

КЛАССИФИКАЦИЯ

Сергей Николенко

СПбГУ — Санкт-Петербург

29 сентября 2017 г.

Random facts:

- 29 сентября 61 г. до н.э. Гней Помпей отметил свой третий триумф за победу над пиратами, конец Митридатовых войн, а также свой 45-й день рождения
- 29 сентября 1922 г. из Петрограда отплыл пароход «Обербургомистр Хакен» с Бердяевым, Франком и Ильиным на борту -- знаменитый «философский пароход»
- 29 сентября 1897 г. в Буэнос-Айресе на премьере спектакля «Креольский суд» было впервые исполнено танго

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).

- 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 стараемся, чтобы было можно).
- Если несколько классов – удобно 1-of- K :

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

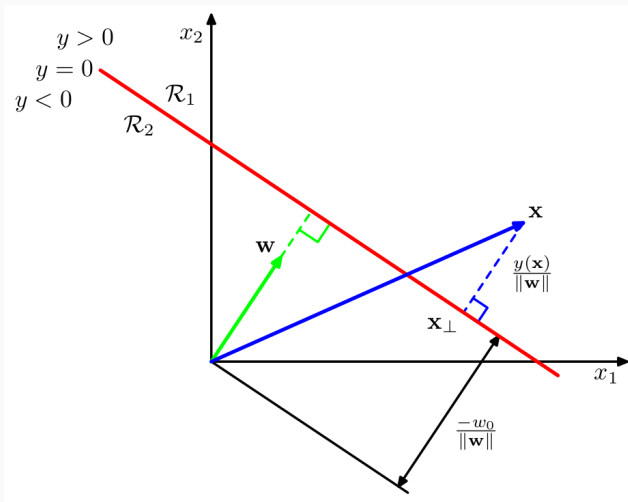
- Тоже можно интерпретировать как вероятности – или пропорционально им.

- Начнём с геометрии: рассмотрим линейную дискриминантную функцию

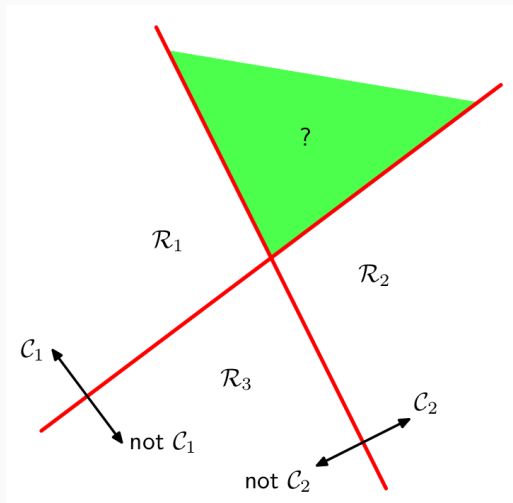
$$y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0.$$

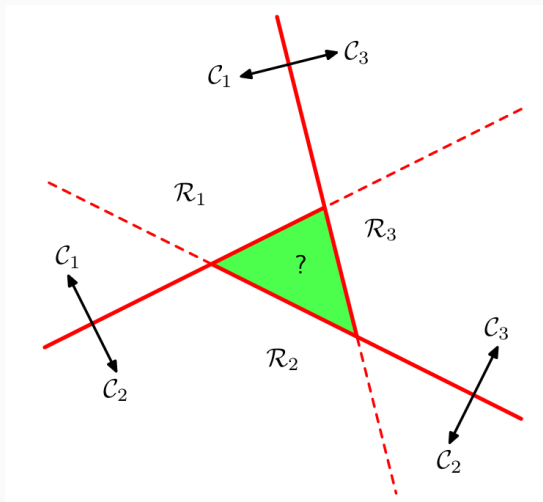
- Это гиперплоскость, и \mathbf{w} – нормаль к ней.
- Расстояние от начала координат до гиперплоскости равно $\frac{-w_0}{\|\mathbf{w}\|}$.
- $y(\mathbf{x})$ связано с расстоянием до гиперплоскости: $d = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$.

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.



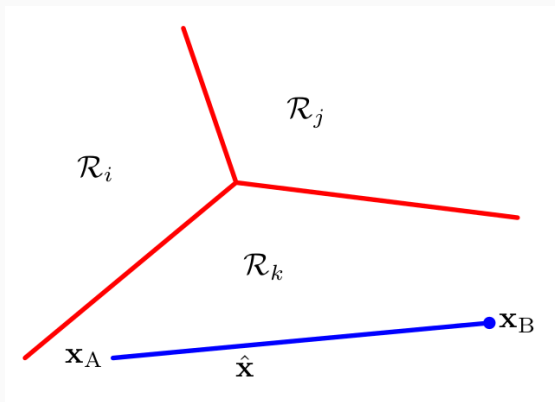


- Лучше рассмотреть единый дискриминант из K линейных функций:

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

- 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

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

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

- Теперь можно найти и дискриминантную функцию:

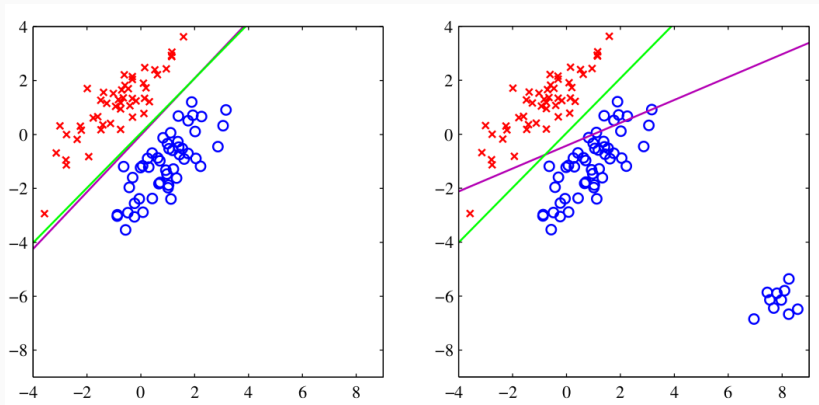
$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^T \mathbf{x} = \mathbf{T}^T (\mathbf{X}^\dagger)^T \mathbf{x}.$$

- Это решение сохраняет линейность.

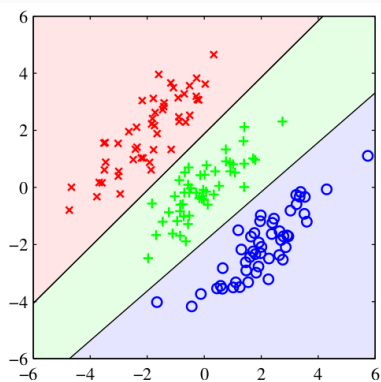
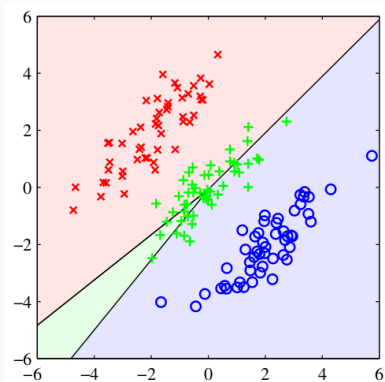
Упражнение. Докажите, что в схеме кодирования 1-of- K предсказания $y_k(\mathbf{x})$ для разных классов при любом \mathbf{x} будут давать в сумме 1. Почему они всё-таки не будут разумными оценками вероятностей?

- 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.

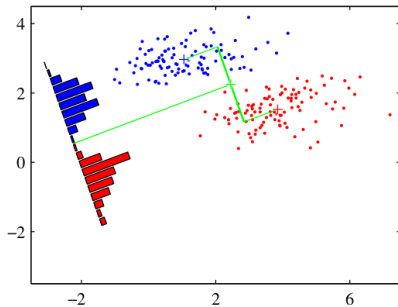
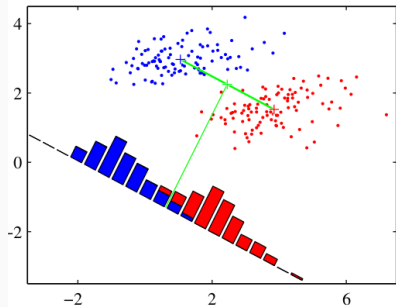
- Рассмотрим два класса \mathcal{C}_1 и \mathcal{C}_2 с N_1 и N_2 точками.
- Первая идея – надо найти серединный перпендикуляр между центрами кластеров

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

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} ...

ЛИНЕЙНЫЙ ДИСКРИМИНАНТ ФИШЕРА

- ...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!