

Обобщающая способность Методы отбора признаков

К. В. Воронцов, А. В. Зухба
vokov@forecsys.ru
a__l@mail.ru

20 февраля 2018 г.

Содержание

1 Методы отбора признаков

- Типы признаков (напоминание)
- Жадные алгоритмы и стратегии перебора
- Стохастический поиск

2 Измерение качества и выбор модели

- Качество регрессии
- Качество классификации
- Выбор модели

Задача статистического (машинного) обучения с учителем

Задача восстановления зависимости $y: X \rightarrow Y$
по точкам *обучающей выборки* (x_i, y_i) , $i = 1, \dots, \ell$.

Дано: $x_i \mapsto (f_1(x_i), \dots, f_n(x_i))$ — объекты обучающей выборки,
 $y_i = y(x_i)$ — правильные ответы, $i = 1, \dots, \ell$:

$$\begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix} \xrightarrow{y} \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}$$

Найти: функцию $a(x)$, способную давать правильные ответы
на *тестовых объектах* $\tilde{x}_i \mapsto (f_1(\tilde{x}_i), \dots, f_n(\tilde{x}_i))$, $i = 1, \dots, k$:

$$\begin{pmatrix} f_1(\tilde{x}_1) & \dots & f_n(\tilde{x}_1) \\ \dots & \dots & \dots \\ f_1(\tilde{x}_k) & \dots & f_n(\tilde{x}_k) \end{pmatrix} \xrightarrow{a?} \begin{pmatrix} a(\tilde{x}_1) \\ \dots \\ a(\tilde{x}_k) \end{pmatrix}$$

Типы признаков и типы задач

Типы признаков, $f_j: X \rightarrow D_j$, в зависимости от D_j :

- $D_j = \{0, 1\}$ — *бинарный* признак;
- $|D_j| < \infty$ — *номинальный* признак;
- D_j упорядочено — *порядковый* признак;
- $D_j = \mathbb{R}$ — *количественный* признак.

Типы задач, в зависимости от Y :

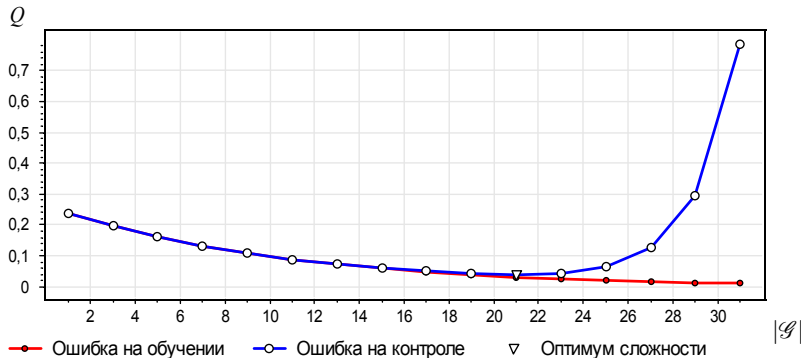
- $Y = \{0, 1\}$ или $Y = \{-1, +1\}$ — *классификация* на 2 класса;
- $Y = \{1, \dots, M\}$ — на M непересекающихся классов;
- $Y = \{0, 1\}^M$ — на M классов, которые могут пересекаться;
- $Y = \mathbb{R}$ — задача восстановления регрессии;
- Y упорядочено — задача ранжирования (learning to rank).

Задача отбора признаков

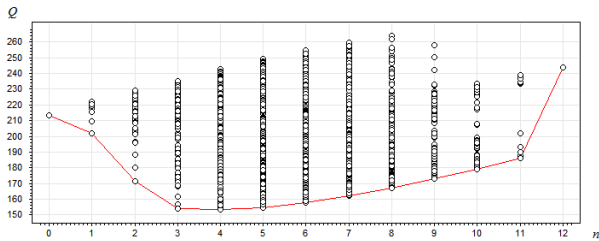
$\mathcal{F} = \{f_j: X \rightarrow D_j: j = 1, \dots, n\}$ — множество признаков;

$\mu_{\mathcal{G}}$ — метод обучения, использующий только признаки $\mathcal{G} \subseteq \mathcal{F}$;

$Q(\mathcal{G}) = Q(\mu_{\mathcal{G}}, X^\ell)$ — какой-либо критерий.



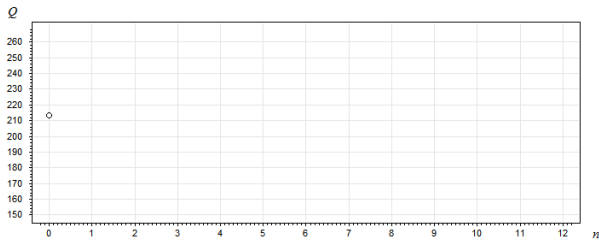
Алгоритм полного перебора (Full Search)



Вход: множество \mathcal{F} , критерий Q , параметр d ;

- 1: $Q^* := Q(\emptyset)$; — инициализация;
- 2: **для всех** $j = 1, \dots, n$, где j — сложность наборов:
- 3: найти лучший набор сложности j :
 $\mathcal{G}_j := \arg \min_{\mathcal{G}: |\mathcal{G}|=j} Q(\mathcal{G})$;
- 4: **если** $Q(\mathcal{G}_j) < Q^*$ **то** $j^* := j$; $Q^* := Q(\mathcal{G}_j)$;
- 5: **если** $j - j^* \geq d$ **то вернуть** \mathcal{G}_{j^*} ;

Алгоритм полного перебора (Full Search)

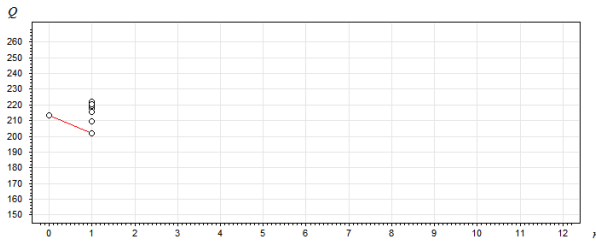


$j = 0$

Вход: множество \mathcal{F} , критерий Q , параметр d ;

- 1: $Q^* := Q(\emptyset)$; — инициализация;
- 2: **для всех** $j = 1, \dots, n$, где j — сложность наборов:
- 3: найти лучший набор сложности j :
 $\mathcal{G}_j := \arg \min_{\mathcal{G}: |\mathcal{G}|=j} Q(\mathcal{G})$;
- 4: **если** $Q(\mathcal{G}_j) < Q^*$ **то** $j^* := j$; $Q^* := Q(\mathcal{G}_j)$;
- 5: **если** $j - j^* \geq d$ **то вернуть** \mathcal{G}_{j^*} ;

Алгоритм полного перебора (Full Search)

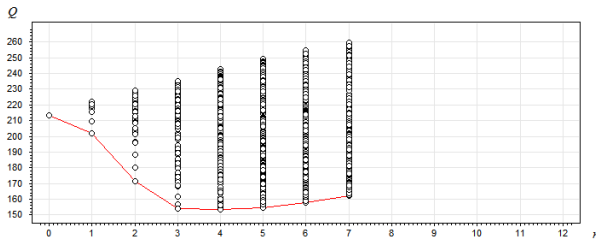


$$j = 1$$
$$j^* = 1$$

Вход: множество \mathcal{F} , критерий Q , параметр d ;

- 1: $Q^* := Q(\emptyset)$; — инициализация;
- 2: **для всех** $j = 1, \dots, n$, где j — сложность наборов:
- 3: найти лучший набор сложности j :
$$\mathcal{G}_j := \arg \min_{\mathcal{G}: |\mathcal{G}|=j} Q(\mathcal{G});$$
- 4: **если** $Q(\mathcal{G}_j) < Q^*$ **то** $j^* := j$; $Q^* := Q(\mathcal{G}_j)$;
- 5: **если** $j - j^* \geq d$ **то вернуть** \mathcal{G}_{j^*} ;

Алгоритм полного перебора (Full Search)



$$j = 7$$
$$j^* = 4$$

Вход: множество \mathcal{F} , критерий Q , параметр d ;

- 1: $Q^* := Q(\emptyset)$; — инициализация;
- 2: **для всех** $j = 1, \dots, n$, где j — сложность наборов:
- 3: найти лучший набор сложности j :
 $\mathcal{G}_j := \arg \min_{\mathcal{G}: |\mathcal{G}|=j} Q(\mathcal{G})$;
- 4: **если** $Q(\mathcal{G}_j) < Q^*$ **то** $j^* := j$; $Q^* := Q(\mathcal{G}_j)$;
- 5: **если** $j - j^* \geq d$ **то вернуть** \mathcal{G}_{j^*} ;

Алгоритм полного перебора (Full Search)

Преимущества:

- простота реализации;
- гарантированный результат;
- полный перебор эффективен, когда
 - информативных признаков не много, $j^* \lesssim 5$;
 - всего признаков не много, $n \lesssim 20..100$.

Недостатки:

- в остальных случаях оооооочень долго — $O(2^n)$;
- чем больше перебирается вариантов, тем больше переобучение (особенно, если лучшие из вариантов существенно различны и одинаково плохи).

Способы устранения:

- эвристические методы сокращённого перебора.

Алгоритм жадного добавления (Add)

Вход: множество \mathcal{F} , критерий Q , параметр d ;

- 1: $\mathcal{G}_0 := \emptyset$; $Q^* := Q(\emptyset)$; — инициализация;
- 2: **для всех** $j = 1, \dots, n$, где j — сложность наборов:
- 3: найти признак, наиболее выгодный для добавления:
$$f^* := \arg \min_{f \in \mathcal{F} \setminus \mathcal{G}_{j-1}} Q(\mathcal{G}_{j-1} \cup \{f\});$$
- 4: добавить этот признак в набор:
$$\mathcal{G}_j := \mathcal{G}_{j-1} \cup \{f^*\};$$
- 5: **если** $Q(\mathcal{G}_j) < Q^*$ **то** $j^* := j$; $Q^* := Q(\mathcal{G}_j)$;
- 6: **если** $j - j^* \geq d$ **то вернуть** \mathcal{G}_{j^*} ;

Алгоритм жадного добавления (Add)

Преимущества:

- работает быстро — $O(n^2)$, точнее $O(n(j^* + d))$;
- возможны быстрые инкрементные алгоритмы, пример — *шаговая регрессия* (step-wise regression).

Недостатки:

- Add склонен включать в набор лишние признаки.

Способы устранения:

- Add-Del — чередование добавлений и удалений (см. далее);
- поиск в ширину (см. ещё далее).

Алгоритм поочерёдного добавления и удаления (Add-Del)

1: $\mathcal{G}_0 := \emptyset$; $Q^* := Q(\emptyset)$; $t := 0$; — инициализация;

2: **повторять**

3: **пока** $|\mathcal{G}_t| < n$ добавлять признаки (Add):

4: $t := t + 1$; — началась следующая итерация;

5: $f^* := \arg \min_{f \in \mathcal{F} \setminus \mathcal{G}_{t-1}} Q(\mathcal{G}_{t-1} \cup \{f\})$; $\mathcal{G}_t := \mathcal{G}_{t-1} \cup \{f^*\}$;

6: **если** $Q(\mathcal{G}_t) < Q^*$ **то** $t^* := t$; $Q^* := Q(\mathcal{G}_t)$;

7: **если** $t - t^* \geq d$ **то прервать цикл**;

8: **пока** $|\mathcal{G}_t| > 0$ удалять признаки (Del):

9: $t := t + 1$; — началась следующая итерация;

10: $f^* := \arg \min_{f \in \mathcal{G}_{t-1}} Q(\mathcal{G}_{t-1} \setminus \{f\})$; $\mathcal{G}_t := \mathcal{G}_{t-1} \setminus \{f^*\}$;

11: **если** $Q(\mathcal{G}_t) < Q^*$ **то** $t^* := t$; $Q^* := Q(\mathcal{G}_t)$;

12: **если** $t - t^* \geq d$ **то прервать цикл**;

13: **пока** значения критерия $Q(\mathcal{G}_{t^*})$ уменьшаются;

14: **вернуть** \mathcal{G}_{t^*} ;

Алгоритм поочерёдного добавления и удаления (Add-Del)

Преимущества:

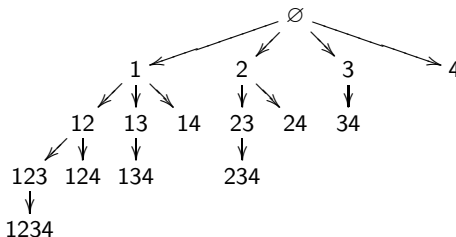
- как правило, лучше, чем Add и Del по отдельности;
- возможны быстрые инкрементные алгоритмы, пример — *шаговая регрессия* (step-wise regression).

Недостатки:

- работает дольше, оптимальность не гарантирует.

Поиск в глубину (DFS, метод ветвей и границ)

Пример: дерево наборов признаков, $n = 4$



Основные идеи:

- нумерация признаков по возрастанию номеров — чтобы избежать повторов при переборе подмножеств;
- если набор \mathcal{S} бесперспективен, то больше не пытаться его наращивать.

Поиск в глубину (DFS, метод ветвей и границ)

Обозначим Q_j^* — значение критерия на самом лучшем наборе мощности j из всех до сих пор просмотренных.

Оценка бесперспективности набора признаков \mathcal{G} :
набор \mathcal{G} не наращивается, если

$$\exists j: \quad Q(\mathcal{G}) \geq \kappa Q_j^* \quad \text{и} \quad |\mathcal{G}| \geq j + d,$$

$d \geq 0$ — целочисленный параметр,

$\kappa \geq 1$ — вещественный параметр.

Чем меньше d и κ , тем сильнее сокращается перебор.

Поиск в глубину (DFS, метод ветвей и границ)

Вход: множество \mathcal{F} , критерий Q , параметры d и κ ;

- 1: **ПРОЦЕДУРА** Нарастить (\mathcal{G});
 - 2: **если** найдётся $j \leq |\mathcal{G}| - d$ такое, что $Q(\mathcal{G}) \geq \kappa Q_j^*$, **то**
 - 3: **выход**;
 - 4: $Q_{|\mathcal{G}|}^* := \min\{Q_{|\mathcal{G}|}^*, Q(\mathcal{G})\}$;
 - 5: **для всех** $f_s \in \mathcal{F}$ таких, что $s > \max\{t \mid f_t \in \mathcal{G}\}$
 Нарастить ($\mathcal{G} \cup \{f_s\}$);
-

- 6: Инициализация массива лучших значений критерия:
 $Q_j^* := Q(\emptyset)$ для всех $j = 1, \dots, n$;
- 7: Упорядочить признаки по убыванию информативности;
- 8: Нарастить (\emptyset);
- 9: **вернуть** \mathcal{G} , для которого $Q(\mathcal{G}) = \min_{j=1, \dots, n} Q_j^*$;

Поиск в ширину (BFS)

Он же *многорядный итерационный алгоритм МГУА*
(МГУА — метод группового учёта аргументов).

Философский принцип *неокончательных решений* Габора:
принимая решения, следует оставлять максимальную свободу
выбора для принятия последующих решений.

Усовершенствуем алгоритм Add:
на каждой j -й итерации будем строить не один набор,
а множество из B_j наборов, называемое j -м *рядом*:

$$R_j = \{\mathcal{G}_j^1, \dots, \mathcal{G}_j^{B_j}\}, \quad \mathcal{G}_j^b \subseteq \mathcal{F}, \quad |\mathcal{G}_j^b| = j, \quad b = 1, \dots, B_j.$$

где $B_j \leq B$ — параметр *ширины поиска*.

Поиск в ширину (BFS)

Вход: множество \mathcal{F} , критерий Q , параметры d, B ;

- 1: первый ряд состоит из всех наборов длины 1:
 $R_1 := \{\{f_1\}, \dots, \{f_n\}\}; \quad Q^* = Q(\emptyset);$
- 2: **для всех** $j = 1, \dots, n$, где j — сложность наборов:
- 3: отсортировать ряд $R_j = \{\mathcal{G}_j^1, \dots, \mathcal{G}_j^{B_j}\}$
по возрастанию критерия: $Q(\mathcal{G}_j^1) \leq \dots \leq Q(\mathcal{G}_j^{B_j});$
- 4: **если** $B_j > B$ **то**
- 5: $R_j := \{\mathcal{G}_j^1, \dots, \mathcal{G}_j^B\};$ — B лучших наборов ряда;
- 6: **если** $Q(\mathcal{G}_j^1) < Q^*$ **то** $j^* := j; \quad Q^* := Q(\mathcal{G}_j^1);$
- 7: **если** $j - j^* \geq d$ **то вернуть** $\mathcal{G}_{j^*}^1;$
- 8: породить следующий ряд:
 $R_{j+1} := \{\mathcal{G} \cup \{f\} \mid \mathcal{G} \in R_j, f \in \mathcal{F} \setminus \mathcal{G}\};$

Поиск в ширину (BFS)

- **Трудоёмкость:**
 $O(Bn^2)$, точнее $O(Bn(j^* + d))$.
- **Проблема дубликатов:**
после сортировки (шаг 3) проверить на совпадение только соседние наборы с равными значениями внутреннего и внешнего критерия.
- **Адаптивный отбор признаков:**
на шаге 8 добавлять к j -му ряду только признаки f с наибольшей информативностью $I_j(f)$:

$$I_j(f) = \sum_{b=1}^{B_j} [f \in \mathcal{G}_j^b].$$

Генетический алгоритм поиска (идея и терминология)

$\mathcal{G} \subseteq \mathcal{F}$ — индивид (в МГУА «модель»);

$R_t := \{\mathcal{G}_t^1, \dots, \mathcal{G}_t^{B_t}\}$ — поколение (в МГУА — «ряд»);

$\beta = (\beta_j)_{j=1}^n$, $\beta_j = [f_j \in \mathcal{G}]$ — хромосома, кодирующая \mathcal{G} ;

Бинарная операция скрещивания $\beta = \beta' \times \beta''$:

$$\beta_j = \begin{cases} \beta'_j, & \text{с вероятностью } 1/2; \\ \beta''_j, & \text{с вероятностью } 1/2; \end{cases}$$

Унарная операция мутации $\beta = \sim \beta'$

$$\beta_j = \begin{cases} 1 - \beta'_j, & \text{с вероятностью } p_m; \\ \beta'_j, & \text{с вероятностью } 1 - p_m; \end{cases}$$

где параметр p_m — вероятность мутации.

Генетический (эволюционный) алгоритм

Вход: множество \mathcal{F} , критерий Q , параметры: d , p_m ,
 B — размер популяции, T — число поколений;

-
- 1: инициализировать случайную популяцию из B наборов:
 $B_1 := B$; $R_1 := \{\mathcal{G}_1^1, \dots, \mathcal{G}_1^{B_1}\}$; $Q^* := Q(\emptyset)$;
 - 2: **для всех** $t = 1, \dots, T$, где t — номер поколения:
 - 3: ранжирование индивидов: $Q(\mathcal{G}_t^1) \leq \dots \leq Q(\mathcal{G}_t^{B_t})$;
 - 4: **если** $B_t > B$ **то**
 - 5: селекция: $R_t := \{\mathcal{G}_t^1, \dots, \mathcal{G}_t^B\}$;
 - 6: **если** $Q(\mathcal{G}_t^1) < Q^*$ **то** $t^* := t$; $Q^* := Q(\mathcal{G}_t^1)$;
 - 7: **если** $t - t^* \geq d$ **то вернуть** $\mathcal{G}_{t^*}^1$;
 - 8: породить $t+1$ -е поколение путём скрещиваний и мутаций:
 $R_{t+1} := \{\sim(\mathcal{G}' \times \mathcal{G}'') \mid \mathcal{G}', \mathcal{G}'' \in R_t\} \cup R_t$;

Эвристики для управления процессом эволюции

- Увеличивать вероятности перехода признаков от более успешного родителя к потомку.
- Накапливать оценки информативности признаков.
Чем более информативен признак, тем выше вероятность его включения в набор во время мутации.
- Применение совокупности критериев качества.
- Скрещивать только лучшие индивиды (элитаризм).
- Переносить лучшие индивиды в следующее поколение.
- В случае стагнации увеличивать вероятность мутаций.
- Параллельно выращивается несколько изолированных популяций (островная модель эволюции).

Генетический (эволюционный) алгоритм

Преимущества:

- it is fun!
- возможность введения различных эвристик;
- решает задачи даже с очень большим числом признаков.

Недостатки:

- относительно медленная сходимость;
- отсутствие теории;
- подбор параметров — непростое искусство;

Случайный поиск — упрощенный генетический алгоритм

Модификация: шаг 8

- породить $t+1$ -е поколение путём многократных *мутаций*:

$$R_{t+1} := \{\sim \mathcal{G}, \dots, \sim \mathcal{G} \mid \mathcal{G} \in R_t\} \cup R_t;$$

Недостатки:

- ничем не лучше ГА;
- очень медленная сходимость.

Способ устранения:

- СПА — случайный поиск с адаптацией.

Основная идея адаптации:

- увеличивать вероятность появления тех признаков, которые часто входят в наилучшие наборы,
- одновременно уменьшать вероятность появления признаков, которые часто входят в наихудшие наборы.

Случайный поиск с адаптацией (СПА)

Вход: множество \mathcal{F} , критерий Q , параметры d, j_0, T, r, h ;

- 1: $p_1 = \dots = p_n := 1/n$; — равные вероятности признаков;
- 2: **для всех** $j = j_0, \dots, n$, где j — сложность наборов:
- 3: **для всех** $t = 1, \dots, T$, где t — номер итерации:
- 4: r случайных наборов признаков из распределения $\{p_1, \dots, p_n\}$:
 $R_{jt} := \{\mathcal{G}_{jt}^1, \dots, \mathcal{G}_{jt}^r\}$, $|\mathcal{G}_{jt}^1| = \dots = |\mathcal{G}_{jt}^r| = j$;
- 5: $\mathcal{G}_{jt}^{\min} := \arg \min_{\mathcal{G} \in R_{jt}} Q(\mathcal{G})$; — лучший из r наборов;
- 6: $\mathcal{G}_{jt}^{\max} := \arg \max_{\mathcal{G} \in R_{jt}} Q(\mathcal{G})$; — худший из r наборов;
- 7: $H := 0$; наказание для всех $f_s \in \mathcal{G}_{jt}^{\max}$:
 $\Delta p_s := \min\{p_s, h\}$; $p_s := p_s - \Delta p_s$; $H := H + \Delta p_s$;
- 8: поощрение для всех $f_s \in \mathcal{G}_{jt}^{\min}$: $p_s := p_s + H/j$;
- 9: $\mathcal{G}_j := \arg \min_{\mathcal{G} \in R_{j1}, \dots, R_{jT}} Q(\mathcal{G})$; — лучший набор сложности j ;
- 10: **если** $Q(\mathcal{G}_j) < Q^*$ **то** $j^* := j$; $Q^* := Q(\mathcal{G}_j)$;
- 11: **если** $j - j^* \geq d$ **то вернуть** \mathcal{G}_{j^*} ;

Случайный поиск с адаптацией (СПА)

Рекомендации по выбору параметров r, T, h :

$T \approx 10..50$ — число итераций;

$r \approx 20..100$ — число наборов, создаваемых на каждой итерации;

$h \approx \frac{1}{r n}$ — скорость адаптации;

Преимущества:

- трудоёмкость порядка $O(Tr(j^* + d))$ операций;
- меньшее число параметров, по сравнению с генетикой;
- довольно быстрая сходимость.

Недостатки:

- при большом числе признаков СПА малоэффективен.

Обучение регрессии — это оптимизация

Задача регрессии, $Y = \mathbb{R}$

- ❶ Выбираем *модель регрессии*, например, линейную:

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

- ❷ Выбираем функцию потерь (например, квадратичную):

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

- ❸ Минимизируем потери *методом наименьших квадратов*:

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

- ❹ Проверяем прогностическую (обобщающую) способность:

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

Loss Functions of dotted forecasts

Mean squared error (MSE):

$$MSE = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2.$$

Mean absolute error (MAE):

$$MAE = \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - y_i|.$$

Mean absolute percentage error (MAPE):

$$MAPE = \frac{100}{N} \sum_{i=1}^N \left| \frac{\hat{y}_i - y_i}{y_i} \right|.$$

Symmetric mean absolute percentage error (SMAPE):

$$SMAPE = \frac{200}{N} \sum_{i=1}^N \left| \frac{\hat{y}_i - y_i}{\hat{y}_i + y_i} \right|.$$

Обучение классификации — тоже оптимизация

Задача классификации, $Y = \{-1, +1\}$

- 1 Выбираем *модель классификации*, например, линейную:

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

- 2 Выбираем функцию потерь (бинарную или её аппроксимацию):

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

- 3 Минимизируем частоту ошибок на обучающей выборке:

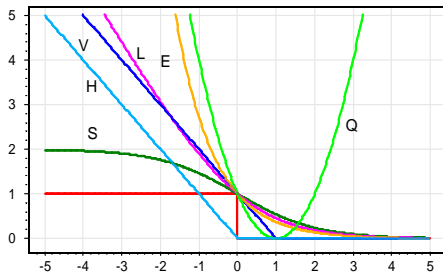
$$Q(w) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i, w) y_i < 0] \leq \frac{1}{\ell} \sum_{i=1}^{\ell} \mathcal{L}(\langle x_i, w \rangle y_i) \rightarrow \min_w$$

- 4 Проверяем прогностическую (обобщающую) способность:

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

Непрерывные аппроксимации пороговой функции потерь

Часто используемые непрерывные функции потерь $\mathcal{L}(M)$:



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

— 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.

Балансировка ошибок I и II рода на примере линейного алгоритма

Задача классификации на два класса, $Y = \{-1, +1\}$;
 λ_y — штраф за ошибку на объекте класса y ;
модель классификации: $a(x, w, w_0) = \text{sign}(f(x, w) - w_0)$.

- $f(x, w) = \langle x, w \rangle$ — не зависит от $\{\lambda_y\}$;
- $w_0 = \ln \frac{\lambda_-}{\lambda_+}$ — зависит только от $\{\lambda_y\}$.

На практике штрафы $\{\lambda_y\}$ могут многократно пересматриваться.

Постановка задачи

- Нужен удобный способ выбора w_0 в зависимости от $\{\lambda_y\}$, не требующий построения w заново.
- Нужна характеристика качества классификатора, инвариантная относительно выбора $\{\lambda_y\}$.

Определение ROC-кривой

ROC — «receiver operating characteristic».

- Каждая точка кривой соответствует некоторому $a(x; w, w_0)$.
- по оси X : доля *ошибочных положительных классификаций* (FPR — false positive rate):

$$\text{FPR}(a, X^\ell) = \frac{\sum_{i=1}^{\ell} [y_i = -1][a(x_i; w, w_0) = +1]}{\sum_{i=1}^{\ell} [y_i = -1]};$$

$1 - \text{FPR}(a)$ называется *специфичностью* алгоритма a .

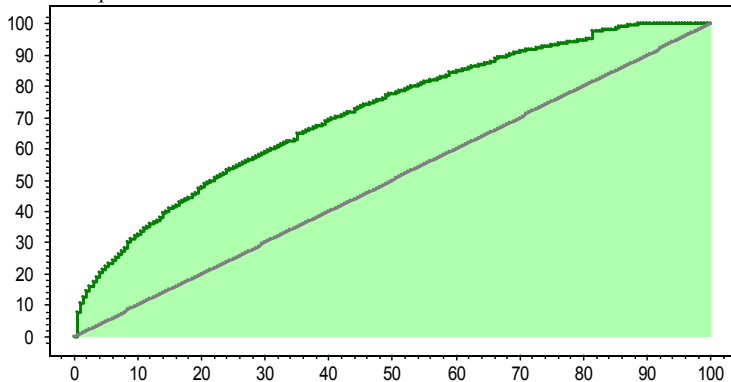
- по оси Y : доля *правильных положительных классификаций* (TPR — true positive rate):

$$\text{TPR}(a, X^\ell) = \frac{\sum_{i=1}^{\ell} [y_i = +1][a(x_i; w, w_0) = +1]}{\sum_{i=1}^{\ell} [y_i = +1]};$$

$\text{TPR}(a)$ называется также *чувствительностью* алгоритма a .

Пример ROC-кривой

TPR , true positive rate, %



FPR , false positive rate, %

AUC, площадь под ROC-кривой

наихудшая ROC-кривая

Алгоритм эффективного построения ROC-кривой

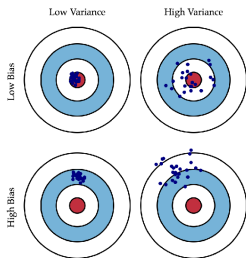
Вход: выборка X^ℓ ; дискриминантная функция $f(x, w)$;

Выход: $\{(FPR_i, TPR_i)\}_{i=0}^\ell$, AUC — площадь под ROC-кривой.

- 1: $\ell_y := \sum_{i=1}^\ell [y_i = y]$, для всех $y \in Y$;
- 2: упорядочить выборку X^ℓ по убыванию значений $f(x_i, w)$;
- 3: поставить первую точку в начало координат:
 $(FPR_0, TPR_0) := (0, 0)$; $AUC := 0$;
- 4: **для** $i := 1, \dots, \ell$
- 5: **если** $y_i = -1$ **то** — сместиться на один шаг вправо:
- 6: $FPR_i := FPR_{i-1} + \frac{1}{\ell_-}$; $TPR_i := TPR_{i-1}$;
 $AUC := AUC + \frac{1}{\ell_-} TPR_i$;
- 7: **иначе** — сместиться на один шаг вверх:
- 8: $FPR_i := FPR_{i-1}$; $TPR_i := TPR_{i-1} + \frac{1}{\ell_+}$;

Разложение ошибки на смещение и разброс

- **шум** — ошибка идеального алгоритма
- **смещение (bias)** — отклонение среднего ответа обученного алгоритма от ответа идеального алгоритма
- **разброс (variance)** — разброс ответов обученных алгоритмов относительно среднего ответа



Квадратичная функция потерь

Задача регрессии: $Y = \mathbb{R}$

Квадратичная функция потерь: $L(y, a) = (a(x) - y)^2$

Вероятностная постановка: $X^\ell = (x_i, y_i)_{i=1}^\ell \sim p(x, y)$

Метод обучения: $\mu: 2^X \rightarrow A$, т.е. выборка \rightarrow алгоритм

Среднеквадратичный риск:

$$R(a) = E_{x,y}(a(x) - y)^2 = \int_X \int_Y (a(x) - y)^2 p(x, y) dx dy$$

Минимум среднеквадратичного риска, «недостижимый идеал»:

$$a^*(x) = E(y|x) = \int_Y y p(y|x) dx$$

Основная мера качества метода обучения μ :

$$Q(\mu) = E_{X^\ell} E_{x,y}(\mu(X^\ell)(x) - y)^2$$

Разложение ошибки на шум, вариацию и смещение

Теорема

В случае квадратичной функции потерь для любого μ

$$\begin{aligned} Q(\mu) = & \underbrace{E_{x,y} (a^*(x) - y)^2}_{\text{шум (noise)}} + \\ & + \underbrace{E_{x,y} (\bar{a}(x) - a^*(x))^2}_{\text{смещение (bias)}} + \\ & + \underbrace{E_{x,y} E_{X^\ell} (\mu(X^\ell)(x) - \bar{a}(x))^2}_{\text{разброс (variance)}}, \end{aligned}$$

$\bar{a}(x) = E_{X^\ell} (\mu(X^\ell)(x))$ — средний ответ обученного алгоритма

Метод k ближайших соседей

Вероятностная модель данных: $p(y|x) = f(x) + \mathcal{N}(0, \sigma^2)$

Метод k ближайших соседей:

$$a(x) = \frac{1}{k} \sum_{j=1}^k y(x^{(j)}),$$

где $x^{(j)}$ — j -й сосед объекта x

$a^*(x) = f(x)$ — истинная зависимость

$\bar{a}(x) = \frac{1}{k} \sum_{j=1}^k f(x^{(j)})$ — средний ответ

Разложение bias-variance:

$$Q(\mu) = \underbrace{\sigma^2}_{\text{шум}} + \underbrace{E_{x,y} \left(\bar{a}(x) - f(x) \right)^2}_{\text{смещение}} + \underbrace{\frac{1}{k} \sigma^2}_{\text{разброс}}$$