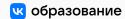
КЛАССИФИКАЦИЯ

Сергей Николенко



Академия больших данных MADE — VK 04 марта 2022 г.

Random facts:

- 04 марта 1733 г. был издан указ Анны Иоанновны «Об учреждении полиции в городах», согласно которому в 23 городах России были созданы полицмейстерские конторы
- 04 марта 1762 г., ровно через 29 лет после этого, император Пётр III подписал указ «О свободной для всех торговле»
- 04 марта 1803 г., ещё через 39 лет, император Александр I издал «Указ о вольных хлебопашцах», по которому землевладельцы получили право освобождать крестьян с обязательным наделением их землёй
- 04 марта 1837 г., ещё через 34 года и через 103 года после первого факта, Михаил Лермонтов был арестован за стихотворение «Смерть поэта»
- 04 марта 1990 г. на 5-летний срок был избран Съезд народных депутатов РСФСР, высший законодательный орган РСФСР; он был распущен указом Бориса Ельцина 21 сентября 1993 г. и разогнан вооружённой силой 4 октября 1993 г.
- 04 марта 2012 г. прошли выборы президента Российской Федерации; Владимир Путин был избран президентом России на третий срок

BIAS-VARIANCE-NOISE DECOMPOSI-

TION

- 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

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

 Давайте подсчитаем ожидаемую ошибку и перепишем её в другой форме:

$$\begin{split} \mathbf{E}[L] &= \mathbf{E}[(y - f(\mathbf{x}))^2] = \mathbf{E}[(y - \mathbf{E}\left[y \mid \mathbf{x}\right] + \mathbf{E}\left[y \mid \mathbf{x}\right] - f(\mathbf{x}))^2] = \\ &= \int \left(f(\mathbf{x}) - \mathbf{E}\left[y \mid \mathbf{x}\right]\right)^2 p(\mathbf{x}) d\mathbf{x} + \int \left(\mathbf{E}\left[y \mid \mathbf{x}\right] - y\right)^2 p(\mathbf{x}, y) d\mathbf{x} dy, \end{split}$$
 notomy yto

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

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

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

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

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

• A noise, $\int \left(\mathbf{E}\left[y\mid\mathbf{x}\right]-y\right)^2p(\mathbf{x},y)d\mathbf{x}dy$, – 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 $\hat{f}(\mathbf{x}) = \mathbf{E}\left[y \mid \mathbf{x}\right]$, 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:
 - \cdot 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 $\left(f(\mathbf{x})-\hat{f}(\mathbf{x})\right)^2$, а теперь будет $\left(f(\mathbf{x};D)-\hat{f}(\mathbf{x})\right)^2$, и его можно усреднить по D, применив такой же трюк:

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

и в ожидании получится...

• ...и в ожидании получится

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

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

BIAS-VARIANCE-NOISE

$${\sf Expected\ loss} = ({\sf bias})^2 + {\sf variance} + {\sf noise},$$

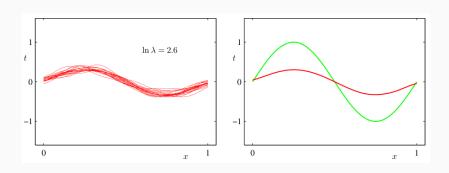
где

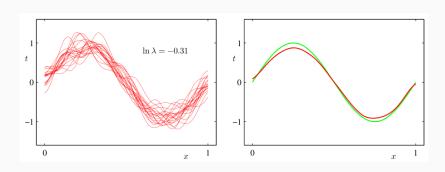
$$\begin{split} (\mathrm{bias})^2 &= \left(\mathbf{E}_D\left[f(\mathbf{x};D)\right] - \hat{f}(\mathbf{x})\right)^2, \\ \mathrm{variance} &= \mathbf{E}_D\left[\left(f(\mathbf{x};D) - \mathbf{E}_D\left[f(\mathbf{x};D)\right]\right)^2\right], \\ \mathrm{noise} &= \int \left(\mathbf{E}\left[y\mid \mathbf{x}\right] - y\right)^2 p(\mathbf{x},y) d\mathbf{x} dy. \end{split}$$

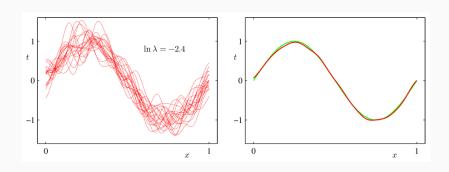
/ı

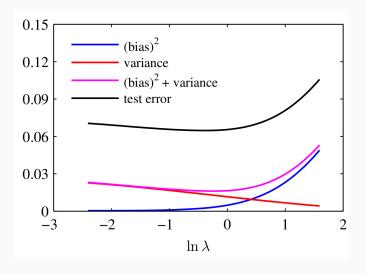
ПРИМЕР

- 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 lpha.
- Будем набрасывать много датасетов и смотреть, что меняется при этом.









Введение в классификацию

Задача классификации

- Теперь классификация: определить вектор ${\bf x}$ в один из K классов 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 $C_1,\,t=1$ соответствует $C_2.$
- Оценку t можно интерпретировать как вероятность (по крайней мере, мы постараемся, чтобы было можно).
- Если несколько классов удобно 1-of-*K*:

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

• Тоже можно интерпретировать как вероятности – или пропорционально им.

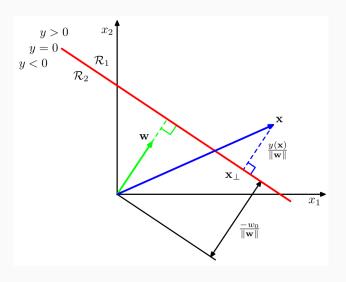
Разделяющая гиперплоскость

• Начнём с геометрии: рассмотрим линейную дискриминантную функцию

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

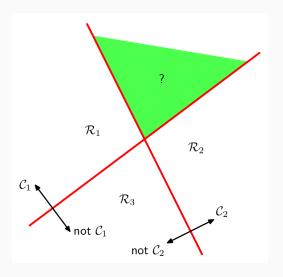
- \cdot Это гиперплоскость, и ${f w}$ нормаль к ней.
- Расстояние от начала координат до гиперплоскости равно $\frac{-w_0}{\|\mathbf{w}\|}$.
- $y(\mathbf{x})$ связано с расстоянием до гиперплоскости: $d=\frac{y(\mathbf{x})}{\|\mathbf{w}\|}.$

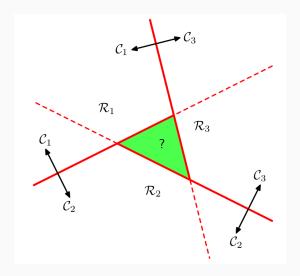
Разделяющая гиперплоскость



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

НЕСКОЛЬКО КЛАССОВ



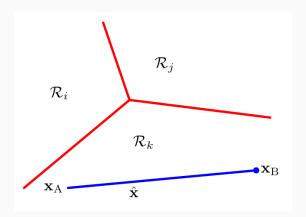


• Лучше рассмотреть единый дискриминант из K линейных функций:

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

- · Классифицировать в C_k , если $y_k(\mathbf{x})$ максимален.
- · Тогда разделяющая поверхность между C_k и C_j будет гиперплоскостью вида $y_k(\mathbf{x})=y_j(\mathbf{x})$:

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



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

МЕТОД НАИМЕНЬШИХ КВАДРАТОВ

• 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 W, оптимизируя сумму квадратов; функция ошибки:

$$E_D(\mathbf{W}) = \frac{1}{2} \mathrm{Tr} \left[\left(\mathbf{X} \mathbf{W} - \mathbf{T} \right)^\top \left(\mathbf{X} \mathbf{W} - \mathbf{T} \right) \right].$$

• Берём производную, решаем...

МЕТОД НАИМЕНЬШИХ КВАДРАТОВ

• ...получается привычное

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

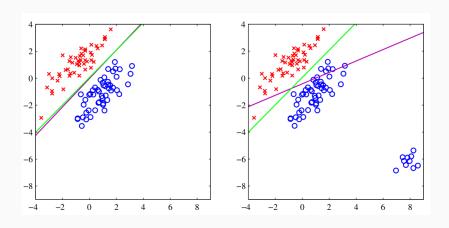
где \mathbf{X}^{\dagger} – псевдообратная Мура-Пенроуза.

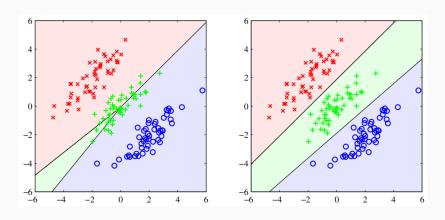
• Теперь можно найти и дискриминантную функцию:

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

МЕТОД НАИМЕНЬШИХ КВАДРАТОВ

- 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.





• 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.

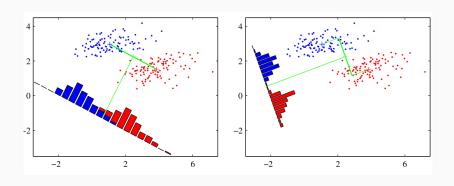
- Другой взгляд на классификацию: в линейном случае мы хотим спроецировать точки в размерность 1 (на нормаль разделяющей гиперплоскости) так, чтобы в этой размерности 1 они хорошо разделялись.
- Т.е. классификация это такой метод радикального сокращения размерности.
- Давайте посмотрим на классификацию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

- Рассмотрим два класса C_1 и C_2 с N_1 и N_2 точками.
- Первая идея надо найти серединный перпендикуляр между центрами кластеров

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{C_1} \mathbf{x}, \text{ if } \mathbf{m}_2 = \frac{1}{N_2} \sum_{C_2} \mathbf{x},$$

т.е. максимизировать $\mathbf{w}^{ op}\left(\mathbf{m}_{2}-\mathbf{m}_{1}
ight)$.

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



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 $y_n = \mathbf{w}^{ op} \mathbf{x}_n$

$$s_1 = \sum_{n \in C_1} \left(y_n - m_1 \right)^2 \text{ if } s_1 = \sum_{n \in C_2} \left(y_n - m_2 \right)^2.$$

• Критерий Фишера:

$$\begin{split} J(\mathbf{w}) &= \frac{\left(m_2 - m_1\right)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}, \text{ где} \\ \mathbf{S}_B &= \left(\mathbf{m}_2 - \mathbf{m}_1\right) \left(\mathbf{m}_2 - \mathbf{m}_1\right)^\top, \\ \mathbf{S}_W &= \sum_{n \in C_1} \left(\mathbf{x}_n - \mathbf{m}_1\right) \left(\mathbf{x}_n - \mathbf{m}_1\right)^\top + \sum_{n \in C_2} \left(\mathbf{x}_n - \mathbf{m}_2\right) \left(\mathbf{x}_n - \mathbf{m}_2\right)^\top. \end{split}$$

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

• Дифференцируя по w...

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

$$\left(\mathbf{w}^{\top}\mathbf{S}_{B}\mathbf{w}\right)\mathbf{S}_{W}\mathbf{w}=\left(\mathbf{w}^{\top}\mathbf{S}_{W}\mathbf{w}\right)\mathbf{S}_{B}\mathbf{w}.$$

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

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

• 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 C_1 выберем целевое значение $\frac{N_1+N_2}{N_1}$, а для класса 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

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

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

$$\begin{split} \mathbf{S}_T &= \sum_n \left(\mathbf{x}_n - \mathbf{m}\right) \left(\mathbf{x}_n - \mathbf{m}\right)^\top, \\ \mathbf{S}_B &= \mathbf{S}_T - \mathbf{S}_W. \end{split}$$

• Обобщить критерий можно разными способами, например:

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

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

$$\begin{split} \mathbf{s}_W &= \sum_{k=1}^K \sum_{n \in C_k} \left(\mathbf{y}_n - \boldsymbol{\mu}_k\right) \left(\mathbf{y}_n - \boldsymbol{\mu}_k\right)^\top, \\ \mathbf{s}_B &= \sum_{k=1}^K N_k \left(\boldsymbol{\mu}_k - \boldsymbol{\mu}\right) \left(\boldsymbol{\mu}_k - \boldsymbol{\mu}\right)^\top, \end{split}$$

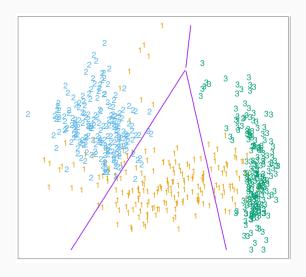
где
$$\mu_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}_n$$
.

LDA и QDA

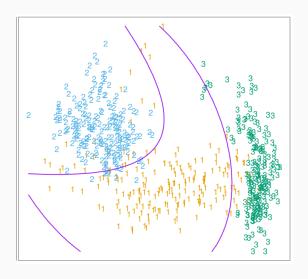
Нелинейные поверхности

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

Нелинейные поверхности



Нелинейные поверхности



ГЕНЕРАТИВНЫЕ МОДЕЛИ

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

$$p(C_1 \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid C_1)p(C_1)}{p(\mathbf{x} \mid C_1)p(C_1) + p(\mathbf{x} \mid C_2)p(C_2)}.$$

ГЕНЕРАТИВНЫЕ МОДЕЛИ

• Перепишем:

$$\begin{split} p(C_1 \mid \mathbf{x}) &= \frac{p(\mathbf{x} \mid C_1) p(C_1)}{p(\mathbf{x} \mid C_1) p(C_1) + p(\mathbf{x} \mid C_2) p(C_2)} = \frac{1}{1 + e^{-a}} = \sigma(a), \end{split}$$
 где
$$a &= \ln \frac{p(\mathbf{x} \mid C_1) p(C_1)}{p(\mathbf{x} \mid C_2) p(C_2)}, \qquad \sigma(a) = \frac{1}{1 + e^{-a}}.$$

ГЕНЕРАТИВНЫЕ МОДЕЛИ

• $\sigma(a)$ – 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.

Несколько классов

• В случае нескольких классов получится

$$p(C_k \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid C_k)p(C_k)}{\sum_j p(\mathbf{x} \mid C_j)p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}}.$$

- · Здесь $a_k = \ln p(\mathbf{x} \mid C_k) p(C_k).$
- · $\frac{e^{a_k}}{\sum_j e^{a_j}}$ нормализованная экспонента, или softmax-функция (сглаженный максимум).

ПРИМЕР

• Давайте рассмотрим гауссовы распределения для классов:

$$p(\mathbf{x} \mid C_k) = N(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}).$$

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

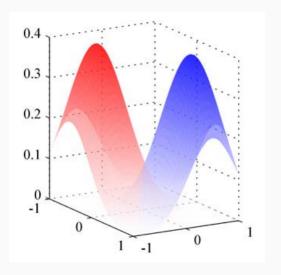
• ...получится

$$p(C_1 \mid \mathbf{x}) = \sigma(\mathbf{w}^{\top}\mathbf{x} + w_0)$$
, где
$$\mathbf{w} = \Sigma^{-1}\left(\mu_1 - \mu_2\right),$$

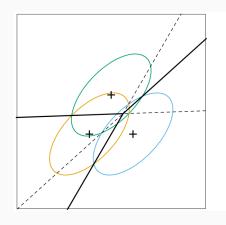
$$w_0 = -\frac{1}{2}\mu_1^{\top}\Sigma^{-1}\mu_1 + \frac{1}{2}\mu_2^{\top}\Sigma^{-1}\mu_2 + \ln\frac{p(C_1)}{p(C_2)}.$$

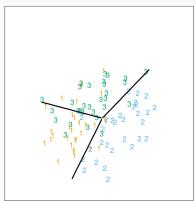
• Т.е. в аргументе сигмоида получается линейная функция от ${\bf x}$. Поверхности уровня – это когда $p(C_1 \mid {\bf x})$ постоянно, т.е. гиперплоскости в пространстве ${\bf x}$. Априорные вероятности $p(C_k)$ просто сдвигают эти гиперплоскости.

Разделяющая гиперплоскость



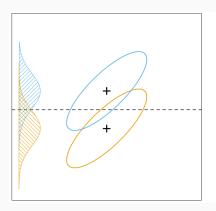
Разделяющая гиперплоскость

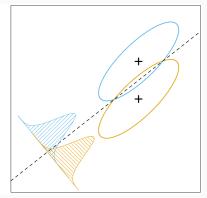




Дискриминант Фишера

Кстати, с дискриминантом Фишера эта разделяющая поверхность отлично сходится.





Несколько классов

• С несколькими классами получится тоже примерно так же:

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

где
$$\pi_k = p(C_k).$$

- 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 $p(\mathbf{x} \mid C_k)$, если даны только данные?
- Можно по методу максимального правдоподобия.
- \cdot 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 C_1 , $t_n=0$ значит C_2 .
- · Обозначим $p(C_1)=\pi$, $p(C_2)=1-\pi$.

• Для одной точки в классе C_1 :

$$p(\mathbf{x}_n, C_1) = p(C_1)p(\mathbf{x}_n \mid C_1) = \pi N(\mathbf{x}_n \mid \mu_1, \Sigma).$$

• В классе C_2 :

$$p(\mathbf{x}_n,C_2) = p(C_2)p(\mathbf{x}_n \mid C_2) = (1-\pi)N(\mathbf{x}_n \mid \mu_2,\Sigma).$$

• Функция правдоподобия:

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

• Максимизируем логарифм правдоподобия. Сначала по π , там останется только

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

и, взяв производную, получим, совершенно неожиданно,

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

 \cdot Теперь по μ_1 ; всё, что зависит от μ_1 :

$$\sum_n t_n \ln N(\mathbf{x}_n \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = -\frac{1}{2} \sum_n t_n \left(\mathbf{x}_n - \boldsymbol{\mu}_1 \right)^\top \boldsymbol{\Sigma}^{-1} \left(\mathbf{x}_n - \boldsymbol{\mu}_1 \right) + C.$$

• 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

$$\begin{split} \hat{\Sigma} &= \frac{N_1}{N_1 + N_2} \mathbf{S}_1 + \frac{N_2}{N_1 + N_2} \mathbf{S}_2, \text{ где} \\ \mathbf{S}_1 &= \frac{1}{N_1} \sum_{n \in C_1} \left(\mathbf{x}_n - \boldsymbol{\mu}_1 \right) \left(\mathbf{x}_n - \boldsymbol{\mu}_1 \right)^\top, \\ \mathbf{S}_2 &= \frac{1}{N_2} \sum_{n \in C_2} \left(\mathbf{x}_n - \boldsymbol{\mu}_2 \right) \left(\mathbf{x}_n - \boldsymbol{\mu}_2 \right)^\top. \end{split}$$

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

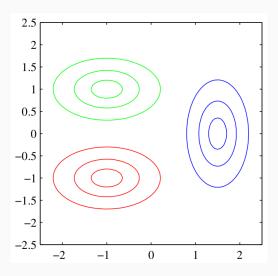
• Это самым прямым образом обобщается на случай нескольких классов. Упражнение. Сделайте это.

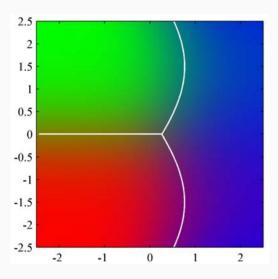
- А вот с разными матрицами ковариаций уже будет по-другому.
- Квадратичные члены не сократятся.
- Разделяющие поверхности станут квадратичными; QDA quadratic discriminant analysis.

· В QDA получится

$$\delta_k(\mathbf{x}) = -\frac{1}{2}\log|\Sigma_k| - \frac{1}{2}\left(\mathbf{x} - \boldsymbol{\mu}_k\right)^{-1}\Sigma_k^{-1}\left(\mathbf{x} - \boldsymbol{\mu}_k\right) + \log\pi_k.$$

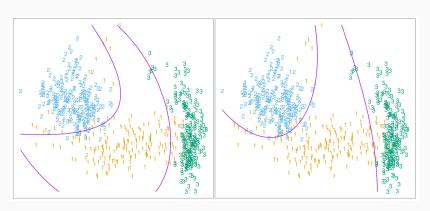
- Разделяющая поверхность между C_i и C_j это $\{\mathbf{x} \mid \delta_i(\mathbf{x}) = \delta_j(\mathbf{x})\}.$
- Оценки максимального правдоподобия такие же, только надо отдельно матрицы ковариаций оценивать.





LDA vs. QDA

Разница между LDA с квадратичными членами и QDA обычно невелика.



LDA vs. QDA

- 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.

LDA vs. QDA

- 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:

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

Снижение ранга в LDA

- 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.

Снижение ранга в LDA

- 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

В прошлый раз

• Итак, мы рассмотрели логистический сигмоид:

$$\begin{split} p(C_1 \mid \mathbf{x}) &= \frac{p(\mathbf{x} \mid C_1) p(C_1)}{p(\mathbf{x} \mid C_1) p(C_1) + p(\mathbf{x} \mid C_2) p(C_2)} = \frac{1}{1 + e^{-a}} = \sigma(a), \end{split}$$
 где $a = \ln \frac{p(\mathbf{x} \mid C_1) p(C_1)}{p(\mathbf{x} \mid C_2) p(C_2)}, \qquad \sigma(a) = \frac{1}{1 + e^{-a}}.$

 Вывели из него LDA и QDA, обучив их методом максимального правдоподобия.

Два класса

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

$$p(C_1 \mid \phi) = y(\phi) = \sigma(\mathbf{w}^\top \phi), \quad p(C_2 \mid \phi) = 1 - p(C_1 \mid \phi).$$

• Логистическая регрессия – это когда мы напрямую оптимизируем ${f w}$.

Два класса

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

$$p(\mathbf{t}\mid\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}, \quad y_n = p(C_1\mid\boldsymbol{\phi}_n).$$

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

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

Два класса

• Пользуясь тем, что $\sigma' = \sigma(1-\sigma)$, берём градиент (похоже на перцептрон):

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n.$$

- Если теперь сделать градиентный спуск, получим как раз разделяющую поверхность.
- Заметим, правда, что если данные действительно разделимы, то может получиться жуткий оверфиттинг: $\|\mathbf{w}\| \to \infty$, и сигмоид превращается в функцию Хевисайда. Надо регуляризовать.

IRLS

- В логистической регрессии не получается замкнутого решения из-за сигмоида.
- Но функция $E(\mathbf{w})$ всё равно выпуклая, и можно воспользоваться методом Ньютона-Рапсона на каждом шаге использовать локальную квадратичную аппроксимацию к функции ошибки:

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \mathbf{H}^{-1} \nabla E(\mathbf{w}),$$

где ${\bf H}$ (Hessian) – матрица вторых производных $E({\bf w}).$

IRLS

• Замечание: давайте применим Ньютона-Рапсона к обычной линейной регрессии с квадратической ошибкой:

$$\begin{split} \nabla E(\mathbf{w}) &= \sum_{n=1}^{N} \left(\mathbf{w}^{\top} \boldsymbol{\phi}_{n} - t_{n} \right) \boldsymbol{\phi}_{n} = \boldsymbol{\Phi}^{\top} \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Phi}^{\top} \mathbf{t}, \\ \nabla \nabla E(\mathbf{w}) &= \sum_{n=1}^{N} \boldsymbol{\phi}_{n} \boldsymbol{\phi}_{n}^{\top} = \boldsymbol{\Phi}^{\top} \boldsymbol{\Phi}, \end{split}$$

и шаг оптимизации будет

$$\begin{split} \mathbf{w}^{\text{new}} &= \mathbf{w}^{\text{old}} - \left(\boldsymbol{\Phi}^{\top} \boldsymbol{\Phi}\right)^{-1} \left[\boldsymbol{\Phi}^{\top} \boldsymbol{\Phi} \mathbf{w}^{\text{old}} - \boldsymbol{\Phi}^{\top} \mathbf{t}\right] = \\ &= \left(\boldsymbol{\Phi}^{\top} \boldsymbol{\Phi}\right)^{-1} \boldsymbol{\Phi}^{\top} \mathbf{t}, \end{split}$$

т.е. мы за один шаг придём к решению.

• Для логистической регрессии:

$$\begin{split} \nabla E(\mathbf{w}) &= \sum_{n=1}^{N} \left(y_n - t_n\right) \phi_n = \boldsymbol{\Phi}^\top \left(\mathbf{y} - \mathbf{t}\right), \\ \mathbf{H} &= \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} y_n (1 - y_n) \phi_n \boldsymbol{\phi}_n^\top = \boldsymbol{\Phi}^\top R \boldsymbol{\Phi} \end{split}$$

для диагональной матрицы R с $R_{nn}=y_n(1-y_n).$

• Формула шага оптимизации:

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \left(\boldsymbol{\Phi}^{\top} \boldsymbol{R} \boldsymbol{\Phi}\right)^{-1} \boldsymbol{\Phi}^{\top} \left(\mathbf{y} - \mathbf{t}\right) = \left(\boldsymbol{\Phi}^{\top} \boldsymbol{R} \boldsymbol{\Phi}\right)^{-1} \boldsymbol{\Phi}^{\top} \boldsymbol{R} \mathbf{z},$$

где
$$\mathbf{z} = \Phi \mathbf{w}^{\text{old}} - R^{-1} (\mathbf{y} - \mathbf{t}).$$

- Получилось как бы решение взвешенной задачи минимизации квадратического отклонения с матрицей весов R.
- · Отсюда название: iterative reweighted least squares (IRLS).

Спасибо!

Спасибо за внимание!