

# Лекция 2

## РАНДОМИЗАЦИЯ И ТЕОРИЯ ОТПЕЧАТКОВ

### 1. ОТПЕЧАТКИ ДЛЯ ИДЕНТИФИКАЦИИ

### 2. КАК РАБОТАЮТ АНТИВИРУСНЫЕ ПРОГРАММЫ

В последнее время компьютерные вирусы стали настоящим “бедствием” для всех, кто пользуется компьютерами. С целью обнаружения и профилактики “вирусов” используются *антивирусные программы*. Одним из методов, которым антивирусные программы обнаруживают вредоносное программное обеспечение, является сканирование файлов для поиска известных вирусов. Анализируя файлы, антивирусная программа обращается к антивирусным базам, составленным производителем программы-антивируса. В случае соответствия какого-либо участка кода просматриваемого файла сигнатуре вируса в базах программа-антивирус может по запросу выполнить одно из действий, включая удаление или лечение файла. Поскольку в базах данных содержатся тысячи известных вирусов, алгоритмы поиска совпадений эталонов и участков кода является одной из актуальных задач.

**2.1. Поиск шаблона в файле.** Задача обнаружения вирусов, внедрившихся в тело файла, называется в компьютерных науках *поиском шаблона в файле*.

Рассмотрим классическую задачу поиска шаблона  $Y$  в файле  $X$ . Мы хотим определить встречается ли шаблон  $Y = y_1y_2 \dots y_m$  в файле  $X = x_1x_2 \dots x_n$ ,  $n > m$ .

Непосредственный способ последовательных сравнений шаблона с участками файла выполняется за  $O((n - m)m)$  операций. Существует сложный детерминированный алгоритм, который выполняется за время  $O(n + m)$  (понятно, что это оптимальный результат). Существует простой рандомизированный алгоритм Карпа–Рабина, который также выполняется за время  $O(n + m)$  и основан на той же идее, что и алгоритм, представленный в разделе 2.

Для простоты считаем, что алфавит бинарный, то есть  $X$  и  $Y$  состоят из нулей и единиц (из *битов*). Пусть  $X(j) = x_j x_{j+1} \dots x_{j+m-1}$  обозначает участок файла  $X$  длины  $m$ , начиная с позиции  $j$ . Отпечаток участка  $X(j)$  мы определяем по формуле

$$\text{FING}(X(j)) = g(x_j \dots x_{j+m-1}) \bmod p,$$

где

$$\begin{aligned} g(x_j \dots x_{j+m-1}) &= x_j 2^{m-1} + x_{j+1} 2^{m-2} + \dots + x_{j+m-2} 2^1 + x_{j+m-1} 2^0 \\ &= \sum_{k=j}^{j+m-1} x_k 2^{m-1+j-k}. \end{aligned}$$

Аналогично определяется и  $\text{FING}(Y)$ .

Ниже приведен алгоритм Карпа–Рабина.

**Input:** Файл  $X$ ,  $|X| = n$ ; шаблон  $Y$ ,  $|Y| = m$ ;  
достаточно большое число  $c$ .

**Output:** Индикатор **найден/не найден** шаблон в файле

1. Вычислить  $k = cm \ln(cm)$ ;
2. Выбрать случайное простое число  $p$  из  $\{1, \dots, k\}$ ;
3. **for**  $j = 1$  **to**  $n - m + 1$  **do**  
    **if**  $\text{FING}(X(j)) = \text{FING}(Y)$  **then** шаблон найден **stop**

**Алгоритм 1.** Рандомизированный поиск шаблона

Если алгоритм заканчивается, а оператор **stop** не срабатывает, то записи  $Y$  действительно нет в файле  $X$ . Ошибка может возникнуть, если алгоритм определит наличие шаблона (ведь два разных числа могут иметь одинаковые остатки от деления на  $p$ ).

Чтобы определить вероятность ошибки, мы рассуждаем так же, как и в разделе 2. Пусть  $j$  фиксировано и  $X(j) \neq Y$ . Тогда абсолютная величина числа  $X(j) - Y$  не превосходит  $2^m$ , значит оно имеет не больше, чем  $m$  простых делителей. Кроме того,  $\text{FING}(X(j)) = \text{FING}(Y)$  тогда и только тогда, когда  $p$  делит  $X(j) - Y$ . Поэтому

$$(1) \quad \mathbf{P}(\text{FING}(X(j)) = \text{FING}(Y)) \leq \frac{m}{\pi(k)}.$$

Если же  $X(j) \neq Y$  для всех  $1 \leq j \leq n - m + 1$ , то

$$(2) \quad \begin{aligned} & \mathbf{P}(\text{алгоритм “находит” шаблон}) \\ &= \mathbf{P}\left(\bigcup_{j=1}^{n-m+1} \{\text{FING}(X(j)) = \text{FING}(Y)\}\right) \\ &\leq \frac{(n-m+1)m}{\pi(k)} \sim \frac{1}{c}, \end{aligned}$$

если выбрать  $k = c n m \ln(c n m)$  (рассуждения при доказательстве эквивалентности  $\sim$  такие же, как и в разделе 2).

Время выполнения алгоритма можно оценить “в лоб”: каждая итерация цикла требует вычисления отпечатка числа, состоящего из  $m$  бит, а это занимает время  $O(m)$ , что означает, что общее время выполнения равно  $O(nm)$ .

Однако, на самом деле, эту оценку можно заметно улуч-

шить, заметив, что

$$\begin{aligned}
 (3) \quad g(X(j+1)) &= \sum_{k=j+1}^{j+m} x_k 2^{m+j-k} \\
 &= x_{j+m} + \sum_{k=j}^{j+m-1} x_k 2^{m+j-k} - x_j 2^m \\
 &= 2(g(X(j)) - 2^{m-1}x_j) + x_{j+m}.
 \end{aligned}$$

Каждая итерация в (3) выполняется за время  $O(1)$  (фактически необходимы 2 сложения и 2 умножения).

Чтобы запустить итерационный процесс, необходимо вычислить  $\text{fing}(X(1))$  и  $\text{fing}(Y)$ . Для этого требуется  $O(m)$  операций. Поэтому весь алгоритм выполняется за время  $O(n+m)$ , как и утверждалось выше.

### 3. ПРОВЕРКА РАВЕНСТВА МАТРИЦ

Вычисление произведения матриц играет ключевую роль во многих вычислениях. Размеры перемножаемых матриц могут быть чрезвычайно велики, поэтому эффективность операции умножения играет ключевую роль в подобных вычислениях. Такую же важную роль играет проверка вычислений, которая сводится к проверке равенства матриц  $AB = C$ .

Пусть заданы три матрицы  $A, B, C$  размера  $n \times n$ , элементами которых являются только 0 или 1. Мы рассматриваем арифметические операции по модулю 2, то есть

$$x + y = \begin{cases} 0, & \text{если } x = y = 0 \text{ или } x = y = 1, \\ 1, & \text{если } x = 0, y = 1 \text{ или } x = 1, y = 0. \end{cases}$$

Как наиболее эффективно проверить матричное равенство

$$AB = C?$$

Простой и одновременно естественный способ состоит из двух шагов:

1. перемножить матрицы  $A$  и  $B$ ;
2. проверить равенство для каждого из элементов матриц  $AB$  и  $C$ .

Для вычисления произведения матриц  $A$  и  $B$  требуется  $O(n^3)$  операций: элементы матрицы  $AB$  вычисляются по формуле

$$(AB)_{ij} = \sum_{k=1}^n a_{ik}b_{kj} \bmod 2, \quad 1 \leq i, j \leq n.$$

Таким образом, для вычисления элемента  $(AB)_{ij}$  требуется  $n$  операций умножения и  $n$  операций сложения. Поскольку матрица размера  $n \times n$  имеет  $n^2$  элементов, то для вычисления произведения  $AB$  действительно требуется  $O(n^3)$  операций. Существует остроумный детерминированный метод, который позволяет осуществить операцию умножения двух матриц за  $O(n^{2.375})$  операций.

Мы рассмотрим рандомизированный алгоритм сравнения матриц  $AB$  и  $C$ , который нуждается в  $O(n^2)$  операциях. Этот алгоритм позволяет утверждать, что равенство  $AB = C$  выполнено лишь с определенной вероятностью, которую можно оценить довольно точно.

**3.1. Рандомизированный алгоритм проверки равенства матриц.** Сначала алгоритм генерирует случайный вектор  $\bar{\mathbf{r}}$  размерности  $n$ , который имеет равномерное распределение в  $\{0, 1\}^n$ . Это означает, что координаты вектора  $\bar{\mathbf{r}}$  равны либо 0, либо 1, причем

$$P(\bar{\mathbf{r}} = \bar{\mathbf{e}}) = \frac{1}{2^n}$$

для любого вектора  $\bar{\mathbf{e}}$ , составленного из 0 и 1.

Затем алгоритм вычисляет последовательно векторы

$$B\bar{\mathbf{r}}, \quad A(B\bar{\mathbf{r}}), \quad C\bar{\mathbf{r}}.$$

Работа алгоритма заканчивается проверкой равенства

$$A(B\bar{\mathbf{r}}) = C\bar{\mathbf{r}}.$$

Если это равенство нарушается, то принимается верное решение о  $AB \neq C$ . В противном же случае алгоритм решает, что  $AB = C$ , а это уже верно не всегда. Поскольку  $\bar{\mathbf{r}}$  — случайный вектор, то ошибочное решение также случайно; его вероятность можно оценить.

Ниже приведена блок-схема данного алгоритма.

**Input:** Матрицы  $A, B, C$  размера  $n \times n$ .  
**Output:**  $d = 1$ , если  $AB = C$ ; или  $d = 0$ , если  $AB \neq C$ .  
 1. Сгенерировать случайный вектор  $\bar{\mathbf{r}} = (r_1, \dots, r_n) \in \{0, 1\}^n$ .  
 2. Вычислить  $B\bar{\mathbf{r}}, A(B\bar{\mathbf{r}}), C\bar{\mathbf{r}}$ .  
 3. Сравнить  $A(B\bar{\mathbf{r}})$  и  $C\bar{\mathbf{r}}$ :  
     if  $A(B\bar{\mathbf{r}}) = C\bar{\mathbf{r}}$  then  $d = 1$  else  $d = 0$ .  
 4. stop.

## Алгоритм 2. Сравнение матриц

**3.2. Время генерирования вектора.** Для генерации вектора  $\bar{\mathbf{r}}$ , имеющего равномерное распределение в  $\{0, 1\}^n$ , необходимо  $O(n)$  операций. Это вытекает из следующего результата.

**Лемма 1.** *Вектор  $\bar{\mathbf{r}}$  имеет равномерное распределение в множестве  $\{0, 1\}^n$  тогда и только тогда, когда его координаты  $r_k, 1 \leq k \leq n$ , независимы и имеют симметричное распределение Бернулли:*

$$P(r_k = 0) = P(r_k = 1) = \frac{1}{2}.$$

Из леммы вытекает, что генерирование вектора  $\bar{\mathbf{r}}$ , имеющего равномерное распределение в  $\{0, 1\}^n$ , сводится к генерированию  $n$  случайных величин Бернулли. Поскольку генерирование одной случайной величины Бернулли требует  $O(1)$  операций, то генерирование вектора  $\bar{\mathbf{r}}$  действительно требует  $O(n)$  операций.

**3.3. Время вычисления произведения матрицы на вектор.** Элементы векторов  $B\bar{\mathbf{r}}$ ,  $A(B\bar{\mathbf{r}})$  и  $C\bar{\mathbf{r}}$  вычисляются по обычной формуле, которую запишем для  $B\bar{\mathbf{r}}$ :

$$(B\bar{\mathbf{r}})_i = \sum_{k=1}^n b_{ik}r_k \bmod 2, \quad 1 \leq i \leq n.$$

Вычисление произведения матрицы размера  $n \times n$  на вектор размерности  $n$  требует  $O(n^2)$  операций, поскольку вычисление каждого элемента произведения требует  $n$  операций умножения и  $n$  операций сложения. Поэтому  $O(n^2)$  операций достаточно для вычисления всех векторов  $B\bar{\mathbf{r}}$ ,  $A(B\bar{\mathbf{r}})$ ,  $C\bar{\mathbf{r}}$ .

**3.4. Общее время.** Сравнение двух векторов требует  $O(n)$  операций. Следовательно, алгоритм заканчивает работу за  $O(n) + O(n^2) + O(n) = O(n^2)$  операций.

**3.5. Вероятность ошибки алгоритма.** Если на самом деле  $AB = C$ , то  $A(B\bar{\mathbf{r}}) = C\bar{\mathbf{r}}$  для любого  $\bar{\mathbf{r}}$ . Если же  $AB \neq C$ , то тем не менее равенство  $A(B\bar{\mathbf{r}}) = C\bar{\mathbf{r}}$  все же может выполняться для некоторых векторов  $\bar{\mathbf{r}}$ . Если такое случается, то алгоритм дает неверный ответ. Вероятность неверного ответа оценивается в следующей теореме.

**Теорема 1.** Если  $AB \neq C$ , а вектор  $\bar{\mathbf{r}}$  имеет равномерное распределение в множестве  $\{0, 1\}^n$ , то

$$P(AB\bar{\mathbf{r}} = C\bar{\mathbf{r}}) \leq \frac{1}{2}.$$

**3.6. Как сделать вероятность ошибки малой.** Улучшить оценку вероятности ошибки, полученную в теореме 1, невозможно. Поэтому кажется, что пользоваться описанным методом крайне опасно при проверке равенства двух матриц: ошибка может наблюдаться при каждой второй проверке. Однако минимальным усилием ситуация существенно улучшается. Именно, повторим описанную процедуру еще раз. При этом, естественно, случайный вектор  $\bar{\mathbf{r}}$  будет другим. Обозначим эти два вектора, возникающие при первом и втором повторе процедуры, через  $\bar{\mathbf{r}}'$  и  $\bar{\mathbf{r}}''$ . Поскольку векторы  $\bar{\mathbf{r}}'$  и  $\bar{\mathbf{r}}''$  независимы, вероятность двойной ошибки можно оценить так:

$$\begin{aligned} P(AB\bar{\mathbf{r}}' = C\bar{\mathbf{r}}', AB\bar{\mathbf{r}}'' = C\bar{\mathbf{r}}'') \\ = P(AB\bar{\mathbf{r}}' = C\bar{\mathbf{r}}')P(AB\bar{\mathbf{r}}'' = C\bar{\mathbf{r}}'') \\ \leq \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}. \end{aligned}$$

Повторяя процедура 10 раз, снижаем вероятность ошибки до  $\frac{1}{2^{10}} \leq \frac{1}{1000}$ . Если необходимо, чтобы вероятность ошибки не превосходила  $\delta > 0$ , то достаточно повторить процедуру  $1 + \lceil \log_2 \delta^{-1} \rceil$  раз. Такое повторение несколько не влияет на оценку  $O(n^2)$  времени выполнения процедуры.

**Доказательства.** Сначала докажем лемму 1, а потом теорему 1.

*Доказательство леммы 1.*  $\triangleleft$  Пусть  $r_1, \dots, r_n$  — независимые случайные величины Бернулли. Докажем, что вектор  $\bar{\mathbf{r}} = (r_1, \dots, r_n)$  имеет равномерное распределение на множестве  $\{0, 1\}^n$ . Выберем произвольный вектор  $\bar{\varepsilon} \in \{0, 1\}^n$ , то есть  $\bar{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)$ , где  $\varepsilon_i \in \{0, 1\}$ ,  $1 \leq i \leq n$ . Тогда в силу независимости

$$\begin{aligned} P(\bar{\mathbf{r}} = \bar{\varepsilon}) &= P(r_1 = \varepsilon_1, \dots, r_n = \varepsilon_n) \\ &= P(r_1 = \varepsilon_1) \dots P(r_n = \varepsilon_n) = \frac{1}{2^n}. \end{aligned}$$



Это и означает, что вектор  $\bar{\mathbf{r}}$  имеет равномерное распределение на множестве  $\{0, 1\}^n$ .

▷ Пусть теперь вектор  $\bar{\mathbf{r}}$  имеет равномерное распределение на множестве  $\{0, 1\}^n$ . Докажем, что его координаты  $r_1, \dots, r_n$  независимы и имеют распределение Бернулли. С целью демонстрации идеи рассмотрим сначала случай  $n = 2$ : для  $\varepsilon_1 \in \{0, 1\}$  имеем

$$\begin{aligned} P(r_1 = \varepsilon_1) &= P(r_1 = \varepsilon_1, \varepsilon_2 = 0) + P(r_1 = \varepsilon_1, \varepsilon_2 = 1) \\ &= \frac{1}{4} + \frac{1}{4} = \frac{1}{2}, \end{aligned}$$

то есть  $r_1$  — случайная величина Бернулли. Аналогично доказывается бернулливость величины  $r_2$ . Теперь их независимость доказать просто:

$$P(r_1 = \varepsilon_1, r_2 = \varepsilon_2) = \frac{1}{4} = \frac{1}{2} \cdot \frac{1}{2} = P(r_1 = \varepsilon_1) \cdot P(r_2 = \varepsilon_2).$$

Доказательство в общем случае совершенно аналогично: если  $\varepsilon_1 \in \{0, 1\}$ , то

$$\begin{aligned} P(r_1 = \varepsilon_1) &= \sum_{\varepsilon_2, \dots, \varepsilon_n \in \{0, 1\}} P(r_1 = \varepsilon_1, \dots, r_n = \varepsilon_n) \\ &= \sum_{\varepsilon_2, \dots, \varepsilon_n \in \{0, 1\}} \frac{1}{2^n} = 2^{n-1} \cdot \frac{1}{2^n} = \frac{1}{2}. \end{aligned}$$

Это и означает, что  $r_1$  — случайная величина Бернулли. Аналогично доказывается бернулливость других величин  $r_2, \dots, r_n$ . Для доказательства независимости  $r_1, \dots, r_k$  выберем  $\varepsilon_1, \dots, \varepsilon_k \in \{0, 1\}$  и запишем

$$\begin{aligned} &P(r_1 = \varepsilon_1, \dots, r_k = \varepsilon_k) \\ &= \sum_{\varepsilon_{k+1}, \dots, \varepsilon_n \in \{0, 1\}} P(r_1 = \varepsilon_1, \dots, r_k = \varepsilon_k, r_{k+1} = \varepsilon_{k+1}, r_n = \varepsilon_n) \\ &= 2^{n-k} \cdot \frac{1}{2^n} = \frac{1}{2^k} = P(r_1 = \varepsilon_1) \dots P(r_k = \varepsilon_k). \end{aligned}$$

Это и означает, что случайные величины  $r_1, \dots, r_k$  независимы. Доказательство независимости других комбинаций из набора  $r_1, \dots, r_n$  абсолютно аналогично.  $\square$

*Доказательство теоремы 1.* Положим  $D \stackrel{\text{def}}{=} AB - C$ . Если  $AB\bar{\mathbf{r}} = C\bar{\mathbf{r}}$ , то  $D\bar{\mathbf{r}} = 0$ . Если при этом  $D \neq 0$ , то у матрицы  $D$  имеется ненулевой элемент. Не теряя общности считаем, что  $d_{11} \neq 0$ . Из  $D\bar{\mathbf{r}} = 0$  вытекает

$$\sum_{j=1}^n d_{1j}r_j = 0,$$

откуда

$$(4) \quad r_1 = -\frac{1}{d_{11}} \sum_{j=2}^n d_{1j}r_j.$$

Чтобы лучше представить дальнейшую оценку вероятности  $P(AB\bar{\mathbf{r}} = C\bar{\mathbf{r}})$ , рассмотрим сначала случай  $n = 2$ . Тогда (4) равносильно

$$(5) \quad r_1 = -\frac{d_{12}}{d_{11}}r_2 = \theta r_2, \quad \text{где} \quad \theta \stackrel{\text{def}}{=} -\frac{d_{12}}{d_{11}}.$$

Поскольку (5) вытекает из равенства  $AB\bar{\mathbf{r}} = C\bar{\mathbf{r}}$ , то

$$\{\omega: AB\bar{\mathbf{r}} = C\bar{\mathbf{r}}\} \subseteq \{\omega: r_1 = \theta r_2\}.$$

Следовательно,

$$\begin{aligned} P(AB\bar{\mathbf{r}} = C\bar{\mathbf{r}}) &= P(AB\bar{\mathbf{r}} = C\bar{\mathbf{r}}, r_2 = 0) + P(AB\bar{\mathbf{r}} = C\bar{\mathbf{r}}, r_2 = 1) \\ &\leq P(r_1 = \theta r_2, r_2 = 0) + P(r_1 = \theta r_2, r_2 = 1) \\ &= P(r_1 = 0, r_2 = 0) + P(r_1 = \theta, r_2 = 1). \end{aligned}$$

В силу независимости величин  $r_1$  и  $r_2$  имеем

$$\mathbb{P}(AB\bar{\Gamma} = C\bar{\Gamma}) \leq \mathbb{P}(r_1 = 0) \cdot \mathbb{P}(r_2 = 0) + \mathbb{P}(r_1 = \theta) \cdot \mathbb{P}(r_2 = 1).$$

Ясно, что  $\mathbb{P}(r_1 = \theta) \leq \frac{1}{2}$ , причем равенство возможно только при  $\theta = 0$  или  $\theta = 1$  (в остальных случаях эта вероятность равна 0). Поэтому первый сомножитель у слагаемых в правой части последнего неравенства не превосходит  $\frac{1}{2}$ , откуда и вытекает

$$\mathbb{P}(AB\bar{\Gamma} = C\bar{\Gamma}) \leq \frac{1}{2}(\mathbb{P}(r_2 = 0) + \mathbb{P}(r_2 = 1)) = \frac{1}{2}.$$

Доказательство для других  $n$  абсолютно аналогично: положим

$$\xi \stackrel{\text{def}}{=} -\frac{1}{d_{11}} \sum_{j=2}^n d_{1j} r_j.$$

Тогда  $\{AB\bar{\Gamma} = C\bar{\Gamma}\} \subseteq \{r_1 = \xi\}$  и из (4) получаем

$$\begin{aligned} & \mathbb{P}(AB\bar{\Gamma} = C\bar{\Gamma}) \\ &= \sum_{\varepsilon_2, \dots, \varepsilon_n \in \{0,1\}} \mathbb{P}(AB\bar{\Gamma} = C\bar{\Gamma}, r_2 = \varepsilon_2, \dots, r_n = \varepsilon_n) \\ &\leq \sum_{\varepsilon_2, \dots, \varepsilon_n \in \{0,1\}} \mathbb{P}(r_1 = \xi, r_2 = \varepsilon_2, \dots, r_n = \varepsilon_n). \end{aligned}$$

На событии  $\{r_1 = \xi, r_2 = \varepsilon_2, \dots, r_n = \varepsilon_n\}$  величина  $\xi$  равна

$$x \stackrel{\text{def}}{=} -\frac{1}{d_{11}} \sum_{j=2}^n d_{1j} \varepsilon_j$$

( $x$  отличается от  $\xi$  тем, что каждая величина  $r_j$ ,  $2 \leq j \leq n$ , в определении  $\xi$  заменяется на число  $\varepsilon_j$ ). В силу независимости случайной величины  $r_1$  и случайного вектора  $(r_2, \dots, r_n)$

оценки можно продолжить следующим образом:

$$\begin{aligned} & \mathbb{P}(AB\bar{\Gamma} = C\bar{\Gamma}) \\ & \leq \sum_{\varepsilon_2, \dots, \varepsilon_n \in \{0,1\}} \mathbb{P}(r_1 = x, r_2 = \varepsilon_2, \dots, r_n = \varepsilon_n) \\ & = \sum_{\varepsilon_2, \dots, \varepsilon_n \in \{0,1\}} \mathbb{P}(r_1 = x) \mathbb{P}(r_2 = \varepsilon_2, \dots, r_n = \varepsilon_n) \\ & \leq \frac{1}{2} \sum_{\varepsilon_2, \dots, \varepsilon_n \in \{0,1\}} \mathbb{P}(r_2 = \varepsilon_2, \dots, r_n = \varepsilon_n) \\ & = \frac{1}{2}. \end{aligned}$$

□