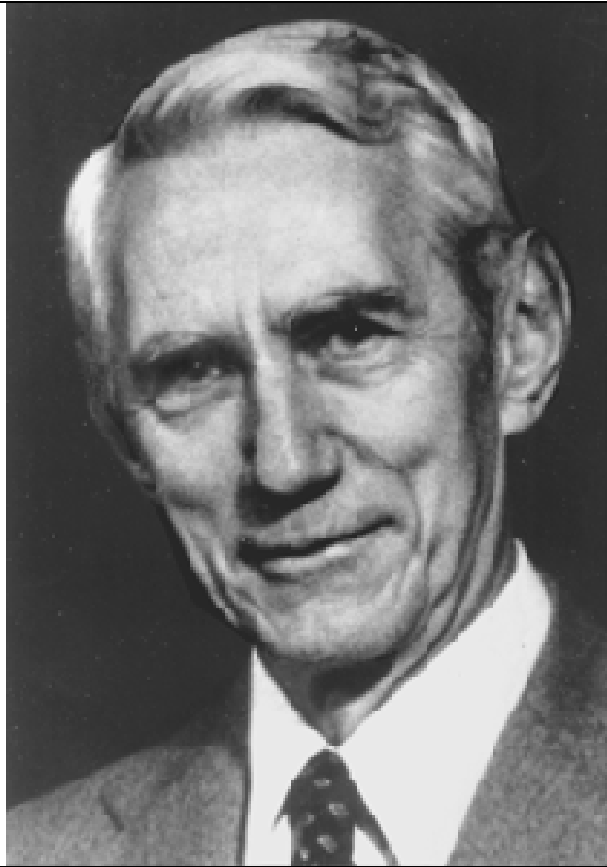


## **Тема 23. Алфавітне й рівномірне кодування. Роздільні коди. Оптимальне кодування**

### **План лекції**

- **Алфавітне й рівномірне кодування**
- **Достатні умови однозначності декодування. Властивості роздільних кодів.**
- **Оптимальні коди. Код Фано — код, близький до оптимального**

Кодування – у широкому розумінні – це перехід від одного способу подання інформації до іншого, який дає змогу відновлювати початкову інформацію. Теорія кодування виникла в 40-х роках ХХ століття після робіт К. Шеннона. У цій теорії досліджуються методи кодування, пов'язані з основними математичними моделями, які відображають істотні риси реальних інформаційних систем.



Клод Елвуд Шеннон — американський електротехнік і математик, «батько теорії інформації». Шеннон відомий тим, що запропонував теорію інформації в науковій статті, опублікованій в 1948 році.

Народився: 30 квітня 1916 р., Мічиган, США

Помер: 24 лютого 2001 р., Медфорд, Массачусетс, США

Нагороди: Медаль пошани IEEE, Національна наукова медаль США в галузі технічних наук

Освіта: Мічиганський університет, Массачусетський технологічний інститут

Теорема про пропускну здатність каналу: Будь-який канал з шумом характеризується максимальною швидкістю передачі повідомлення, ця межа названа в честь Шеннона. При передачі інформації зі швидкостями, які перевищують цю межу, відбувається невідворотне спотворення даних, але знизу до цієї межі можна наближатися з необхідною точністю, забезпечуючи задано малу ймовірність помилки передачі інформації в зашумленому каналі.

## Алфавітне й рівномірне кодування

Без утрати загальності можна сформулювати задачу кодування так. Нехай задано алфавіт  $A = \{a_1, \dots, a_r\}$  зі скінченної кількості символів, які називають *буквами*. Скінченну послідовність букв  $a_{i_1} a_{i_2} \dots a_{i_l}$  алфавіту  $A$  називають *словом* у цьому алфавіті. Для позначення слів ми будемо використовувати малі грецькі букви:  $\alpha = a_{i_1} a_{i_2} \dots a_{i_l}$ . Число  $l$  – кількість букв у слові  $\alpha$  – називають *довжиною* слова  $\alpha$  та позначають  $l(\alpha)$ . Множину всіх слів у алфавіті  $A$  позначають як  $A^*$ . Порожнє слово позначають  $\lambda$ ; зазначимо, що  $\lambda \notin A$ ,  $\lambda \in A^*$ ,  $l(\lambda)=0$ .

Якщо слово  $\alpha$  має вигляд  $\alpha = \alpha_1 \alpha_2$ , то  $\alpha_1$  називають *початком* (*префіксом*) слова  $\alpha$ , а  $\alpha_2$  – *закінченням* (*постфіксом*) слова  $\alpha$ . Якщо при цьому  $\alpha_1 \neq \lambda$  (відповідно  $\alpha_2 \neq \lambda$ ), то  $\alpha_1$  (відповідно  $\alpha_2$ ) називають *власним початком* (відповідно *власним закінченням*) слова  $\alpha$ .

Нехай  $S$  – підмножина множини  $A^*$ :  $S \subset A^*$ ; елементи множини  $S$  називають *повідомленнями*. Нехай також задано скінченний алфавіт  $B = \{b_1, \dots, b_q\}$ . Як  $\beta$  позначимо слово в алфавіті  $B$  і як  $B^*$  – множину всіх слів в алфавіті  $B$ . Задамо відображення  $f$ , яке кожному слову  $\alpha$ ,  $\alpha \in S$ , ставить у відповідність слово  $\beta = f(\alpha)$ ,  $\beta \in B^*$ . Слово  $\beta$  називають *кодом повідомлення*  $\alpha$ . Перехід від слова  $\alpha$  до його коду називають *кодуванням*.

Якщо  $|B| = q$ , то кодування називають *q-ковим*. Найчастіше використовують алфавіт  $B = \{0, 1\}$ , тобто *двійкове кодування*. Саме його ми й розглядатимемо в наступних підрозділах, випускаючи слово «двійкове».

Відображення  $f$  задають якимось алгоритмом. Є два способи задати відображення  $f$ .

**1. Алфавітне кодування.** Алфавітне кодування задають *схемою* (або *таблицею кодів*)  $\sigma$ :

$$\begin{aligned} a_1 &\rightarrow \beta_1, \\ a_2 &\rightarrow \beta_2, \\ &\dots\dots\dots \\ a_r &\rightarrow \beta_r, \end{aligned}$$

де  $a_i \in A$ ,  $\beta_i \in B^*$ ,  $i = 1, \dots, r$ .

Схема  $\sigma$  задає відповідність між буквами алфавіту  $A$  та деякими словами в алфавіті  $B$ . Вона визначає алфавітне кодування так: кожному слову  $\alpha = a_{i_1} a_{i_2} \dots a_{i_l}$  із множини  $S$  ставиться у відповідність слово  $\beta = \beta_{i_1} \beta_{i_2} \dots \beta_{i_l}$ . Це слово  $\beta$  називають *кодом* повідомлення  $\alpha$ . Слова  $\beta_1, \dots, \beta_r$  (див. схему  $\sigma$ ) називають *елементарними кодами*.

Схеми алфавітного кодування використовують, зокрема, для побудови оптимальних кодів, або кодів з мінімальною надлишковістю.

**2. Рівномірне кодування.** *Рівномірне (або блокове) кодування з параметрами  $k$  та  $n$  визначають так. Повідомлення  $\alpha$  розбивають на блоки довжиною  $k$ :*

$$\alpha = (x_1 \dots x_k) (x_{k+1} \dots x_{2k}) \dots (x_{mk+1} \dots x_{mk+s}),$$

*де  $s \leq k$  (останній блок може бути коротшим, у такому разі спосіб його кодування спеціально обумовлюють),  $x_i \in A$ ,  $i=1, \dots, mk+s$ .*

*Блоки довжиною  $k$  розглядають як „букви” якогось алфавіту (таких блоків є, очевидно,  $r^k$ , бо алфавіт  $A$  складається з  $r$  букв) і кодують словами в алфавіті  $B$  довжиною  $n$  за схемою рівномірного кодування  $\sigma_{k,n}$ :*

$$\begin{aligned} \alpha_1 &\rightarrow \beta_1, \\ \alpha_2 &\rightarrow \beta_2, \\ \alpha_3 &\rightarrow \beta_3, \\ &\dots\dots\dots \\ \alpha_{r^k} &\rightarrow \beta_{r^k}. \end{aligned}$$

На практиці використовують значення  $r=2$ , тобто двійкові блоки довжиною  $k$  кодують двійковими словами довжиною  $n$ , де  $n > k$ . Такі схеми рівномірного кодування використовують для побудови кодів, стійких до перешкод у каналах зв'язку. Головними формами стійкості коду до перешкод є здатність *виявляти деяку кількість помилок*, а також сильніша властивість – *виправляти помилки*. Це досягається введенням певної надлишковості.

*Надлишковістю схеми  $\sigma_{k,n}$  на символ повідомлення називають величину  $R = \frac{n-k}{k} = \frac{n}{k} - 1$ .*

Отже, ми описали дві можливості задати відображення  $f: S \rightarrow B^*$ . Однозначне декодування можливе, якщо існує обернене відображення  $f^{-1}$ .

## Достатні умови однозначності декодування. Властивості роздільних кодів

Розглянемо схему алфавітного кодування  $\sigma$  та різні слова, складені з елементарних кодів. Схему  $\sigma$  називають *роздільною*, якщо з того що

$$\beta_{i_1} \dots \beta_{i_k} = \beta_{j_1} \dots \beta_{j_l}$$

випливає, що  $k = l$  та  $i_t = j_t$  для кожного  $t = 1, \dots, k$ , тобто будь-яке слово, складене з елементарних кодів, можна єдиним способом розкласти на елементарні коди. Очевидно, що алфавітне кодування з роздільною схемою вможливорює однозначне декодування.

Схему  $\sigma$  називають *префіксною*, якщо для будь-яких  $i$  та  $j$  ( $i, j = 1, \dots, r$ ;  $i \neq j$ ) елементарний код  $\beta_i$  не є префіксом елементарного коду  $\beta_j$ .

### Теорема 1. Префіксна схема є роздільною.

**Доведення.** Від протилежного. Оскільки схема  $\sigma$  префіксна, то всі елементарні коди в  $\sigma$  попарно різні, тобто  $\beta_i \neq \beta_j$ , якщо  $i \neq j$ . Нехай кодування зі схемою  $\sigma$  не роздільне. Тоді існує таке слово  $\beta \in B^*$ , що можуть бути два розклади на елементарні коди:

$$\beta = \beta_{i_1} \dots \beta_{i_k}, \beta = \beta_{j_1} \dots \beta_{j_l}.$$

Нехай  $\beta_{i_1} = \beta_{j_1}, \dots, \beta_{i_{p-1}} = \beta_{j_{p-1}}$  але  $\beta_{i_p} \neq \beta_{j_p}$ . У такому разі одне зі слів  $\beta_{i_p}, \beta_{j_p}$  являє собою префікс іншого, а це суперечить тому, що схема  $\sigma$  префіксна.

Властивість префіксності достатня, але не необхідна умова роздільності схеми.

**Приклад.** Нехай  $A=\{a, b\}$ ,  $B=\{0, 1\}$ , схема  $\sigma: a \rightarrow 0, b \rightarrow 01$ . Ця схема не префіксна, але роздільна. Справді, перед кожним входженням 1 в слові  $\beta_2$  безпосередньо є 0. Це дає змогу виділити всі пари (01). Частина слова, що залишилась, складатиметься із символів 0.

Нехай  $\beta = y_1 y_2 \dots y_n$  – слово з  $B^*$ . Позначимо як  $\tilde{\beta}$  слово, одержане оберненням слова  $\beta$ , тобто  $\tilde{\beta} = y_n y_{n-1} \dots y_1$ . Позначимо як  $\tilde{\sigma}$  схему

$$\begin{aligned} a_1 &\rightarrow \tilde{\beta}_1, \\ a_2 &\rightarrow \tilde{\beta}_2, \\ &\dots\dots\dots \\ a_r &\rightarrow \tilde{\beta}_r. \end{aligned}$$

**Приклад.** Візьмемо схему  $\sigma$  з попереднього прикладу. Тоді  $\tilde{\sigma}$  має такий вигляд:  $a \rightarrow 0, b \rightarrow 10$ . Схема  $\tilde{\sigma}$  префіксна, тому за теоремою 1 – роздільна.

*Зауваження.* Схеми  $\sigma$  та  $\tilde{\sigma}$  водночас або роздільні, або ні.

Це зауваження дає змогу посилити теорему 1.

**Теорема 2.** Якщо або схема  $\sigma$ , або схема  $\tilde{\sigma}$  префіксна, то обидві схеми  $\sigma$  та  $\tilde{\sigma}$  роздільні.

Можна навести приклад такої роздільної схеми  $\sigma$ , що ні  $\sigma$ , ні  $\tilde{\sigma}$  не префіксні (див. задачу 2 на стор. 297 підручника). Отже, теорема 2 також дає достатню, але не необхідну умову роздільності схеми.

Нехай задано схему алфавітного кодування  $\sigma: a_i \rightarrow \beta_i, a_i \in A, \beta_i \in B^*, i = 1, \dots, r$ . Для того, щоб схема алфавітного кодування була роздільною, необхідно, щоб довжини елементарних кодів задовольняли *нерівність Мак-Міллана* (В. McMillan).

**Теорема 3 (нерівність Мак-Міллана).** Якщо схема алфавітного кодування  $\sigma$  роздільна, то

$$\sum_{i=1}^r \frac{1}{2^{l_i}} \leq 1, \text{ де } l_i = l(\beta_i).$$

*Зауваження 1.* Зазначимо, що нерівність Мак-Міллана – це необхідна умова роздільності схеми алфавітного кодування: якщо схема роздільна, то нерівність виконується. Якщо ж нерівність виконується, то з цього НЕ випливає, що схема роздільна: є нероздільні схеми, для яких нерівність Мак-Міллана виконується, наприклад:

$$a_1 \rightarrow 1, \quad a_2 \rightarrow 01, \quad a_3 \rightarrow 10.$$

Справді, повідомлення 1101 можна розкласти на елементарні коди двома способами: 1/1/01, що дасть при декодуванні  $a_1 a_1 a_2$ , і 1/10/1, що дасть  $a_1 a_3 a_1$ .

*Зауваження 2.* У теоремі 3 розглянуто нерівність Мак-Міллана для двійкового кодування. Для загального випадку  $q$ -кового кодування вона має вигляд:

$$\sum_{i=1}^r \frac{1}{q^{l_i}} \leq 1, \text{ де } l_i = l(\beta_i).$$

Розглянемо ще один важливий факт.

**Теорема 4.** Якщо числа  $l_1, \dots, l_r$  задовольняють нерівність

$$\sum_{i=1}^r \frac{1}{2^{l_i}} \leq 1 \text{ (нерівність Мак-Міллана),}$$

то існує префіксна схема алфавітного кодування  $\sigma$ :

$$a_1 \rightarrow \beta_1,$$

.....

$$a_r \rightarrow \beta_r,$$

де  $l_1 = l(\beta_1), \dots, l_r = l(\beta_r)$ .

**Наслідок 1.** Нерівність Мак-Міллана – необхідна й достатня умова існування алфавітного кодування з префіксною схемою та довжинами елементарних кодів, що дорівнюють  $l_1, \dots, l_r$ .

**Наслідок 2.** Якщо існує роздільна схема алфавітного кодування із заданими довжинами елементарних кодів, то існує й префіксна схема з тими самими довжинами елементарних кодів.

Для доведення потрібно спочатку застосувати теорему 3, а потім – теорему 4.



## Оптимальні коди. Код Фано – код, близький до оптимального

Досі термін «код» було використано в загальноприйнятому розумінні (як код повідомлення). Однак у теорії кодування, а ще раніше в техніці слово «код» трактують також і як *множину елементарних кодів*. Починаючи з цього місця, будемо використовувати слово «код» у двох значеннях. З контексту буде зрозуміло, яке саме з них ми маємо на увазі.

Нехай задано алфавіт  $A = \{a_1, \dots, a_r\}$  і ймовірності  $P = (p_1, \dots, p_r)$  появи букв у повідомленні; тут  $p_i$  – ймовірність появи букви  $a_i$ . Не зупиняючись на строгому означенні ймовірності, зауважимо лише, що ймовірність деякої букви можна уявити як частину випадків, у яких вона з'являється, від загальної кількості букв у повідомленні. Так, якщо  $p_1 = 0.25$ , то це означає, що з, наприклад, 1000 переданих букв близько 250 разів з'явиться буква  $a_1$ . Не втрачаючи загальності, можна вважати, що

$$p_1 \geq p_2 \geq \dots \geq p_r > 0,$$

тобто можна одразу вилучити букви, які не можуть з'явитись у повідомленні, і впорядкувати букви за спаданням ймовірностей їх появи. Крім того,  $p_1 + p_2 + \dots + p_r = 1$ . Кортеж  $P = (p_1, \dots, p_r)$  називають *розподілом ймовірностей* (появи букв алфавіту  $A$ ); це певна характеристика алфавіту повідомлень.

Для кожної роздільної схеми алфавітного кодування  $\sigma$  величину

$$l_{сер}^{\sigma}(P) = \sum_{i=1}^r p_i l_i,$$

де  $l_i = l(\beta_i)$ ,  $i = 1, \dots, r$ , і називають *середньою довжиною* кодування за схемою  $\sigma$  для розподілу ймовірностей  $P$ .

Можна побачити, що величина  $l_{сер}^{\sigma}$  ( $l_{сер}^{\sigma} > 1$ ) показує, у скільки разів зросте середня довжина слова в разі кодування зі схемою  $\sigma$ .

Можна довести, що мінімальне значення

$$l_* = \min_{\sigma} l_{сер}^{\sigma}$$

обов'язково досягається на якійсь схемі алфавітного кодування  $\sigma$ .

Коди, визначені схемою  $\sigma$  з  $l_{сер} = l_*$ , називають *кодами з мінімальною надлишковістю* (або *оптимальними кодами*) для розподілу ймовірностей  $P$ . За наслідком 2 з теорем 3 та 4 існує алфавітне кодування з префіксною схемою, яке дає оптимальні коди. У зв'язку з цим, будуючи оптимальні коди, можна обмежитися лише префіксними схемами.

Проста ідея побудови коду, близького до оптимального, належить американському вченому в галузі інформатики Р. Фано (R. Fano). Сформулюємо алгоритм кодування за методом Фано.

1. Упорядковуємо букви алфавіту  $A$  за спаданням імовірностей їх появи в повідомленні.
2. Розбиваємо множину букв, записаних у зазначеному порядку, на дві послідовні (тобто без перестановок букв) частини так, щоб сумарні ймовірності кожної з них були якомога близькими одна до одної. Кожній букві з першої частини приписуємо символ 0, другої – символ 1. Далі те саме робимо із кожною частиною, якщо вона містить принаймні дві букви. Процедуру продовжуємо доти, доки всю множину не буде розбито на окремі букви.



### **Роберт Фано**

11 листопада 1917 р. (м. Турин, Італія – 13 липня 2016 р.(м. Нейплз, Флорида, США)

Освіта: Массачусетський технологічний інститут

У 1939 р переїхав у США. У 1961 р. під керівництвом його та Фернандо Корбатто було створено систему, яка підтримувала режим розподілу часу (Compatible Time-Sharing System, CTSS). [Автор коду Фано](#). У 1963 р. очолив проект MAC і залишався на цій посаді до 1968 року. [Фото 30.05.2014 зі святкування 50 річниці проекту MAC](#).

**Приклад.** Нехай задано розподіл імовірностей  $P=(0.4, 0.15, 0.15, 0.15, 0.15)$ . Побудуємо код за методом Фано. Розв'язок подано в таблиці 1.

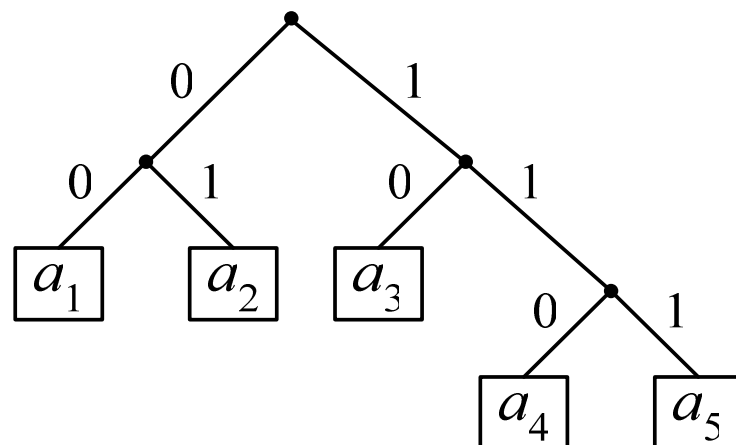
Середня довжина коду дорівнює

$$l_{\text{сер}}^F = 0.4 \times 2 + 0.15 \times 2 + 0.15 \times 2 + 0.15 \times 3 + 0.15 \times 3 = 2.3,$$

тут верхній індекс  $F$  означає, що код отримано методом Фано.

**Таблиця 1**

Букви алфавіту $A$	Імовірності появи букв	Розбиття множини букв			Елементарні коди
$a_1$	0.4	0	0		00
$a_2$	0.15		1		01
$a_3$	0.15	1	0		10
$a_4$	0.15		1	0	110
$a_5$	0.15			1	111



**Рис. 1**

Алгоритм Фано має просту інтерпретацію за допомогою бінарного дерева. Від кореня відходять два ребра, ліве позначено символом 0, а праве – символом 1. Ці два ребра відповідають розбиттю множини всіх букв на дві майже рівноймовірні частини, одній з яких поставлено у відповідність символ 0, а другій – символ 1. Ребра, що виходять із вершин наступного рівня, відповідають розбиттю одержаних частин знову на дві майже рівноймовірні послідовні частини. Цей процес продовжують доти, доки множини букв не буде розбито на окремі букви. Кожний листок дерева відповідає певному елементарному коду. Щоб виписати цей код, потрібно пройти шлях від кореня до відповідного листка.

На рисунку 1. зображено кодове дерево, отримане методом Фано для розподілу ймовірностей із останнього прикладу.

Зазначимо, що незалежно від способу кодування кожному бінарному дереву відповідає набір двійкових елементарних кодів. У такому разі дерево називають *кодовим*. Якщо елементарні коди відповідають листкам кодового дерева, то відповідна схема алфавітного кодування є префіксною (отже, забезпечується однозначність декодування).