

# СТАТИСТИЧЕСКАЯ ТЕОРИЯ ПРИНЯТИЯ РЕШЕНИЙ

---

Сергей Николенко

TRA Robotics — Санкт-Петербург  
22 апреля 2018 г.

---

*Random facts:*

Есть человек, и день его рождения  
Сама природа празднует у нас.  
Неодолимой силой пробужденья  
Она крушит последний лед и наст...

Есть человек, он зорек и отважен,  
Он воплощенье светлого всего,  
Мы вечно совершенствуемся, даже  
Не называя имени его.

Сергей Смирнов

# СТАТИСТИЧЕСКАЯ ТЕОРИЯ ПРИНЯТИЯ РЕШЕНИЙ

---

- Сейчас мы попытаемся понять, что же на самом деле происходит в этих методах.
- Начнём с обычной регрессии – непрерывный вещественный вход  $\mathbf{x} \in \mathbb{R}^p$ , непрерывный вещественный выход  $y \in \mathbb{R}$ ; у них есть некоторое совместное распределение  $p(\mathbf{x}, y)$ .
- Мы хотим найти функцию  $f(\mathbf{x})$ , которая лучше всего предсказывает  $y$ .

- Введём *функцию потерь* (loss function)  $L(y, f(\mathbf{x}))$ , которая наказывает за ошибки; естественно взять квадратичную функцию потерь

$$L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2.$$

- Тогда каждому  $f$  можно сопоставить *ожидаемую ошибку предсказания* (expected prediction error):

$$\text{EPE}(f) = \mathbb{E}(y - f(\mathbf{x}))^2 = \int \int (y - f(\mathbf{x}))^2 p(\mathbf{x}, y) d\mathbf{x} dy.$$

- И теперь самая хорошая функция предсказания  $\hat{f}$  – это та, которая минимизирует  $\text{EPE}(f)$ .

- Это можно переписать как

$$\text{EPE}(f) = E_{\mathbf{x}} E_{y|\mathbf{x}} [(y - f(\mathbf{x}))^2 | \mathbf{x}],$$

и, значит, можно теперь минимизировать EPE поточечно:

$$\hat{f}(\mathbf{x}) = \arg \min_c E_{y|\mathbf{x}'} [(y - c)^2 | \mathbf{x}' = \mathbf{x}],$$

а это можно решить и получить

$$\hat{f}(\mathbf{x}) = E_{y|\mathbf{x}'}(y | \mathbf{x}' = \mathbf{x}).$$

- Это решение называется *функцией регрессии* и является наилучшим предсказанием  $y$  в любой точке  $\mathbf{x}$ .

- Теперь мы можем понять, что такое  $k$ -NN.
- Давайте оценим это ожидание:

$$f(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}'}(y \mid \mathbf{x}' = \mathbf{x}).$$

- Оценка ожидания – это среднее всех  $y$  с данным  $\mathbf{x}$ . Конечно, у нас таких нету, поэтому мы приближаем это среднее как

$$\hat{f}(\mathbf{x}) = \text{Average}[y_i \mid \mathbf{x}_i \in N_k(\mathbf{x})].$$

- Это сразу два приближения: ожидание через среднее и среднее в точке через среднее в ближних точках.
- Иначе говоря,  $k$ -NN предполагает, что в окрестности  $\mathbf{x}$  функция  $y(\mathbf{x})$  не сильно меняется, а лучше всего – она кусочно-постоянна.

- А линейная регрессия – это модельный подход, мы предполагаем, что функция регрессии линейна от своих аргументов:

$$f(\mathbf{x}) \approx \mathbf{x}^T \mathbf{w}.$$

- Теперь мы не берём условие по  $\mathbf{x}$ , как в  $k$ -NN, а просто собираем много значений для разных  $\mathbf{x}$  и обучаем модель.

# КЛАССИФИКАЦИЯ

- То же самое можно и с задачей классификации сделать. Пусть у нас переменная  $g$  с  $K$  возможными значениями  $g_1, \dots, g_k$  предсказывается.
- Введём функцию потери, равную 1 за каждый неверный ответ. Получим

$$\text{EPE} = \mathbb{E} [L(g, \hat{g}(\mathbf{x}))].$$

- Перепишем как раньше:

$$\text{EPE} = \mathbb{E}_{\mathbf{x}} \sum_{k=1}^K L(g_k, \hat{g}(\mathbf{x})) p(g_k | \mathbf{x}).$$

- Опять достаточно оптимизировать поточечно:

$$\hat{g}(\mathbf{x}) = \arg \min_g \sum_{k=1}^K L(g_k, \hat{g}(\mathbf{x})) p(g_k | \mathbf{x}).$$



- Опять достаточно оптимизировать поточечно:

$$\hat{g}(\mathbf{x}) = \arg \min_g \sum_{k=1}^K L(g_k, \hat{g}(\mathbf{x})) p(g_k | \mathbf{x}).$$

- Для 0-1 функции потери это упрощается до

$$\hat{g}(\mathbf{x}) = \arg \min_g [1 - p(g | \mathbf{x})], \text{ т.е.}$$

$$\hat{g}(\mathbf{x}) = g_k, \text{ если } p(g_k | \mathbf{x}) = \max_g p(g | \mathbf{x}).$$

- Это называется *оптимальным байесовским классификатором*; если модель известна, то его обычно можно построить.

- Рассмотрим совместное распределение  $p(y, \mathbf{x})$  и квадратичную функцию потерь  $L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$ .
- Мы знаем, что тогда оптимальная оценка – это функция регрессии

$$\hat{f}(\mathbf{x}) = \mathbb{E}[y | \mathbf{x}] = \int yp(y | \mathbf{x})dx.$$

- Давайте подсчитаем ожидаемую ошибку и перепишем её в другой форме:

$$\begin{aligned} E[L] &= E[(y - f(\mathbf{x}))^2] = E[(y - E[y | \mathbf{x}] + E[y | \mathbf{x}] - f(\mathbf{x}))^2] = \\ &= \int (f(\mathbf{x}) - E[y | \mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} + \int (E[y | \mathbf{x}] - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy, \end{aligned}$$

потому что

$$\int (f(\mathbf{x}) - E[y | \mathbf{x}]) (E[y | \mathbf{x}] - y) p(\mathbf{x}, y) d\mathbf{x} dy = 0.$$

- Эта форма записи – разложение на bias-variance и noise:

$$E[L] = \int (f(\mathbf{x}) - E[y | \mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} + \int (E[y | \mathbf{x}] - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy,$$

- Отсюда, кстати, тоже сразу видно, что от  $f(\mathbf{x})$  зависит только первый член, и он минимизируется, когда

$$f(\mathbf{x}) = \hat{f}(\mathbf{x}) = E[y | \mathbf{x}].$$

- А noise,  $\int (E[y | \mathbf{x}] - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy$ , – это просто свойство данных, дисперсия шума.

- Если бы у нас был всемогущий компьютер и неограниченный датасет, мы бы, конечно, на этом и закончили, посчитали бы  $\hat{f}(\mathbf{x}) = \mathbb{E}[y \mid \mathbf{x}]$ , и всё.
- Однако жизнь – борьба, и у нас есть только ограниченный датасет из  $N$  точек. Предположим, что этот датасет берётся по распределению  $p(\mathbf{x}, y)$  – т.е. фактически рассмотрим много-много экспериментов такого вида:
  - взяли датасет  $D$  из  $N$  точек по распределению  $p(\mathbf{x}, y)$ ;
  - подсчитали нашу чудо-регрессию;
  - получили новую функцию предсказания  $f(\mathbf{x}; D)$ .
- Разные датасеты будут приводить к разным функциям предсказания...

- ...а потому давайте усредним теперь по датасетам.
- Наш первый член в ожидаемой ошибке выглядел как  $(f(\mathbf{x}) - \hat{f}(\mathbf{x}))^2$ , а теперь будет  $(f(\mathbf{x}; D) - \hat{f}(\mathbf{x}))^2$ , и его можно усреднить по  $D$ , применив такой же трюк:

$$\begin{aligned} & (f(\mathbf{x}; D) - \hat{f}(\mathbf{x}))^2 \\ &= (f(\mathbf{x}; D) - \mathbb{E}_D[f(\mathbf{x}; D)] + \mathbb{E}_D[f(\mathbf{x}; D)] - \hat{f}(\mathbf{x}))^2 \\ &= (f(\mathbf{x}; D) - \mathbb{E}_D[f(\mathbf{x}; D)])^2 + (\mathbb{E}_D[f(\mathbf{x}; D)] - \hat{f}(\mathbf{x}))^2 + 2(\dots)(\dots), \end{aligned}$$

и в ожидании получится...

- ...и в ожидании получится

$$\begin{aligned} \mathbb{E}_D \left[ (f(\mathbf{x}; D) - \hat{f}(\mathbf{x}))^2 \right] &= \\ &= \mathbb{E}_D \left[ (f(\mathbf{x}; D) - \mathbb{E}_D [f(\mathbf{x}; D)])^2 \right] + \left( \mathbb{E}_D [f(\mathbf{x}; D)] - \hat{f}(\mathbf{x}) \right)^2. \end{aligned}$$

- Разложили на дисперсию  $\mathbb{E}_D \left[ (f(\mathbf{x}; D) - \mathbb{E}_D [f(\mathbf{x}; D)])^2 \right]$  и квадрат систематической ошибки  $\left( \mathbb{E}_D [f(\mathbf{x}; D)] - \hat{f}(\mathbf{x}) \right)^2$ ; это и есть bias-variance decomposition.

Expected loss = (bias)<sup>2</sup> + variance + noise,

где

$$(\text{bias})^2 = \left( \mathbb{E}_D [f(\mathbf{x}; D)] - \hat{f}(\mathbf{x}) \right)^2,$$

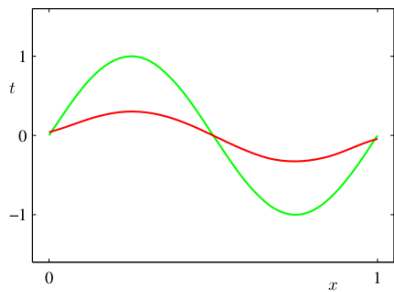
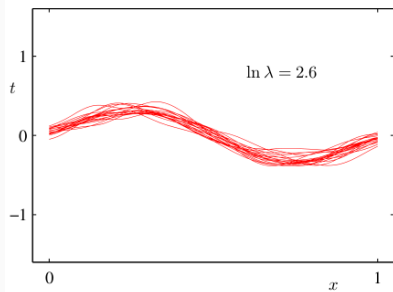
$$\text{variance} = \mathbb{E}_D \left[ (f(\mathbf{x}; D) - \mathbb{E}_D [f(\mathbf{x}; D)])^2 \right],$$

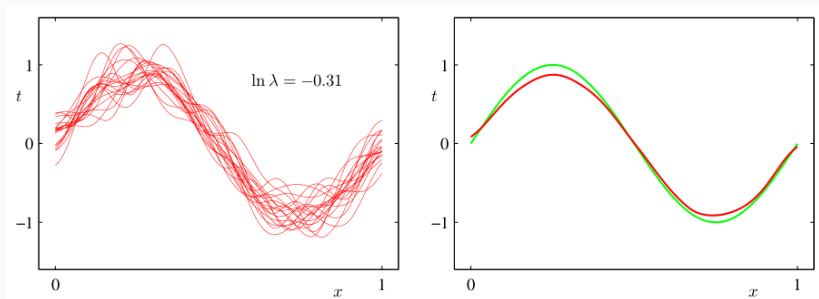
$$\text{noise} = \int (\mathbb{E}[y | \mathbf{x}] - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy.$$



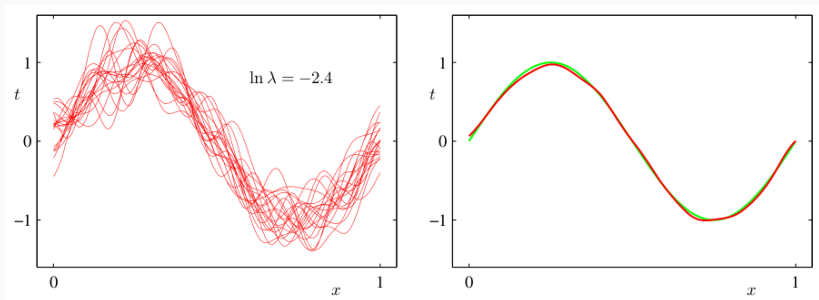
- Теперь давайте посмотрим на пример: опять та же синусоида, опять приближаем её линейной регрессией с полиномиальными признаками (максимальным их числом).
- И мы регуляризуем эту регрессию с параметром  $\alpha$ .
- Будем набрасывать много датасетов и смотреть, что меняется при этом.

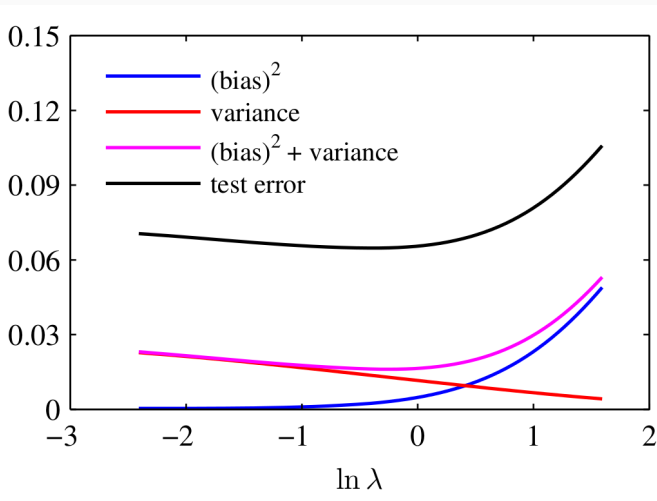
# РЕГУЛЯРИЗАТОР И BIAS-VARIANCE





# РЕГУЛЯРИЗАТОР И BIAS-VARIANCE





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  $\mu_N = \beta \Sigma_N \Phi^\top \mathbf{t}$ ):

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

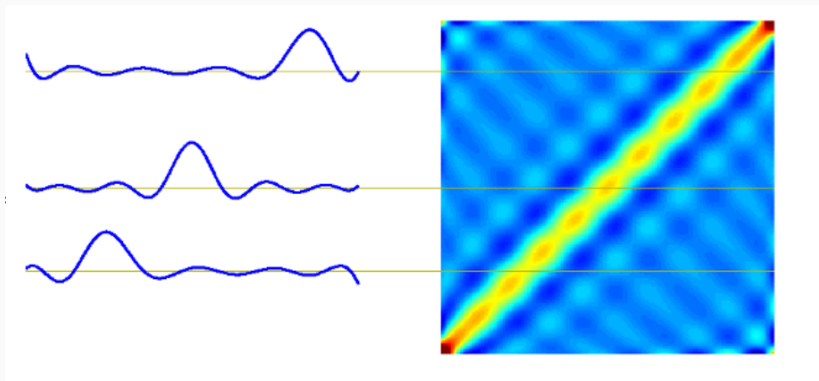
- $y(\mathbf{x}, \mu_N) = \sum_{n=1}^N \beta \phi(\mathbf{x})^\top \Sigma_N \phi(\mathbf{x}_n) t_n$ .
- Это значит, что предсказание можно переписать как

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

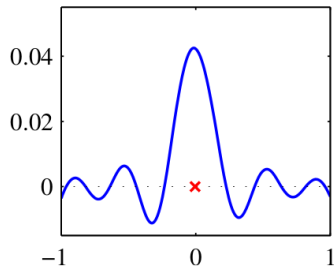
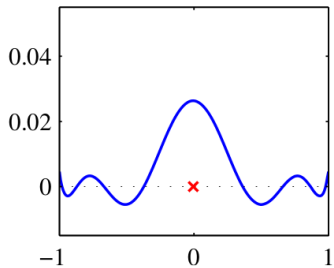
- Т.е. мы предсказываем следующую точку как линейную комбинацию значений в известных точках.
- Функция  $k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^\top \Sigma_N \phi(\mathbf{x}')$  называется *эквивалентным ядром* (equivalent kernel).



# 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

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

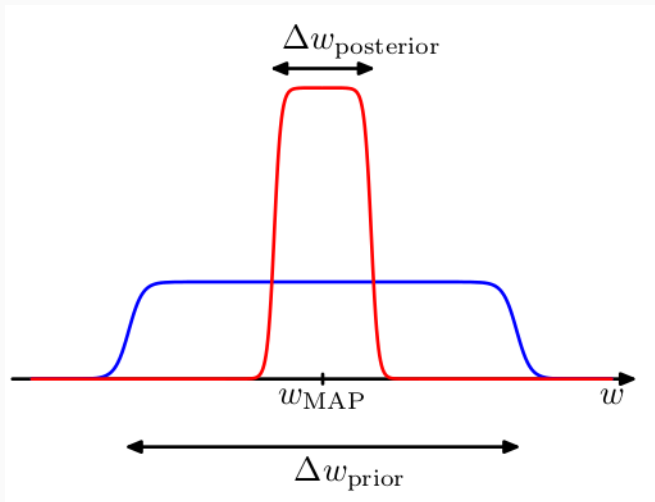
- Т.е. это вероятность сгенерировать  $D$ , если выбрать параметры модели по её априорному распределению, а потом накидывать данные.
- Это, кстати, в точности знаменатель из теоремы Байеса:

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

- 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 $p(D)$



## ПРИБЛИЖЕНИЕ $p(D)$

- Тогда получится

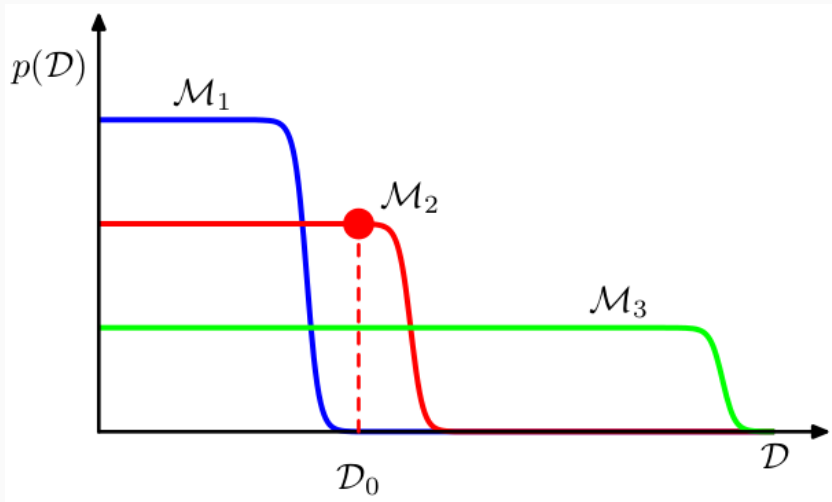
$$p(D) = \int p(D | w)p(w)dw \approx p(D | w_{\text{MAP}}) \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}},$$
$$\ln p(D) \approx \ln p(D | w_{\text{MAP}}) + \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right).$$

- 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- $K$ :

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

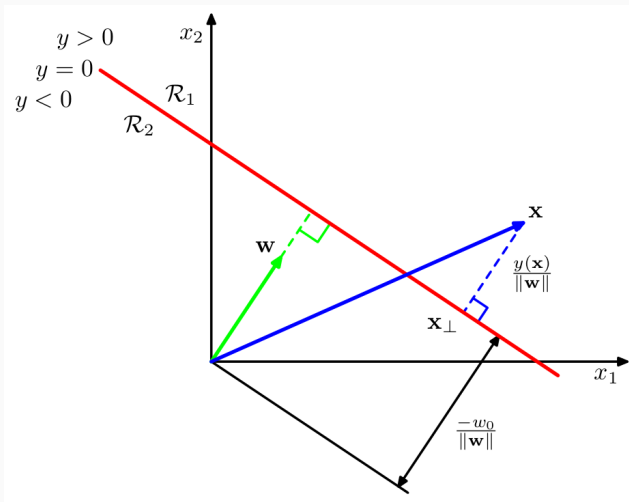
- Тоже можно интерпретировать как вероятности – или пропорционально им.

- Начнём с геометрии: рассмотрим линейную дискриминантную функцию

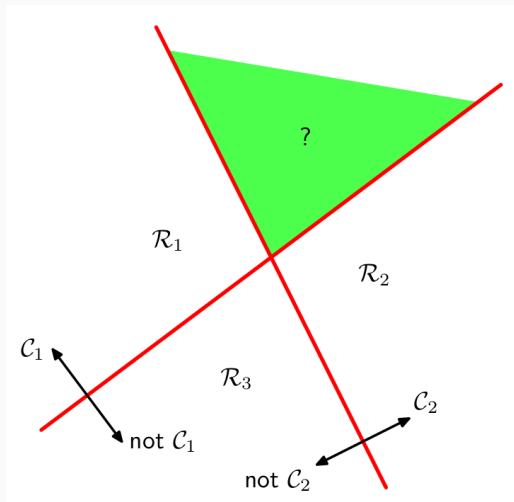
$$y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0.$$

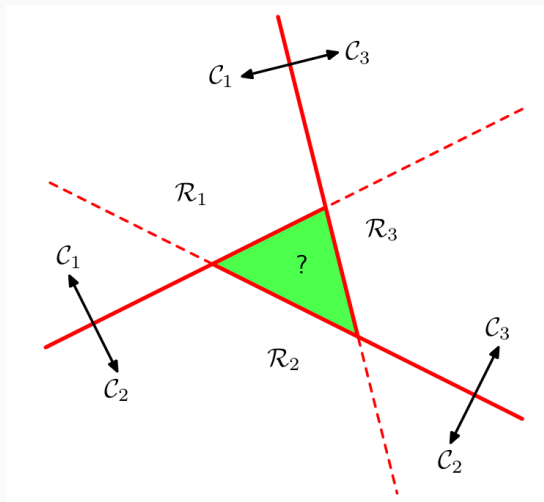
- Это гиперплоскость, и  $\mathbf{w}$  – нормаль к ней.
- Расстояние от начала координат до гиперплоскости равно  $\frac{-w_0}{\|\mathbf{w}\|}$ .
- $y(\mathbf{x})$  связано с расстоянием до гиперплоскости:  $d = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$ .

# 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.



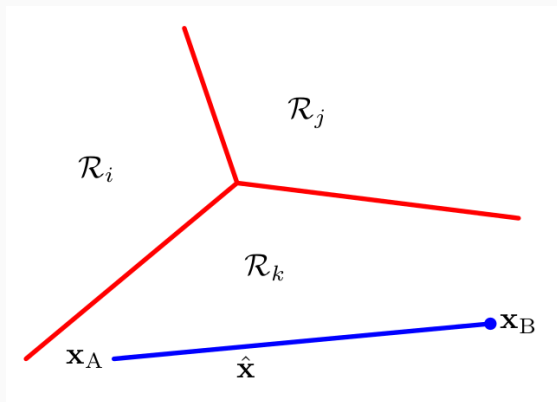


- Лучше рассмотреть единый дискриминант из  $K$  линейных функций:

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

- 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:

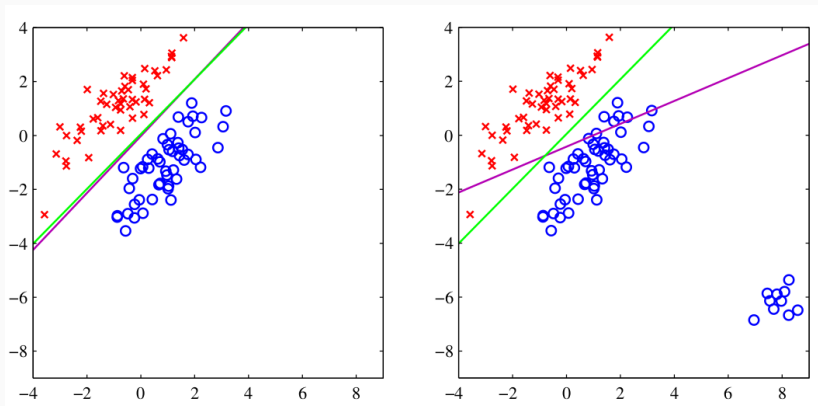
$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^\top \mathbf{x} = \mathbf{T}^\top (\mathbf{X}^\dagger)^\top \mathbf{x}.$$

- Это решение сохраняет линейность.

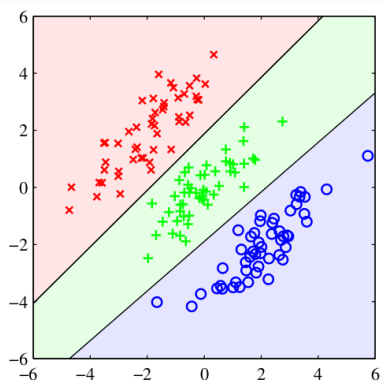
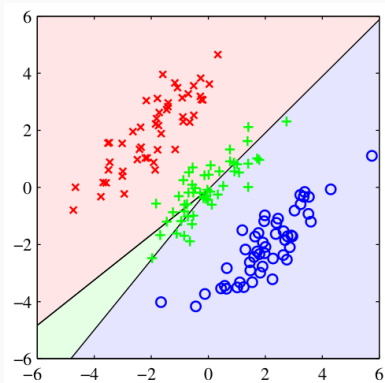
**Упражнение.** Докажите, что в схеме кодирования 1-of- $K$  предсказания  $y_k(\mathbf{x})$  для разных классов при любом  $\mathbf{x}$  будут давать в сумме 1. Почему они всё-таки не будут разумными оценками вероятностей?

- 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.



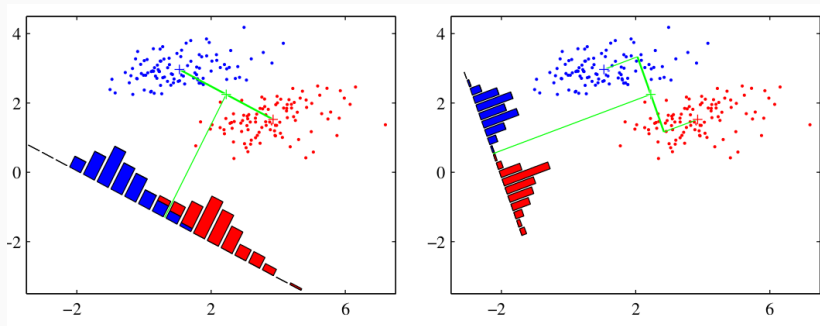
- Рассмотрим два класса  $\mathcal{C}_1$  и  $\mathcal{C}_2$  с  $N_1$  и  $N_2$  точками.
- Первая идея – надо найти серединный перпендикуляр между центрами кластеров

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

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



Спасибо за внимание!