

Отчет по практическому заданию №2 "Применение линейных моделей для определения токсичности комментария".

Логистическая регрессия и градиентный спуск.

Содержание

1	Введение	2
2	Эксперименты	2
2.1	Исследование поведения градиентного спуска	2
2.1.1	Параметр размера шага step_alpha	2
2.1.2	Параметр размера шага step_beta	3
2.1.3	Начальное приближение w_0	4
2.2	Исследование поведения стохастического градиентного спуска	5
2.2.1	Параметр размера шага step_alpha	5
2.2.2	Параметр размера шага step_beta	6
2.2.3	Размер подвыборки batch_size	7
2.2.4	Начальное приближение w_0	8
2.3	Сравнение градиентного спуска и стохастического градиентного спуска	9
2.4	Лемматизация и удаление стоп-слов	10
2.5	Сравнение представлений BagOfWords и TF-IDF с различными параметрами	10
3	Применение лучших алгоритмов с каждого эксперимента к тестовой выборке	10

1 Введение

В данном документе представлен отчет о проделанных экспериментах по практическому заданию №2, анализ результатов. Краткое описание задания: необходимо реализовать линейный классификатор с произвольной функцией потерь.

2 Эксперименты

В этом блоке приведены все обязательные эксперименты, которые изложены в формулировке задания. Все эксперименты проводились на упрощенном датасете (рассматривается задача бинарной классификации) из соревнования **Toxic Comment Classification Challenge**, в котором нужно определить токсичность комментария.

Стандартный дизайн эксперимента:

- Оценка качества и подбор параметров модели проводились на каждой эпохе с помощью отложенной тренировочной выборки (30%). Все графики ниже построены по значениям ассигасу, посчитанным на отложенной выборке.
- В тренировочную выборку был добавлен признак, состоящий из всех единиц, который позволяет учитывать смещение (**bias**). Было решено не использовать смещение в $L2$ -регуляризации, чтобы даже при плохом выборе коэффициента регуляризации решающая гиперплоскость не вырождалась в 0.
- В стохастическом градиентном спуске проверяется критерий останова на каждой эпохе (не итерации).

2.1 Исследование поведения градиентного спуска

Обновления весов модели при использовании градиентного спуска происходит по следующей формуле:

$$w_t = w_{t-1} - \frac{\alpha}{t^\beta} \times \frac{1}{N} \times \sum_{i=1}^N \nabla_w \mathcal{L}(x_i, y_i | w_{t-1}), \quad (1)$$

где t - номер итерации, β - **step_beta**, $\nabla_w \mathcal{L}(x_i, y_i | w_{t-1})$ - градиент функции потерь.

2.1.1 Параметр размера шага **step_alpha**

Параметр **step_alpha** (α) используется в градиентном спуске при обновлении весов в формуле 1. Рассмотрим следующие зависимости при разных значениях параметра **step_alpha**:

1. зависимость значения функции потерь от реального времени работы метода
2. зависимость точности (ассигасу) от реального времени работы метода
3. зависимость значения функции потерь от итерации метода
4. зависимость точности (ассигасу) от итерации метода

Соответствующие графики приведены на: рис. 1, 2, 3, 4.

Рис. 1: Зависимость значения функции потерь от реального времени работы градиентного спуска

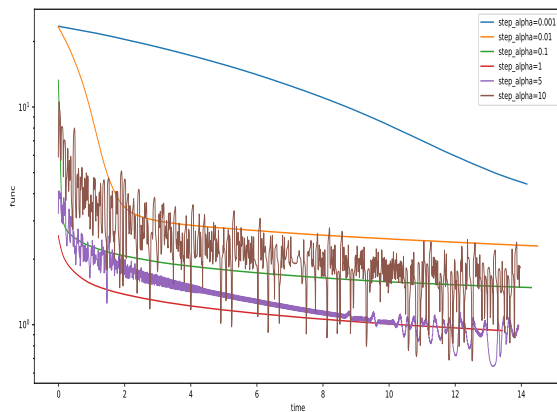


Рис. 2: Зависимость значения точности (ассигасу) от реального времени работы градиентного спуска

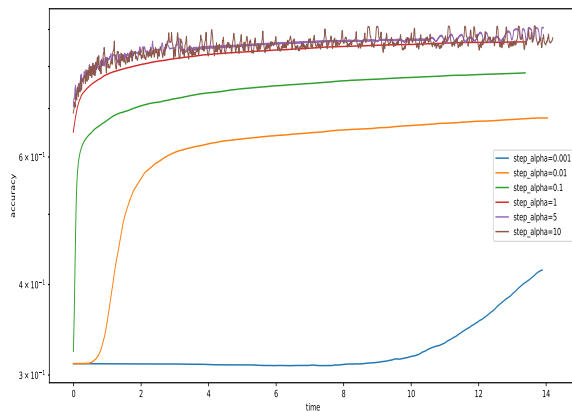


Рис. 3: Зависимость значения функции потерь от итерации метода градиентного спуска

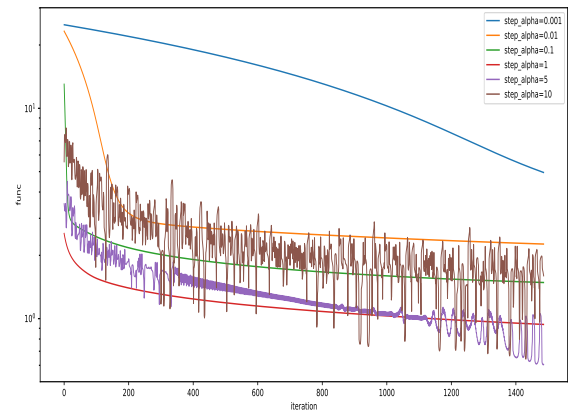
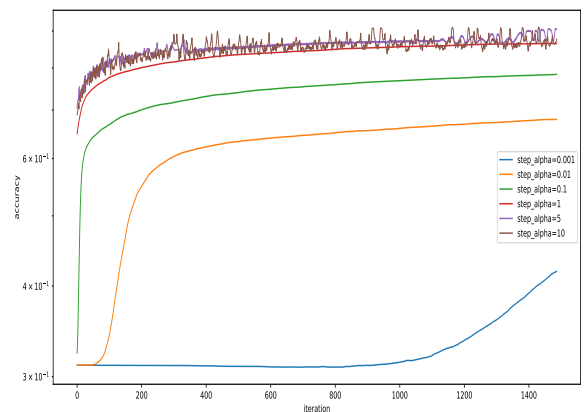


Рис. 4: Зависимость значения точности (ассигасу) от итерации метода градиентного спуска



Из графиков видно, что при значениях α , близких к нулю алгоритму нужно больше времени для сходимости. С другой стороны, если значения слишком большие, то алгоритм становится крайне не стабильным.

2.1.2 Параметр размера шара step_beta

Параметр **step_beta** (β) используется в градиентном спуске при обновлении весов в формуле 1. Аналогично предыдущему пункту рассмотрим зависимости из 2.1.1 при разных значениях параметра **step_beta** и проанализируем соответствующие графики, представленные на рис. 5, 6, 7, 8.

Рис. 5: Зависимость значения функции потерь от реального времени работы градиентного спуска

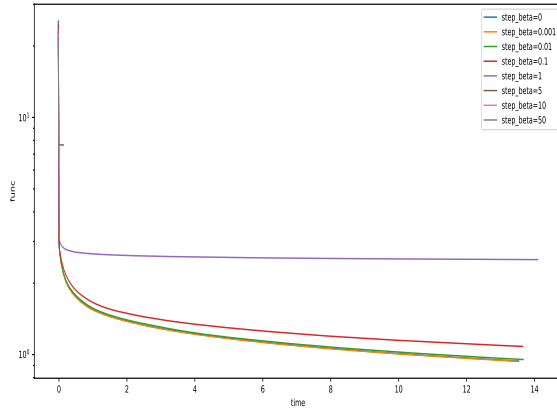


Рис. 6: Зависимость значения точности (ассигасу) от реального времени работы градиентного спуска

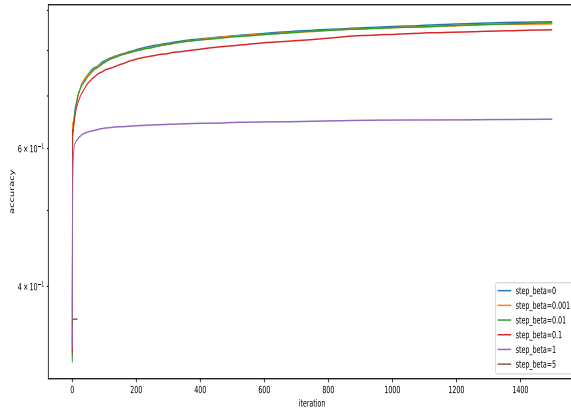


Рис. 7: Зависимость значения функции потерь от итерации метода градиентного спуска

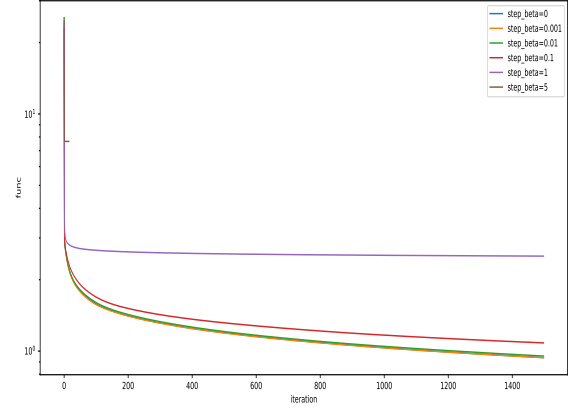
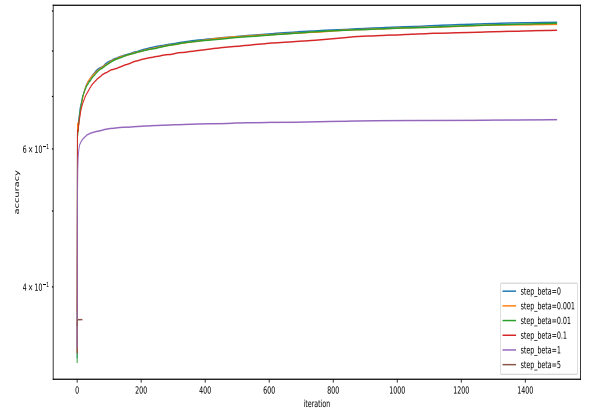


Рис. 8: Зависимость значения точности (ассигасу) от итерации метода градиентного спуска



Из графиков можно заметить, что значения **step_beta**, близкие к 0, приводят к одинаковому качеству. С увеличением же параметра **step_beta** значение точности уменьшается. При **step_beta = 5** изменение функции потерь так мало, что критерий останова срабатывает до первых 200 итераций.

2.1.3 Начальное приближение w_0

Начальное приближение нужно для инициализации весов модели. В данной работе были рассмотрены следующие варианты задания w_0 :

- нулевой вектор
- вектор с координатами из $U(0, 1)$
- вектор с координатами из $U(100, 500)$
- вектор с координатами из $U(1000, 5000)$
- вектор с координатами из $U(10000, 50000)$
- вектор с координатами из $N(0, 1)$
- вектор с координатами из $N(0.5, 0.5)$

Графики зависимостей 2.1.1 представлены на рис. 9, 10, 11, 12.

Рис. 9: Зависимость значения функции потерь от реального времени работы стохастического градиентного спуска

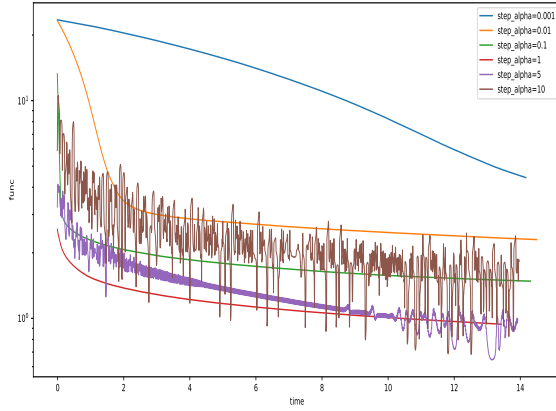


Рис. 10: Зависимость значения точности (ассигасу) от реального времени работы стохастического градиентного спуска

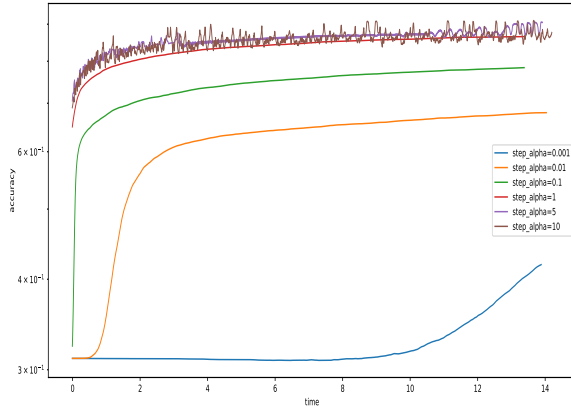


Рис. 11: Зависимость значения функции потерь от итерации метода стохастического градиентного спуска

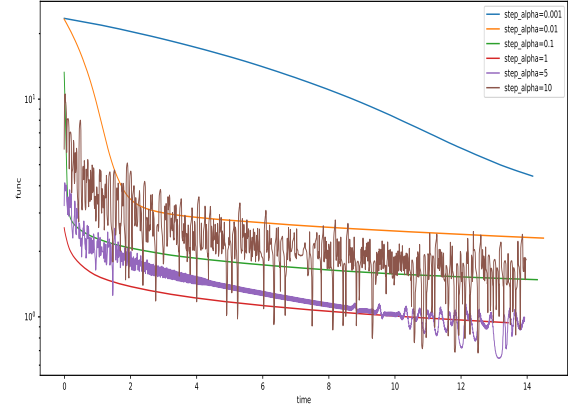
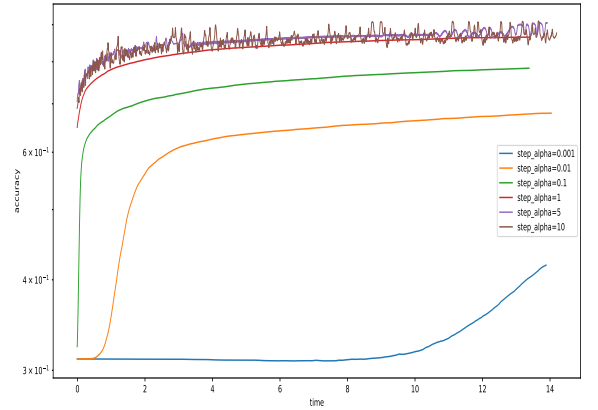


Рис. 12: Зависимость значения точности (ассигасу) от итерации метода стохастического градиентного спуска



2.2 Исследование поведения стохастического градиентного спуска

Обновления весов модели при использовании стохастического градиентного спуска происходит по следующей формуле:

$$w_t = w_{t-1} - \frac{\alpha}{t^\beta} \times \frac{1}{|I|} \times \sum_{i \in I} \nabla_w \mathcal{L}(x_i, y_i | w_{t-1}), \quad (2)$$

где t - номер итерации, β - **step_beta**, I - некоторое подмножество индексов тренировочной выборки, $\nabla_w \mathcal{L}(x_i, y_i | w_{t-1})$ - градиент функции потерь.

2.2.1 Параметр размера шара step_alpha

Параметр **step_alpha** (α) используется в стохастическом градиентном спуске при обновлении весов в формуле 2. Рассмотрим следующие зависимости при разных значениях параметра **step_alpha**:

1. зависимость значения функции потерь от реального времени работы метода
2. зависимость точности (ассигасу) от реального времени работы метода
3. зависимость значения функции потерь от эпохи метода
4. зависимость точности (ассигасу) от эпохи метода

Соответствующие графики приведены на: рис. 17, 18, 19, 20.

Рис. 13: Зависимость значения функции потерь от реального времени работы градиентного спуска

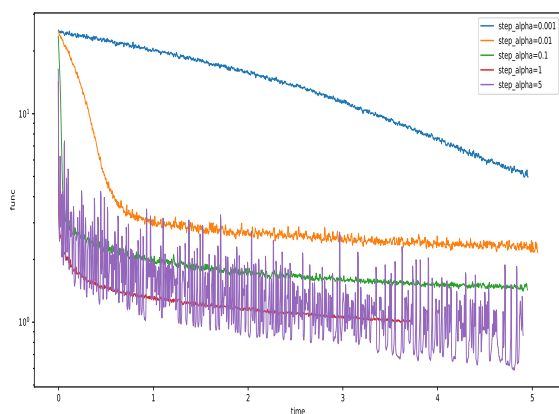


Рис. 14: Зависимость значения точности (ассигасу) от реального времени работы градиентного спуска

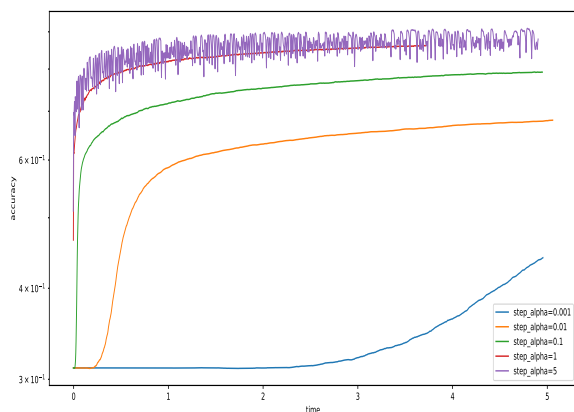


Рис. 15: Зависимость значения функции потерь от эпохи метода градиентного спуска

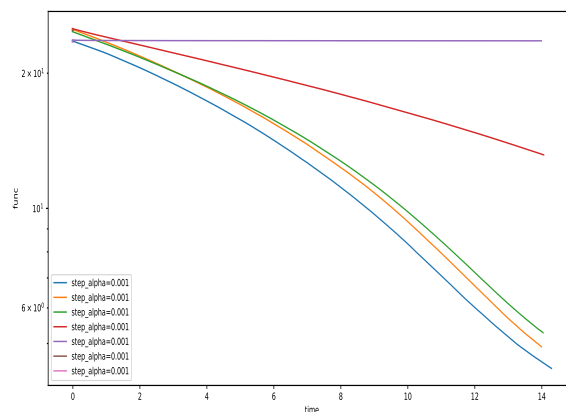
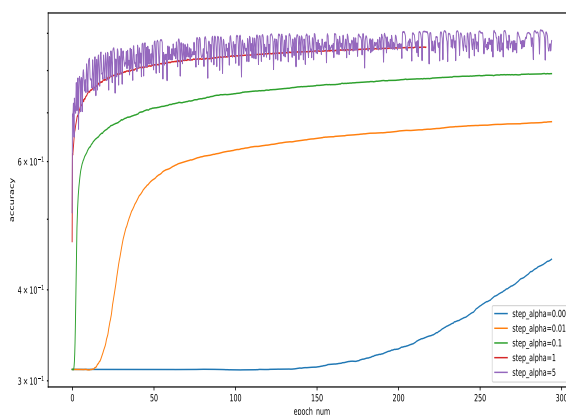


Рис. 16: Зависимость значения точности (ассигасу) от эпохи метода градиентного спуска



На графиках просматривается аналогичная ситуация, что и с градиентным спуском: при значениях **step_alpha**, близких к 0 возникает эффект недообучения, а при больших - появляется нестабильность кривой обучения, но есть точки, в которых достигается наивысшая точность. В таком случае можно производить сохранение весов модели на итерации со значением наилучшей точности. Возможно, такой метод поможет справиться с нестабильностью.

2.2.2 Параметр размера шага **step_beta**

Параметр **step_beta** (β) используется в градиентном спуске при обновлении весов в формуле 1. Аналогично предыдущему пункту рассмотрим зависимости из 2.2.1 при разных значениях параметра **step_beta** и проанализируем соответствующие графики, представленные на рис. ??, ??, ??, ??.

Рис. 17: Зависимость значения функции потерь от реального времени работы градиентного спуска

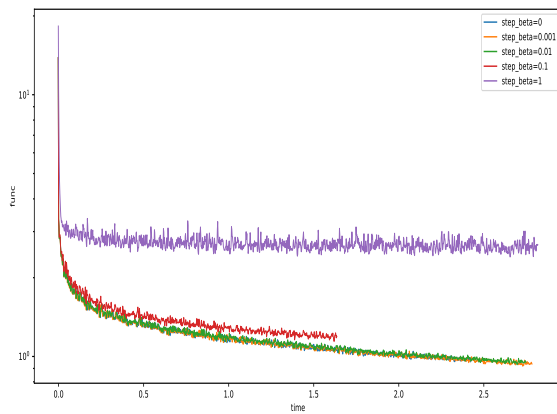


Рис. 18: Зависимость значения точности (ассигасу) от реального времени работы градиентного спуска

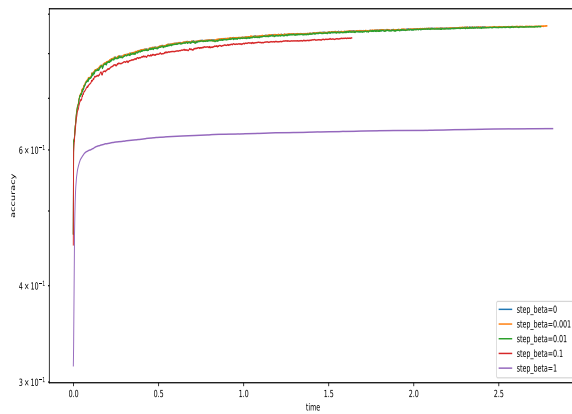


Рис. 19: Зависимость значения функции потерь от эпохи метода градиентного спуска

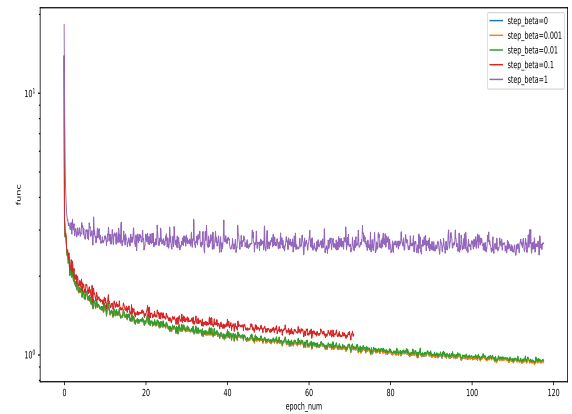
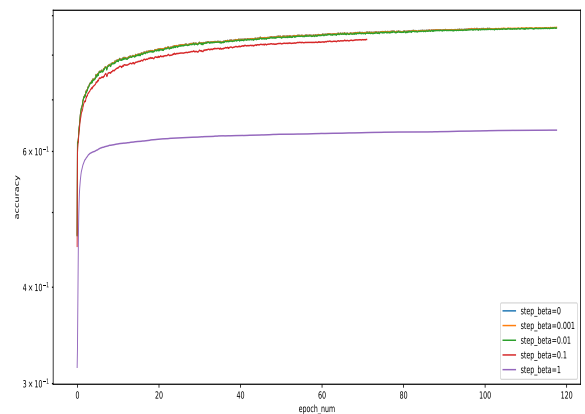


Рис. 20: Зависимость значения точности (ассигасу) от эпохи метода градиентного спуска



2.2.3 Размер подвыборки batch_size

Размер подвыборки определяет количество элементов тренировочной выборки, которые будут использованы для подсчета градиента.

Соответствующие графики приведены на рис. 21, 22, 23, 24.

Рис. 21: Зависимость значения функции потерь от реального времени работы градиентного спуска

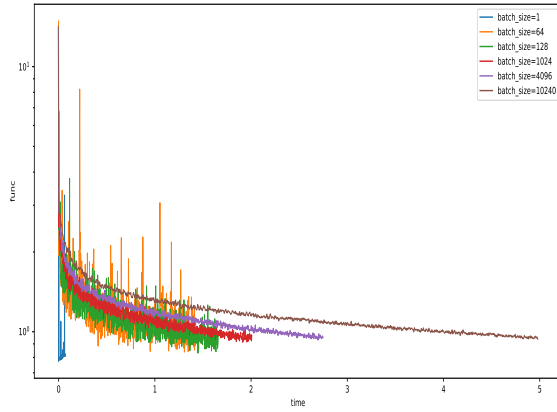


Рис. 22: Зависимость значения точности (ассигасу) от реального времени работы градиентного спуска

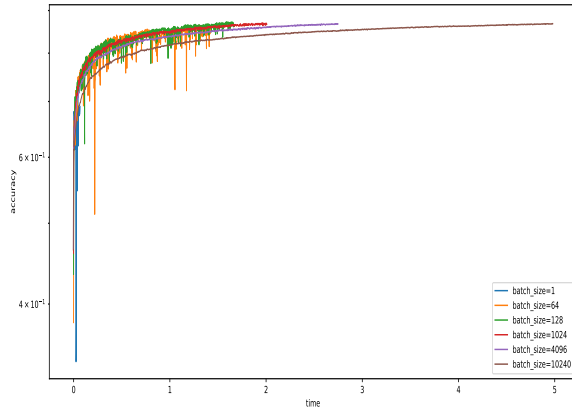


Рис. 23: Зависимость значения функции потерь от эпохи метода градиентного спуска

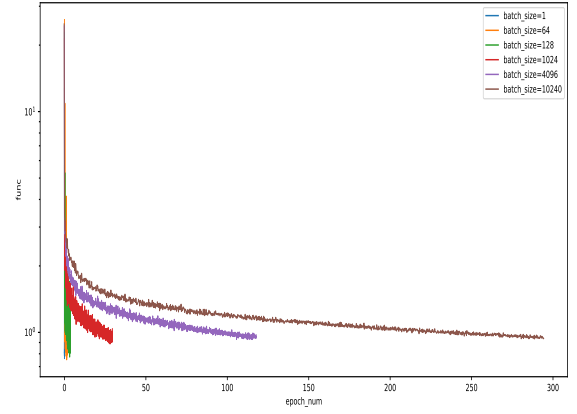
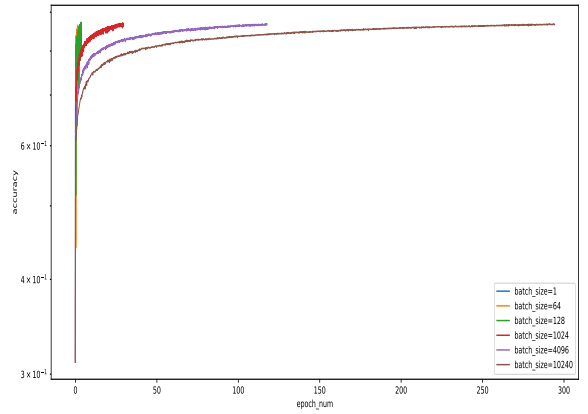


Рис. 24: Зависимость значения точности (ассигасу) от эпохи метода градиентного спуска



2.2.4 Начальное приближение w_0

Начальное приближение нужно для инициализации весов модели. В данной работе были рассмотрены следующие варианты задания w_0 :

- нулевой вектор
- вектор с координатами из $U(0, 1)$
- вектор с координатами из $U(100, 500)$
- вектор с координатами из $U(1000, 5000)$
- вектор с координатами из $U(10000, 50000)$
- вектор с координатами из $N(0, 1)$
- вектор с координатами из $N(0.5, 0.5)$

Графики зависимостей 2.2.1 представлены на рис. 25, 26, 27, 28.

Рис. 25: Зависимость значения функции потерь от реального времени работы градиентного спуска

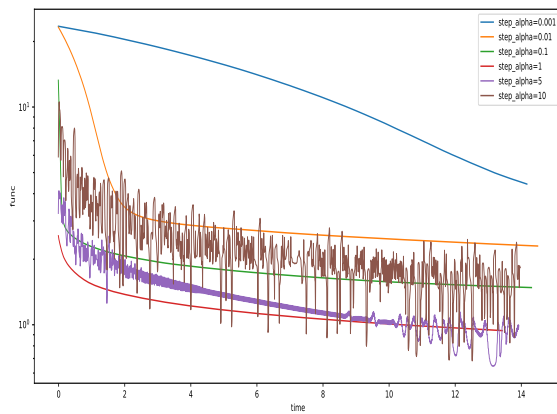


Рис. 27: Зависимость значения функции потерь от эпохи метода градиентного спуска

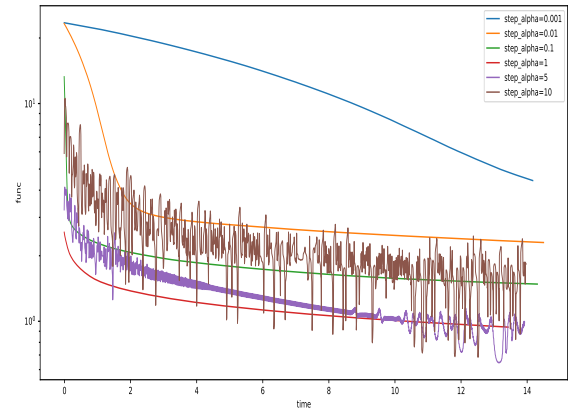


Рис. 26: Зависимость значения точности (ассигасу) от реального времени работы градиентного спуска

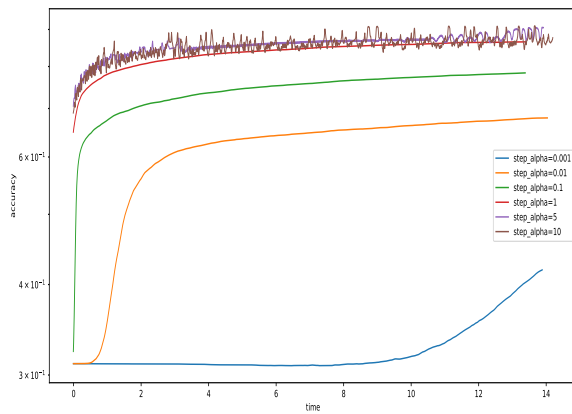
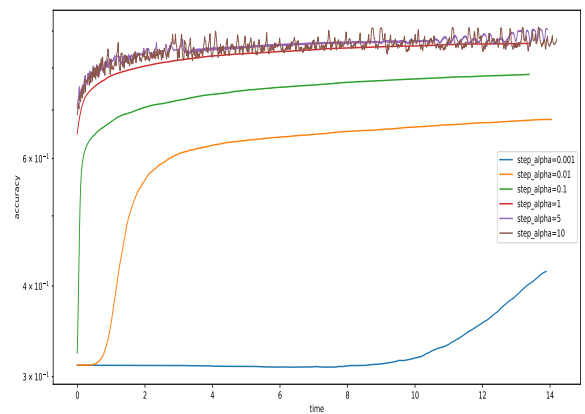


Рис. 28: Зависимость значения точности (ассигасу) от эпохи метода градиентного спуска



2.3 Сравнение градиентного спуска и стохастического градиентного спуска

В данном разделе проведено сравнение методов по трем характеристикам:

- время сходимости метода
- точность (ассигасу)
- значения функции потерь

Результаты экспериментов приведены на рис.

- 2.4 Лемматизация и удаление стоп-слов
- 2.5 Сравнение представлений BagOfWords и TF-IDF с различными параметрами
- 3 Применение лучших алгоритмов с каждого эксперимента к тестовой выборке