

Методы машинного обучения

Линейные модели классификации

Воронцов Константин Вячеславович

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

вопросы к лектору: k.vorontsov@iai.msu.ru

материалы курса:

github.com/MSU-ML-COURSE/ML-COURSE-24-25

орг.вопросы по курсу: ml.cmc@mail.ru

1 Минимизация эмпирического риска

- Регрессия и классификация
- Понятия отступа в задачах классификации
- Многоклассовая классификация

2 Метод стохастического градиента

- Градиентный метод оптимизации
- Метод стохастической аппроксимации
- Ускоренные градиентные методы

3 Эвристики

- Инициализация весов и порядок объектов
- Выбор величины градиентного шага
- Регуляризация и отбор признаков

Оптимизационная задача обучения регрессии

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

- ❶ Модель регрессии — *линейная* с параметром $w \in \mathbb{R}^n$:

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

- ❷ Функция потерь — *квадратичная*:

$$\mathcal{L}(w, x) = (a(x, w) - y(x))^2$$

- ❸ Метод обучения — *метод наименьших квадратов*:

$$Q(w) = \sum_{i=1}^{\ell} \mathcal{L}(w, x_i) = \sum_{i=1}^{\ell} (a(x_i, w) - y_i)^2 \rightarrow \min_w$$

- ❹ Проверка по тестовой выборке $X^k = (\tilde{x}_i, \tilde{y}_i)_{i=1}^k$:

$$\tilde{Q}(w) = \frac{1}{k} \sum_{i=1}^k (a(\tilde{x}_i, w) - \tilde{y}_i)^2$$

Оптимизационная задача обучения бинарной классификации

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

- ❶ Модель классификации — *линейная* с параметром $w \in \mathbb{R}^n$:

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

- ❷ Функция потерь — бинарная, заменяем её **аппроксимацией**:

$$\mathcal{L}(w, x) = [a(x, w)y(x) < 0] = [\langle x, w \rangle y(x) < 0] \leq L(\langle x, w \rangle y(x))$$

- ❸ Метод обучения — *минимизация эмпирического риска*:

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

- ❹ Проверка по тестовой выборке $X^k = (\tilde{x}_i, \tilde{y}_i)_{i=1}^k$:

$$\tilde{Q}(w) = \frac{1}{k} \sum_{i=1}^k [\langle \tilde{x}_i, w \rangle \tilde{y}_i < 0]$$

Бинарный разделяющий классификатор (margin-based classifier)

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

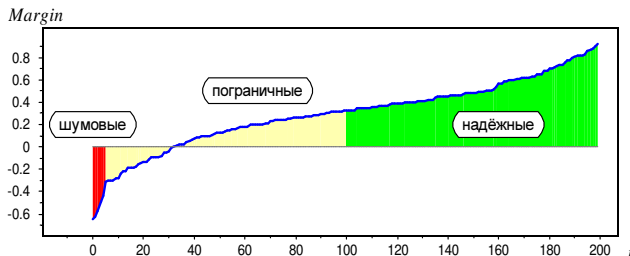
$g(x, w)$ — разделяющая (дискриминантная) функция

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

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

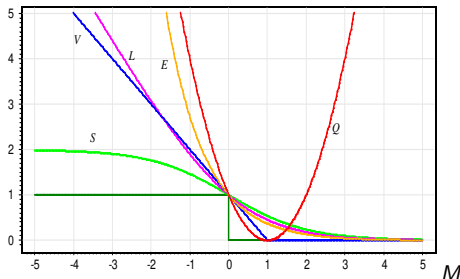
$M_i(w) < 0 \iff$ алгоритм $a(x, w)$ ошибается на x_i

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



Непрерывные аппроксимации пороговой функции потерь

Часто используемые непрерывные функции потерь $L(M)$:



$$V(M) = (1 - M)_+$$

— кусочно-линейная (SVM);

$$H(M) = (-M)_+$$

— кусочно-линейная (Hebb's rule);

$$L(M) = \log_2(1 + e^{-M})$$

— логарифмическая (LR);

$$Q(M) = (1 - M)^2$$

— квадратичная (FLD);

$$S(M) = 2(1 + e^M)^{-1}$$

— сигмоидная (ANN);

$$E(M) = e^{-M}$$

— экспоненциальная (AdaBoost);

$[M < 0]$

— пороговая функция потерь.

Линейный классификатор — математическая модель нейрона

Линейная модель нейрона МакКаллока-Питтса [1943]:

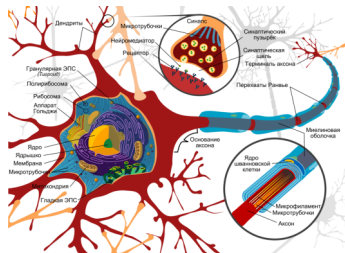
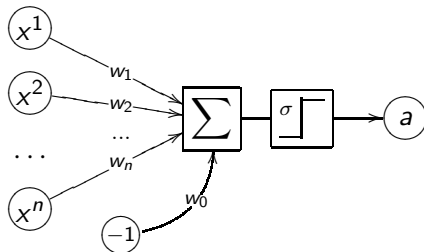
$$a(x, w) = \sigma(\langle w, x \rangle) = \sigma\left(\sum_{j=1}^n w_j f_j(x) - w_0\right),$$

$\sigma(z)$ — функция активации (например, sign),

w_j — весовые коэффициенты синаптических связей,

w_0 — порог активации,

$w, x \in \mathbb{R}^{n+1}$, если ввести константный признак $f_0(x) \equiv -1$



Задача многоклассовой классификации (multiclass classification)

Обучающая выборка: $X^\ell = (x_i, y_i)_{i=1}^\ell$, $x_i \in \mathbb{R}^n$, $y_i = y(x_i) \in Y$

- 1 Модель классификации — *линейная*, $w = (w_y : y \in Y)$:

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

- 2 Функция потерь — бинарная, заменяем её аппроксимацией:

$$\mathcal{L}(w, x) = \sum_{z \neq y(x)} [\langle x, w_{y(x)} \rangle < \langle x, w_z \rangle] \leq \sum_{z \neq y(x)} L(\langle x, w_{y(x)} - w_z \rangle)$$

- 3 Метод обучения — *минимизация эмпирического риска*:

$$Q(w) = \sum_{i=1}^{\ell} \sum_{z \neq y_i} L(\langle x_i, w_{y_i} - w_z \rangle) \rightarrow \min_w$$

- 4 Проверка по тестовой выборке $X^k = (\tilde{x}_i, \tilde{y}_i)_{i=1}^k$

Многоклассовый разделяющий классификатор

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

$g_y(x, w_y)$ — дискриминантная функция класса $y \in Y$

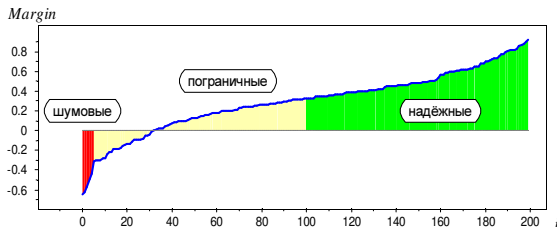
x : $g_y(x, w_y) = g_z(x, w_z)$ — разделяющая поверхность между y, z

$M_{iy}(w) = g_{y_i}(x_i, w_{y_i}) - g_y(x_i, w_y)$ — отступ объекта x_i по классу y

$M_i(w) = \min_{y \neq y_i} M_{iy}(w)$ — отступ (margin) объекта x_i

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

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



Градиентный метод оптимизации

Минимизация эмпирического риска:

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

Метод *градиентного спуска*:

$w^{(0)}$:= начальное приближение;

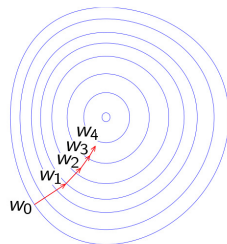
$$w^{(t+1)} := w^{(t)} - h \cdot \nabla Q(w^{(t)}), \quad \nabla Q(w) = \left(\frac{\partial Q(w)}{\partial w_j} \right)_{j=0}^n,$$

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

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

Идея ускорения сходимости:

брать объекты x_i по одному и сразу обновлять вектор весов



Алгоритм SG (Stochastic Gradient)

$$Q = \sum_{i=1}^{\ell} \mathcal{L}(w, x_i) \rightarrow \min_w \text{ — минимизация эмпирического риска}$$

Вход: выборка X^ℓ , темп обучения h , темп забывания λ ;

Выход: вектор весов w ;

- 1 инициализировать веса $w_j, j = 0, \dots, n$;
- 2 инициализировать оценку функционала:
 $Q := \text{среднее } \mathcal{L}(w, x_i) \text{ по случайному подмножеству } \{x_i\}$;
- 3 **повторять**
 - 4 | выбрать объект x_i из X^ℓ случайным образом;
 - 5 | вычислить потерю: $\varepsilon_i := \mathcal{L}(w, x_i)$;
 - 6 | сделать градиентный шаг: $w := w - h \nabla \mathcal{L}(w, x_i)$;
 - 7 | оценить функционал: $Q := \lambda \varepsilon_i + (1 - \lambda)Q$;
- 8 **пока** значение Q и/или веса w не сойдутся;

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

Откуда взялась 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:

$$\bar{Q}_m = \frac{1}{m}\varepsilon_m + \frac{1}{m}\varepsilon_{m-1} + \frac{1}{m}\varepsilon_{m-2} + \dots$$

$$\bar{Q}_m = \frac{1}{m}\varepsilon_m + \left(1 - \frac{1}{m}\right)\bar{Q}_{m-1}$$

Экспоненциальное скользящее среднее:

$$\bar{Q}_m = \lambda\varepsilon_m + (1 - \lambda)\lambda\varepsilon_{m-1} + (1 - \lambda)^2\lambda\varepsilon_{m-2} + \dots$$

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

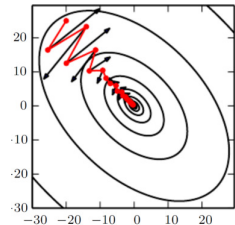
Параметр λ (порядка $\frac{1}{m}$) — темп забывания предыстории ряда.

Метод накопления инерции (momentum)

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

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

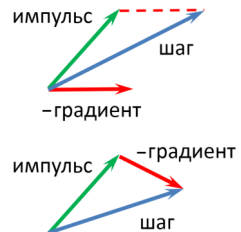
$$w := w - hv$$



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

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

$$w := w - hv$$



Варианты инициализации весов

- ❶ $w_j := 0$ для всех $j = 0, \dots, n$;
- ❷ небольшие случайные значения:
 $w_j := \text{random} \left(-\frac{1}{2n}, \frac{1}{2n} \right)$;
- ❸ $w_j := \frac{\langle y, f_j \rangle}{\langle f_j, f_j \rangle}$, $f_j = (f_j(x_i))_{i=1}^\ell$ — вектор значений признака.

Эта оценка w оптимальна, если

- 1) функция потерь квадратична и
- 2) признаки некоррелированы, $\langle f_j, f_k \rangle = 0$, $j \neq k$.

- ❹ обучение по небольшой случайной подвыборке объектов;
- ❺ мультистарт: многократные запуски из разных случайных начальных приближений и выбор лучшего решения.

Варианты порядка предъявления объектов

Возможны варианты:

- ❶ *перетасовка объектов (shuffling)*:
попеременно брать объекты из разных классов;
- ❷ чаще брать объекты, на которых ошибка больше:
чем меньше M_i , тем больше вероятность взять объект;
- ❸ чаще брать объекты, на которых уверенность меньше:
чем меньше $|M_i|$, тем больше вероятность взять объект;
- ❹ вообще не брать «хорошие» объекты, у которых $M_i > \mu_+$
(при этом немного ускоряется сходимость);
- ❺ вообще не брать объекты-«выбросы», у которых $M_i < \mu_-$
(при этом может улучшиться качество классификации);

Параметры μ_+ , μ_- придётся подбирать.

Варианты выбора градиентного шага

- ❶ сходимость гарантируется (для выпуклых функций) при

$$h_t \rightarrow 0, \quad \sum_{t=1}^{\infty} h_t = \infty, \quad \sum_{t=1}^{\infty} h_t^2 < \infty,$$

в частности можно положить $h_t = 1/t$

- ❷ *метод скорейшего градиентного спуска* основан на поиске оптимального *адаптивного шага* h^* :

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

При квадратичной функции потерь $h^* = \|x_i\|^{-2}$

- ❸ пробные случайные шаги для «выбивания» итерационного процесса из локальных экстремумов
- ❹ метод Левенберга-Марквардта (сходимость второго порядка)

Диагональный метод Левенберга-Марквардта

Метод Ньютона-Рафсона, $\mathcal{L}(w, x_i) \equiv L(\langle w, x_i \rangle y_i)$:

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

где $\mathcal{L}''(w, x_i) = \left(\frac{\partial^2 \mathcal{L}(w, x_i)}{\partial w_j \partial w_{j'}} \right)$ — гессиан, $n \times n$ -матрица

Эвристика. Считаем, что гессиан диагонален:

$$w_j := w_j - h \max \left\{ \mu, \frac{\partial^2 \mathcal{L}(w, x_i)}{\partial w_j^2} \right\}^{-1} \frac{\partial \mathcal{L}(w, x_i)}{\partial w_j},$$

$h > 0$ — темп обучения, можно полагать $h = 1$,

$\mu > 0$ — параметр, предотвращающий обнуление знаменателя.

Вблизи минимума сходимость второго порядка с темпом h

Вдали от минимума сходимость первого порядка с темпом $\frac{h}{\mu}$

Мультиколлинеарность и переобучение в линейных моделях

Причины — те же, что и в линейных моделях регрессии:

- объектов меньше, чем признаков, либо
- признаки линейно зависимы (мультиколлинеарны):
если $\exists u \in \mathbb{R}^n: \forall x_i \in X^\ell \langle u, x_i \rangle = 0$, то решение
не единственно и не устойчиво: $\forall \gamma \in \mathbb{R} \langle w + \gamma u, x_i \rangle = \langle w, x_i \rangle$

Проявления мультиколлинеарности:

- слишком большие веса $|w_j|$ разных знаков
- вес w_j не интерпретируется как важность признака f_j
- переобучение: $Q(w^*, X^\ell) \ll Q(w^*, X^k)$

Способы устранения мультиколлинеарности и переобучения:

- регуляризация: $\|w\| \rightarrow \min$;
- отбор признаков: $f_1, \dots, f_n \rightarrow f_{j_1}, \dots, f_{j_m}, \quad m \ll n$.
- преобразование признаков: $f_1, \dots, f_n \rightarrow g_1, \dots, g_m, \quad m \ll n$;

L_2 -регуляризация (сокращение весов, weight decay)

Штраф за увеличение нормы вектора весов:

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

Градиент:

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

Модификация градиентного шага:

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

Методы подбора коэффициента регуляризации τ :

- hold-out или скользящий контроль
- стохастическая адаптация

Негладкие регуляризаторы для отбора и группировки признаков

Общий вид регуляризаторов (μ — параметр селективности):

$$\sum_{i=1}^{\ell} \mathcal{L}(w, x_i) + \sum_{j=1}^n R_{\mu}(\alpha_j) \rightarrow \min_{\alpha}$$

Регуляризаторы с эффектами отбора и группировки признаков:

LASSO (L_1): $R_{\mu}(\alpha) = \mu|\alpha|$

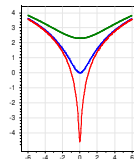
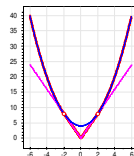
Elastic Net: $R_{\mu}(\alpha) = \mu|\alpha| + \tau\alpha^2$

Support Feature Machine (SFM):

$$R_{\mu}(\alpha) = \begin{cases} 2\mu|\alpha|, & |\alpha| \leq \mu; \\ \mu^2 + \alpha^2, & |\alpha| \geq \mu; \end{cases}$$

Relevance Feature Machine (RFM):

$$R_{\mu}(\alpha) = \ln(\mu\alpha^2 + 1)$$



Резюме в конце лекции

- Метод стохастического градиента (SG)
 - подходит для любых моделей и функций потерь
 - подходит для обучения по большим данным
- *Аппроксимация пороговой функции потерь $L(M)$*
позволяет использовать градиентную оптимизацию
- Функции $L(M)$, штрафующие за приближение к границе классов, увеличивают зазор между классами, благодаря чему повышается надёжность классификации
- *Регуляризация* снижает переобучение, возникающее в линейных моделях из-за мультиколлинеарности
- *Недостаток*: подбор эвристик является искусством (не забыть про сходимость, застревание, переобучение,...)