

Національний технічний університет України

«Київський політехнічний інститут»

Фізико-технічний інститут

Криптографія

Комп'ютерний практикум №1

Експериментальна оцінка ентропії на символ джерела відкритого тексту

Виконали:

студенти групи ФБ-32

Грабовецький Микита

Драбок Алла

Київ - 2025

Мета роботи:

Засвоєння понять ентропії на символ джерела та його надлишковості, вивчення та порівняння різних моделей джерела відкритого тексту для наближеного визначення ентропії, набуття практичних навичок щодо оцінки ентропії на символ джерела.

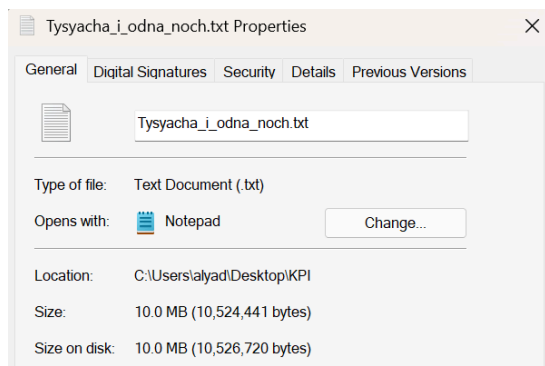
Порядок виконання роботи:

1. Написати програми для підрахунку частот букв і частот біграм в тексті, а також підрахунку H_1 та H_2 за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення H_1 та H_2 на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення H_1 та H_2 на тому ж тексті, в якому вилучено всі пробіли.
2. За допомогою програми CoolPinkProgram оцінити значення $H^{(10)}$, $H^{(20)}$, $H^{(30)}$.
3. Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела.

Варіант: 13

Хід роботи:

Ми обрали книгу, яку будемо аналізувати в нашій програмі, а саме: «Тисяча і одна ніч». Як бачимо, вона важить 10 мб, але будемо використовувати лише перші 2 мегабайти (прописано в коді). Кількість літер після фільтрації складає 1975786.



Наша програма замінює всі непотрібні символи на нижні підкреслення (таким чином ми позначити пробіли), а всі подвійні нижні підкреслення (тобто пробіли) видаляє. Замінює «ё» на «е», «ъ» на «ь» і «<» на «_» (знову ж, таким чином позначили пробіл).

Підраховує частоти літер, біграм, що перетинаються та не перетинаються. Топ-30 літер та топ-10 біграм будуть показані нижче, однак у нас також створюється таблиця Excel, де зберігаються абсолютно всі можливі результати (додана у репозиторій). Також обчислюємо ентропії H_1 та H_2 на тексті з та без пробілів.

Топ-30 літер:

'_': 0.18351
'о': 0.08441
'а': 0.07528
'и': 0.06751
'е': 0.06434
'н': 0.04734
'т': 0.04695
'л': 0.04586
'с': 0.04331
'в': 0.03350
'р': 0.03288
'д': 0.02742
'к': 0.02636
'м': 0.02566
'у': 0.02451
'п': 0.02113
'я': 0.01799
'ь': 0.01661
'з': 0.01538
'б': 0.01459
'г': 0.01421
'ы': 0.01332
'ч': 0.01061
'й': 0.00854
'х': 0.00850
'ш': 0.00823
'ж': 0.00669
'ю': 0.00545
'ц': 0.00474
'э': 0.00241

Топ-10 біграм (перетин):

'и_': 0.03079
'_и': 0.02214
'а_': 0.02111
'о_': 0.01829
'_с': 0.01704
'_п': 0.01563
'е_': 0.01451
'_н': 0.01433
'_о': 0.01366
'_в': 0.01335

Топ-10 біграм (без перетину):

'и_': 0.03080
'_и': 0.02218
'а_': 0.02102
'о_': 0.01824
'_с': 0.01711
'_п': 0.01560
'е_': 0.01453
'_н': 0.01426
'_о': 0.01379
'_в': 0.01325

Результаты для $H^{(20)}$

Произвольная часть текста:
арушать_его_ту_же_и

Использованные буквы:

Порядок n-граммы:
5 символов
10 символов
15 символов
20 символов
25 символов
30 символов
35 символов
40 символов
45 символов
50 символов

Введенный символ:
Символ по счету:
Номер эксперимента: 50

Неравенство для энтропии:
 $2,97654261593089 < H < 3,49995660038497$

Двоичная таблица угаданных символов:
00000000100000000000000000000000
00000100000000000000000000000000
00000000000000100000000000000000
00000000100000000000000000000000
10000000000000000000000000000000

Поле ввода символов:
Продолжить Другой

Вероятности:
 $q[1] = 0,3061224$
 $q[2] = 0,0816326$
 $q[3] = 0,0612244$
 $q[4] = 0$
 $q[5] = 0,0408163$
 $q[6] = 0,0612244$
 $q[7] = 0$
 $q[8] = 0,0612244$
 $q[9] = 0,0816326$
 $q[10] = 0,020408$
 $q[11] = 0,061224$
 $q[12] = 0,020408$
 $q[13] = 0$
 $q[14] = 0,040816$
 $q[15] = 0,040816$
 $q[16] = 0,040816$
 $q[17] = 0$
 $q[18] = 0,020408$
 $q[19] = 0$
 $q[20] = 0$
 $q[21] = 0$
 $q[22] = 0$
 $q[23] = 0,040816$
 $q[24] = 0$
 $q[25] = 0$
 $q[26] = 0,020408$
 $q[27] = 0$
 $q[28] = 0$
 $q[29] = 0$
 $q[30] = 0$
 $q[31] = 0$
 $q[32] = 0$

Строка состояния:

Результаты для $H^{(30)}$

Произвольная часть текста:
оры_не_имеют_никакого_значени

Использованные буквы:
_, о, а, и, е, н, т, л, с,

Порядок n-граммы:
5 символов
10 символов
15 символов
20 символов
25 символов
30 символов
35 символов
40 символов
45 символов
50 символов

Введенный символ: с
Символ по счету: 9
Номер эксперимента: 50

Неравенство для энтропии:
 $3,04745268039148 < H < 3,68674567447369$

Двоичная таблица угаданных символов:
10000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000

Поле ввода символов:
Продолжить Другой

Вероятности:
 $q[1] = 0,2653061$
 $q[2] = 0,0612244$
 $q[3] = 0,0816326$
 $q[4] = 0,1020408$
 $q[5] = 0,0408163$
 $q[6] = 0,0408163$
 $q[7] = 0,0408163$
 $q[8] = 0,0816326$
 $q[9] = 0$
 $q[10] = 0,040816$
 $q[11] = 0,040816$
 $q[12] = 0$
 $q[13] = 0$
 $q[14] = 0,040816$
 $q[15] = 0,020408$
 $q[16] = 0,020408$
 $q[17] = 0,020408$
 $q[18] = 0$
 $q[19] = 0$
 $q[20] = 0$
 $q[21] = 0,040816$
 $q[22] = 0$
 $q[23] = 0,020408$
 $q[24] = 0,020408$
 $q[25] = 0$
 $q[26] = 0$
 $q[27] = 0$
 $q[28] = 0,020408$
 $q[29] = 0$
 $q[30] = 0$
 $q[31] = 0$
 $q[32] = 0$

Строка состояния:
Вы не угадали. Введите другую букву

Аналіз результатів:

Спочатку проаналізуємо наші значення **ентропій**:

- $H_0 = 5$ біт/символ - це максимальна ентропія, яка відповідає ситуації, коли всі 32 символи алфавіту з'являються з однаковою ймовірністю. Однак це дуже ідеалізована ситуація, а мова має свої особливості.
- $H_1 \approx 4.32$ - значення вже менше, оскільки частота появи літер нерівномірна (наприклад, літери 'о', 'а', 'е' зустрічаються набагато частіше, ніж 'ф', 'щ' або 'э'), тож враховуючи це можна трошки краще передбачити наступний символ.
- $H_2 \approx 3.83$ - це значення ще нижче, тому що модель біграм враховує залежність між сусідніми літерами (наприклад, після приголосної літери ймовірність появи голосної набагато вища, ніж іншої приголосної), а також існують стійкі буквосполучення та майже неможливі. Тому знання попередньої літери суттєво звужує коло можливих наступних літер, що ще більше знижує ентропію.

Маємо: $H_0 > H_1 > H_2$.

Тепер проаналізуємо $H(10)$, $H(20)$ та $H(30)$ з програми **CoolPinkProgram**. По ідеї, зі збільшенням кількості символів у послідовності діапазон значень ентропії мав би зміщуватись вниз, тобто мало б бути $H(10) > H(20) > H(30)$. Однак фактично вийшли такі результати:

- $3.09 < H(10) < 3.43$
- $2.98 < H(20) < 3.5$
- $3.05 < H(30) < 3.69$

Це є аномалією, адже фактично виходить, що у нас $H(30) > H(10) > H(20)$. Найімовірніше, що це похибка, спричинена специфікою тексту, а також кількістю наших експериментів. Для більш точних значень краще було б взяти більшу кількість тестів.

Тож підсумуючи, маємо: $H_0 > H_1 > H_2 > H(30) > H(10) > H(20)$.

Отже, чим більший порядок моделі (чим довший ланцюжок попередніх символів ми враховуємо), тим точніше ми можемо передбачити наступний символ. Це відбувається тому, що моделі вищого порядку враховують не тільки частоти літер і біграм, але й складніші мовні патерни, синтаксис та семантику. Тому з ростом n ентропія H_n зменшується.

Тепер порівняємо наші результати в розрізі **тексту з та без пробілів**.

- H_1 з пробілами ≈ 4.32
- H_1 без пробілів ≈ 4.44

- H_2 з пробілами ≈ 3.83
- H_2 без пробілів ≈ 4.05

Як бачимо, значення ентропій у тексті без пробілів більше. Але цьому є логічне пояснення. У таблиці із частотами букв чітко видно, що “_” (тобто пробіл) є найчастішим. І коли ми його видаляємо, то виходить, що розподіл частот літер також змінюється, а від нього залежить ентропія. Тобто видалення одного дуже частого символу зробило розподіл ймовірностей решти літер трохи більш рівномірним, що і призвело до збільшення ентропії. Тож якщо казати простими словами, то коли зникає найбільш передбачуваний елемент, система в цілому стає менш передбачуваною.

Тепер перейдемо до аналізу **надлишковості**.

Надлишковість (R) — це показник, який демонструє, наскільки мова є "стискаємою" або передбачуваною. Шукаємо за формулою: $1 - (H_n/H_0)$. З наших результатів бачимо, що чим менша ентропія, тим більша надлишковість. Що є логічним, адже ми ми виявляємо все більше закономірностей у мові. Тож конкретні цифри:

- При $R_1 \approx 13.7\%$ ми враховуємо лише частоту літер.
- При $R_2 \approx 23.4\%$ ми вже враховуємо зв'язки між парами літер, що розкриває більшу частину структури мови.

При більших n надлишковість буде ще більшою, оскільки модель враховуватиме слова, фрази та граматичні конструкції, роблячи текст ще більш передбачуваним.

За даними з інтернету надлишковість української мови оцінюється у приблизно 70%. Результати нашої програми показали 13.7% та 23.4% для літер та біграм відповідно для тексту із пробілами й 11.2% та 18.9% для тексту без пробілів. CoolPinkProgram безпосередньо не показує значення надлишковості, однак ми можемо його порахувати самостійно за нашою формулою, взявши $H(10)$, яке є середнім для нашого випадку. Тож:

$$H(10) \text{ середнє} = (3.09 + 3.43) / 2 = 3.26 \text{ біт/секунду}$$

$$R = 1 - (3.26/5) = 1 - 0.652 = 0.348 = 34.8\%$$

Як бачимо, це значення вже ближче, але все ще далеко від реального.

Значення з нашої програми на пайтоні найнижчі, адже наші моделі дуже прості. Вони не враховують існування слів, граматики, синтаксису, сенсу речень, що є основними

джерелами надлишковості в мові. Значення з програми CoolPinkProgram вже краще, адже програма використовує моделі значно вищого порядку, які здатні вловити більш складні статистичні закономірності тексту. Але цього все ще недостатньо для розкриття повної надлишковості мови. Отже, чим складніша модель (чим більше контексту вона враховує), тим повніше вона розкриває справжню надлишковість мови.

Висновки:

У межах лабораторної роботи було проведено аналіз тексту шляхом розрахунку частот символів та біграм, ентропій та надлишковості. Ми виявили, що зі збільшенням порядку моделі (від H_0 до H_2 і вище), ентропія джерела послідовно зменшується. А значне зниження ентропії від $H_1 \approx 4.32$ до $H_2 \approx 3.83$ біт/символ довело наявність зв'язків між сусідніми літерами. Надлишковість мови, що є показником її передбачуваності, відповідно зросла (від $R_1 \approx 13.7\%$ до $R_2 \approx 23.4\%$), підтверджуючи виявлення мовних патернів. Видалення найбільш частого символу (пробілу) збільшило H_1 та H_2 , що свідчить про те, що його присутність робила текст більш передбачуваним. Експериментальні результати з високопорядковими n -грамами (до $n=30$) дали оцінку надлишковості близько 35%, що є ближчим до справжніх лінгвістичних показників.