

1. Какие другие параметры сети Вы предлагаете использовать и почему не использовали?

Ответ: параметры, которые я менял – выходной размер LSTM-слоя и процент прореживания в слое Dropout. Во всех случаях, когда были изменены начальные значения, результаты оказывались немного хуже.

2. Что такое бутстреп?

Ответ: в статистике и анализе данных бутстрепом называют статистическую процедуру, основанную на выборке с замещением для определения точности (смещения) выборочных оценок дисперсии, среднего, стандартного отклонения, доверительных интервалов и других структурных характеристик совокупности.

В основе идеи бутстрепа лежит оценка структурных характеристик генеральной совокупности на основе перевыборки (resampling) из выборки. Иными словами, перевыборка по отношению к выборке рассматривается как выборка по отношению к генеральной совокупности.

Алгоритм работы метода следующий:

1. Из генеральной совокупности формируется случайная выборка из $N(t)$ наблюдений (например, если требуется определить среднюю сумму чека посетителя супермаркета, будем оценивать ее на основе выборки из 1 000 клиентов).
2. К выборке применяется случайная перевыборка с возвратом (псевдовыборка) того же объема, но в которую некоторые наблюдения могут попасть несколько раз, а другие не попасть совсем. Например, если выборка содержала 5 значений (1, 2, 3, 4, 5), то результатом перевыборки может быть (2, 2, 4, 5, 5). Затем вычисляется ее среднее.
3. Процедура перевыборки повторяется достаточно много раз (несколько десятков, сотен или даже тысяч), и для каждого случая вычисляется среднее.
4. Из полученного набора средних значений вычисляется среднее и рассматривается как среднее всей генеральной совокупности.

Важнейшим преимуществом бутстрепа являются:

- простота реализации;
- отсутствие необходимости гипотез о параметрах распределения данных;
- возможность оценивания многих статистических характеристик (среднего, дисперсии, стандартного отклонения, доверительных интервалов, квантилей, коэффициентов корреляции и др.).

К недостатку метода можно отнести использование малореалистичного предположения о независимости перевыборок и значительные вычислительные затраты при их многократном построении.

3. В чем заключается смысл ансамблирования комитетом большинства?

Ответ: В задачах классификации простейший пример ансамбля – комитет большинства: $a(x) = \text{mode}(b_1(x), \dots, b_n(x))$, где mode – мода.

Если рассмотреть большее число алгоритмов, то по неравенству Хёфдинга ошибка комитета большинства

$$\sum_{t=0}^{\lfloor n/2 \rfloor} C_n^t (1-p)^t p^{n-t} \leq e^{-\frac{1}{2}n(2p-1)^2}$$

, где p – вероятность ошибки каждого алгоритма, т.е. она экспоненциально убывает с ростом числа базовых алгоритмов.