

0.1 Колмогоровская сложность моделей

Одним из фундаментальных способов определить сложность произвольного математического объекта является колмогоровская сложность. Ниже представлено формальное определение колмогоровской сложности и основные ее свойства.

Определение 1. Пусть задано вычислимое частично определенное отображение из множества бинарных слов в себя:

$$T : \{0, 1\}^* \rightarrow \{0, 1\}^*.$$

Колмогоровской сложностью бинарной строки x назовем минимальную длину описания относительно T :

$$K_T(x) = \min_{f \in \{0, 1\}^*} \{|f| : T(f) = x\},$$

Заметим, что колмогоровская сложность зависит от отображения T . В [1] доказано, что колмогоровская сложность $K_T(x)$ при двух отображениях T_1, T_2 отличается лишь на некоторую константу, не зависящих от строки x . Поэтому для дальнейшего изложения зафиксируем некоторое отображение T и положим $K(x) = K_T(x)$.

Обобщим понятие колмогоровской сложности на случай двух бинарных строк.

Определение 2. Пусть задано вычислимое и частично определенное отображение из декартового произведения двух множеств бинарных слов в себя:

$$T : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{0, 1\}^*.$$

Условной колмогоровской сложностью бинарной строки y при условии x назовем минимальную длину описания относительно T :

$$K_T(y|x) = \min_{f \in \{0, 1\}^*} \{|f| : T(f, y) = x\},$$

Аналогично простой колмогоровской сложности, зафиксируем некоторое отображение T и положим $K_T(y|x) = K(y|x)$.

Рассмотрим некоторые свойства условной колмогоровской сложности.

Оценка условной Колмогоровской сложности [1]

$$K(x, y) \leq K(x) + K(y|x) + O(\log K(x, y)).$$

Количество информации в паре x, y симметрично с точностью до константы:

$$I(x : y) = I(y : x) + O(\log K(x, y)),$$

где величина $I(x : y) = K(y) - K(y|x)$ задает количество информации в x об объекте y .

Отметим, что схожими свойствами обладает взаимная информация и энтропия, определения которых даны ниже.

Определение 3. Пусть задана дискретная случайная величина x с вероятностным распределением p , принимающая значения x_1, \dots, x_n . Энтропией распределения случайной величины x назовем:

$$H(x) = - \sum_{i=1}^n p(x = x_i) \log p(x = x_i).$$

Оценка условной энтропии

$$H(x, y) = H(x) + K(y|x).$$

Определение 4. Взаимной информацией I двух случайных величин x, y назовем следующее выражение:

$$I(x, y) = H(x) - H(x|y), \quad H(x) = - \sum_i p_x(x_i) \log p_x(x_i)$$

Взаимная информация симметрична:

$$I(x, y) = I(y, x).$$

Таким образом, свойства энтропии и колмогоровской сложности, а также количества информации $I(x : y)$ и взаимной информации, во многом совпадают. Докажем теорему о связи колмогоровской сложности и энтропии распределения, подытоживающую связь этих двух математических объектов.

Теорема 1. [1] Пусть задана некоторая строка x длины n с частотами $p = (p_0, 1 - p_0)$ появлений нулей и единиц в строке. Тогда

$$K(x) \leq H(x) + O(\log m).$$

Неравенство обращается в равенство для большинства строк x длины n .

Доказательство. Всего слов, которые можно получить с использованием заданных частот:

$$C = \frac{m!}{(p_0 m)!((1 - p_0)m)!}.$$

Т.к. количество таких слов конечно, то их можно пронумеровать и построить отображение, выдающее строку x по ее порядковому номеру. Таким образом, условная колмогоровская сложность ограничена сверху:

$$K(x|C, p_0) \leq \log C + O(1).$$

Воспользуемся формулой Стирлинга:

$$n! = \sqrt{(2\pi + o(1))n} \frac{n^n}{e}.$$

И получим оценку:

$$C \leq 2^{nH(x) + O(\log n)}, \quad K(x|C, p_0) \leq nH(x) + O(\log n).$$

Для того, чтобы избавиться от условия в $K(x|C, p_0)$ потребуется $O(\log n)$ бит для описания чисел $p_0 n, n$.

TODO: Поскольку слов с более короткими описаниями меньше, чем C , то для большинства слов будет достигаться предложенная оценка. \square

Частным случаем колмогоровской сложности является префиксная колмогоровская сложность. Эта сложность задается машиной Тьюринга специального вида, имеющей две ленты: однонаправленную ленту для чтения и двунаправленную рабочую ленту. Будем полагать что машина Тьюринга T останавливается на p с выводом x : $T(p) = x$, если вся запись p осталась слева от читающей каретки, x осталась слева от пишущей каретки и T остановлена. Колмогоровскую сложность относительно префиксных машин Тьюринга назовем префиксной колмогоровской сложностью KP . Отметим, что префиксная колмогоровская сложность является частным случаем колмогоровской сложности, а потому:

$$K(x) \leq KP(x).$$

Теорема (без доказательств) 1. Пусть задана вычислимая функция p_x вероятности на множестве бинарных строк. Тогда

$$0 \leq E_{p_x} KP(x) - H(x) \leq KP(p_x) + O(1),$$

где $K(p_x)$ определяется как минимальная длина программы для префиксной машины Тьюринга, вычисляющей p_x .

Таким образом, существует связь между энтропией и колмогоровской сложностью (как для обычного варианта сложности, так и для префиксной колмогоровской сложности): для простых распределений в смысле сложности функции p_x энтропия будет приближать математическим ожиданием префиксной колмогоровской сложности.

0.2 Колмогоровская сложность и принцип минимальной длины описания

Рассмотрим задачу выбора модели для заданной выборки. Будем полагать что заданная выборка описывается в виде некоторой бинарной строки x . В дальнейшем будем отождествлять выборку и ее бинарное описание x .

Задачу выбора модели для выборки можно рассматривать как задачу нахождения колмогоровской сложности для выборки. В случае, если модель является дискриминативной, то вместо колмогоровской сложности можно использовать условную колмогоровскую сложность. В общем виде колмогоровская сложность и префиксная сложность невычислимы [1], поэтому рассмотрим упрощенный подход к выбору модели: вместо сложности строки x будем искать некоторое множество S , в которое входит x , и чья префиксная сложность описания невелика. Таким образом, мы сможем найти “хорошую” машину Тьюринга не для конкретной строки, а для некоторого семейства строк (или выборок), обладающих общими свойствами или регулярностью.

Определение 5. Сложностью конечного множества S назовем следующей величину:

$$K(S) = \min_{f \in \{0,1\}^*} \{|f| : T_i(f) \text{ перечисляет все элементы множества } S\}.$$

Вместо задачи нахождения минимальной сложности для выборки x будем искать множество S , которое описывается некоторой машиной Тьюринга, и в которое входит заданная строка x . Приведем формулу для оценки разности между сложностью выборки x и множества S , в которое входит данная выборка.

Теорема (без доказательств) 2. Для любого $x \in S$ справедливо неравенство [3]:

$$K(x) \leq K(S) + \log |S| + O(1).$$

На практике задача выбора модели подразумевает, что мы можем выбрать модель, которая описывает выборку (или множество выборок) S неидеально, а с некоторым допустимым уровнем потери информации. Тогда задача выбора модели для заданной выборки ставится следующим образом:

$$\arg \min_S \{\log |S| + K(S) : x \in S, K(x) \leq \alpha\}, \quad (0.1)$$

где α — максимально допустимая сложность множества S .

Заметим, что решение задачи выбора модели в приведенном выше виде является вычислимой, то есть можно предложить алгоритм, вычисляющий данную задачу. Приведем схему данного алгоритма:

```

Положим  $\hat{p}, \hat{S}$  неопределенными;
for all  $S, p : T(p) = S, |p| \leq \alpha$  do
  if  $\hat{S}$  неопределен или  $|p| + \log(S) \leq \hat{p} + \log \hat{S}$ , then
     $\hat{p}, \hat{S} = p, S$ 
  end if
end for
return  $\hat{S}$ 

```

Т.к. множество всех программ с длиной менее α является конечным, то алгоритм остановится, а потому вычислим. По построению он также доставляет решение оптимизационной задачи (??).

Теорема (без доказательств) 3. Обозначим за $S^*, K(S^*) \leq \alpha$ множество, доставляющее минимум следующей функции:

$$\delta(S) = \log |S| - K(x|S).$$

Тогда $\delta(S^*) - \delta \hat{S} = O(\log |x|)$.

Таким образом, пара (\hat{p}, \hat{S}) доставляет минимуму суммы $|p| + \log(S)$ при ограничениях на длину программы $|p| \leq \alpha$. Первое слагаемое данной суммы — это длина программы, то есть сложность описания множества \hat{S} . Сама сумма $|p| + \log(S)$ характеризует длину кода, требуемого для двухэтапного описания строки x , где на первом этапе мы описываем множество S , а на втором этапе находим x во множестве S .

Описанный метод получения кода для описания выборки x соответствует *Принципу минимальной длины описания* или MDL: для заданной выборки требуется найти минимальный код, который описывает выборку (возможно, с некоторой наперед заданной допустимой ошибкой).

0.3 Вероятностная интерпретация минимальной длины описания

Рассмотрим подробнее задачу кодирования выборки. Задачу можно рассматривать как проблему передачу информации от кодировщика декодировщику. Задана выборка $\mathbf{X}, x \in \mathbf{X}$. Кодировщик кодирует информацию о выборке \mathbf{X} с помощью некоторого кода \mathbf{f} и передает ее декодировщику. Декодировщик декодирует код $\mathbf{f}(\mathbf{X})$, полученный от кодировщика и восстанавливает исходную выборку \mathbf{X} (возможно, с некоторой потерей информации, если это оговорено заранее). Допустим, для кодирования информации о выборке доступно несколько методов кодирования, при этом для разных объектов выборки минимальный по количеству информации метод будет отличаться. Тогда кодировщик должен передать не только информацию о выборке, но и информацию о самом методе кодирования. Проблему выбора оптимального способа кодирования из нескольких можно рассмотреть при помощи вероятностной интерпретации минимальной длины описания.

Теорема (без доказательств) 4. Пусть задан конечное или счетное множество Z с введенной на нем вероятностью P . Тогда существует код C , такой что для любого $x \in Z$ длина его описания L будет равняться $L(C(x)) = -\log P(x)$.

Таким образом, вместо методов кодирования можно рассматривать задачу выбора функции вероятности (правдоподобия?), введенной на признаковом описании выборки. Задача формулируется следующим образом:

$$\log p(\mathbf{X}|\mathbf{f}) + L(\mathbf{f}),$$

где $L(\mathbf{f})$ — длина описания вероятностной модели \mathbf{f} .

Одним из критериев качества вероятностного кодирования с помощью смеси кодов относительно фиксированного кодирования является регрет (TODO):

$$R(x) = -\log P(x) + \min_{\mathbf{f}} (\log P(x|\mathbf{f})).$$

Регрет характеризует разницу между длиной рассматриваемого $\log P(x)$ кода для x в сравнении с наилучшим кодом из некоторого множества \mathbf{f} .

Пусть модель \mathbf{f} зависит от вектора параметра \mathbf{w} :

$$\log p(x|\mathbf{f}) = \log p(x|\mathbf{w}(x)),$$

где $\mathbf{w}(x)$ — оптимальные параметры для объекта x .

Тогда регрет выглядит следующим образом:

$$R(P, x) = -\log P(x) + \min_{\mathbf{w}} (\log P(x|\mathbf{w}(x))).$$

Для всей выборки регрет определяется следующим образом:

$$R(\mathbf{X}) = \max_{x \in \mathbf{X}} (-\log P(x) + \min_{\mathbf{w}} (\log P(x|\mathbf{w}(x)))).$$

Теорема 2 ((Штарьков, 1987) ??,??). Пусть величина

$$\log \sum_x p(x|\mathbf{w}(x))$$

конечна. Тогда следующее выражение доставляет единственный минимум регрета $R(\mathbf{X})$:

$$\frac{p(x|\mathbf{w}(x))}{\sum_{x' \in \mathbf{X}} p(x'|\mathbf{w}(x'))}.$$

Доказательство. Рассмотрим выражение регрета для данной вероятностной меры:

$$-\log P(x) + \min_{\mathbf{w}} (\log P(x|\mathbf{w}(w))) = \log \sum_x p(x|\mathbf{w}(x)).$$

Выражение регрета не зависит от объекта выборки x . Заметим, что для любых двух отличных распределений p_1, p_2 существует хотя бы один элемент, такой что

$$p_1(x) < p_2(x).$$

Действительно, пусть это неверно: $\forall x p_1(x) \geq p_2(x), \exists x' : p_1(x') > p_2(x')$. Просуммируем вероятности p_1 и p_2 и воспользуемся тем, что сумма вероятностей по всем объектам будет равна единице:

$$1 = \sum p_1 > \sum p_2 = 1,$$

приходим к противоречию.

Тогда для любого распределения на выборке существует элемент x' :

$$\frac{p(x'|\mathbf{w}(x'))}{\sum_x p(\mathbf{w}|w(x))} > p(x').$$

Тогда

$$R(P, x) \geq -\log p(x') + \log p(x'|\mathbf{w}(x')) > \frac{p(x'|\mathbf{w}(x'))}{\sum_x p(x|\mathbf{w}(x))} = R(x).$$

□

Альтернативным методом выбора модели является двусвязный байесовский вывод. На *первом уровне* байесовского вывода находится апостериорное распределение параметров:

$$\hat{\mathbf{w}} = \arg \max \frac{p(\mathbf{X}|\mathbf{w}, \mathbf{f})p(\mathbf{w}|\mathbf{f})}{p(\mathbf{X}|\mathbf{f})},$$

где $p(\mathbf{X}|\mathbf{w}, \mathbf{f})$ — правдоподобие выборки при условии модели \mathbf{f} с параметрами \mathbf{w} , $p(\mathbf{w}|\mathbf{f})$ — априорное распределение параметров при условии модели \mathbf{f} . Оно задается на основе наших предположений о природе выборки и о способах ее порождения. Знаменатель $p(\mathbf{w}|\mathbf{f})p(\mathbf{X}|\mathbf{f})$ называется *evidence* или *обоснованность модели* и определяет, насколько хорошо модель способна описать выборку в среднем по всем возможным значениям параметров.

На *втором уровне* байесовского вывода осуществляется выбор модели на основе обоснованности модели:

$$\mathbf{f} = \arg \max p(\mathbf{w}|\mathbf{f})p(\mathbf{X}|\mathbf{f}) = \int_{\mathbf{w}} p(\mathbf{X}|\mathbf{w}, \mathbf{f})p(\mathbf{w}|\mathbf{f})d\mathbf{w}$$

Эта величина также выступает байесовской интерпретацией минимальной длины. Для подтверждения этого факта докажем следующую теорему.

Теорема 3. Пусть правдоподобие $p(\mathbf{X}|\mathbf{w}, \mathbf{f})$ соответствует экспоненциальному семейству распределений, т.е.

$$p(x|\mathbf{w}, \mathbf{f}) = h(x)g(\boldsymbol{\eta})\exp(\boldsymbol{\eta} \cdot \mathbf{T}(x)),$$

где h, g, \mathbf{T} — некоторые функции, $\boldsymbol{\eta}$ — некоторый параметр распределения.

Пусть в качестве априорного распределения выступает распределение Джеффри:

$$p(\mathbf{w}|\mathbf{f}) = \frac{\sqrt{I}}{\int_w \sqrt{I(w)}}, \text{ где } I \text{ — определить матрицы Фишера:}$$

$$I(\mathbf{w}) = \det\left\{\mathbb{E}_w \frac{-\partial^2}{\partial w_i \partial w_j} \log p(\mathbf{X}|\mathbf{w}, \mathbf{f})\right\}_{ij}.$$

Тогда регрет можно аппроксимировать следующим выражением:

$$\frac{k}{2} \log \frac{n}{2\pi} - \log p(\mathbf{X}|\mathbf{w}, \mathbf{f}) + \log \sqrt{I(\mathbf{w})}.$$

Доказательство. Воспользуемся аппроксимацией Лапласа для упрощения формулы обоснованности модели. Разложим $\log p(\mathbf{X}, \mathbf{w}|\mathbf{f}) = \log p(\mathbf{X}|\mathbf{w}, \mathbf{f})p(\mathbf{w}|\mathbf{f})$ в точке локального максимума \mathbf{w}_0 в ряд Тейлора:

$$\log p(\mathbf{X}, \mathbf{w}|\mathbf{f}) \approx \log p(\mathbf{X}, \mathbf{w}_0|\mathbf{f}) - \frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^\top \mathbf{A}(\mathbf{w} - \mathbf{w}_0),$$

$$A_{ij} = -\frac{\partial^2}{\partial w_i \partial w_j} \log p(\mathbf{X}, \mathbf{w}|\mathbf{f})|_{\mathbf{w}=\mathbf{w}_0}.$$

Полученное распределение представимо в виде ненормированного гауссового распределения. для такой аппроксимации плотности вероятности запишем нормирующий коэффициент:

$$\log p(\mathbf{X}, \mathbf{w}|\mathbf{f}) \approx \log p(\mathbf{X}|\mathbf{w}_0, \mathbf{f}) + p(\mathbf{w}_0|\mathbf{f}) - \log \sqrt{\frac{(2\pi)^k}{\det \mathbf{A}}}.$$

Подставляя в полученную формулу распределение Джеффри получим:

$$p(\mathbf{X}, \mathbf{w}|\mathbf{f}) \approx p(\mathbf{X}|\mathbf{w}_0, \mathbf{f}) - \log \sqrt{\frac{(2\pi)^k}{2}}.$$

TODO: там еще информация Фишера + экспоненциальное семейство. Тогда регрет будет равен:

$$R(P, x) \approx \log \sqrt{\frac{(2\pi)^k}{2}}.$$

□

Теорема (без доказательств) 5. При количество параметров, стремящемся к бесконечности, регрет оптимальной модели отличается от байесовской оценки на константу:

$$R(P, x) = \frac{k}{2} \log \frac{n}{2\pi} + \log \int_{\mathbf{w}} \sqrt{I(\mathbf{w})} d\mathbf{w} + o(1).$$

TODO: вывод, что эти вещи совпадают.

Список литературы

- [1] Успенский В., Шень А., Верещагин Н. Колмогоровская сложность и алгоритмическая случайность. – Litres, 2017
- [2] Grunwald P., Vitányi P. Shannon information and Kolmogorov complexity //arXiv preprint cs/0410002. – 2004.
- [3] Vereshchagin N. K., Vitányi P. M. B. Kolmogorov’s structure functions and model selection //IEEE Transactions on Information Theory. – 2004. – Т. 50. – №. 12. – С. 3265-3290.
- [4] Grunwald P. A tutorial introduction to the minimum description length principle //arXiv preprint math/0406077. – 2004.
- [5] Штарьков Ю. М. Универсальное последовательное кодирование отдельных сообщений //Проблемы передачи информации. – 1987. – Т. 23. – №. 3. – С. 3-17.