

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ
АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ ИТМО»

ФАКУЛЬТЕТ ТЕХНОЛОГИЙ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА



Отчёт по лабораторной работе № 2
“Hashcoding против спама”
по дисциплине “Алгоритмы и структуры данных”
Семестр 2

Выполнил студент:
Мавров Артём Николаевич
гр. J3113
ИСУ 466574
Отчет сдан: 24.04.2025

Санкт-Петербург
2025

Содержание

1	Введение	3
2	Теоретическая часть	4
2.1	Murmur hash 3	4
2.2	Feature hashing	4
2.3	Логическая регрессия	5
3	Практическая часть	6
3.1	Murmur hashing 3	6
3.2	Feature hashing	7
3.3	Логическая регрессия	7
3.4	Тестирование	9
4	Заключение	10

Введение

В данной лабораторной работе предстоит:

- Реализовать алгоритм хеширования “Murmur hash 3”
- Преобразовать строки в вектора с помощью “Feature hashing”
- Реализовать логическую регрессию для определения спама в датасете
- Обучить и протестировать модель.

Весь код реализован на C++. Репозиторий с кодом можно найти по [ссылке](#).

Теоретическая часть

Murmur hash 3

MurmurHash3 — это некриптографическая хеш-функция, разработанная для быстрого хеширования с хорошими распределительными свойствами. Она основана на принципах перемешивания и эффекта лавины.

Определение: Пусть $x \in \{0, 1\}^*$ — входная последовательность байтов. MurmurHash3 отображает x в фиксированную длину $h(x) \in \{0, 1\}^n$ по формуле:

$$h(x) = \text{Mix}(\text{BlockMix}(x), \text{Seed})$$

Обоснование:

- **Перемешивание:** Каждый блок входных данных умножается на константу и циклически сдвигается:

$$\text{BlockMix}(x_i) = \text{ROTL}(x_i \cdot c_1, r_1) \cdot c_2$$

Это обеспечивает диффузию, аналогичную универсальному хешированию.

- **Эффект лавины:** Малое изменение входа приводит к значительному изменению выхода:

$$\forall x, x' : d_H(x, x') = 1 \Rightarrow d_H(h(x), h(x')) \approx \frac{n}{2}$$

где d_H — расстояние Хэмминга.

- **Снижение смещённости:** Финальные операции дополнительно перемешивают биты для устранения статистических зависимостей:

$$h \leftarrow h \oplus (h \gg r) \cdot c_3$$

MurmurHash3 не является криптостойким, но обеспечивает высокую скорость и равномерное распределение для общих задач хеширования.

Feature hashing

Feature hashing — это метод отображения категориальных признаков в фиксированное пространство меньшей размерности с помощью хеш-функции.

Определение: Пусть входной вектор x имеет категориальные признаки, каждый из которых представлен парой (f_j, v_j) , где f_j — имя признака, v_j — его значение. Тогда хеш-преобразование:

$$\phi_i(x) = \sum_{j:h(f_j)=i} \xi(f_j) \cdot v_j, \quad i = 1, \dots, d$$

где:

- $h : \mathcal{F} \rightarrow \{1, \dots, d\}$ — хеш-функция, определяющая индекс в выходном пространстве размерности d ;
- $\xi : \mathcal{F} \rightarrow \{-1, +1\}$ — случайная знаковая хеш-функция, уменьшающая корреляцию;
- $\phi(x) \in \mathbb{R}^d$ — выходной вектор признаков фиксированной размерности.

Логическая регрессия

Логистическая регрессия моделирует вероятность бинарного события. Пусть:

- $y \in \{0, 1\}$ - целевая переменная
- $\mathbf{x} \in \mathbb{R}^n$ - вектор признаков
- $\mathbf{w} \in \mathbb{R}^n$ - веса модели
- $b \in \mathbb{R}$ - смещение

Линейная комбинация:

$$z = \mathbf{w}^T \mathbf{x} + b$$

Вероятность вычисляется через сигмоидную функцию:

$$P(y = 1|\mathbf{x}) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Функция потерь Кросс-энтропия для m примеров:

$$J(\mathbf{w}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \sigma(z^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(z^{(i)}))]$$

Оптимизация Градиентный спуск для обновления весов:

$$\mathbf{w} := \mathbf{w} - \alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

где α - скорость обучения.

Практическая часть

Murmur hashing 3

Функция реализована согласно описанию из теоретической части.

```
1 uint32_t rotl32(uint32_t x, int r) {
2     return (x << r) | (x >> (32 - r));
3 }
4
5 uint32_t fmix32(uint32_t h) {
6     h ^= h >> 16;
7     h *= 0x85ebca6b;
8     h ^= h >> 13;
9     h *= 0xc2b2ae35;
10    h ^= h >> 16;
11    return h;
12 }
13
14 void MurmurHash3 ( const void * key, int len, uint32_t seed, void *
15 out ) {
16     uint32_t h1 = seed;
17
18     const uint32_t c1 = 0xcc9e2d51;
19     const uint32_t c2 = 0x1b873593;
20
21     const uint8_t* data = (const uint8_t*)key;
22     const int nblocks = len / 4;
23
24     const uint32_t* blocks = (const uint32_t*)(data);
25     for (int i = 0; i < nblocks; i++) {
26         uint32_t k1 = blocks[i];
27         k1 *= c1;
28         k1 = rotl32(k1, 15);
29         k1 *= c2;
30         h1 ^= k1;
31         h1 = rotl32(h1, 13);
32         h1 = h1 * 5 + 0xe6546b64;
33     }
34
35     const uint8_t* tail = (const uint8_t*)(data + nblocks * 4);
36     uint32_t k1 = 0;
37     switch (len & 3) {
38         case 3: k1 ^= tail[2] << 16; break;
39         case 2: k1 ^= tail[1] << 8; break;
40         case 1: k1 ^= tail[0];
41                 k1 *= c1; k1 = rotl32(k1, 15); k1 *= c2; h1 ^= k1;
42     }
43
44     h1 ^= len;
45     h1 = fmix32(h1);
46     std::memcpy(out, &h1, sizeof(h1));
47 }
```

Feature hashing

В ходе тестирования было принято убрать L-2 нормализацию, так как она негативно сказывалась на метриках Recall и Accuracy. Причиной такого нетипичного поведения с большой долей вероятности является не корректность датасета.

Также для улучшения метрик была попытка добавить n-граммы, но это также оказало негативный эффект.

```
1 vector<double> text_to_features(const std::string& text) {
2     vector<double> feats(HASH_DIM, 0);
3
4     string low_text = text;
5     transform(low_text.begin(), low_text.end(), low_text.begin(), ::
        tolower);
6
7     istringstream iss(low_text);
8     string word;
9     while (iss >> word) {
10         uint32_t hash = 0;
11         MurmurHash3(&word, word.size(), 42, &hash);
12         feats[hash % HASH_DIM] += 1;
13     }
14
15     return feats;
16 }
```

Логическая регрессия

Для логической регрессии была выбрана функция $\sigma(z) = \frac{1}{1+e^{-z}}$. Обучения осуществляется функцией train, evaluate используется для подсчёта метрик.

```
1 class LogisticRegression {
2 public:
3     size_t dim = HASH_DIM;
4     double lr = 0.05;
5     double reg_lambda = 0.05;
6     double class_weight_0 = 1.0; // (ham)
7     double class_weight_1 = 1.0; // (spam)
8     LogisticRegression()
9     : weights(dim, 0.0),
10       learning_rate(lr),
11       lambda(reg_lambda),
12       epochs(2000) {}
13
14     void train(const std::vector<Sample>& trainData, const std:::
        vector<Sample>& validData) {
15         for (int epoch = 0; epoch < epochs; ++epoch) {
16             vector<double> gradients(dim, 0.0);
17             double loss = 0.0;
18
19             for (const auto& sample : trainData) {
20                 double pred = sigmoid(dot_product(weights, sample.
                    features));
21                 double error = pred - sample.label;
```

```

22         if (pred == 0) {
23             pred += 1e-9;
24         } else if (pred == 1) {
25             pred -= 1e-9;
26         }
27         loss += -sample.label * log(pred) - (1 - sample.label
28             ) * log(1 - pred);
29
30         for (size_t i = 0; i < dim; ++i) {
31             double class_weight = (sample.label == 1) ?
32                 class_weight_1 : class_weight_0;
33             gradients[i] += class_weight * error * sample.
34                 features[i];
35         }
36     }
37
38     loss /= trainData.size();
39     for (size_t i = 0; i < dim; ++i) {
40         weights[i] -= learning_rate * (gradients[i] /
41             trainData.size() + lambda * weights[i]);
42     }
43
44     if (epoch % 50 == 0) {
45         learning_rate *= 1.0;
46     }
47 }
48
49 std::vector<double> evaluate(const std::vector<Sample>& data)
50 const {
51     double TP = 0.0, TN = 0.0, FP = 0.0, FN = 0.0;
52
53     for (const auto& sample : data) {
54         int prediction = predict(sample.features);
55
56         if (sample.label == 1 && prediction == 1) {
57             TP += 1.0;
58         } else if (sample.label == 0 && prediction == 0) {
59             TN += 1.0;
60         } else if (sample.label == 0 && prediction == 1) {
61             FP += 1.0;
62         } else if (sample.label == 1 && prediction == 0) {
63             FN += 1.0;
64         }
65     }
66
67     return {TP, TN, FP, FN};
68 }
69
70 int predict(const FEATURE& feats) const {
71     return (sigmoid(dot_product(weights, feats)) >= 0.5) ? 1 : 0;
72 }

```



```

70 private:
71     std::vector<double> weights;
72     double learning_rate;
73     double lambda;
74     int epochs;
75
76     static double dot_product(const FEATURE& w, const FEATURE& x) {
77         double res = 0.0;
78         for (size_t i = 0; i < w.size(); ++i) {
79             res += w[i] * x[i];
80         }
81         return res;
82     }
83
84     static double sigmoid(double z) {
85         return 1.0 / (1.0 + exp(-z));
86     }
87 };

```

Тестирование

В ходе тестирования было выявлено, что L-2 нормализация негативно сказывается на метриках, n-граммы так-же не оказывают положительного влияния. Эпохи так-же не дают положительный результат. Также был понижен learning rate для градиентного спуска.

Учитывая эти доработки метрика *Accuracy* ≈ 0.65 , *Recall* ≈ 0.97 .

Заключение

По итогам лабораторной работы были получены и закреплены теоретические знания хэш-функций, устройства логической регрессии. Реализована хэш-функция Murmur hashing 3, преобразованы строки в вектора при помощи Feature hashing, реализован класс логической регрессии.

Все функции протестированы и грамотно работают. Из недочётов можно выявить невысокие показатели матрицы Ассигасу.

Для улучшения метрик, можно:

- Продолжать тестировать модель с разным набором констант, для подбора более лучших.
- Попробовать использовать такие методы как: метод максимального правдоподобия, метод Ньютона-Рафсона, кросс-валидация.