

# КЛАССИФИКАТОРЫ

---

Сергей Николенко

uData School — Киев

8 июня 2018 г.

---

## *Random facts:*

- 8 июня 1648 г. Богдан Хмельницкий написал письмо Алексею Михайловичу с просьбой принять Гетманщину в подданство Русского царства
- 8 июня 1783 г. произошло извержение вулкана Лаки в Исландии, одно из крупнейших в истории; погиб каждый пятый житель страны, три четверти домашнего скота, 75% территории оказались покрыты лавой; два года после извержения известны как Móðuharðindin («бедствия в тумане»)
- 8 июня 1824 г. Ноа Кашинг из Квебека запатентовал стиральную машину, а 8 июня 1860 г. Айс Макгаффи из Чикаго запатентовал пылесос

ЭКВИВАЛЕНТНОЕ ЯДРО И  
СРАВНЕНИЕ МОДЕЛЕЙ

---

- Вспомним наши байесовские предсказания:

$$p(t | \mathbf{t}, \alpha, \beta) = \mathcal{N}(t | \mu_N^\top \phi(\mathbf{x}), \sigma_N^2),$$

$$\text{где } \sigma_N^2 = \frac{1}{\beta} + \phi(\mathbf{x})^\top \Sigma_N \phi(\mathbf{x}).$$

- Давайте перепишем среднее апостериорного распределения в другой форме (вспомним, что  $\mu_N = \beta \Sigma_N \Phi^\top \mathbf{t}$ ):

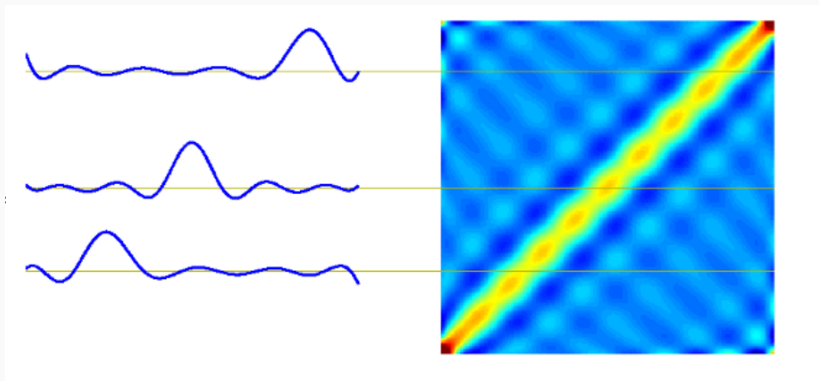
$$\begin{aligned} y(\mathbf{x}, \mu_N) &= \mu_N^\top \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^\top \Sigma_N \Phi^\top \mathbf{t} = \\ &= \sum_{n=1}^N \beta \phi(\mathbf{x})^\top \Sigma_N \phi(\mathbf{x}_n) t_n. \end{aligned}$$

- $y(\mathbf{x}, \mu_N) = \sum_{n=1}^N \beta \phi(\mathbf{x})^\top \Sigma_N \phi(\mathbf{x}_n) t_n$ .
- Это значит, что предсказание можно переписать как

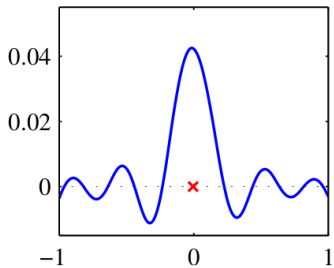
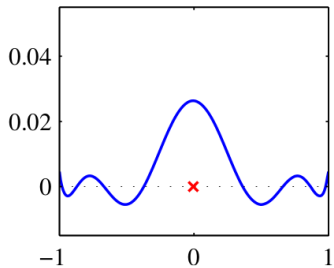
$$y(\mathbf{x}, \mu_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n.$$

- Т.е. мы предсказываем следующую точку как линейную комбинацию значений в известных точках.
- Функция  $k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^\top \Sigma_N \phi(\mathbf{x}')$  называется *эквивалентным ядром* (equivalent kernel).

# ЭКВИВАЛЕНТНОЕ ЯДРО



# ЭКВИВАЛЕНТНОЕ ЯДРО



## Выводы про эквивалентное ядро

- Эквивалентное ядро  $k(\mathbf{x}, \mathbf{x}')$  локализовано вокруг  $\mathbf{x}$  как функция  $\mathbf{x}'$ , т.е. каждая точка оказывает наибольшее влияние около себя и затухает потом.
- Можно было бы с самого начала просто определить ядро и предсказывать через него, безо всяких базисных функций  $\phi$  – такой подход мы ещё будем рассматривать.

**Упражнение.** Докажите, что  $\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1$ .

- Мы говорили о том, что при увеличении числа параметров модели возникает оверфиттинг.
- Как этого избежать? Как сравнить модели с разным числом параметров?
- Теория байесовского вывода предлагает такой выход: давайте будем не точечные оценки параметров модели рассматривать, а тоже интегрировать по параметрам модели.



- Пусть мы хотим сравнить модели из множества  $\{\mathcal{M}_i\}_{i=1}^L$ .
- Модель – это распределение вероятностей над данными  $D$ .
- По тестовому набору  $D$  можно оценить апостериорное распределение

$$p(\mathcal{M}_i | D) \propto p(\mathcal{M}_i)p(D | \mathcal{M}_i).$$

- Если знать апостериорное распределение, то можно сделать предсказание:

$$p(t | \mathbf{x}, D) = \sum_{i=1}^L p(t | \mathbf{x}, \mathcal{M}_i, D)p(\mathcal{M}_i | D).$$

- *Model selection* (выбор модели) – это когда мы приближаем предсказание, выбирая просто самую (апостериорно) вероятную модель.

- Если модель определена параметрически, через  $\mathbf{w}$ , то

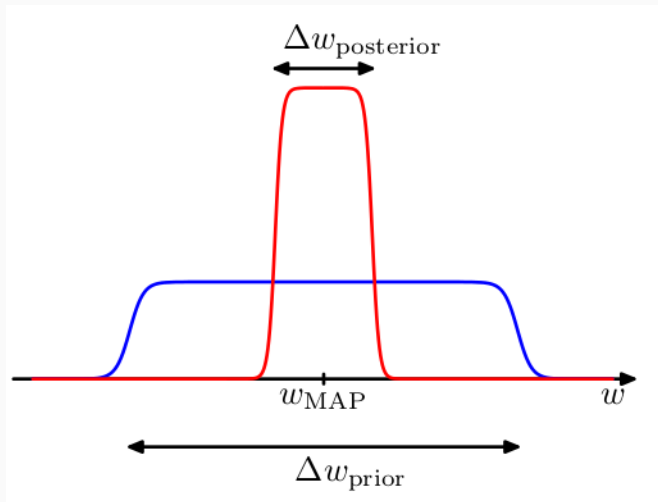
$$p(D | \mathcal{M}_i) = \int p(D | \mathbf{w}, \mathcal{M}_i)p(\mathbf{w} | \mathcal{M}_i)d\mathbf{w}.$$

- Т.е. это вероятность сгенерировать  $D$ , если выбрать параметры модели по её априорному распределению, а потом накидывать данные.
- Это, кстати, в точности знаменатель из теоремы Байеса:

$$p(\mathbf{w} | \mathcal{M}_i, D) = \frac{p(D | \mathbf{w}, \mathcal{M}_i)p(\mathbf{w} | \mathcal{M}_i)}{p(D | \mathcal{M}_i)}.$$

- Предположим, что у модели один параметр  $w$ , а апостериорное распределение – это острый пик вокруг  $w_{\text{MAP}}$  шириной  $\Delta w_{\text{posterior}}$ .
- Тогда можно приблизить  $p(D) = \int p(D | w)p(w)dw$  как значение в максимуме, умноженное на ширину.
- Предположим ещё, что априорное распределение тоже плоское,  $p(w) = \frac{1}{\Delta w_{\text{prior}}}$ .

# ПРИБЛИЖЕНИЕ $p(D)$



- Тогда получится

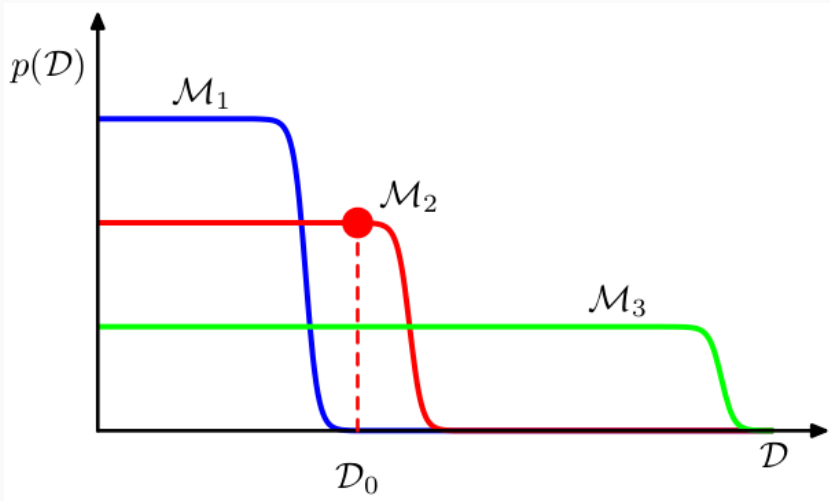
$$p(D) = \int p(D | w)p(w)dw \approx p(D | w_{\text{MAP}}) \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}},$$
$$\ln p(D) \approx \ln p(D | w_{\text{MAP}}) + \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right).$$

- Это значит, что мы добавляем штраф за «слишком узкое» апостериорное распределение – то есть в точности штраф за оверфиттинг!
- Для модели из  $M$  параметров, если предположить, что у них одинаковые  $\Delta w_{\text{posterior}}$ , получим

$$\ln p(D) \approx \ln p(D | w_{\text{MAP}}) + M \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right).$$

- Другими словами: давайте посмотрим, какие датасеты может генерировать та или иная модель.
- Простая модель (e.g., линейная) генерирует похожие датасеты, «мало» разных датасетов, у неё высокая  $p(D | \mathcal{M})$ .
- Сложная модель (e.g., многочлен девятой степени) генерирует «много» разных датасетов, у неё низкая  $p(D | \mathcal{M})$ .
- Но сложная может хорошо выразить датасеты, которые не может выразить простая; поэтому в сумме надо выбирать «среднюю».

# ПРИБЛИЖЕНИЕ $p(D)$





- Sanity check: тут какие-то штрафы мы навводили; будет ли истинный правильный ответ  $p(D | \mathcal{M}_{\text{true}})$  всегда оптимальным в этом смысле?
- Конечно, для конкретного датасета может так повезти, что не будет.
- Но если усреднить по всем датасетам, выбранным по  $p(D | \mathcal{M}_{\text{true}})$ ...

- ...то получится

$$E \left[ \ln \frac{p(D | \mathcal{M}_{\text{true}})}{p(D | \mathcal{M})} \right] = \int p(D | \mathcal{M}_{\text{true}}) \ln \frac{p(D | \mathcal{M}_{\text{true}})}{p(D | \mathcal{M})} dD.$$

- Это называется *расстоянием Кульбака-Лейблера* (Kullback-Leibler divergence) между распределениями  $p(D | \mathcal{M}_{\text{true}})$  и  $p(D | \mathcal{M})$ .

**Упражнение.** Докажите, что расстояние Кульбака-Лейблера всегда неотрицательно, т.е. что  $p(D | \mathcal{M}_{\text{true}}) \geq p(D | \mathcal{M})$  для любой  $\mathcal{M}$ .

## ВВЕДЕНИЕ В КЛАССИФИКАЦИЮ

---

- Теперь классификация: определить вектор  $\mathbf{x}$  в один из  $K$  классов  $\mathcal{C}_k$ .
- В итоге у нас так или иначе всё пространство разобьётся на эти классы.
- Т.е. на самом деле мы ищем *разделяющую поверхность* (decision surface, decision boundary).

## ЗАДАЧА КЛАССИФИКАЦИИ

- Как кодировать? Бинарная задача – очень естественно, переменная  $t$ ,  $t = 0$  соответствует  $\mathcal{C}_1$ ,  $t = 1$  соответствует  $\mathcal{C}_2$ .
- Оценку  $t$  можно интерпретировать как вероятность (по крайней мере, мы постараемся, чтобы было можно).
- Если несколько классов – удобно 1-of- $K$ :

$$\mathbf{t} = (0, \dots, 0, 1, 0, \dots)^\top.$$

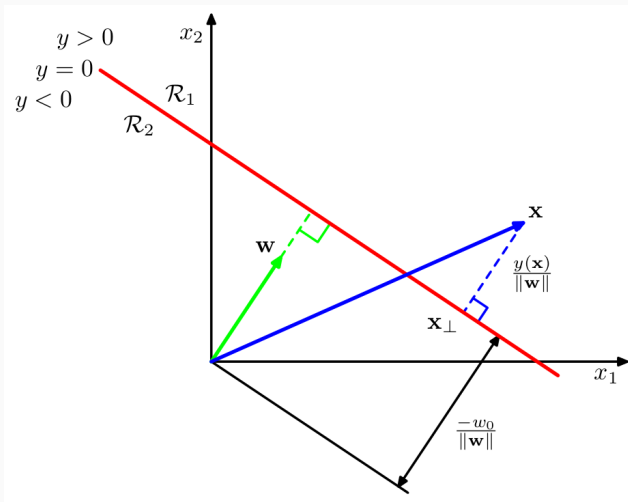
- Тоже можно интерпретировать как вероятности – или пропорционально им.

- Начнём с геометрии: рассмотрим линейную дискриминантную функцию

$$y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0.$$

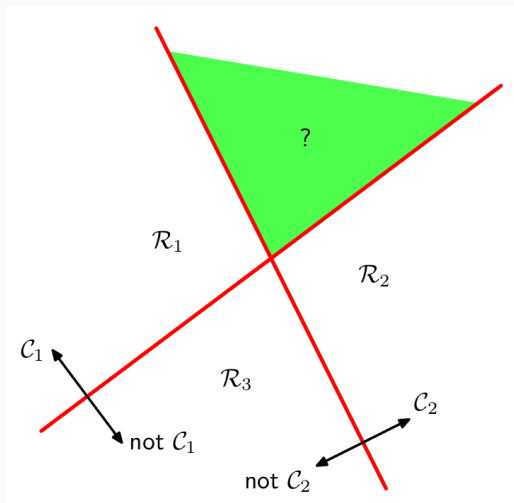
- Это гиперплоскость, и  $\mathbf{w}$  – нормаль к ней.
- Расстояние от начала координат до гиперплоскости равно  $\frac{-w_0}{\|\mathbf{w}\|}$ .
- $y(\mathbf{x})$  связано с расстоянием до гиперплоскости:  $d = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$ .

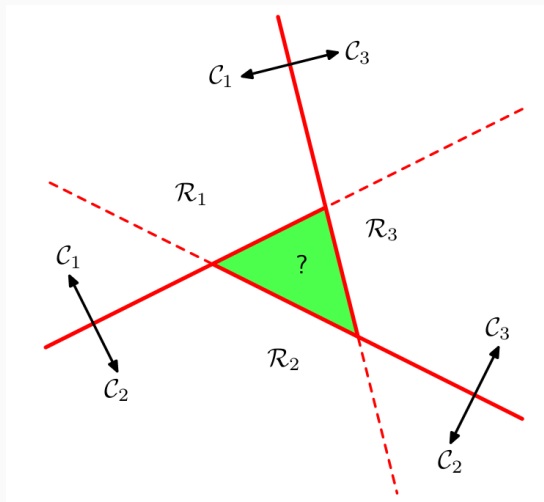
# РАЗДЕЛЯЮЩАЯ ГИПЕРПЛОСКОСТЬ



- С несколькими классами выходит задача.
- Можно рассмотреть  $K$  поверхностей вида «один против всех».
- Можно –  $\binom{K}{2}$  поверхностей вида «каждый против каждого».
- Но всё это как-то нехорошо.





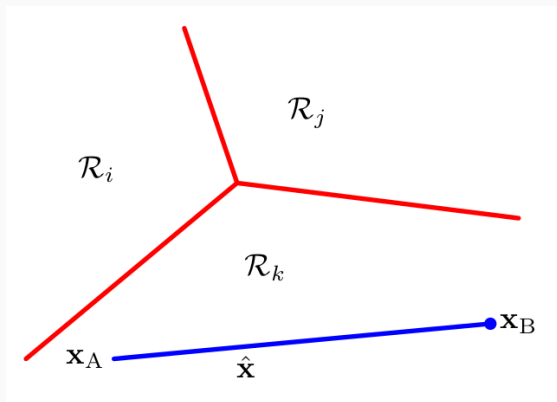


- Лучше рассмотреть единый дискриминант из  $K$  линейных функций:

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + w_{k0}.$$

- Классифицировать в  $\mathcal{C}_k$ , если  $y_k(\mathbf{x})$  – максимален.
- Тогда разделяющая поверхность между  $\mathcal{C}_k$  и  $\mathcal{C}_j$  будет гиперплоскостью вида  $y_k(\mathbf{x}) = y_j(\mathbf{x})$ :

$$(\mathbf{w}_k - \mathbf{w}_j)^\top \mathbf{x} + (w_{k0} - w_{j0}).$$



**Упражнение.** Докажите, что области, соответствующие классам, при таком подходе всегда односвязные и выпуклые.

- Мы снова можем воспользоваться методом наименьших квадратов: запишем  $y_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + w_{k0}$  вместе (спрятав свободный член) как

$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^\top \mathbf{x}.$$

- Можно найти  $\mathbf{W}$ , оптимизируя сумму квадратов; функция ошибки:

$$E_D(\mathbf{W}) = \frac{1}{2} \text{Tr} [(\mathbf{XW} - \mathbf{T})^\top (\mathbf{XW} - \mathbf{T})].$$

- Берём производную, решаем...

- ...получается привычное

$$\mathbf{W} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{T} = \mathbf{X}^\dagger \mathbf{T},$$

где  $\mathbf{X}^\dagger$  – псевдообратная Мура-Пенроуза.

- Теперь можно найти и дискриминантную функцию:

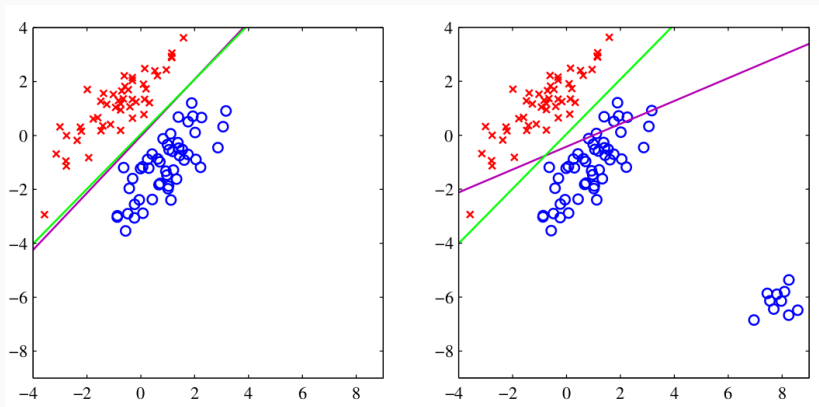
$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^\top \mathbf{x} = \mathbf{T}^\top (\mathbf{X}^\dagger)^\top \mathbf{x}.$$

- Это решение сохраняет линейность.

**Упражнение.** Докажите, что в схеме кодирования 1-of- $K$  предсказания  $y_k(\mathbf{x})$  для разных классов при любом  $\mathbf{x}$  будут давать в сумме 1. Почему они всё-таки не будут разумными оценками вероятностей?

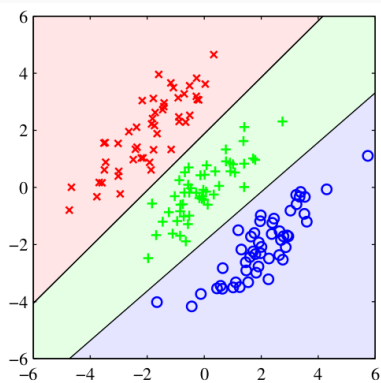
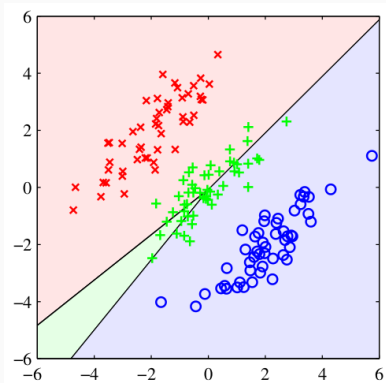
- Проблемы наименьших квадратов:
  - outliers плохо обрабатываются;
  - «слишком правильные» предсказания добавляют штраф.

# ПРОБЛЕМЫ НАИМЕНЬШИХ КВАДРАТОВ





# ПРОБЛЕМЫ НАИМЕНЬШИХ КВАДРАТОВ



- Почему так? Почему наименьшие квадраты так плохо работают?

- Почему так? Почему наименьшие квадраты так плохо работают?
- Они предполагают гауссовское распределение ошибки.
- Но, конечно, распределение у бинарных векторов далеко не гауссово.

- Другой взгляд на классификацию: в линейном случае мы хотим спроецировать точки в размерность 1 (на нормаль разделяющей гиперплоскости) так, чтобы в этой размерности 1 они хорошо разделялись.
- Т.е. классификация – это такой метод радикального сокращения размерности.
- Давайте посмотрим на классификацию с этих позиций и попробуем добиться оптимальности в каком-то смысле.

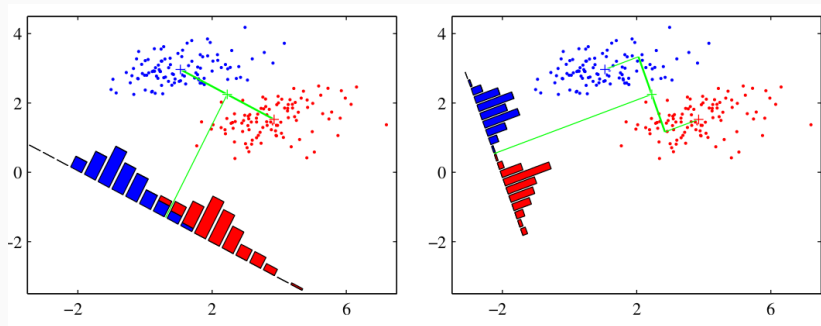
- Рассмотрим два класса  $\mathcal{C}_1$  и  $\mathcal{C}_2$  с  $N_1$  и  $N_2$  точками.
- Первая идея – надо найти серединный перпендикуляр между центрами кластеров

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{\mathcal{C}_1} \mathbf{x}, \text{ и } \mathbf{m}_2 = \frac{1}{N_2} \sum_{\mathcal{C}_2} \mathbf{x},$$

т.е. максимизировать  $\mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1)$ .

- Надо ещё добавить ограничение  $\|\mathbf{w}\| = 1$ , но всё равно не ахти как работает.

# ЛИНЕЙНЫЙ ДИСКРИМИНАНТ ФИШЕРА



Чем левая картинка хуже правой?

- Слева больше дисперсия каждого кластера.
- Идея: минимизировать перекрытие классов, оптимизируя и проекцию расстояния, и дисперсию.
- Выборочные дисперсии в проекции: для  $y_n = \mathbf{w}^\top \mathbf{x}_n$

$$s_1 = \sum_{n \in \mathcal{C}_1} (y_n - m_1)^2 \quad \text{и} \quad s_2 = \sum_{n \in \mathcal{C}_2} (y_n - m_2)^2.$$

- Критерий Фишера:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}, \text{ где}$$

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^\top,$$

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1) (\mathbf{x}_n - \mathbf{m}_1)^\top + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2) (\mathbf{x}_n - \mathbf{m}_2)^\top.$$

(between-class covariance и within-class covariance).

- Дифференцируя по  $\mathbf{w}$ ...



- ...получим, что  $J(\mathbf{w})$  максимален при

$$(\mathbf{w}^\top \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^\top \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}.$$

- Т.к.  $\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top$ ,  $\mathbf{S}_B \mathbf{w}$  всё равно будет в направлении  $\mathbf{m}_2 - \mathbf{m}_1$ , а длина  $\mathbf{w}$  нас не интересует.
- Поэтому получается

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1).$$

- В итоге мы выбрали направление проекции, и осталось только разделить данные на этой проекции.

- Любопытно, что дискриминант Фишера тоже можно получить из наименьших квадратов.
- Давайте для класса  $\mathcal{C}_1$  выберем целевое значение  $\frac{N_1+N_2}{N_1}$ , а для класса  $\mathcal{C}_2$  возьмём  $-\frac{N_1+N_2}{N_2}$ .

**Упражнение.** Докажите, что при таких целевых значениях наименьшие квадраты – это дискриминант Фишера.

- А что будет с несколькими классами? Рассмотрим  $\mathbf{y} = \mathbf{W}^\top \mathbf{x}$ , обобщим внутреннюю дисперсию как

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k = \sum_{k=1}^K \sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k) (\mathbf{x}_n - \mathbf{m}_k)^\top.$$

- Чтобы обобщить внешнюю (межклассовую) дисперсию, просто возьмём остаток полной дисперсии

$$\mathbf{S}_T = \sum_n (\mathbf{x}_n - \mathbf{m}) (\mathbf{x}_n - \mathbf{m})^\top,$$

$$\mathbf{S}_B = \mathbf{S}_T - \mathbf{S}_W.$$

- Обобщить критерий можно разными способами, например:

$$J(\mathbf{W}) = \text{Tr} [\mathbf{s}_W^{-1} \mathbf{s}_B],$$

где  $\mathbf{s}$  – ковариации в пространстве проекций на  $\mathbf{y}$ :

$$\mathbf{s}_W = \sum_{k=1}^K \sum_{n \in \mathcal{C}_k} (\mathbf{y}_n - \mu_k) (\mathbf{y}_n - \mu_k)^\top,$$

$$\mathbf{s}_B = \sum_{k=1}^K N_k (\mu_k - \mu) (\mu_k - \mu)^\top,$$

где  $\mu_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{y}_n$ .

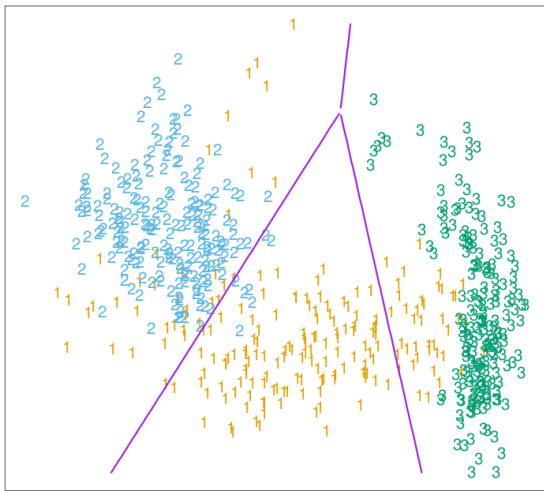
## LDA и QDA

---

- В прошлый раз мы рассмотрели задачу классификации.
- Построили разделяющую гиперплоскость методом наименьших квадратов.
- И методом линейного дискриминанта Фишера.
- А потом научились обучать перцептрон и доказали сходимость метода.

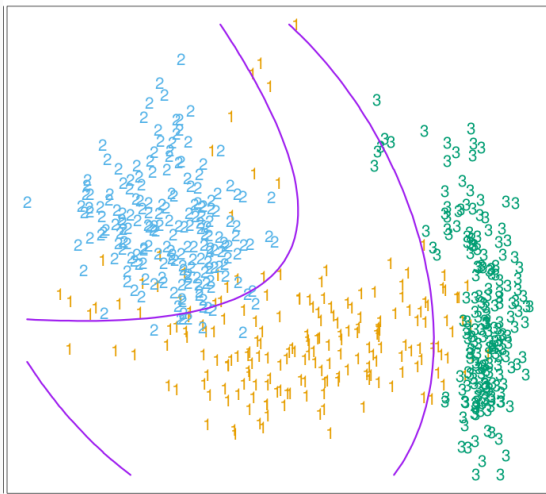
- Мы учились проводить разделяющие гиперплоскости.
- Но как же нелинейные поверхности?
- Можно делать нелинейные из линейных, увеличивая размерность.

# НЕЛИНЕЙНЫЕ ПОВЕРХНОСТИ





# НЕЛИНЕЙНЫЕ ПОВЕРХНОСТИ



- Теперь классификация через генеративные модели: давайте каждому классу сопоставим плотность  $p(\mathbf{x} | \mathcal{C}_k)$ , найдём априорные распределения  $p(\mathcal{C}_k)$ , будем искать  $p(\mathcal{C}_k | \mathbf{x})$  по теореме Байеса.
- Для двух классов:

$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x} | \mathcal{C}_2)p(\mathcal{C}_2)}.$$

- Перепишем:

$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x} | \mathcal{C}_2)p(\mathcal{C}_2)} = \frac{1}{1 + e^{-a}} = \sigma(a),$$

где

$$a = \ln \frac{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_2)p(\mathcal{C}_2)}, \quad \sigma(a) = \frac{1}{1 + e^{-a}}.$$

- $\sigma(a)$  – логистический сигмоид:

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

- $\sigma(-a) = 1 - \sigma(a)$ .
- $a = \ln\left(\frac{\sigma}{1-\sigma}\right)$  – логит-функция.

**Упражнение.** Докажите эти свойства.

- В случае нескольких классов получится

$$p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j)p(\mathcal{C}_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}}.$$

- Здесь  $a_k = \ln p(\mathbf{x} | \mathcal{C}_k)p(\mathcal{C}_k)$ .
- $\frac{e^{a_k}}{\sum_j e^{a_j}}$  – нормализованная экспонента, или softmax-функция (сглаженный максимум).

- Давайте рассмотрим гауссовы распределения для классов:

$$p(\mathbf{x} | \mathcal{C}_k) = \mathcal{N}(\mathbf{x} | \mu_k, \Sigma).$$

- Сначала пусть  $\Sigma$  у всех одинаковые, а классов всего два.
- Посчитаем логистический сигмоид...

- ...получится

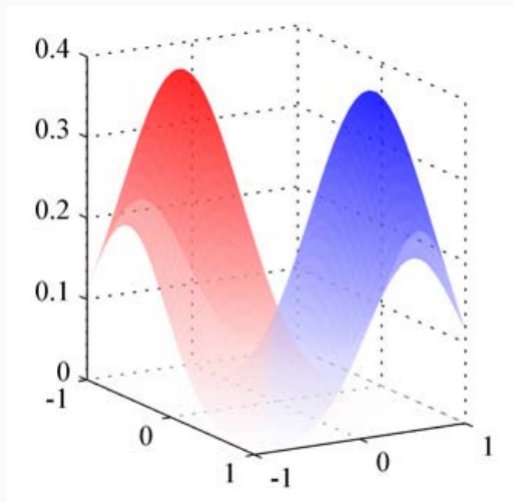
$$p(\mathcal{C}_1 | \mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + w_0), \text{ где}$$

$$\mathbf{w} = \Sigma^{-1} (\mu_1 - \mu_2),$$

$$w_0 = -\frac{1}{2}\mu_1^\top \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^\top \Sigma^{-1} \mu_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}.$$

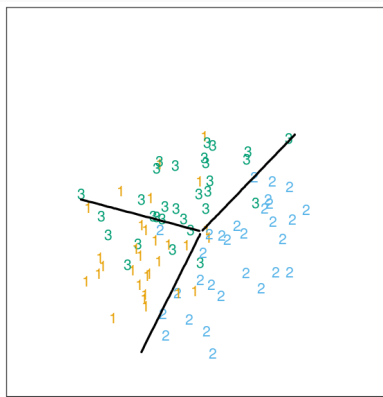
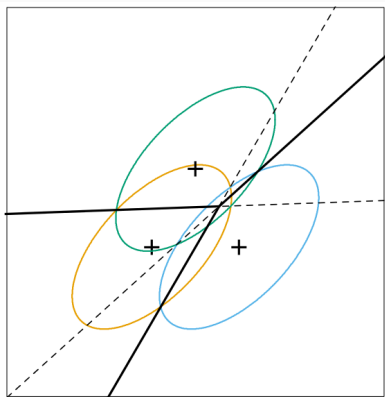
- Т.е. в аргументе сигмоида получается линейная функция от  $\mathbf{x}$ . Поверхности уровня – это когда  $p(\mathcal{C}_1 | \mathbf{x})$  постоянно, т.е. гиперплоскости в пространстве  $\mathbf{x}$ . Априорные вероятности  $p(\mathcal{C}_k)$  просто сдвигают эти гиперплоскости.

# РАЗДЕЛЯЮЩАЯ ГИПЕРПЛОСКОСТЬ



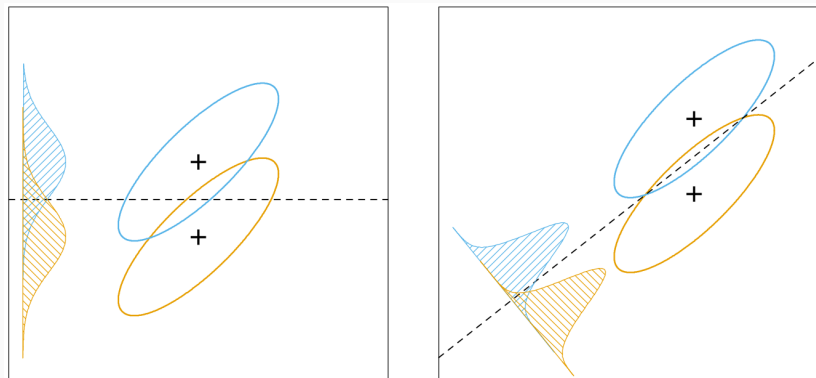


# РАЗДЕЛЯЮЩАЯ ГИПЕРПЛОСКОСТЬ



## ДИСКРИМИНАНТ ФИШЕРА

Кстати, с дискриминантом Фишера эта разделяющая поверхность отлично сходится.



- С несколькими классами получится тоже примерно так же:

$$\delta_k(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \ln \pi_k,$$

где  $\pi_k = p(\mathcal{C}_k)$ .

- Получились линейные  $\delta_k(\mathbf{x})$ , и опять разделяющие поверхности линейные (тут разделяющие поверхности – когда две максимальных вероятности равны).
- Этот метод называется LDA – linear discriminant analysis.

- Как оценить распределения  $p(\mathbf{x} | \mathcal{C}_k)$ , если даны только данные?
- Можно по методу максимального правдоподобия.
- Опять рассмотрим тот же пример: два класса, гауссианы с одинаковой матрицей ковариаций, и есть  $D = \{\mathbf{x}_n, t_n\}_{n=1}^N$ , где  $t_n = 1$  значит  $\mathcal{C}_1$ ,  $t_n = 0$  значит  $\mathcal{C}_2$ .
- Обозначим  $p(\mathcal{C}_1) = \pi$ ,  $p(\mathcal{C}_2) = 1 - \pi$ .

- Для одной точки в классе  $\mathcal{C}_1$ :

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1)p(\mathbf{x}_n | \mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma).$$

- В классе  $\mathcal{C}_2$ :

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2)p(\mathbf{x}_n | \mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \mu_2, \Sigma).$$

- Функция правдоподобия:

$$\begin{aligned} p(\mathbf{t} | \pi, \mu_1, \mu_2, \Sigma) &= \\ &= \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma)]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \mu_2, \Sigma)]^{1-t_n}. \end{aligned}$$

- Максимизируем логарифм правдоподобия. Сначала по  $\pi$ , там останется только

$$\sum_{n=1}^N [t_n \ln \pi + (1 - t_n) \ln(1 - \pi)],$$

и, взяв производную, получим, совершенно неожиданно,

$$\hat{\pi} = \frac{N_1}{N_1 + N_2}.$$

- Теперь по  $\mu_1$ ; всё, что зависит от  $\mu_1$ :

$$\sum_n t_n \ln \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma) = -\frac{1}{2} \sum_n t_n (\mathbf{x}_n - \mu_1)^\top \Sigma^{-1} (\mathbf{x}_n - \mu_1) + C.$$

- Берём производную, и получается, опять внезапно,

$$\hat{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n.$$

- Аналогично,

$$\hat{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n.$$

- Для матрицы ковариаций придётся постараться; в результате получится

$$\hat{\Sigma} = \frac{N_1}{N_1 + N_2} \mathbf{S}_1 + \frac{N_2}{N_1 + N_2} \mathbf{S}_2, \text{ где}$$
$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mu_1) (\mathbf{x}_n - \mu_1)^\top,$$
$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mu_2) (\mathbf{x}_n - \mu_2)^\top.$$

- Тоже совершенно неожиданно: взвешенное среднее оценок для двух матриц ковариаций.



- Это самым прямым образом обобщается на случай

нескольких классов.

**Упражнение.** Сделайте это.

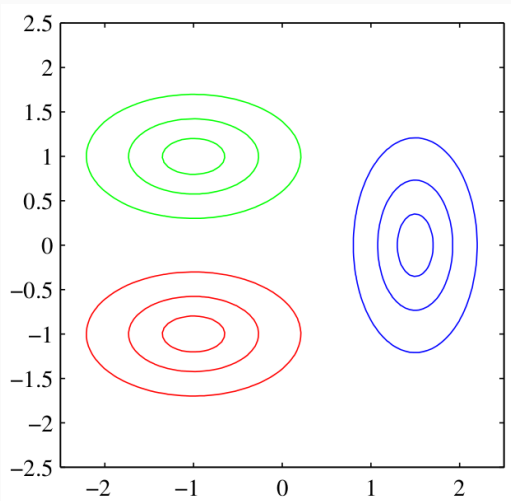
- А вот с разными матрицами ковариаций уже будет по-другому.
- Квадратичные члены не сократятся.
- Разделяющие поверхности станут квадратичными; QDA – quadratic discriminant analysis.

- В QDA получится

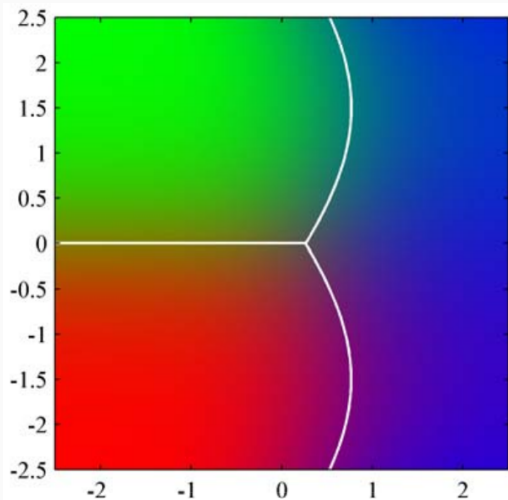
$$\delta_k(\mathbf{x}) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k) + \log \pi_k.$$

- Разделяющая поверхность между  $\mathcal{C}_i$  и  $\mathcal{C}_j$  – это  $\{\mathbf{x} \mid \delta_i(\mathbf{x}) = \delta_j(\mathbf{x})\}$ .
- Оценки максимального правдоподобия такие же, только надо отдельно матрицы ковариаций оценивать.

## РАЗНЫЕ МАТРИЦЫ КОВАРИАЦИИ

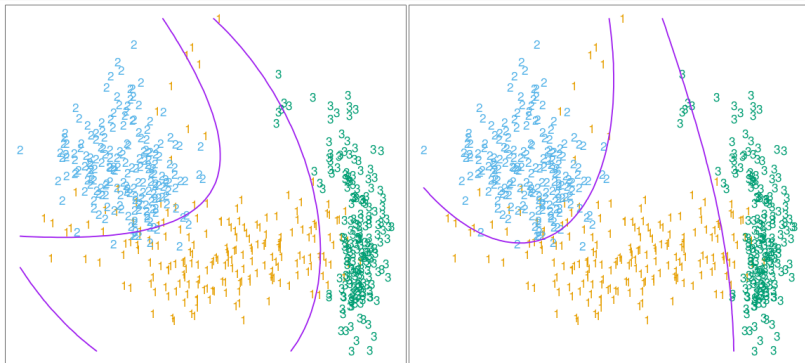


## РАЗНЫЕ МАТРИЦЫ КОВАРИАЦИЙ



# LDA vs. QDA

Разница между LDA с квадратичными членами и QDA обычно невелика.



- LDA и QDA неплохо работают на практике. Часто это первая идея в классификации.
- Число параметров:
  - у LDA  $(K - 1)(d + 1)$  параметр: по  $d + 1$  на каждую разницу вида  $\delta_k(\mathbf{x}) - \delta_K(\mathbf{x})$ ;
  - у QDA  $(K - 1)(d(d + 3)/2 + 1)$  параметр, но он выглядит гораздо лучше своих лет.

- Почему хорошо работают?
- Скорее всего, потому, что линейные и квадратичные оценки достаточно стабильны: даже если bias относительно большой (как будет, если данные всё-таки не гауссианами порождены), variance будет маленькой.



- Компромисс между LDA и QDA – регуляризованный дискриминантный анализ, RDA.
- Стянем ковариации каждого класса к общей матрице ковариаций:

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma},$$

где  $\hat{\Sigma}_k$  – оценка из QDA,  $\hat{\Sigma}$  – оценка из LDA.

- Или стянем к единичной матрице:

$$\hat{\Sigma}_k(\gamma) = \gamma \hat{\Sigma}_k + (1 - \gamma) \hat{\sigma}^2 \mathbf{I}.$$

- Предположим, что размерность  $d$  больше, чем число классов  $K$ .
- Тогда центроиды классов  $\hat{\mu}_k$  лежат в подпространстве размерности  $\leq K - 1$ .
- И когда мы определяем ближайший центроид, нам достаточно считать расстояния только в этом подпространстве.
- Таким образом, можно сократить ранг задачи.

- Куда именно проецировать? Не обязательно само подпространство, порождённое центроидами, будет оптимальным.
- Это мы уже проходили: для размерности 1 это линейный дискриминант Фишера.
- Это он и есть: оптимальное подпространство будет там, где межклассовая дисперсия максимальна по отношению к внутриклассовой.

# ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ

---

- Итак, мы рассмотрели логистический сигмоид:

$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x} | \mathcal{C}_2)p(\mathcal{C}_2)} = \frac{1}{1 + e^{-a}} = \sigma(a),$$

$$\text{где } a = \ln \frac{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_2)p(\mathcal{C}_2)}, \quad \sigma(a) = \frac{1}{1 + e^{-a}}.$$

- Вывели из него LDA и QDA, обучили их методом максимального правдоподобия, а потом отвлеклись на naïve Bayes.

- Два класса, и апостериорное распределение – логистический сигмоид на линейной функции:

$$p(\mathcal{C}_1 | \phi) = y(\phi) = \sigma(\mathbf{w}^T \phi), \quad p(\mathcal{C}_2 | \phi) = 1 - p(\mathcal{C}_1 | \phi).$$

- *Логистическая регрессия* – это когда мы напрямую оптимизируем  $\mathbf{w}$ .

- Для датасета  $\{\phi_n, t_n\}$ ,  $t_n \in \{0, 1\}$ ,  $\phi_n = \phi(\mathbf{x}_n)$ :

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}, \quad y_n = p(\mathcal{C}_1 | \phi_n).$$

- Ищем параметры максимального правдоподобия, минимизируя  $-\ln p(\mathbf{t} | \mathbf{w})$ :

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)].$$

Спасибо за внимание!