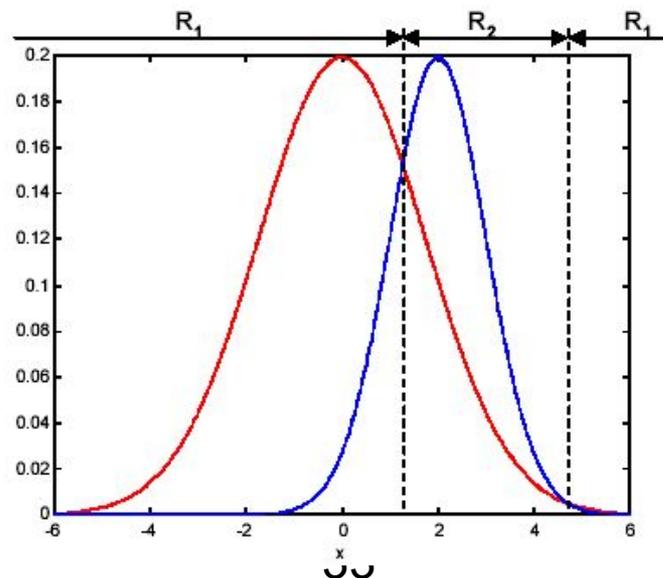
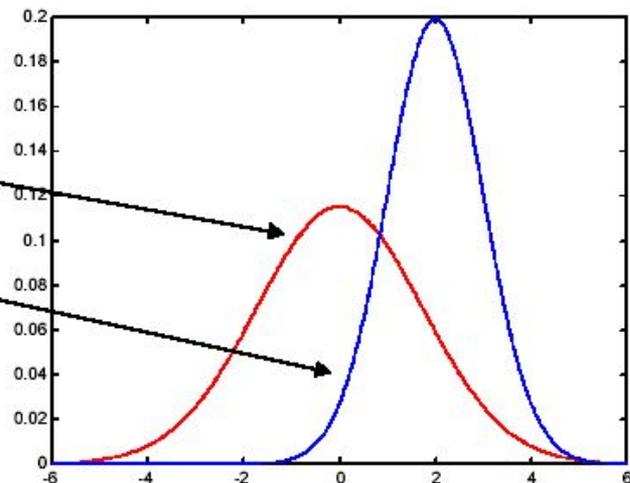
$$P(x|\omega_2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-2)^2}$$

- Пусть  $P(\omega_1) = P(\omega_2) = 0.5$ ,  
 $C_{11} = C_{22} = 0$ ,  $C_{12} = 1$ ,  $C_{21} = 3^{1/2}$
- Построить решающее правило, минимизирующее вероятность ошибки

$$\Lambda(x) = \frac{\frac{1}{\sqrt{2\pi}\sqrt{3}} e^{-\frac{1x^2}{2 \cdot 3}}}{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-2)^2}} > \frac{1}{\sqrt{3}}$$

$$\frac{e^{-\frac{1x^2}{2 \cdot 3}}}{e^{-\frac{1}{2}(x-2)^2}} > 1 \quad -\frac{1}{2 \cdot 3} x^2 + \frac{1}{2}(x-2)^2 > 0$$

$$2x^2 - 12x + 12 > 0 \Rightarrow x = 4.73, 1.27$$



# Вариации критерия отношения правдоподобия

- Решающее правило КОП, минимизирующее байесовский риск обычно называется **критерием Байеса**

$$\Lambda(x) = \frac{P(x | \omega_1) > \frac{(C_{12} - C_{22}) P[\omega_2]}{P(x | \omega_2) < \frac{(C_{21} - C_{11}) P[\omega_1]}{\omega_2}}}{\omega_1} \quad \text{Критерий Байеса}$$

- В частном случае, при минимизации вероятности ошибки, решающее правило называется **максимальным апостериорным критерием (МАК)**

$$C_{ij} = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases} \Rightarrow \Lambda(x) = \frac{P(x | \omega_1) > \frac{P(\omega_2)}{P(x | \omega_2) < \frac{P(\omega_1)}{\omega_2}}}{\omega_1} \Leftrightarrow \frac{P(\omega_1 | x) > 1}{P(\omega_2 | x) < \omega_2} \quad \text{Максимальный апостериорный критерий}$$

- В случае равных априорных вероятностей  $P(\omega_i) = 1/2$  решающее правило называется **критерием максимального правдоподобия**

$$C_{ij} = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases} \Rightarrow \Lambda(x) = \frac{P(x | \omega_1) > 1}{P(x | \omega_2) < \omega_2} \quad \text{Критерий максимального правдоподобия}$$

$$P(\omega_i) = \frac{1}{C} \quad \forall i$$

# Правило минимизации $P[\text{error}]$ для многоклассовых задач

- Правило минимизации вероятности ошибки  $P[\text{error}]$  легко обобщается на случай многих классов
  - Для простоты вводится дополнительная вероятность – вероятность правильного выбора класса:

$$P[\text{error}] = 1 - P[\text{correct}]$$

- Вероятность корректного распознавания

$$P[\text{correct}] = \sum_{i=1}^C P(\omega_i) \int_{R_i} P(x | \omega_i) dx$$

- Задача минимизации  $P[\text{error}]$  эквивалентна максимизации  $P[\text{correct}]$
- $P[\text{correct}]$  в терминах апостериорных вероятностей:

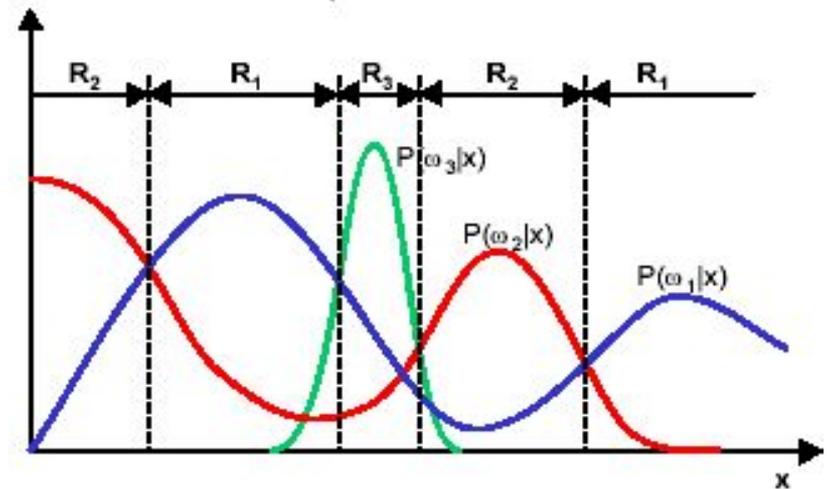
$$P[\text{correct}] = \sum_{i=1}^C P(\omega_i) \int_{R_i} P(x | \omega_i) dx = \sum_{i=1}^C \int_{R_i} P(x | \omega_i) P(\omega_i) dx = \sum_{i=1}^C \underbrace{\int_{R_i} P(\omega_i | x) P(x) dx}_{\mathfrak{A}_i}$$

# Правило минимизации $P[\text{error}]$ для многоклассовых задач

- Правило минимизации вероятности ошибки  $P[\text{error}]$  легко обобщается на случай многих классов
  - $P[\text{correct}]$  в терминах апостериорных вероятностей:

$$P[\text{correct}] = \sum_{i=1}^C P(\omega_i) \int_{R_i} P(x|\omega_i) dx = \sum_{i=1}^C \int_{R_i} P(x|\omega_i) P(\omega_i) dx = \sum_{i=1}^C \underbrace{\int_{R_i} P(\omega_i|x) P(x) dx}_{J_i}$$

- Для максимизации  $P[\text{correct}]$  нужно максимизировать каждый из интегралов  $J_i$ . Интегралы максимизируются выбором классов  $\omega_i$ , дающих максимум  $P[\omega_i|x]$ , т.е. определением области  $R_i$ , в которой  $P[\omega_i|x]$  максимальна.



- Таким образом, правило минимизации  $P[\text{error}]$  – максимальный апостериорный критерий

# Минимизация байесовского риска

## для многоклассовых задач

### ■ Новые обозначения

- $\alpha_i$  – решение о выборе класса  $\omega_i$
- $\alpha(x)$  – общее правило выбора, отображающее вектор  $x$  на классы  $\omega_i$ :  $\alpha(x) \in \{\alpha_1, \alpha_2, \dots, \alpha_C\}$

$$P[\text{correct}] = \sum_{i=1}^C P(\omega_i) \int_{R_i} P(x | \omega_i) dx = \sum_{i=1}^C \int_{R_i} P(x | \omega_i) P(\omega_i) dx = \sum_{i=1}^C \underbrace{\int_{R_i} P(\omega_i | x) P(x) dx}_{S_i}$$

- Условный риск  $R(\alpha_i | x)$  отнесения вектора  $x$  к классу  $\omega_i$  есть

$$\mathfrak{R}(\alpha(x) \rightarrow \alpha_i) = \mathfrak{R}(\alpha_i | x) = \sum_{j=1}^C C_{ij} P(\omega_j | x)$$

- Байесовский риск, связанный с решением по правилу  $\alpha(x)$  есть

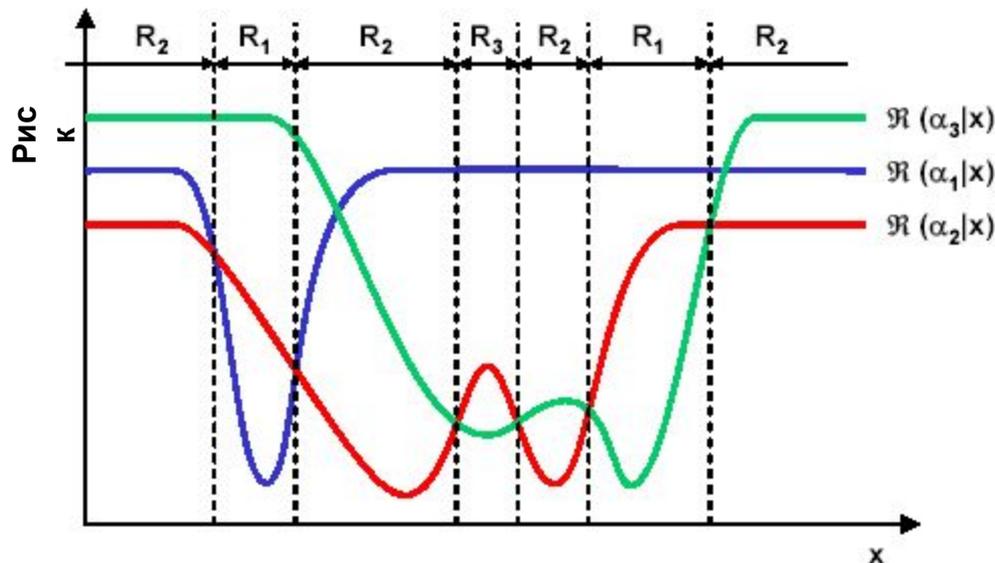
$$\mathfrak{R}(\alpha(x)) = \int \mathfrak{R}(\alpha(x) | x) P(x) dx$$

# Минимизация байесовского риска для многоклассовых задач

- Байесовский риск, связанный с решающим правилом  $\alpha(x)$  есть

$$\mathfrak{R}(\alpha(x)) = \int \mathfrak{R}(\alpha(x) | x) P(x) dx$$

- Для минимизации этого выражения необходимо минимизировать условный риск  $R(\alpha(x)|x)$  в каждой точке пространства  $x$ , что эквивалентно выбору класса  $\omega_i$ , такого, что  $R(\alpha_i|x)$  минимален

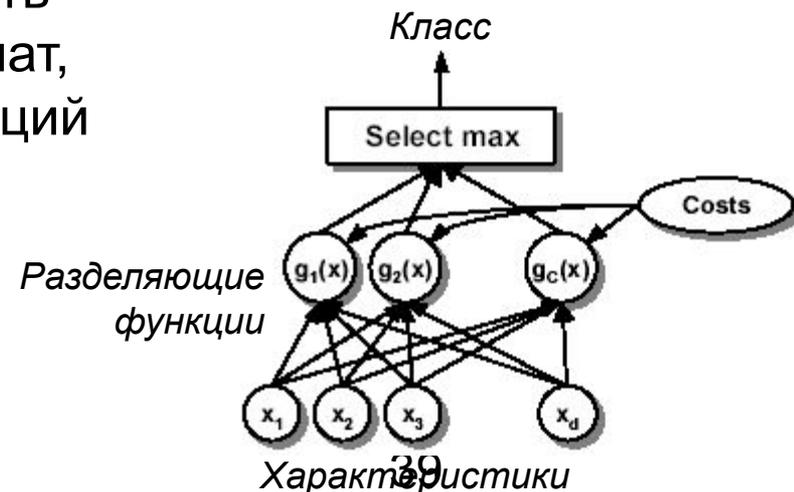


# Разделяющие функции

- Все решающие правила, рассмотренные в этой лекции, имеют одинаковую структуру
  - В каждой точке  $x$  пространства признаков выбрать класс  $\omega_i$  такой, что минимизируется (максимизируется) некоторая мера  $g_i(x)$
- Эта структура может быть формализована с помощью множества разделяющих функций  $g_i(x)$ ,  $i=1, \dots, C$  и следующего решающего правила

“отнести  $x$  к классу  $\omega_i$ , если  $g_i(x) > g_j(x) \quad \forall j \neq i$ ”

- Тогда решающее правило может быть визуализировано как сеть или автомат, вычисляющий  $C$  разделяющих функций и выбирающий класс, соответствующий наибольшему значению



# Разделяющие функции

- Основные критерии как разделяющие функции

Критерий	Разделяющая функция
Байеса	$g_i(x) = -\mathcal{R}(\alpha_i x)$
Максимильный апостериорный	$g_i(x) = P(\omega_i x)$
Максимального правдоподобия	$g_i(x) = P(x \omega_i)$

# Байесовские классификаторы для нормально распределенных классов

- Для случая нормально распределенных классов критерий минимизации ошибки (максимизации апостериорной вероятности, МАВ)

“отнести  $x$  к классу  $\omega_i$ , если  $g_i(x) > g_j(x) \quad \forall j \neq i$ ”

где  $g_i(x) = P(\omega_i | x)$ , может быть существенно упрощен

# Байесовские классификаторы для нормально распределенных классов

- Выражения для гауссовых плотностей в общем виде

- Плотность многомерного нормального распределения

$$f_x(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right]$$

- Преобразованная разделяющая функция критерия МАВ (по формуле Байеса)

$$g_i(x) = P(\omega_i | x) = \frac{P(x | \omega_i)P(\omega_i)}{P(x)} = \frac{1}{(2\pi)^{n/2} |\Sigma_i|^{1/2}} \exp\left[-\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i)\right] P(\omega_i) \frac{1}{P(x)}$$

- Исключая постоянные члены, получим

$$g_i(x) = |\Sigma_i|^{-1/2} \exp\left[-\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i)\right] P(\omega_i)$$

- После логарифмирования

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i) - \frac{1}{2} \log(|\Sigma_i|) + \log(P(\omega_i))$$

- Это выражение называется **квадратичной разделяющей функцией**

# Случай 1: $\Sigma_i = \sigma^2 I$

- Случай возникает, когда характеристики образцов статистически независимы и имеют одинаковую вариацию по всем классам

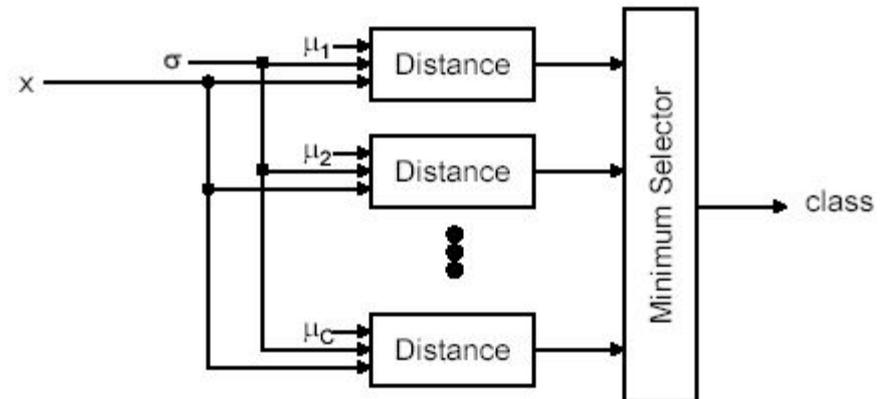
$$g_i(x) = -\frac{1}{2\sigma^2}(-2\mu_i^T x + \mu_i^T \mu_i) + \log(P(\omega_i)) = w_i^T x + w_{i0}$$

$$\text{where } \begin{cases} w_i = \frac{\mu_i}{\sigma^2} \\ w_{i0} = -\frac{1}{2\sigma^2} \mu_i^T \mu_i + \log(P(\omega_i)) \end{cases}$$

- Т.к.  $g_i(x)$  линейны, границы классов – гиперплоскости
- Если априорные вероятности равны, то

$$g_i(x) = -\frac{1}{2\sigma^2} (x - \mu_i)^T (x - \mu_i)$$

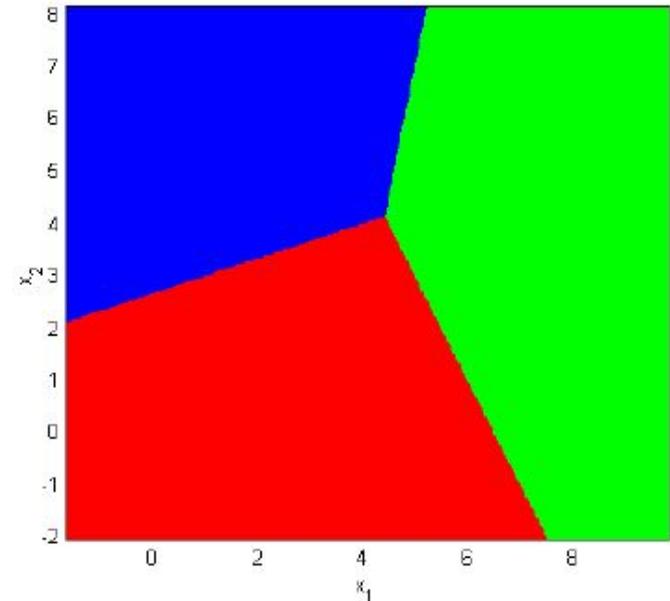
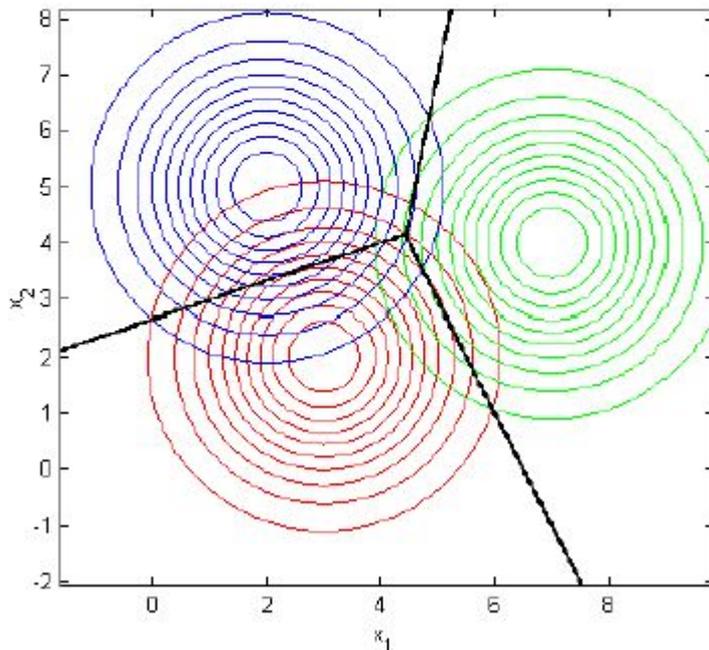
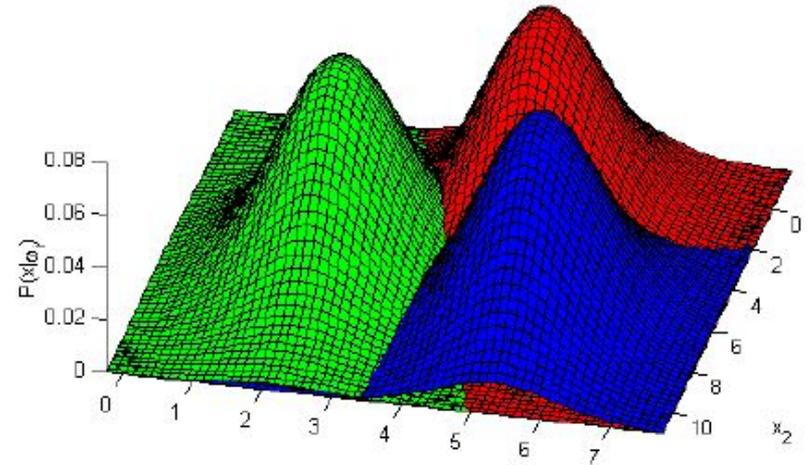
- Это – классификатор по минимальному расстоянию (по ближайшему среднему)
- ГМТ с постоянной вероятностью – гиперсферы
- Для единичной дисперсии, расстояние становится евклидовым



# Случай 1: $\Sigma_i = \sigma^2 I$ , пример

- Двумерная задача, три класса со следующими средними и ковариациями:

$$\begin{aligned} \mu_1 &= \begin{bmatrix} 3 & 2 \end{bmatrix}^T & \mu_2 &= \begin{bmatrix} 7 & 4 \end{bmatrix}^T & \mu_3 &= \begin{bmatrix} 2 & 5 \end{bmatrix}^T \\ \Sigma_1 &= \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} & \Sigma_2 &= \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} & \Sigma_3 &= \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \end{aligned}$$



## Случай 2: $\Sigma_i = \Sigma$ ( $\Sigma$ – диагональная)

- Классы по-прежнему имеют одинаковую матрицу ковариации, но характеристики имеют разные дисперсии

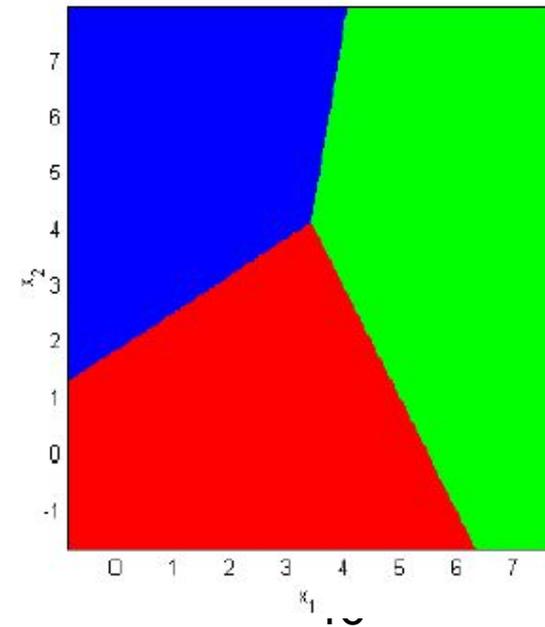
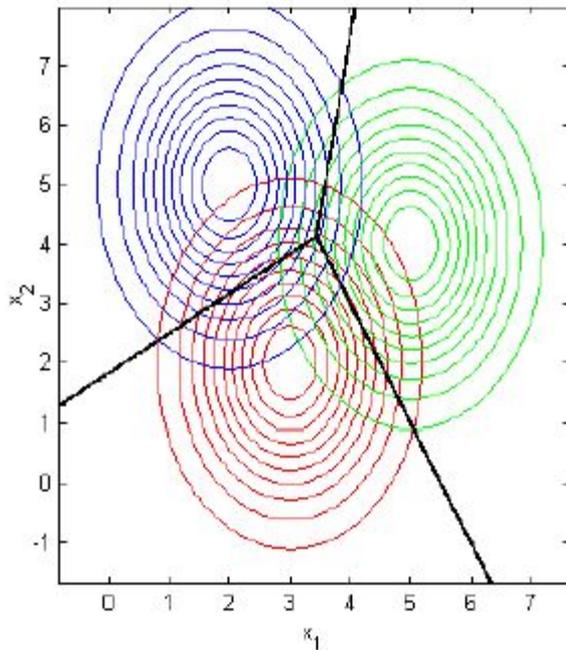
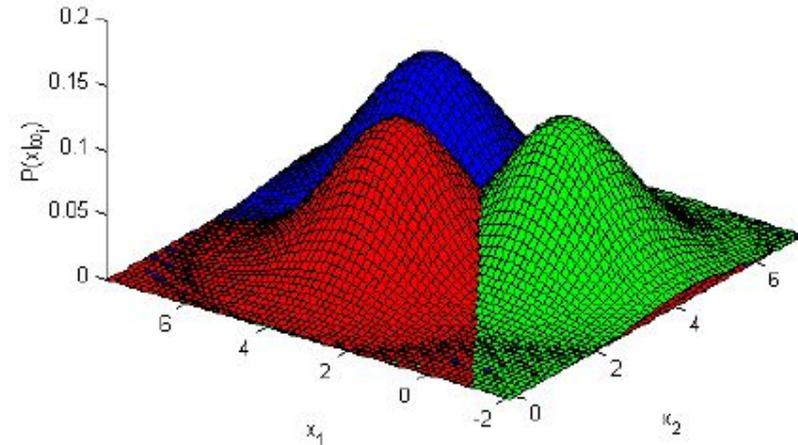
$$g_i(x) = -\frac{1}{2} \sum_{k=1}^N \frac{2x[k]\mu_i[k] + \mu_i[k]^2}{\sigma_k^2} - \frac{1}{2} \log \prod_{k=1}^N \sigma_k^2 + \log(P(\omega_i))$$

- Функция линейна, границы классов – гиперплоскости
- ГМТ точек с одинаковой вероятностью – гиперэллипсоиды
- Единственное отличие от предыдущего случая – нормализация осей

# Случай 2: $\Sigma_i = \Sigma$ ( $\Sigma$ – диагональная), пример

- Двумерная задача, три класса со следующими средними и ковариациями:

$$\begin{aligned} \mu_1 &= \begin{bmatrix} 3 & 2 \end{bmatrix}^T & \mu_2 &= \begin{bmatrix} 5 & 4 \end{bmatrix}^T & \mu_3 &= \begin{bmatrix} 2 & 5 \end{bmatrix}^T \\ \Sigma_1 &= \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} & \Sigma_2 &= \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} & \Sigma_3 &= \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \end{aligned}$$



## Случай 3: $\Sigma_i = \Sigma$ ( $\Sigma$ – не диагональная)

- Классы по-прежнему имеют одинаковую матрицу ковариации, но матрица – не диагональная

$$g_i(x) = w_i^T x + w_{i0}$$

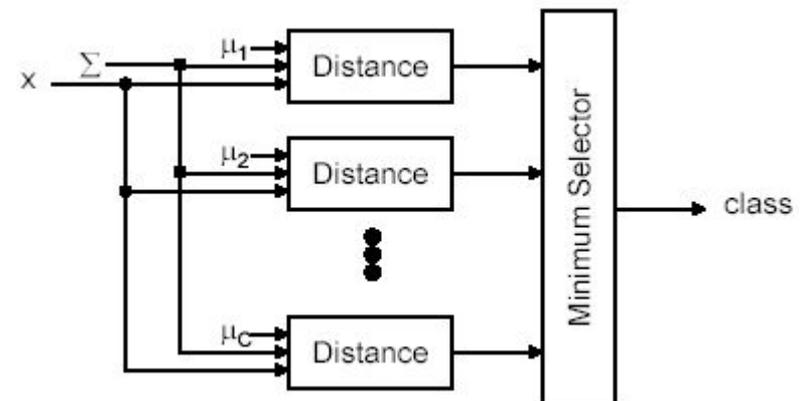
$$\text{where } \begin{cases} w_i = \Sigma^{-1} \mu_i \\ w_{i0} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \log P(\omega_i) \end{cases}$$

- Функция линейна, границы классов – гиперплоскости
- ГМТ точек с одинаковой вероятностью – гиперэллипсоиды, натянутые на собственные вектора матрицы ковариации
- При равных априорных вероятностях – классификатор по минимальному расстоянию Махаланобиса:

$$g_i(x) = -\frac{1}{2} (x - \mu_i)^T \Sigma^{-1} (x - \mu_i)$$

**Расстояние  
Махаланобиса**

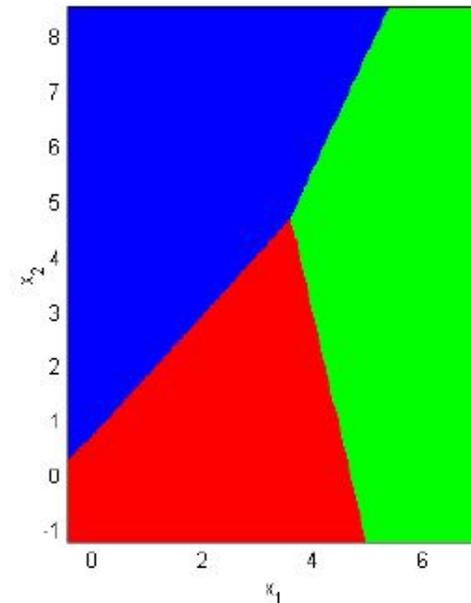
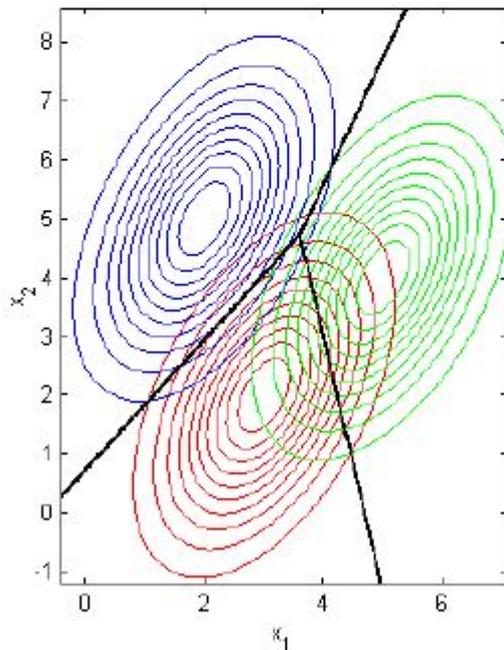
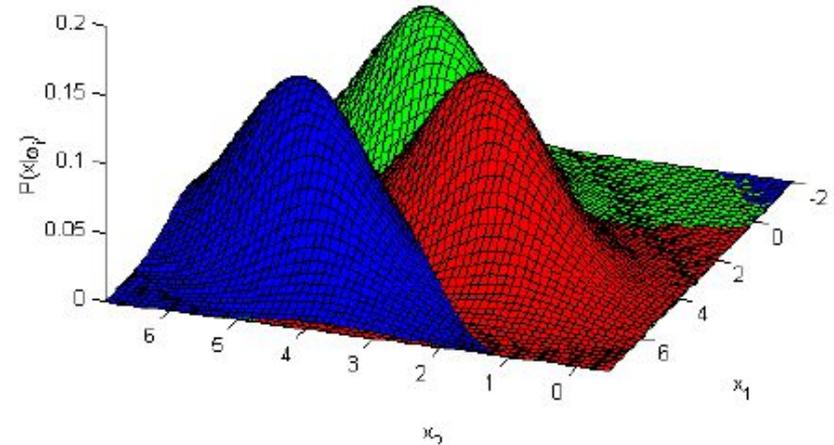
$$\|x - y\|_{\Sigma^{-1}}^2 = (x - y)^T \Sigma^{-1} (x - y)$$



# Случай 3: $\Sigma_i = \Sigma$ ( $\Sigma$ – не диагональная), пример

- Двумерная задача, три класса со следующими средними и ковариациями:

$$\begin{aligned} \mu_1 &= [3 \ 2]^T & \mu_2 &= [5 \ 4]^T & \mu_3 &= [2 \ 5]^T \\ \Sigma_1 &= \begin{bmatrix} 1 & 0.7 \\ 0.7 & 2 \end{bmatrix} & \Sigma_2 &= \begin{bmatrix} 1 & 0.7 \\ 0.7 & 2 \end{bmatrix} & \Sigma_3 &= \begin{bmatrix} 1 & 0.7 \\ 0.7 & 2 \end{bmatrix} \end{aligned}$$



## Случай 4: $\Sigma_i = \sigma_i^2 I$

- Классы имеют разные матрицы ковариации, которые пропорциональны единичной матрице

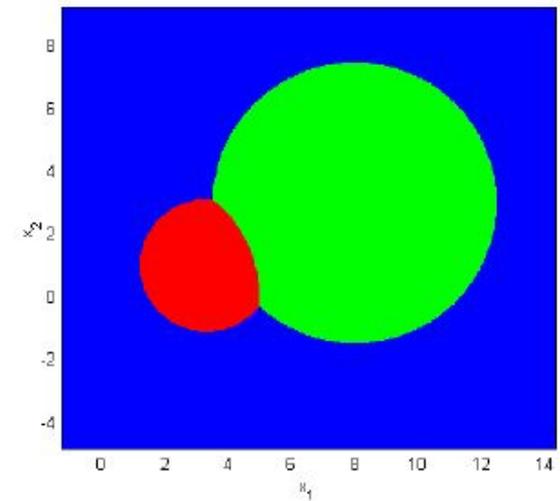
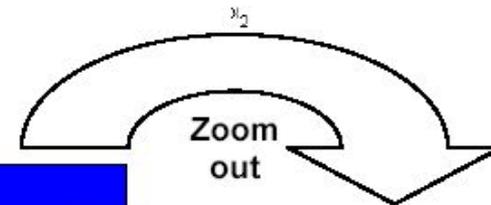
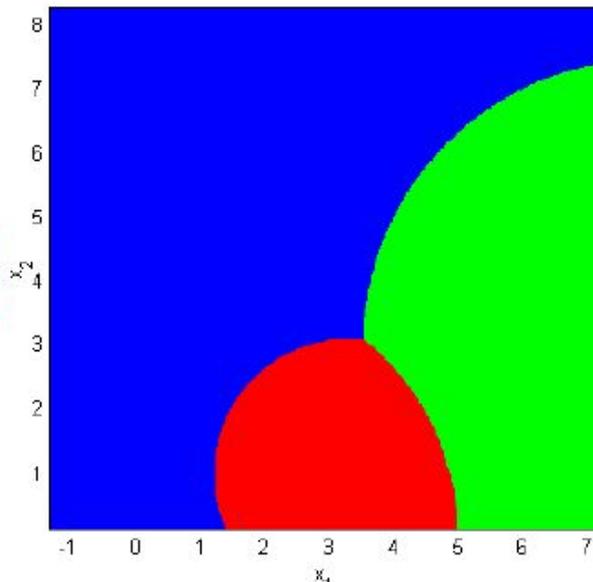
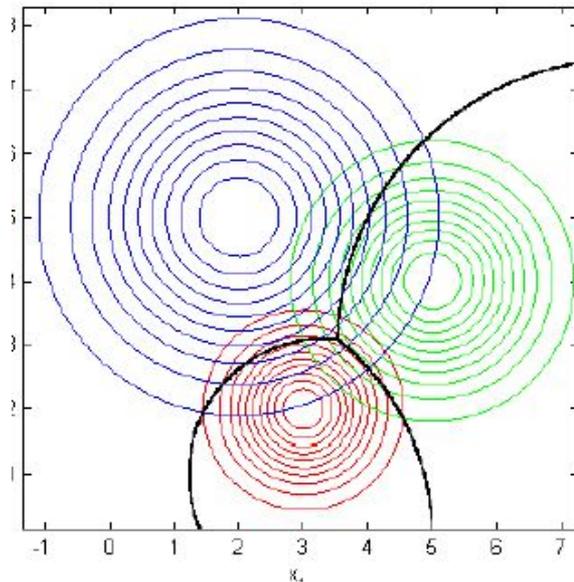
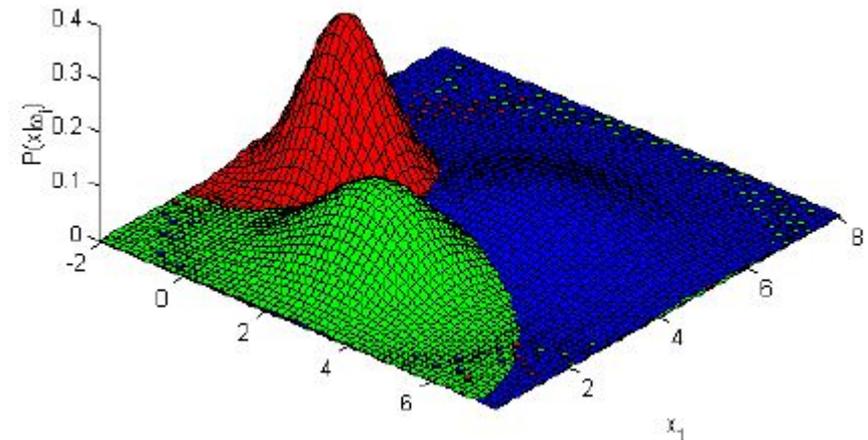
$$\begin{aligned}g_i(x) &= -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i) - \frac{1}{2} \log(|\Sigma_i|) + \log(P(\omega_i)) = \\ &= -\frac{1}{2}(x - \mu_i)^T \sigma_i^{-2}(x - \mu_i) - \frac{1}{2} N \log(\sigma_i^2) + \log(P(\omega_i))\end{aligned}$$

- Границы классов – гиперэллипсоиды
- ГМТ точек с одинаковой вероятностью – гиперсферы, определяемые осями пространства характеристик

# Случай 4: $\Sigma_i = \sigma_i^2 I$ , пример

- Двумерная задача, три класса со следующими средними и ковариациями:

$$\begin{aligned} \mu_1 &= [3 \ 2]^T & \mu_2 &= [5 \ 4]^T & \mu_3 &= [2 \ 5]^T \\ \Sigma_1 &= \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} & \Sigma_2 &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} & \Sigma_3 &= \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \end{aligned}$$



## Случай 5: $\Sigma_i \neq \Sigma_j$ (общий случай)

- Для общего случая разделяющая функция выглядит так:

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i) - \frac{1}{2} \log(|\Sigma_i|) + \log(P(\omega_i))$$

- Или, после преобразований:

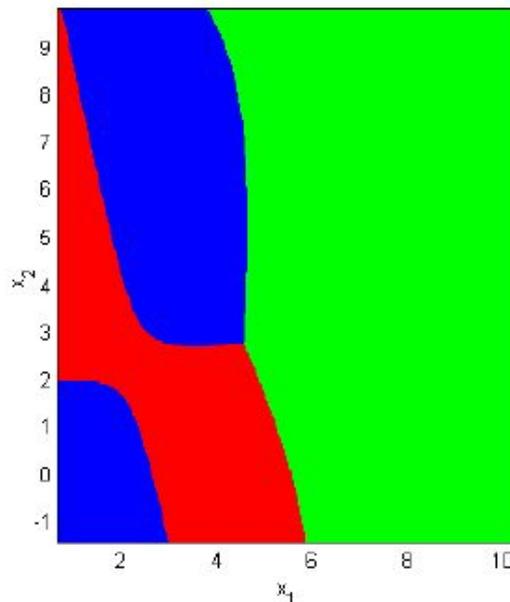
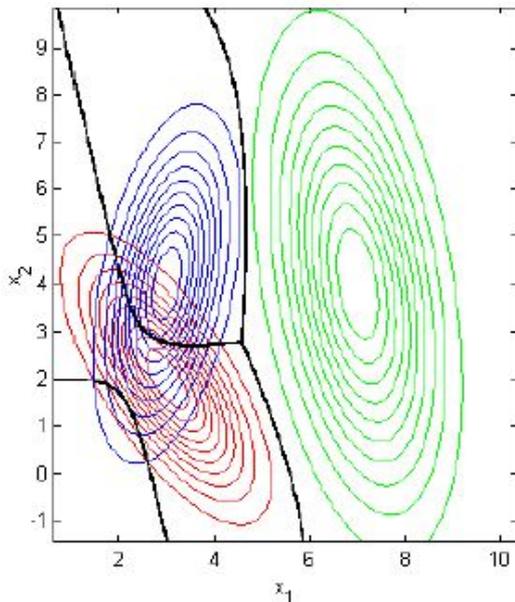
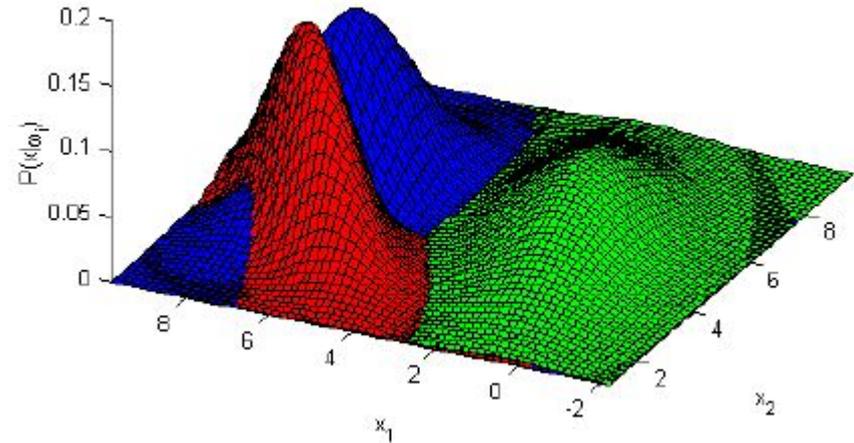
$$g_i(x) = x^T W_i x + w_i^T x + w_{i0}$$
$$\text{where } \begin{cases} W_i = -\frac{1}{2} \Sigma_i^{-1} \\ w_i = \Sigma_i^{-1} \mu_i \\ w_{i0} = -\frac{1}{2} \mu_i^T \Sigma_i^{-1} \mu_i - \frac{1}{2} \log(|\Sigma_i|) + \log(P(\omega_i)) \end{cases}$$

- Границы классов – гиперэллипсоиды и гиперпараболоиды
- ГМТ точек с одинаковой вероятностью – гиперэллипсоиды, определяемые собственными векторами матриц  $\Sigma_i$  для каждого из классов в отдельности

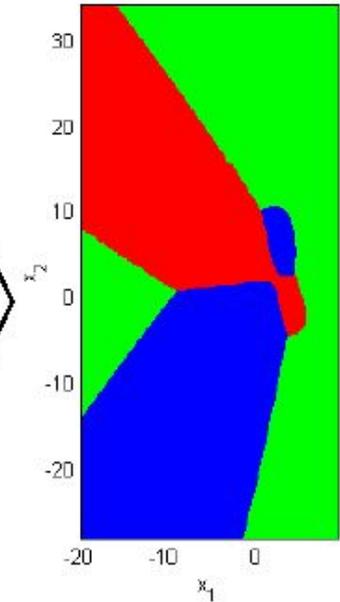
# Случай 5: $\Sigma_i \neq \Sigma_j$ (общий случай), пример

- Двумерная задача, три класса со следующими средними и ковариациями:

$$\begin{aligned} \mu_1 &= [3 \ 2]^T & \mu_2 &= [5 \ 4]^T & \mu_3 &= [2 \ 5]^T \\ \Sigma_1 &= \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} & \Sigma_2 &= \begin{bmatrix} 1 & -1 \\ -1 & 7 \end{bmatrix} & \Sigma_3 &= \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 3 \end{bmatrix} \end{aligned}$$



Zoom out



# Численный пример

- Построить линейную разделяющую функцию для трехмерной двухклассовой задачи распознавания по следующим данным:

$$\mu_1 = [0 \ 0 \ 0]^T; \quad \mu_2 = [1 \ 1 \ 1]^T; \quad \Sigma_1 = \Sigma_2 = \begin{bmatrix} 1/4 & 0 & 0 \\ 0 & 1/4 & 0 \\ 0 & 0 & 1/4 \end{bmatrix}; \quad p(\omega_2) = 2p(\omega_1)$$

- Решение

$$g_i(x) = -\frac{1}{2\sigma^2}(x - \mu_i)^T(x - \mu_i) + \log P(\omega_i) = -\frac{1}{2} \begin{bmatrix} x - \mu_x \\ y - \mu_y \\ z - \mu_z \end{bmatrix}^T \begin{bmatrix} 4 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} x - \mu_x \\ y - \mu_y \\ z - \mu_z \end{bmatrix} + \log P(\omega_i)$$

$$g_1(x) = -\frac{1}{2} \begin{bmatrix} x-0 \\ y-0 \\ z-0 \end{bmatrix}^T \begin{bmatrix} 4 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} x-0 \\ y-0 \\ z-0 \end{bmatrix} + \log \frac{1}{3}; \quad g_2(x) = -\frac{1}{2} \begin{bmatrix} x-1 \\ y-1 \\ z-1 \end{bmatrix}^T \begin{bmatrix} 4 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} x-1 \\ y-1 \\ z-1 \end{bmatrix} + \log \frac{2}{3}$$

$$g_1(x) \underset{\omega_2}{>} g_2(x) \underset{\omega_1}{\Rightarrow} -2(x^2 + y^2 + z^2) + \log \frac{1}{3} \underset{\omega_2}{>} -2((x-1)^2 + (y-1)^2 + (z-1)^2) + \log \frac{2}{3}$$

$$\boxed{x + y + z \underset{\omega_1}{>} \frac{6 - \log 2}{4} = 1.32}$$

- Классифицировать тестовый образец  $x_t = [0.1 \ 0.7 \ 0.8]^T$

$$0.1 + 0.7 + 0.8 = 1.6 \underset{\omega_1}{>} \underset{\omega_2}{1.32} \Rightarrow x_u \in \omega_2$$

# Некоторые заключения

- Байесовские классификаторы для нормально распределенных классов – квадратичные функции
- Байесовские классификаторы для нормально распределенных классов с одинаковыми ковариационными матрицами – линейные функции
- Классификатор по минимальному расстоянию Махалобиса оптимален по Байесу, если
  - классы нормально распределены **и**
  - матрицы ковариаций равны **и**
  - априорные вероятности равны
- Классификатор по минимальному евклидову расстоянию оптимален по Байесу, если
  - классы нормально распределены **и**
  - матрицы ковариаций равны и пропорциональны единичной матрице **и**
  - априорные вероятности равны