## SVM, ЯДЕРНЫЕ МЕТОДЫ И RVM

Сергей Николенко СПбГУ— Санкт-Петербург 29 февраля 2024 г.





#### Random facts:

- 29 февраля в Шотландии с 1288 г. был объявлен днём, когда женщина может предложить брак мужчине; нужно было, чтобы «всякая дама, которая идёт свататься, надевала нижнюю сорочку из багряной фланели и чтобы кромка её была хорошо видна, иначе мужчине придётся заплатить за неё штраф» в 1 фунт
- 29 февраля 1504 г. Христофор Колумб сумел обмануть аборигенов Ямайки, отказывавшихся поставлять европейцам еду; Колумб знал, что в этот день будет лунное затмение, и объявил, что за такое поведение испанский бог лишит индейцев Луны; в итоге индейцы возобновили поставки, а Колумб вернул им Луну
- за 29 февраля 1712 г. в Швеции последовало 30 февраля; чтобы вернуться к юлианскому календарю после пропуска високосных дней (собирались постепенно перейти на григорианский, начали в 1700, но потом началась Северная война и стало не до того), Швеция в 1712 г. добавила к февралю два лишних дня вместо одного
- 29 февраля 1880 г. было закончено строительство тоннеля Готард через Альпы
- 29 февраля 1888 г. Герман Холлерит для ускорения обработки данных переписи населения изобрёл первую электрическую счётную машину — табулятор, работавший на перфокартах; с 1924 г. фирма Холлерита стала называться IBM

# линейной классификации \_\_\_\_\_\_

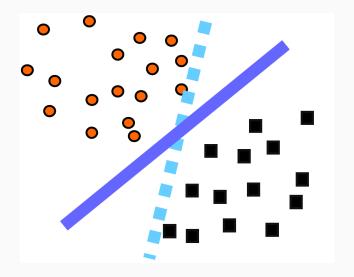
SVM и задача

#### Постановка задачи

- Метод опорных векторов решает задачу классификации.
- Каждый элемент данных точка в n-мерном пространстве  $\mathbb{R}^n$ .
- Формально: есть точки  $x_i$ , i=1..m, у точек есть метки  $y_i=\pm 1.$
- $\cdot$  Мы интересуемся: можно ли разделить данные (n-1)-мерной гиперплоскостью, а также хотим найти эту гиперплоскость.
- Это всё?

#### Постановка задачи

- Нет, ещё хочется научиться разделять этой гиперплоскостью как можно лучше.
- То есть желательно, чтобы два разделённых класса лежали как можно дальше от гиперплоскости.
- Практическое соображение: тогда от небольших возмущений в гиперплоскости ничего не испортится.

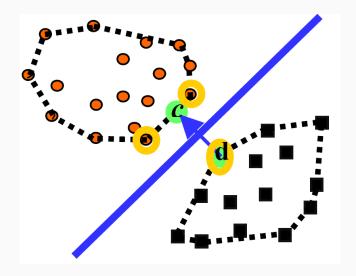


#### Выпуклые оболочки

- Один подход: найти две ближайшие точки в выпуклых оболочках данных, а затем провести разделяющую гиперплоскость через середину отрезка.
- Формально это превращается в задачу квадратичной оптимизации:

$$\min_{\alpha}\left\{||c-d||^2, \text{ где } c=\sum_{y_i=1}\alpha_ix_i, d=\sum_{y_i=-1}\alpha_ix_i\right\}$$
 при условии 
$$\sum_{y_i=1}\alpha_i=\sum_{y_i=-1}\alpha_i=1, \alpha_i\geq 0.$$

• Эту задачу можно решать общими оптимизационными алгоритмами.



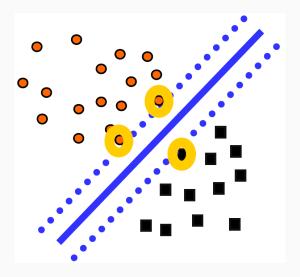
- Другой подход: максимизировать зазор (margin) между двумя параллельными опорными плоскостями, затем провести им параллельную на равных расстояниях от них.
- Гиперплоскость называется опорной для множества точек X, если все точки из X лежат под одну сторону от этой гиперплоскости.
- $\cdot$  Формально: расстояние от точки до гиперплоскости  $y(\mathbf{x})=\mathbf{w}^{ op}\mathbf{x}+w_0=0$  равно  $\frac{|y(\mathbf{x})|}{\|\mathbf{w}\|}.$

- Расстояние от точки до гиперплоскости  $y(\mathbf{x}) = \mathbf{w}^{\top}\mathbf{x} + w_0 = 0$  равно  $\frac{|y(\mathbf{x})|}{\|\mathbf{w}\|}$ .
- \* Все точки классифицированы правильно:  $t_n y(\mathbf{x}_n) > 0$   $(t_n \in \{-1,1\}).$
- И мы хотим найти

$$\begin{split} \arg\max_{\mathbf{w},w_0} \min_n \frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} &= \\ &= \arg\max_{\mathbf{w},w_0} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n \left[ t_n(\mathbf{w}^\top \mathbf{x}_n + w_0) \right] \right\}. \end{split}$$

- $\cdot \, rg \max_{\mathbf{w}, w_0} \left\{ rac{1}{\|\mathbf{w}\|} \min_n \left[ t_n(\mathbf{w}^ op \mathbf{x}_n + w_0) 
  ight] 
  ight\}$  . Сложно.
- $\cdot$  Но если перенормировать  $\mathbf{w}$ , гиперплоскость не изменится.
- · Давайте перенормируем так, чтобы  $\min_n \left[t_n(\mathbf{w}^{\intercal}\mathbf{x}_n + w_0)\right] = 1.$

# ПРИМЕР



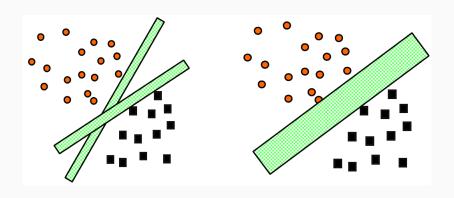
• Получается тоже задача квадратичного программирования:

$$\min_{\vec{w},b} \left\{ \frac{1}{2} ||\mathbf{w}||^2 \right\} \text{ при условии } t_n(\mathbf{w}^\top \mathbf{x}_n + w_0) \geq 1.$$

#### Результаты

- Результаты получаются хорошие. Такой подход позволяет находить устойчивые решения, что во многом решает проблемы с оверфиттингом и позволяет лучше предсказывать дальнейшую классификацию.
- В каком-то смысле в решениях с «толстыми» гиперплоскостями между данными содержится больше информации, чем в «тонких», потому что «толстых» меньше.
- Это всё можно сформулировать и доказать (позже).

# ПРИМЕР



- Напомним, что такое дуальные задачи.
- Прямая задача оптимизации:

$$\min \{f(x)\}$$
 при условии  $h(x) = 0, \ g(x) \le 0, \ x \in X.$ 

• Для дуальной задачи вводим параметры  $\lambda$ , соответствующие равенствам, и  $\mu$ , соответствующие неравенствам.

• Прямая задача оптимизации:

$$\min \{f(x)\}$$
 при условии  $h(x) = 0, \ g(x) \le 0, \ x \in X.$ 

• Дуальная задача оптимизации:

$$\min\left\{\phi(\lambda,\mu)\right\} \text{ при условии } \mu \geq 0,$$
 где  $\phi(\lambda,\mu) = \inf_{x \in X} \left\{f(x) + \lambda^\top h(x) + \mu^\top g(x)\right\}.$ 

- Тогда, если  $(\bar{\lambda},\bar{\mu})$  – допустимое решение дуальной задачи, а  $\bar{x}$  – допустимое решение прямой, то

$$\begin{split} \phi(\bar{\lambda},\bar{\mu}) &= \inf_{x \in X} \left\{ f(x) + \bar{\lambda}^\top h(x) + \bar{\mu}^\top g(x) \right\} \leq \\ &\leq f(\bar{x}) + \bar{\lambda}^\top h(\bar{x}) + \bar{\mu}^\top g(\bar{x}) \leq f(\bar{x}). \end{split}$$

• Это называется слабой дуальностью (только  $\leq$ ), но во многих случаях достигается и равенство.

• Для линейного программирования прямая задача:

$$\min c^{\intercal}x$$
 при условии  $Ax=b,\;x\in X=\{x\leq 0\}.$ 

• Тогда дуальная задача получается так:

$$\begin{split} \phi(\lambda) &= \inf_{x \geq 0} \left\{ c^\top x + \lambda^\top (b - Ax) \right\} = \\ &= \lambda^\top b + \inf_{x \geq 0} \left\{ (c^\top - \lambda^\top A) x \right\} = \\ &= \begin{cases} \lambda^\top b, & \text{если } c^\top - \lambda^\top A \geq 0, \\ -\infty & \text{в противном случае.} \end{cases} \end{split}$$

• Для линейного программирования прямая задача:

$$\min \left\{ c^\top x \right\} \text{ при условии } Ax = b, \ x \in X = \{x \leq 0\}.$$

• Дуальная задача:

 $\max\left\{b^{\top}\lambda\right\}$  при условии  $A^{\top}\lambda\leq c,\;\lambda$  не ограничены.

• Для квадратичного программирования прямая задача:

$$\min\left\{\frac{1}{2}x^\top Qx + c^\top x\right\} \text{ при условии } Ax \leq b,$$

где Q – положительно полуопределённая матрица (т.е.  $x^{\top}Qx \geq 0$  всегда).

• Дуальная задача (проверьте):

$$\max\left\{\frac{1}{2}\mu^\top D\mu + \mu^\top d - \frac{1}{2}c^\top Q^{-1}c\right\} \text{ при условии } c \geq 0,$$

где  $D = -AQ^{-1}A^{\top}$  (отрицательно определённая матрица),  $d = -b - AQ^{-1}c$ .

### Дуальная задача к SVM

• В случае SVM надо ввести множители Лагранжа:

$$L(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_n \alpha_n \left[ t_n(\mathbf{w}^\top \mathbf{x}_n + w_0) - 1 \right], \ \ \alpha_n \geq 0.$$

• Берём производные по  ${f w}$  и  $w_0$ , приравниваем нулю, получаем

$$\mathbf{w} = \sum_{n} \alpha_{n} t_{n} \mathbf{x}_{n},$$
 
$$0 = \sum_{n} \alpha_{n} t_{n}.$$

### Дуальная задача к SVM

• Подставляя в  $L(\mathbf{w}, w_0, \alpha)$ , получим

$$L(\alpha) = \sum_n \alpha_n - \frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m t_n t_m \left( \mathbf{x}_n^\top \mathbf{x}_m \right)$$
 при условии  $\alpha_n \geq 0, \sum_n \alpha_n t_n = 0.$ 

· Это дуальная задача, которая обычно в SVM и используется.

## Предсказание и ККТ

- А для предсказания потом надо посмотреть на знак  $y(\mathbf{x})$ :

$$y(\mathbf{x}) = \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}^{\top} \mathbf{x}_n + w_0.$$

· Получилось, что предсказания зависят от всех точек  $\mathbf{x}_n$ ...

## ПРЕДСКАЗАНИЕ И ККТ

· ...но нет. :) Условия ККТ (Karush-Kuhn-Tucker):

$$\begin{split} \alpha_n &\geq 0, \\ t_n y(\mathbf{x}_n) - 1 &\geq 0, \\ \alpha_n \left( t_n y(\mathbf{x}_n) - 1 \right) &= 0. \end{split}$$

 $\cdot$  Т.е. реально предсказание зависит от небольшого числа опорных векторов, для которых  $t_n y(\mathbf{x}_n) = 1$  (они находятся собственно на границе разделяющей поверхности).

#### Постановка задачи

- Все эти методы работают, когда данные действительно линейно разделимы.
- А что делать, когда их всё-таки немножко не получается разделить?
- Первый вопрос: что делать для первого метода, метода выпуклых оболочек?

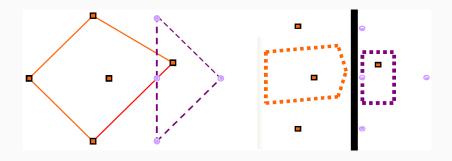
### Редуцированные выпуклые оболочки

• Вместо обычных выпуклых оболочек можно рассматривать редуцированные (reduced), у которых коэффициенты ограничены не 1, а сильнее:

$$c = \sum_{y_i = 1} \alpha_i x_i, \quad 0 \le \alpha_i \le D.$$

- $\cdot$  Тогда для достаточно малых D редуцированные выпуклые оболочки не будут пересекаться.
- И мы будем искать оптимальную гиперплоскость между редуцированными выпуклыми оболочками.

# ПРИМЕР



#### Для метода опорных векторов

• Естественно, для метода опорных векторов тоже надо что-то изменить. Что?

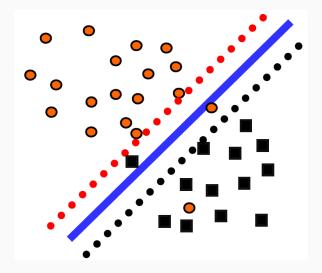
#### Для метода опорных векторов

- Естественно, для метода опорных векторов тоже надо что-то изменить. Что?
- Мы просто добавим в оптимизирующуюся функцию неотрицательную ошибку (slack):

$$\min_{\vec{w},w_0} \left\{ ||\vec{w}||^2 + C \sum_{i=1}^m z_i \right\}$$
 при условии  $t_i(\vec{w}\cdot\vec{x}_i-w_0)+z_i \geq 1.$ 

• Это прямая задача...

# ПРИМЕР



#### Дуальная переформулировка

• ...а вот дуальная:

$$\min_{\alpha} \left\{ \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{m} t_i t_j \alpha_i \alpha_j \left( \vec{x}_i \cdot \vec{x}_j \right) - \sum_{i=1}^{m} \alpha_i, \right.$$
 где  $\sum_{i=1}^{m} t_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C. \right\}$ 

- · Эта формулировка чаще всего используется в теории SVM.
- Единственное отличие от линейно разделимого случая верхняя граница C на  $\alpha_j$ , т.е. на влияние каждой точки.

#### Итого

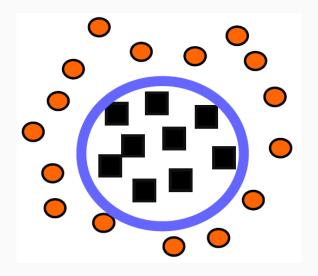
- Метод опорных векторов отлично подходит для линейной классификации.
- Решая задачу квадратичного программирования, мы получаем параметры оптимальной гиперплоскости.
- Точно так же, как и в дуальном случае, если бы мы просто искали середину между выпуклыми оболочками.

SVM и разделение нелинейными
ФУНКЦИЯМИ

#### Введение

- Часто бывает нужно разделять данные не только линейными функциями.
- Что делать в таком случае?

# ПРИМЕР



#### Введение

- Часто бывает нужно разделять данные не только линейными функциями.
- Классический метод: развернуть нелинейную классификацию в пространство большей размерности (feature space), а там запустить линейный классификатор.
- Для этого просто нужно для каждого монома нужной степени ввести новую переменную.

#### ПРИМЕР

• Чтобы в двумерном пространстве [r,s] решить задачу классификации квадратичной функцией, надо перейти в пятимерное пространство:

$$[r,s] \longrightarrow [r,s,rs,r^2,s^2].$$

• Или формальнее; определим  $\theta:\mathbb{R}^2\to\mathbb{R}^5$ :  $\theta(r,s)=(r,s,rs,r^2,s^2)$ . Вектор в  $\mathbb{R}^5$  теперь соответствует квадратичной кривой общего положения в  $\mathbb{R}^2$ , а функция классификации выглядит как

$$f(\vec{x}) = \mathrm{sign}(\theta(\vec{w}) \cdot \theta(\vec{x}) - b).$$

• Если решить задачу линейного разделения в этом новом пространстве, тем самым решится задача квадратичного разделения в исходном.

# Проблемы с классическим подходом

- Во-первых, количество переменных растёт экспоненциально.
- Во-вторых, по большому счёту теряются преимущества того, что гиперплоскость именно оптимальная; например, оверфиттинг опять становится проблемой.
- Важное замечание: концептуально мы задачу уже решили. Остались технические сложности: как обращаться с гигантской размерностью. Но в них-то всё и дело.

#### Основная идея и схема работы SVM

- Тривиальная схема алгоритма классификации такова:
  - входной вектор  $\vec{x}$  трасформируется во входной вектор в feature space (большой размерности);
  - в этом большом пространстве мы вычисляем опорные векторы, решаем задачу разделения;
  - потом по этой задаче классифицируем входной вектор.
- Это нереально, потому что пространство слишком большой размерности.

#### Основная идея и схема работы SVM

- Оказывается, кое-какие шаги здесь можно переставить. Вот так:
  - опорные векторы вычисляются в исходном пространстве малой размерности;
  - там же они перемножаются (сейчас увидим, что это значит);
  - и только потом мы делаем нелинейную трансформацию того, что получится;
  - потом по этой задаче классифицируем входной вектор.
- Осталось теперь объяснить, что всё это значит. :)

#### Постановка задачи

• Напомним, что наша задача поставлена следующим образом:

$$\min_{\alpha} \left\{ \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{m} y_i y_j \alpha_i \alpha_j \left( \vec{x}_i \cdot \vec{x}_j \right) - \sum_{i=1}^{m} \alpha_i, \right.$$
 где 
$$\sum_{i=1}^{m} y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C. \right\}$$

#### Постановка задачи

• Мы теперь хотим ввести некое отображение  $\theta:\mathbb{R}^n o \mathbb{R}^N$ , N>n. Получится:

$$\min_{\alpha} \left\{ \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m y_i y_j \alpha_i \alpha_j \left( \theta(\vec{x}_i) \cdot \theta(\vec{x}_j) \right) - \sum_{i=1}^m \alpha_i, \right.$$
 где 
$$\sum_{i=1}^m y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C. \right\}$$

- Придётся немножко вспомнить (или изучить) функциональный анализ.
- Мы хотим обобщить понятие *скалярного произведения*; давайте введём новую функцию, которая (минуя трансформацию) будет сразу вычислять скалярное произведение векторов в feature space:

$$k(\vec{u},\vec{v}) := \theta(\vec{u}) \cdot \theta(\vec{v}).$$

• Первый результат: любая симметрическая функция  $k(\vec{u}, \vec{v}) \in L_2$  представляется в виде

$$k(\vec{u}, \vec{v}) = \sum_{i=1}^{\infty} \lambda_i \theta_i(\vec{u}) \cdot \theta_i(\vec{v}),$$

где  $\lambda_i \in \mathbb{R}$  — собственные числа, а  $\theta_i$  — собственные векторы интегрального оператора с ядром k, т.е.

$$\int k(\vec{u}, \vec{v}) \theta_i(\vec{u}) \mathrm{d}\vec{u} = \lambda_i \theta_i(\vec{v}).$$

• Чтобы k задавало скалярное произведение, достаточно, чтобы все собственные числа были положительными. А собственные числа положительны тогда и только тогда, когда (теорема Мерсера)

$$\int \int k(\vec{u},\vec{v})g(\vec{u})g(\vec{v})\mathrm{d}\vec{u}\mathrm{d}\vec{v} > 0$$

для всех g таких, что  $\int g^2(\vec{u})\mathrm{d}\vec{u} < \infty$ .

• Вот, собственно и всё. Теперь мы можем вместо подсчёта  $\theta(\vec{u}) \cdot \theta(\vec{v})$  в задаче квадратичного программирования просто использовать подходящее sдро  $k(\vec{u}, \vec{v})$ .

• Итого задача наша выглядит так:

$$\min_{\alpha} \left\{ \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m y_i y_j \alpha_i \alpha_j k(\vec{x}_i, \vec{x}_j) - \sum_{i=1}^m \alpha_i, \right.$$
 где 
$$\sum_{i=1}^m y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C. \right\}$$

- Просто меняя ядро k, мы можем вычислять самые разнообразные разделяющие поверхности.
- Условия на то, чтобы k была подходящим ядром, задаются теоремой Мерсера.

• Рассмотрим ядро

$$k(\vec{u}, \vec{v}) = (\vec{u} \cdot \vec{v})^2.$$

• Какое пространство ему соответствует?

• После выкладок получается:

$$\begin{split} k(\vec{u},\vec{v}) &= (\vec{u}\cdot\vec{v})^2 = \\ &= \left(u_1^2,u_2^2,\sqrt{2}u_1u_2\right)\cdot \left(v_1^2,v_2^2,\sqrt{2}v_1v_2\right). \end{split}$$

• Иначе говоря, линейная поверхность в новом пространстве соответствует квадратичной поверхности в исходном (эллипс, например).

- Естественное обобщение: ядро  $k(\vec{u}, \vec{v}) = (\vec{u} \cdot \vec{v})^d$  задаёт пространство, оси которого соответствуют всем однородным мономам степени d.
- А как сделать пространство, соответствующее произвольной полиномиальной поверхности, не обязательно однородной?

 $\cdot$  Поверхность, описывающаяся полиномом степени d:

$$k(\vec{u},\vec{v}) = (\vec{u}\cdot\vec{v}+1)^d.$$

• Тогда линейная разделимость в feature space в точности соответствует полиномиальной разделимости в базовом пространстве.

· Нормальное распределение (radial basis function):

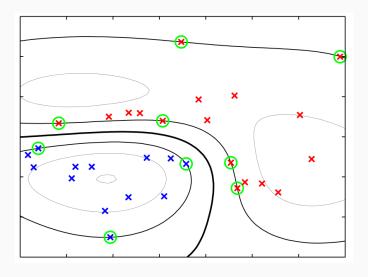
$$k(\vec{u}, \vec{v}) = e^{-\frac{||\vec{u} - \vec{v}||^2}{2\sigma}}.$$

• Двухуровневая нейронная сеть:

$$k(\vec{u},\vec{v}) = o(\eta \vec{u} \cdot \vec{v} + c),$$

где o — сигмоид.

# ПРИМЕР



#### SVM и эмпирический риск

- $\cdot$  Ещё один взгляд на SVM какая вообще задача у любой классификации?
- Мы хотим минимизировать эмпирический риск, то есть число неправильных ответов:

$$\sum_n \left[ y_i \neq t_i \right] \to \ \min_{\mathbf{w}}.$$

• И если функция линейная с параметрами  ${f w}$ ,  $w_0$ , то это эквивалентно

$$\sum_{n} \left[ t_i \left( \mathbf{x}_n^\top \mathbf{w} - w_0 \right) < 0 \right] \to \min_{\mathbf{w}}.$$

- · Величину  $M_i = \mathbf{x}_n^{ op} \mathbf{w} w_0$  назовём отступом (margin).
- Оптимизировать напрямую сложно...

#### SVM и эмпирический риск

• ...поэтому заменим на оценку сверху:

$$\sum_n \left[ M_i < 0 \right] \leq \sum_n \left( 1 - M_i \right) \to \ \min_{\mathbf{w}}.$$

• А потом ещё добавим регуляризатор для стабильности:

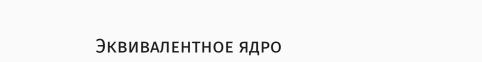
$$\sum_n \left[ M_i < 0 \right] \leq \sum_n \left( 1 - M_i \right) + \frac{1}{2C} \| \mathbf{w} \|^2 \to \min_{\mathbf{w}}.$$

· И это снова получилась задача SVM!

#### **РЕЗЮМЕ**

- Вот какой получается в итоге алгоритм.
  - 1. Выбрать параметр C, от которого зависит акцент на минимизации ошибки или на максимизации зазора.
  - 2. Выбрать ядро и параметры ядра, которые у него, возможно, есть.
  - 3. Решить задачу квадратичного программирования.
  - 4. По полученным значениям опорных векторов определить  $w_0$  (как именно?).
  - 5. Новые точки классифицировать как

$$f(\vec{x}) = \mathrm{sign}(\sum_i y_i \alpha_i k(\vec{x}, \vec{x}_i) - w_0).$$



• Вспомним наши байесовские предсказания:

$$\begin{split} p(t\mid\mathbf{t},\alpha,\beta) &= N(t\mid\mu_N^\top\phi(\mathbf{x}),\sigma_N^2), \end{split}$$
 где  $\sigma_N^2 &= \frac{1}{\beta} + \phi(\mathbf{x})^\top\Sigma_N\phi(\mathbf{x}).$ 

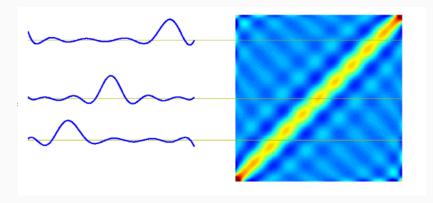
· Давайте перепишем среднее апостериорного распределения в другой форме (вспомним, что  $\mu_N = \beta \Sigma_N \Phi^{\top} \mathbf{t}$ ):

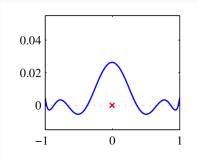
$$\begin{split} y(\mathbf{x}, \boldsymbol{\mu}_N) &= \boldsymbol{\mu}_N^\top \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^\top \boldsymbol{\Sigma}_N \boldsymbol{\Phi}^\top \mathbf{t} = \\ &= \sum_{n=1}^N \beta \phi(\mathbf{x})^\top \boldsymbol{\Sigma}_N \phi(\mathbf{x}_n) t_n. \end{split}$$

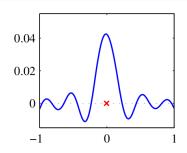
- ·  $y(\mathbf{x}, \mu_N) = \sum_{n=1}^{N} \beta \phi(\mathbf{x})^{\top} \Sigma_N \phi(\mathbf{x}_n) t_n$ .
- Это значит, что предсказание можно переписать как

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

- 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).







#### 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.

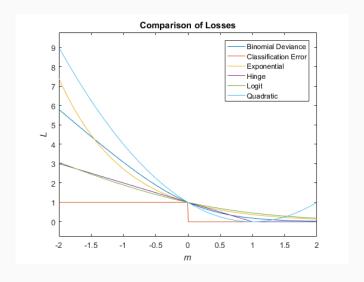
**Упражнение.** Докажите, что  $\sum_{n=1}^{N} k(\mathbf{x}, \mathbf{x}_n) = 1$ .

# RVM

# Функции потерь в классификации

- Сначала ещё один важный другой взгляд: разные методы классификации отличаются друг от друга тем, какую функцию ошибки они оптимизируют.
- У классификации проблема с «правильной» функцией ошибки, то есть ошибкой собственно классификации:
  - она и не везде дифференцируема,
  - и производная её никому не нужна.
- Давайте посмотрим на разные функции потерь (loss functions); мы уже несколько видели, ещё кое-что осталось.

# Функции потерь в классификации



#### Постановка задачи

- · SVM отличный метод. Но и у него есть недостатки.
  - 1. Выходы SVM решения, а апостериорные вероятности непонятно как получить.
  - 2. SVM для двух классов, обобщить на несколько проблематично.
  - 3. Есть параметр C (или  $\nu$ , или ещё вдобавок  $\epsilon$ ), который надо подбирать.
  - 4. Предсказания линейные комбинации ядер, которым необходимо быть положительно определёнными и которые центрированы на точках из датасета.
- Сейчас мы рассмотрим байесовский аналог SVM relevance vector machines (RVM).

#### RVM для РЕГРЕССИИ

- · RVM удобнее сразу формулировать для регрессии.
- Вспомним обычную нашу линейную модель:

$$p(t\mid\mathbf{x},\mathbf{w},eta)=N(t\mid y(\mathbf{x}),eta^{-1}),$$
 где 
$$y(\mathbf{x})=\sum_{i=1}^M w_i\phi_i(\mathbf{x})=\mathbf{w}^{ op}\phi(\mathbf{x}).$$

#### RVM для РЕГРЕССИИ

- RVM это вариант такой модели, который старается работать как SVM.
- Рассмотрим

$$y(\mathbf{x}) = \sum_{n=1}^{N} w_n k(\mathbf{x}, \mathbf{x}_n) + b.$$

• Т.е. мы сразу ищем решение в форме линейной комбинации значений ядра (вспомним «эквивалентное ядро» для линейной регрессии), но, в отличие от SVM, теперь ядро никак не ограничивается.

#### RVM для РЕГРЕССИИ

· Для N наблюдений вектора  ${f x}$  (обозначим через  ${f X}$ ) со значениями  ${f t}$  получим правдоподобие

$$p(\mathbf{t}\mid \mathbf{X}, \mathbf{w}, \boldsymbol{\beta}) = \prod_{n=1}^N p(t_n\mid \mathbf{x}_n, \mathbf{w}, \boldsymbol{\beta}^{-1}).$$

• Априорное распределение тоже будет нормальное, но вместо единого гиперпараметра для всех весов мы введём отдельный гиперпараметр для каждого:

$$p(\mathbf{w}\mid\alpha) = \prod_{i=1}^M N(w_i\mid 0, \alpha_i^{-1}).$$

• Отдельные гиперпараметры:

$$p(\mathbf{w} \mid \alpha) = \prod_{i=1}^{M} N(w_i \mid 0, \alpha_i^{-1}).$$

- Идея здесь в том, что при максимизации апостериорной вероятности большая часть  $\alpha_i$  просто уйдёт на бесконечность, и соответствующие веса будут нулевыми.
- Сейчас увидим, как это получается.

• Апостериорное распределение нам знакомо:

$$p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \alpha, \beta) = N(\mathbf{w} \mid \mathbf{m}, \Sigma),$$
 где

$$\mathbf{m} = \beta \Sigma \Phi^{\top} \mathbf{t},$$
  
$$\Sigma = (\mathbf{A} + \beta \Phi^{\top} \Phi)^{-1},$$

где  ${f A}={
m diag}(lpha_1,\dots,lpha_M)$ , а  $\Phi$  в нашем случае – это  ${f K}$ , симметрическая матрица с элементами  $k({f x}_n,{f x}_m).$ 

• Как найти  $\alpha$  и  $\beta$ ? Нужно максимизировать маргинальное правдоподобие датасета

$$p(\mathbf{t} \mid \mathbf{X}, \alpha, \beta) = \int p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} \mid \alpha) d\mathbf{w}.$$

• Это свёртка двух гауссианов, тоже гауссиан:

$$\begin{split} & \ln p(\mathbf{t}\mid \mathbf{X}, \alpha, \beta) = \ln N(\mathbf{t}\mid 0, \mathbf{C}) = \\ & = -\frac{1}{2}\left[N\ln(2\pi) + \ln |\mathbf{C}| + \mathbf{t}^{\top}\mathbf{C}^{-1}\mathbf{t}\right], \text{ где } \mathbf{C} = \beta^{-1}\mathbf{I} + \Phi \mathbf{A}^{-1}\Phi^{\top}. \end{split}$$

• Как это оптимизировать?

• Можно подсчитать производные и получить

$$\begin{split} \alpha_i &= \frac{\gamma_i}{m_i^2}, \\ \beta^{-1} &= \frac{\|\mathbf{t} - \Phi \mathbf{m}\|^2}{N - \sum_i \gamma_i}, \end{split}$$

где 
$$\gamma_i = 1 - \alpha_i \Sigma_{ii}$$
.

• Теперь можно просто итеративно пересчитывать  $\alpha,\beta$  из  $\mathbf{m},\Sigma$ , потом наоборот, потом опять наоборот, и до сходимости.

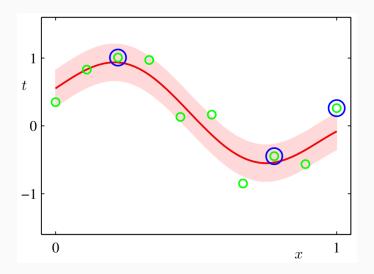
#### RVM для регрессии

- В результате получается обычно, что большинство  $\alpha_i$  неограниченно растут, и соответствующие веса можно считать нулевыми.
- · Оставшиеся называются relevance vectors, их обычно мало.
- Если теперь мы найдём  $\alpha^*, \beta^*$ , то предсказывать в новых точках можно как

$$\begin{split} p(t \mid \mathbf{x}, \mathbf{X}, \mathbf{t}, \alpha^*, \beta^*) &= \int p(t \mid \mathbf{x}, \mathbf{w}, \beta^*) p(\mathbf{w} \mid \mathbf{X}, \mathbf{t}, \alpha^*, \beta^*) d\mathbf{w} = \\ &= N(t \mid \mathbf{m}^\top \phi(\mathbf{x}), \sigma^2(\mathbf{x})), \end{split}$$

где 
$$\sigma^2(\mathbf{x}) = (\beta^*)^{-1} + \phi(\mathbf{x})^{\top} \Sigma \phi(\mathbf{x}).$$

# RVM для РЕГРЕССИИ



• Можно сделать то же самое и для классификации. Рассмотрим классификацию с двумя классами,  $t \in \{0,1\}$ :

$$y(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^{\top} \phi(\mathbf{x})).$$

• И добавим сюда, опять же, априорное распределение с разными  $\alpha_i$  для каждого веса:

$$p(\mathbf{w}\mid\alpha) = \prod_{i=1}^M N(w_i\mid 0, \alpha_i^{-1}).$$

• Идея: инициализируем  $\alpha$ , считаем лапласовское приближение к апостериорному распределению, максимизируем, получаем новое  $\alpha$ , и т.д.

• Апостериорное распределение:

$$\begin{split} & \ln p(\mathbf{w} \mid \mathbf{t}, \alpha) = \ln \left( p(\mathbf{t} \mid \mathbf{w}) p(\mathbf{w} \mid \alpha) \right) - \ln p(\mathbf{t} \mid \alpha) = \\ & = \sum_{n=1}^{N} \left[ t_n \ln y_n + (1 - t_n) \ln (1 - y_n) \right] - \frac{1}{2} \mathbf{w}^{\top} \mathbf{A} \mathbf{w} + \text{const.} \end{split}$$

 Мы уже обсуждали, как его максимизировать – IRLS; для этого подсчитаем

$$\nabla \ln p(\mathbf{w} \mid \mathbf{t}, \alpha) = \Phi^{\top}(\mathbf{t} - \mathbf{y}) - \mathbf{A}\mathbf{w},$$
$$\nabla \nabla \ln p(\mathbf{w} \mid \mathbf{t}, \alpha) = -\left(\Phi^{\top} \mathbf{B} \Phi + \mathbf{A}\right),$$

где  ${f B}$  – диагональная матрица с элементами  $b_n=y_n(1-y_n).$ 

• Лапласовское приближение получится из  $\nabla \ln p(\mathbf{w} \mid \mathbf{t}, lpha)$ , и получится

$$\begin{split} \mathbf{w}^* &= \mathbf{A}^{-1} \boldsymbol{\Phi}^\top \left( \mathbf{t} - \mathbf{y} \right), \\ \boldsymbol{\Sigma} &= \left( \boldsymbol{\Phi}^\top \mathbf{B} \boldsymbol{\Phi} + \mathbf{A} \right)^{-1}, \end{split}$$

и распределение для предсказания получится

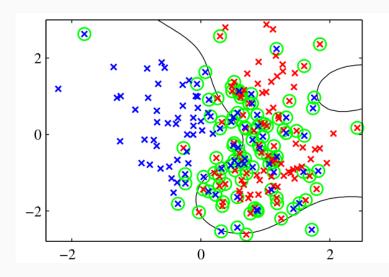
$$\begin{split} p(\mathbf{t} \mid \alpha) &= \int p(\mathbf{t} \mid \mathbf{w}) p(\mathbf{w} \mid \alpha) d\mathbf{w} \approx \\ &\approx p(\mathbf{t} \mid \mathbf{w}^*) p(\mathbf{w}^* \mid \alpha) (2\pi)^{M/2} |\Sigma|^{1/2}. \end{split}$$

- $\begin{aligned} & \boldsymbol{\cdot} \ p(\mathbf{t} \mid \boldsymbol{\alpha}) = \int p(\mathbf{t} \mid \mathbf{w}) p(\mathbf{w} \mid \boldsymbol{\alpha}) d\mathbf{w} \approx \\ & p(\mathbf{t} \mid \mathbf{w}^*) p(\mathbf{w}^* \mid \boldsymbol{\alpha}) (2\pi)^{M/2} |\boldsymbol{\Sigma}|^{1/2}. \end{aligned}$
- Теперь мы оптимизируем это по  $\alpha$ : берём производную, получаем

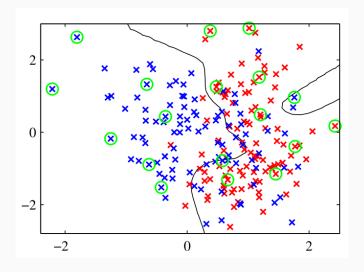
$$-\frac{1}{2}(w_i^*)^2+\frac{1}{2\alpha_i}-\frac{1}{2}\Sigma_{ii}=0, \text{ т.е.}$$
 
$$\alpha_i=\frac{\gamma_i}{(w_i^*)^2}, \ \gamma_i=1-\alpha_i\Sigma_{ii}.$$

• Т.е. формула получилась точно такая же, как в случае регрессии.

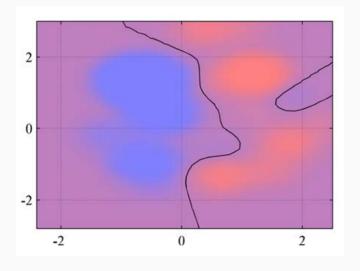
# Было: SVM



## Стало: RVM



## Стало: RVM



### RVM для нескольких классов

 На несколько классов теперь обобщается естественным образом:

$$a_k = \mathbf{w}_k^{\intercal} \mathbf{x}, \quad y_k(\mathbf{x}) = \frac{e^{a_k}}{\sum_{j} e^{a_j}}.$$

• И дальше всё то же самое.

#### Сравнение SVM и RVM

- RVM как-то получилось лучше со всех сторон.
- Главный минус в RVM обучение гораздо дольше (хотя есть алгоритмы и побыстрее, чем мы рассматривали, но всё равно дольше).
- Но даже это не то чтобы минус, потому что в SVM нужна кросс-валидация для подбора параметров, т.е. на самом деле обучение SVM дольше, чем кажется.
- Есть и (даже более важный) плюс с точки зрения скорости в RVM гораздо быстрее применение модели к новым точкам, потому что опорных векторов гораздо меньше.

• В RVM для регрессии получается правдоподобие

$$\begin{split} & \ln p(\mathbf{t} \mid \mathbf{X}, \alpha, \beta) = \ln N(\mathbf{t} \mid 0, \mathbf{C}) = \\ & = -\frac{1}{2} \left[ N \ln(2\pi) + \ln |\mathbf{C}| + \mathbf{t}^{\top} \mathbf{C}^{-1} \mathbf{t} \right], \text{ где } \mathbf{C} = \beta^{-1} \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^{\top}. \end{split}$$

- Выделим вклад в  ${f C}$  одного компонента  $lpha_i$ :

$$\begin{split} \mathbf{C} &= \beta^{-1}\mathbf{I} + \sum_{j \neq i} \alpha_j^{-1} \boldsymbol{\varphi}_j \boldsymbol{\varphi}_j^\top + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^\top = \\ &= \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^\top, \end{split}$$

где  $\varphi_i$  – i-я строка  $\Phi$  ( $\phi_n$  был n-м столбцом).

$$\cdot \ \mathbf{C} = \mathbf{C}_{-i} + \alpha_i^{-1} \varphi_i \varphi_i^{\top}.$$

• Верны следующие тождества для определителя и обратной матрицы:

$$\begin{split} |\mathbf{C}| &= |\mathbf{C}_{-i}||1 + \alpha_i^{-1} \boldsymbol{\varphi}_i^{\intercal} \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i|, \\ \mathbf{C}^{-1} &= \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^{\intercal} \mathbf{C}_{-i}^{-1}}{\alpha_i + \boldsymbol{\varphi}_i^{\intercal} \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i}. \end{split}$$

Упражнение. Докажите это.

 $\cdot$  Значит, L(lpha) можно переписать в виде

$$L(\alpha)=L(\alpha_{-i})+\lambda(\alpha_i), \ \text{где}$$
 
$$\lambda(\alpha_i)=\frac{1}{2}\left[\ln\alpha_i-\ln(\alpha_i+s_i)+\frac{q_i^2}{\alpha_i+s_i}\right].$$

• Здесь

$$\begin{array}{ll} s_i &= \boldsymbol{\varphi}_i^\top \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i & \quad \text{sparsity } \boldsymbol{\varphi}_i \\ q_i &= \boldsymbol{\varphi}_i^\top \mathbf{C}_{-i}^{-1} \mathbf{t} & \quad \text{quality } \boldsymbol{\varphi}_i. \end{array}$$

- ·  $s_i = \varphi_i^{\intercal} \mathbf{C}_{-i}^{-1} \varphi_i$ ,  $q_i = \varphi_i^{\intercal} \mathbf{C}_{-i}^{-1} \mathbf{t}$ .
- Sparsity то, насколько  $\varphi_i$  перекрывается с остальными векторами модели.
- Quality то, насколько  $\varphi_i$  сонаправлен с ошибкой между  ${f t}$  и  ${f y}_{-i}$  (ошибкой модели без  $\varphi_i$ ).
- Чем больше sparsity и чем меньше quality, тем более вероятно, что этот базисный вектор из модели исключат (т.е.  $\alpha_i \to \infty$ ).

· 
$$\lambda(\alpha_i) = \frac{1}{2} \left[ \ln \alpha_i - \ln(\alpha_i + s_i) + \frac{q_i^2}{\alpha_i + s_i} \right].$$

• Возьмём производную, приравняем нулю, получим (т.к.  $lpha_i \geq 0$ )

$$\alpha_i = \begin{cases} \infty, & q_i^2 \leq s_i, \\ \frac{s_i^2}{q_i^2 - s_i}, & q_i^2 > s_i. \end{cases}$$

• Как мы и ожидали.

- И алгоритм теперь получается такой:
  - 1. инициализировать eta,  $\varphi_{\text{1}}$ ,  $lpha_{1}=s_{1}^{2}/(q_{1}^{2}-s_{1})$ , остальные  $lpha_{j}=\infty$ ;
  - 2. вычислить  $\Sigma$ ,  $\mathbf{m}$ ,  $q_i$  и  $s_i$  для всех i;
  - 3. выбрать i, проапдейтить  $\alpha_i$ , проапдейтить  $\beta$ ;
  - 4. goto 2 и так пока не сойдётся.

# Спасибо за внимание!



