

- Введение в машинное обучение •

Минимизация эмпирического риска, линейные модели и регуляризация

Воронцов Константин Вячеславович

`k.v.vorontsov@phystech.edu`

`http://www.MachineLearning.ru/wiki?title=User:Vokov`

Этот курс доступен на странице вики-ресурса

`http://www.MachineLearning.ru/wiki`

«Введение в машинное обучение (курс лекций, К.В.Воронцов)»

Машинное обучение — автоматизация научного метода

Задача машинного обучения начинается с её постановки

- ДНК задачи — Дано, Найти, Критерий

Основные понятия машинного обучения:

- объект, ответ, признак
- модель, функция потерь, эмпирический риск
- алгоритм обучения, недообучение, переобучение

Три главных принципа машинного обучения:

- эмпирическая индукция Фрэнсиса Бэкона
- минимизация (+регуляризация) эмпирического риска
- обучаемая векторизация данных (Deep Learning)

1 Градиентная оптимизация в машинном обучении

- Оптимизационная постановка задачи
- Метод стохастического градиента
- Ускорение сходимости и другие эвристики

2 Основные типы задач обучения с учителем

- Задачи регрессии
- Задачи классификации
- Задачи ранжирования

3 Линейные модели

- Линейный классификатор и логистическая регрессия
- Мультиколлинеарность и регуляризация
- Метод опорных векторов

Общая постановка большинства задач машинного обучения

Дано: X — пространство объектов

$X^\ell = \{x_1, \dots, x_\ell\} \subset X$ — обучающая выборка (training sample)

$a(x, w)$, $a: X \times W \rightarrow Y$ — параметрическая модель, гипотеза

Найти $w \in W \subseteq \mathbb{R}^N$ — вектор параметров модели $a(x, w)$

Критерий \min регуляризованного эмпирического риска
(Regularized ERM — Empirical Risk Minimization):

$$Q(w, X^\ell) = \frac{1}{\ell} \sum_{i=1}^{\ell} \mathcal{L}(w, x_i) + \tau \mathcal{R}(w) \rightarrow \min_w$$

$\mathcal{L}(w, x)$ — функция потерь (loss function),

тем больше, чем хуже ответ модели $a(x, w)$ на объекте x

$\mathcal{R}(w)$ — регуляризатор, не прецедентные требования к модели

Вопрос: зачем нужен коэффициент регуляризации τ ?

Градиентный метод минимизации эмпирического риска

$$Q(w) = \frac{1}{\ell} \sum_{i=1}^{\ell} \mathcal{L}(w, x_i) + \tau \mathcal{R}(w) \rightarrow \min_w$$

Метод *градиентного спуска*:

$w^{(0)}$:= начальное приближение;

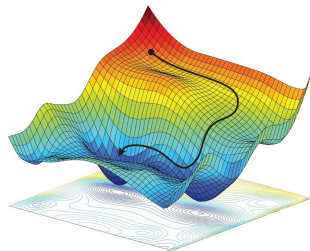
$$w^{(t+1)} := w^{(t)} - h \nabla Q(w^{(t)})$$

где $\nabla Q(w) = \left(\frac{\partial Q(w)}{\partial w_j} \right)_{j=1}^N$ — *вектор градиента*,

h — *градиентный шаг*, называемый также *темпом обучения*

$$w^{(t+1)} := w^{(t)} - h \left(\frac{1}{\ell} \sum_{i=1}^{\ell} \nabla \mathcal{L}(w^{(t)}, x_i) + \tau \nabla \mathcal{R}(w^{(t)}) \right)$$

Идея ускорения сходимости: брать случайное подмножество слагаемых, или вообще один объект, чаще обновляя вектор w



Метод стохастического градиента SG (Stochastic Gradient)

$$Q(w) = \frac{1}{\ell} \sum_{i=1}^{\ell} \mathcal{L}(w, x_i) + \tau \mathcal{R}(w) \rightarrow \min_w$$

Вход: выборка X^ℓ , параметры h, τ, λ ;

Выход: вектор параметров w ;

инициализировать параметры $w_j, j = 1, \dots, N$;

инициализировать оценку $Q(w)$ по небольшой подвыборке;

повторять

 объект x_i выбрать из X^ℓ случайным образом;

 потеря: $\mathcal{L}_i := \mathcal{L}(w, x_i)$;

 градиентный шаг: $w := w - h \nabla \mathcal{L}(w, x_i) - h\tau \nabla \mathcal{R}(w)$;

 рекуррентная оценка критерия: $Q := \lambda \mathcal{L}_i + (1 - \lambda)Q$;

пока значение Q и/или параметры w не сойдутся;

В: какие бывают условия остановки в методах оптимизации?

H. Robbins, S. Monro. A stochastic approximation method. 1951.

Откуда взялась такая рекуррентная оценка функционала?

Проблема: вычисление оценки Q по всей выборке x_1, \dots, x_ℓ намного дольше градиентного шага по одному объекту x_i .

Решение: использовать приближённую рекуррентную формулу.

Среднее арифметическое:

$$\bar{Q}_m = \frac{1}{m}\mathcal{L}_m + \frac{1}{m}\mathcal{L}_{m-1} + \frac{1}{m}\mathcal{L}_{m-2} + \dots$$

$$\bar{Q}_m = \frac{1}{m}\mathcal{L}_m + (1 - \frac{1}{m})\bar{Q}_{m-1}$$

Экспоненциальное скользящее среднее (ЭСС):

$$\bar{Q}_m = \lambda\mathcal{L}_m + (1 - \lambda)\lambda\mathcal{L}_{m-1} + (1 - \lambda)^2\lambda\mathcal{L}_{m-2} + \dots$$

$$\bar{Q}_m = \lambda\mathcal{L}_m + (1 - \lambda)\bar{Q}_{m-1}$$

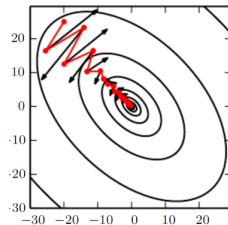
Параметр λ (порядка $\frac{1}{m}$) — *темп забывания* предыстории ряда.

Метод накопления инерции (momentum)

Momentum — экспоненциальное скользящее среднее градиента по последним $\approx \frac{1}{1-\gamma}$ итерациям [Б.Т.Поляк, 1964]:

$$v := \gamma v + (1-\gamma) \nabla \mathcal{L}(w, x_i)$$

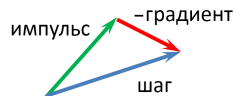
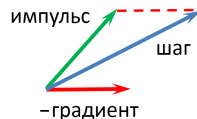
$$w := w - hv$$



NAG (Nesterov's accelerated gradient) — стохастический градиент с инерцией [Ю.Е.Нестеров, 1983]:

$$v := \gamma v + (1-\gamma) \nabla \mathcal{L}(w - h\gamma v, x_i)$$

$$w := w - hv$$



Эвристики (универсальные, не зависящие от вида модели)

- 1 **инициализация w :** случайная, быстрое приближённое решение по малой подвыборке, их комбинация, ...
- 2 **порядок предъявления объектов:**
чем больше \mathcal{L}_i , тем выше вероятность снова выбрать x_i
- 3 **выбор градиентного шага:** *метод скорейшего спуска*
основан на поиске оптимального *адаптивного шага* h^* :

$$\mathcal{L}(w - h\nabla\mathcal{L}(w, x_i), x_i) \rightarrow \min_h$$

- 4 **методы второго порядка:** *метод Ньютона-Рафсона*
 $w := w - h(\mathcal{L}''(w, x_i))^{-1}\nabla\mathcal{L}(w, x_i)$, \mathcal{L}'' — матрица Гессе
диагональный метод Левенберга-Марквардта:
обнулить все недиагональные элементы матрицы Гессе
- 5 **мультистарт:** многократные запуски из разных случайных начальных приближений и выбор лучшего решения

Метод SG: достоинства и недостатки

Достоинства:

- 1 *универсальность*: для любых $a(x, w)$, $\mathcal{L}(w, x)$, $\mathcal{R}(w)$
- 2 *простота*: относительно легко реализуется
- 3 *поточность*: возможность обучения на потоке данных
- 4 *подходит для обработки больших данных*:
 - можно получить неплохое решение, успев обработать лишь малую часть обучающей выборки
 - часто оказывается быстрее и лучше более сложных и ресурсоёмких методов второго порядка

Недостатки:

- 1 подбор комплекса эвристик является искусством
(не забыть про переобучение, застревание, расходимость)

Задача обучения модели регрессии

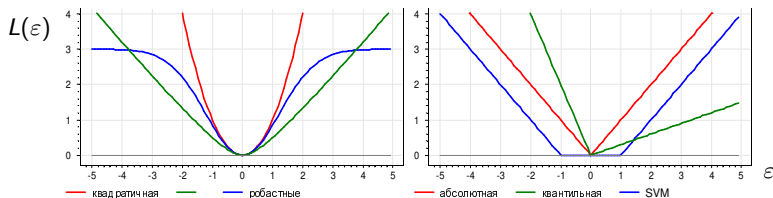
Дано: обучающая выборка $(x_i, y_i)_{i=1}^{\ell}$ с ответами $y_i \in \mathbb{R}$

Найти: вектор параметров w модели регрессии $a(x, w)$

Критерий: минимум эмпирического риска (ERM)

$$Q(w) = \sum_{i=1}^{\ell} L(a(x_i, w) - y_i) \rightarrow \min_w$$

Унимодальные функции потерь $L(\varepsilon)$ от ошибки $\varepsilon = a(x, w) - y$:



В: какие свойства придаёт модели каждая из этих функций?

Пример. Задача прогнозирования спроса объёмов продаж

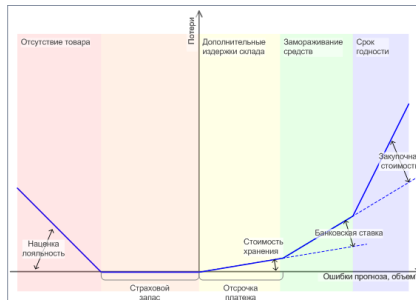
Объект — тройка ⟨товар, магазин, день⟩.

Примеры признаков:

- бинарные: выходной день, праздник, промоакция, и т. д.
- количественные: объёмы продаж в предшествующие дни.

Особенности задачи:

- функция потерь кусочно линейная, не квадратичная, не симметричная
- разреженные данные: если запасы нулевые, то и продажи нулевые



Задача обучения модели бинарной классификации

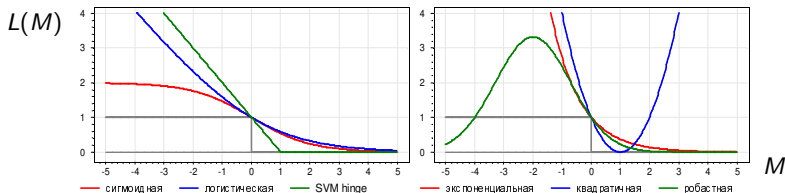
Дано: обучающая выборка $(x_i, y_i)_{i=1}^{\ell}$, $y_i \in \{-1, +1\}$

Найти: вектор w модели классификации $a(x, w) = \text{sign } g(x, w)$

Критерий \min аппроксимированного эмпирического риска:

$$\sum_{i=1}^{\ell} [g(x_i, w)y_i < 0] \leq \sum_{i=1}^{\ell} L(g(x_i, w)y_i) \rightarrow \min_w$$

Убывающие функции потерь $L(M)$ от отступа $M = g(x, w)y$:



В: какие свойства придаёт модели каждая из этих функций?

Бинарный разделяющий классификатор (margin-based classifier)

Бинарный классификатор: $a(x, w) = \text{sign } g(x, w)$, $Y = \{-1, +1\}$

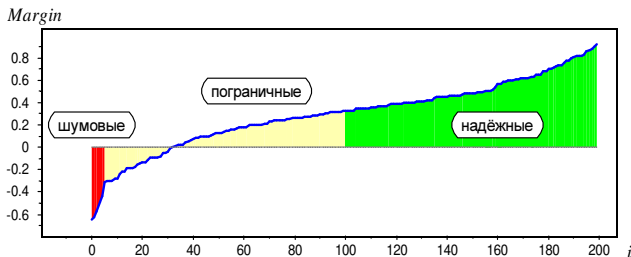
$g(x, w)$ — разделяющая (дискриминантная) функция

$\{x: g(x, w) = 0\}$ — разделяющая поверхность между классами

$M_i(w) = g(x_i, w)y_i$ — отступ (margin) объекта x_i

$M_i(w) < 0 \iff$ модель $a(x, w)$ ошибается на x_i

Ранжирование объектов по возрастанию отступов $M_i(w)$:



Задача обучения модели многоклассовой классификации

Дано: обучающая выборка $(x_i, y_i)_{i=1}^{\ell}$, $y_i \in Y$, $|Y| < \infty$

Найти: вектор параметров $w = (w_y : y \in Y)$ модели:

$$a(x, w) = \arg \max_{y \in Y} g_y(x, w_y)$$

Критерий «каждый против всех остальных» (One-vs.-All):
вводится $M_{iy}(w)$ — отступ объекта x_i относительно класса y

$$\sum_{i=1}^{\ell} \sum_{z \neq y_i} \underbrace{[g_{y_i}(x_i, w_{y_i}) - g_z(x_i, w_z)]}_{M_{iz}(w)} < 0 \leq \sum_{i=1}^{\ell} \sum_{z \neq y_i} L(M_{iz}(w)) \rightarrow \min_w$$

$M_i(w) = \min_{z \neq y_i} M_{iz}(w)$ — **многоклассовый отступ** объекта x_i
(годится для оценивания, но критерий на нём строить нельзя)

Задача обучения модели ранжирования

Дано: обучающая выборка (x_1, \dots, x_ℓ) ,

$i \prec j$ — отношение « x_j лучше, чем x_i » между объектами из X^ℓ

Найти: параметры w модели ранжирования $a(x, w)$, которая аппроксимирует отношение порядка, заданное на выборке:

$$i \prec j \Rightarrow a(x_i, w) < a(x_j, w)$$

Критерий: число неверно ранжированных пар объектов

$$\begin{aligned} Q(w) &= \sum_{i \prec j} \underbrace{[a(x_j, w) - a(x_i, w) < 0]}_{M_{ij}(w)} \leq \\ &\leq \sum_{i \prec j} L(a(x_j, w) - a(x_i, w)) \rightarrow \min_w \end{aligned}$$

где $L(M)$ — убывающая функция *парного отступа* $M_{ij}(w)$

Идея **SG**: градиентные шаги по парам объектов x_i, x_j : $i \prec j$

Примеры задач ранжирования

Ранжирование (Learning to Rank, LtR, L2R, LETOR) нужно в системах человеко-машинного принятия решений, когда **машина** формирует ранжированный список вариантов, **человек** (ответственно) выбирает и принимает решение

- ранжирование выдачи поисковой системы
- ранжирование рекомендаций пользователям (книги, фильмы, музыка, товары интернет-магазина, и т.п.)
- ранжирование вариантов автоматического завершения запроса (Query Auto Completion, auto-suggest)
- ранжирование возможных ответов в диалоговых системах (Question Answering Systems)
- ранжирование вариантов перевода в системах машинного перевода (Machine Translation)

Логистическая регрессия для бинарной классификации

Дано: обучающая выборка $(x_i, y_i)_{i=1}^{\ell}$, $x_i \in \mathbb{R}^n$, $y_i \in \{-1, +1\}$

Найти: параметры $w \in \mathbb{R}^n$ линейной модели классификации

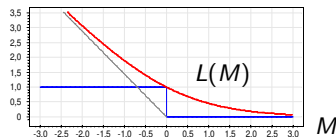
$$a(x, w) = \text{sign} \sum_{j=1}^n w_j f_j(x) = \text{sign} \langle w, x \rangle$$

Критерий \min аппроксимированного эмпирического риска:

$$\sum_{i=1}^{\ell} \ln(1 + \exp(-\underbrace{\langle w, x_i \rangle}_{\text{margin}} y_i)) \rightarrow \min_w$$

логарифмическая функция потерь:

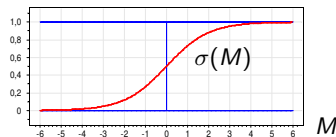
$$L(M) = \ln(1 + e^{-M})$$



Условные вероятности классов:

$$P(y|x, w) = \sigma(\langle w, x \rangle y)$$

$\sigma(M) = \frac{1}{1+e^{-M}}$ — сигмоидная функция



Мультиколлинеарность и переобучение в линейных моделях

$a(x, w) = \langle w, x \rangle$ — линейная модель регрессии

$a(x, w) = \text{sign} \langle w, x \rangle$ — линейная модель классификации

Мультиколлинеарность — линейная зависимость признаков:

- $\exists u \in \mathbb{R}^n: \forall x \in X \quad \langle u, x \rangle = 0$
- неединственность решения: $\forall \gamma \in \mathbb{R} \quad \langle w, x \rangle = \langle w + \gamma u, x \rangle$

Мультиколлинеарность приводит к переобучению:

- неустойчивость: слишком большие веса $|w_j|$ разных знаков
- переобучение: $Q(w^*, X^\ell) \ll Q(w^*, X^k)$

Как уменьшить переобучение:

- регуляризация $\|w\| \rightarrow \min$ (сокращение весов, weight decay)
- отбор признаков: $f_1, \dots, f_n \rightarrow f_{j_1}, \dots, f_{j_m}, \quad m \ll n$
- преобразование признаков: $f_1, \dots, f_n \rightarrow g_1, \dots, g_m, \quad m \ll n$

Регуляризация по L_2 -норме (сокращение весов, weight decay)

Штраф за увеличение нормы вектора весов:

$$\mathcal{L}_\tau(w, x_i) = \mathcal{L}(w, x_i) + \frac{\tau}{2} \|w\|^2 = \mathcal{L}(w, x_i) + \frac{\tau}{2} \sum_{j=1}^n w_j^2 \rightarrow \min_w.$$

Градиент:

$$\nabla \mathcal{L}_\tau(w, x_i) = \nabla \mathcal{L}(w, x_i) + \tau w.$$

Модификация градиентного шага в методе SG:

$$w := w(1 - h\tau) - h\nabla \mathcal{L}(w, x_i).$$

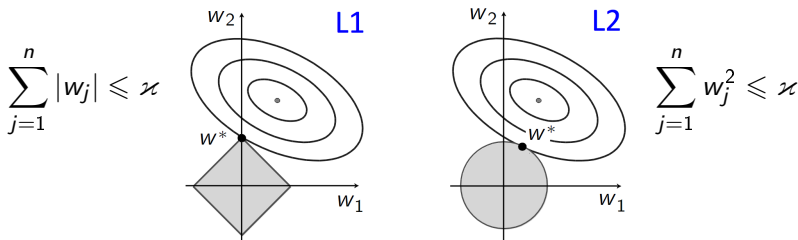
Методы подбора коэффициента регуляризации τ :

- 1 hold-out или скользящий контроль
- 2 стохастическая адаптация по сетке значений $\{\tau_k\}$

Регуляризация по L_p -норме ($p = 1, 2, \dots$) для отбора признаков

LASSO — Least Absolute Shrinkage and Selection Operator

$$Q(w) + \tau \sum_{j=1}^n |w_j|^p \rightarrow \min_w \iff \begin{cases} Q(w) \rightarrow \min_w; \\ \sum_{j=1}^n |w_j|^p \leq \kappa; \end{cases}$$

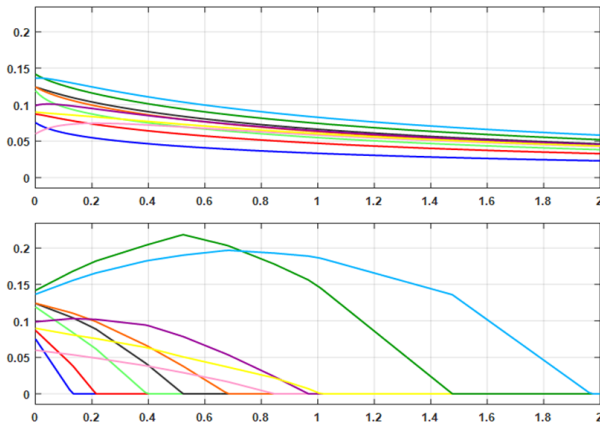


B: как доказать эквивалентность этих двух постановок?

T.Hastie, R.Tibshirani, J.Friedman. The Elements of Statistical Learning. 2017.

Сравнение L_2 и L_1 регуляризации

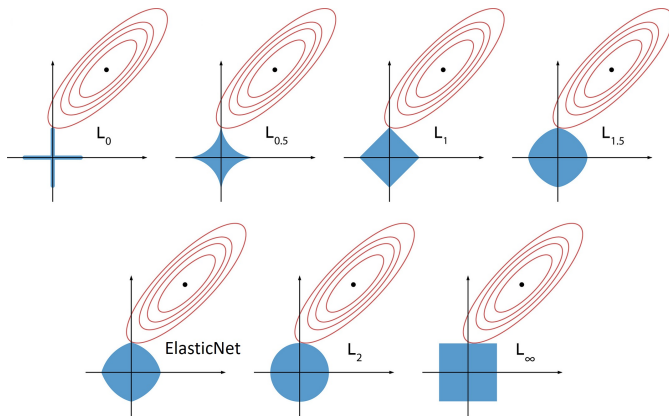
Зависимость весов w_j от параметра селективности τ



LASSO усиливает отбор признаков по мере увеличения τ

Геометрическая интерпретация отбора признаков

Сравнение регуляризаторов по различным L_p -нормам:



Негладкость регуляризатора приводит к отбору признаков

Негладкие регуляризаторы для отбора и группировки признаков

Общий вид регуляризаторов (τ — параметр селективности):

$$\sum_{i=1}^{\ell} \mathcal{L}(w, x_i) + \sum_{j=1}^n R_{\tau}(w_j) \rightarrow \min_w.$$

LASSO (L_1): $R_{\tau}(w) = \tau|w|$

— слишком агрессивно отбрасывает признаки

Elastic Net: $R_{\tau}(w) = \tau|w| + \lambda w^2$

— группирует информативные признаки

Support Feature Machine (SFM):

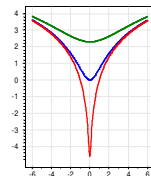
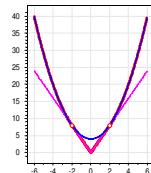
$$R_{\tau}(w) = [|w| \leq \tau] \cdot 2\tau|w| + [|w| > \tau] \cdot (\tau^2 + w^2)$$

— не группирует неинформативные признаки

Relevance Feature Machine (RFM):

$$R_{\tau}(w) = \ln(\tau w^2 + 1)$$

— лучше находит совместно значимые признаки



Метод опорных векторов (Support Vector Machine, SVM)

Дано: обучающая выборка $(x_i, y_i)_{i=1}^{\ell}$, $x_i \in \mathbb{R}^n$, $y_i \in \{-1, +1\}$

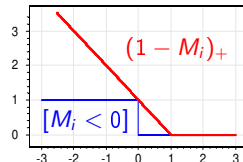
Найти: параметры $w \in \mathbb{R}^n$, $w_0 \in \mathbb{R}$ линейной модели

$$a(x; w, w_0) = \text{sign}(\langle x, w \rangle - w_0)$$

Критерий аппроксимированного и регуляризованного ERM:

$$\sum_{i=1}^{\ell} \left(1 - \underbrace{(\langle x_i, w \rangle - w_0) y_i}_{\text{margin } M_i(w, w_0)} \right)_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}$$

- Аппроксимация штрафует объекты за приближение к границе классов, увеличивая зазор между классами
- Регуляризация штрафует неустойчивые решения в случае мультиколлинеарности

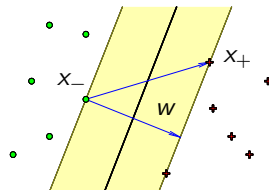


В.Н.Вавник, А.Я.Лернер. Узнавание образов при помощи обобщенных портретов. 1963
С.С.Солтис, В.Вавник. Support vector networks. 1995.

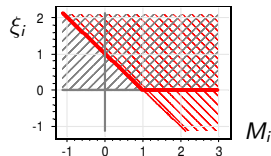
Основные свойства метода опорных векторов

- Геометрия:** ширина зазора (разделяющей полосы) между классами максимальна, разделяющая гиперплоскость проходит посередине полосы.
- Эквивалентная** гладкая задача выпуклого программирования с переменными $\xi_i = (1 - M_i(w, w_0))_+$:

$$\begin{cases} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \rightarrow \min_{w, w_0, \xi} \\ \xi_i \geq 1 - M_i(w, w_0), \quad i = 1, \dots, \ell \\ \xi_i \geq 0, \quad i = 1, \dots, \ell \end{cases}$$



ξ_i — slack variable,
 штраф за нарушение
 разделяющей полосы



Основные свойства метода опорных векторов

- 3 Двойственная задача относительно переменных λ_i :

$$\begin{cases} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle - 2 \sum_{i=1}^{\ell} \lambda_i \rightarrow \min_{\lambda}; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0; \quad 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell \end{cases}$$

- 4 Опорные векторы — объекты x_i : $M_i \leq 1$, $\lambda_i > 0$,
через которые линейно выражается вектор решения w :

$$w = \sum_{i=1}^{\ell} \lambda_i y_i x_i, \quad w_0 = \langle w, x_i \rangle - y_i \text{ для } i: M_i = 1$$

- 5 Разреженность решения: опорными векторами является,
как правило, небольшая доля объектов выборки

S.Fine, K.Scheinberg. INCAS: An incremental active set method for SVM. 2002.

J.Platt. Fast training support vector machines using sequential minimal optimization. 1999.

Основные свойства метода опорных векторов

- 6 Скалярные произведения $f_i(x) = \langle x, x_i \rangle$ с опорными x_i образуют новые признаки для линейного классификатора:

$$a(x; w, w_0) = \text{sign} \left(\sum_{i=1}^{\ell} \lambda_i y_i \langle x, x_i \rangle - w_0 \right)$$

- 7 Подстановка ядра (kernel trick) в двойственной задаче и в классификаторе делает модель нелинейной:

$$a(x; w, w_0) = \text{sign} \left(\sum_{i=1}^{\ell} \lambda_i y_i K(x, x_i) - w_0 \right)$$

- 8 Ядро — скалярное произведение $K(x, x') = \langle \psi(x), \psi(x') \rangle_H$ в «спрямляющем» гильбертовом пространстве H

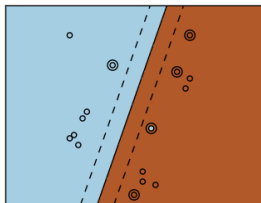
Классификация с различными ядрами

Гиперплоскость в спрямляющем пространстве соответствует нелинейной разделяющей поверхности в исходном.

Примеры с различными ядрами $K(x, x')$

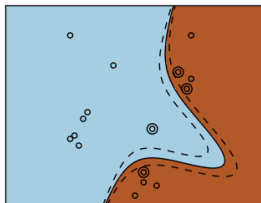
линейное

$$\langle x, x' \rangle$$



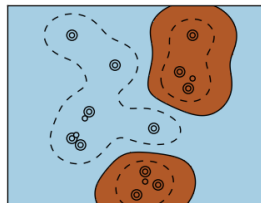
полиномиальное

$$(\langle x, x' \rangle + 1)^d, \quad d=3$$



гауссовское (RBF)

$$\exp(-\gamma \|x - x'\|^2)$$



На практике C и K подбирают под задачу перебором

Пример из Python SkLearn: <http://scikit-learn.org/dev>

Задача обучения SVM-регрессии

Дано: обучающая выборка
 $(x_i, y_i)_{i=1}^{\ell}$, $x_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$

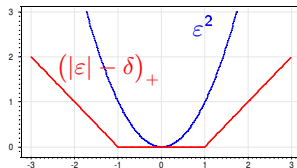
Найти: параметры w и w_0
линейной модели регрессии
 $a(x, w) = \langle w, x \rangle - w_0$

Критерий регуляризованного эмпирического риска:

$$\sum_{i=1}^{\ell} (|\langle w, x_i \rangle - w_0 - y_i| - \delta)_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}$$

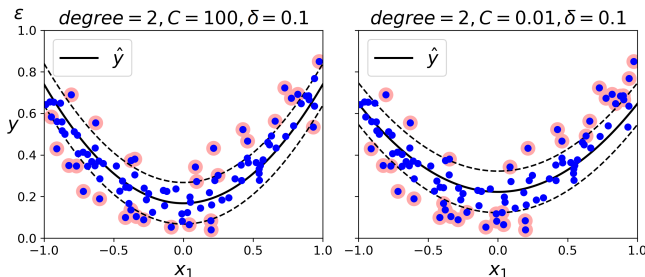
Свойства, общие с SVM-классификацией:

- задача выпуклая, решение единственно
- решение разреженное, выражается через опорные векторы
- возможна замена $\langle x, x_i \rangle$ нелинейным ядром $K(x, x_i)$



SVM-регрессия. Пример 1

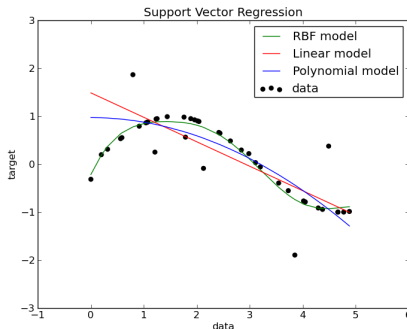
SVM-регрессия с полиномиальным ядром степени 2:



- Выделены опорные векторы
- Типично для SVM: результат обучения слабо зависит от коэффициента C , который подбирают по грубой сетке

SVM-регрессия. Пример 2

Сравнение SVM-регрессии с гауссовским (RBF) ядром, линейной и полиномиальной регрессией:



- SVM-регрессия игнорирует выбросы (робастность)
- Удачный выбор ядра имеет значение!

Резюме в конце лекции

- Метод стохастического градиента (SG) для ERM:
 - подходит для любых моделей и функций потерь
 - подходит для обучения по большим данным
- *Аппроксимация пороговой функции потерь $L(M)$*
 - непрерывная оптимизация для дискретных критериев
 - общий приём для классификации и ранжирования
- *Регуляризация* — в линейных моделях снижает переобучение, возникающее из-за мультиколлинеарности
- *Логистическая регрессия* — метод классификации, оценивающий условные вероятности классов $P(y|x)$
- *Метод опорных векторов* — изящно обобщается ядрами для нелинейной классификации и регрессии
- *Негладкость функции потерь* \Rightarrow отбор объектов (SVM)
- *Негладкость регуляризатора* \Rightarrow отбор признаков (LASSO)