

УДК 004.77

Хасан Али Аль-Абабнех

кандидат технических наук

Hassan Ali al-Ababneh

Ph.D.

ОЦЕНКА ПРОИЗВОДИТЕЛЬНОСТИ ВЕБ-ОРИЕНТИРОВАННЫХ КОМПЬЮТЕРНЫХ СИСТЕМ С ИСПОЛЬЗОВАНИЕМ СТАТИСТИЧЕСКИХ ГИПОТЕЗ

Аннотация: Описаны результаты оценки производительности при передаче веб-ориентированными системами различного рода контента (текстовая аудио- и видеоинформация, архивированные данные). Выполнен статистический анализ полученных результатов, проведена сравнительная характеристика полученных зависимостей с аналогичными, полученными ранее, на фоне обновления архитектурных характеристик серверных компьютерных систем.

Ключевые слова: статистические гипотезы, веб-ориентированные компьютерные системы, контент, нагрузка.

Summary: The results of the performance evaluation for the transmission of various types of content by web-based systems (text audio and video information, archived data) are described. The statistical analysis of the obtained results is performed, the comparative characteristics of the obtained dependences are compared with those obtained earlier, against the background of updating the architectural characteristics of the server computer systems.

Keywords: statistical hypotheses, web-oriented computer systems, content, load.

Постановка проблемы. На сегодняшний день, благодаря широкому распространению веб-приложений и росту популярности веб-ресурсов, практически все компьютерные системы также являются веб-ориентированными. Такие системы определенно серверные и при работе с веб-ресурсами обладают соответствующим «запасом прочности» производительности и отказоустойчивости.

Не смотря на экспоненциальный рост технических характеристик серверов и серверных систем, вопросы производительности продолжают оставаться актуальными. При этом статические параметры компьютерной системы, такие как скорость обработки и передачи информации, объемы памяти для ее хранения, не являются определяющими для подбора оптимальных параметров системы. Для определения производительности компьютерных систем необходимо учитывать целый ряд дополнительных параметров, таких как, например, адаптивность поведения, возможность тонкой настройки (программно-аппаратной) для, например, минимизации затрат памяти, конфигурирование серверов различных типов, и пр. Однако и этого не всегда достаточно, поскольку в данном вопросе важно оценить нагрузочную способность системы в зависимости от плотности контента [1].

Анализ последних исследований и публикаций. В работах [2, 3] рассмотрены различные подходы к решению проблемы обеспечения заданной производительности веб-ориентированных компьютерных систем. При этом было отмечено, что для успешных решений задач производительности, необходимо выделить достаточные условия для реализации численных моделей инфраструктуры с использованием адекватных моделей возможных рабочих нагрузок.

Формулировка целей статьи (постановка задачи) заключается в определении целесообразности применения статистических гипотез при оценке производительности компьютерных систем.

Изложение основного материала. Оценки производительности веб-ориентированных компьютерных систем принято представлять в виде данных (утверждений), характеризующих свойства распределения наблюдаемых в эксперименте случайных величин. Это и есть принцип статистических гипотез, которые делят на следующие виды: однородности, если имеется две или более выборки случайных величин; независимости, если имеется выборка многомерной случайной величины; случайности, если есть предположения о наличии в последовательности наблюдений систематических изменений; о виде распределения, если есть предположения о законе распределения случайной величины (рис. 1).

Проверка статистической гипотезы состоит в том, чтобы сформулировать такое правило, которое позволило бы по результатам проведенных наблюдений принять или отклонить гипотезу. Правило, согласно которому гипотеза принимается или отвергается, называется критерием проверки статистической гипотезы [4, с. 257–264].

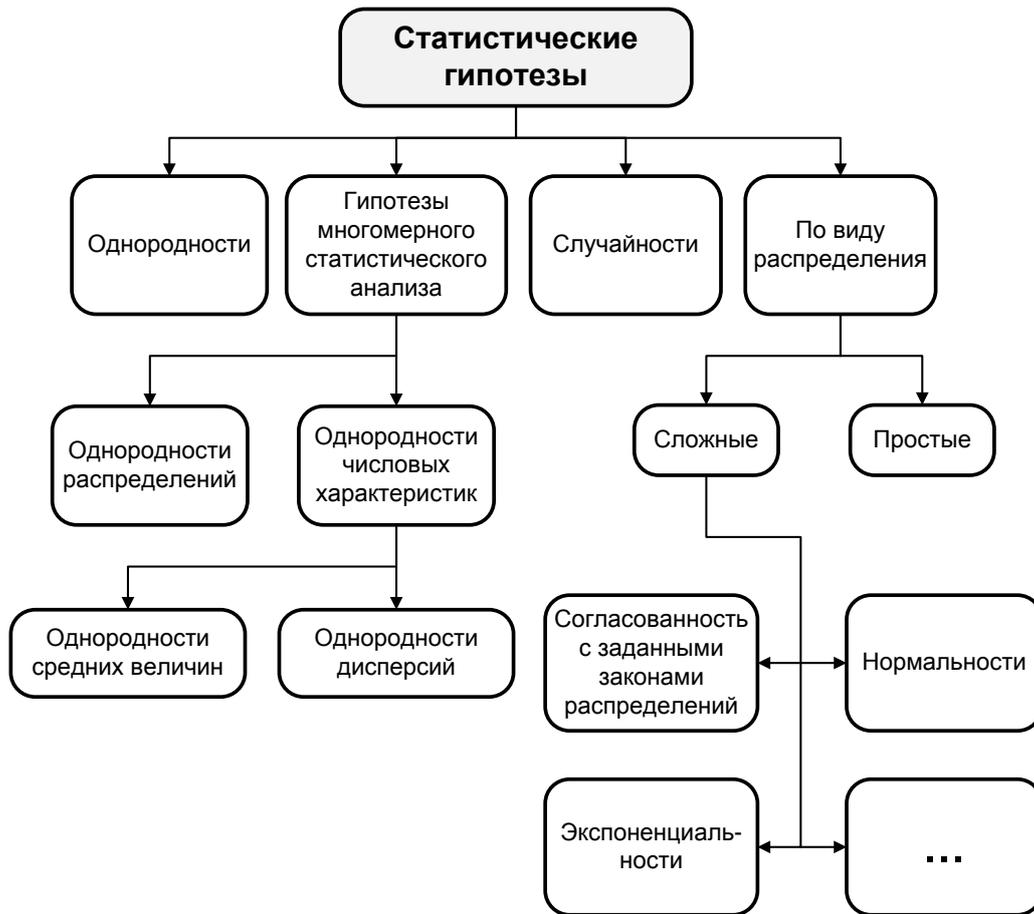


Рис.1 Классификация статистических гипотез

При планировании архитектуры серверной системы важную роль играет вид используемого контента.

Для исследования зависимости объема файла и объема страницы сайта от вида контента были выбраны наиболее типичные виды контента:

- текстовые;
- тексто-графические (в том числе на основе формата pdf);
- графические (преимущественно на основе jpg);
- аудио (mp3 и wma);
- видео (преимущественно на основе flv).

В качестве инструмента для проведения экспериментов был выбран Интернет-браузер Opera, позволяющий анализировать состояние кэш-памяти браузера в процессе исследования с целью сбора необходимой статистической информации.

Сбор статистических параметров проводился по двум основным показателям:

1. Объем информационного файла (носитель информации).
2. Общий объем страницы информационного сайта (в совокупности со всеми сопутствующими файлами).

Объектом для исследований текстового формата была выбрана электронная библиотека Lib.ru (Библиотека Машкова). В ходе исследований оценивался объем отдельно текстового файла и совокупный объем страницы с ресурсом. Результаты оценки производительности при передаче веб-ориентированными системами текстовой информации представлены на рис. 2, 3.

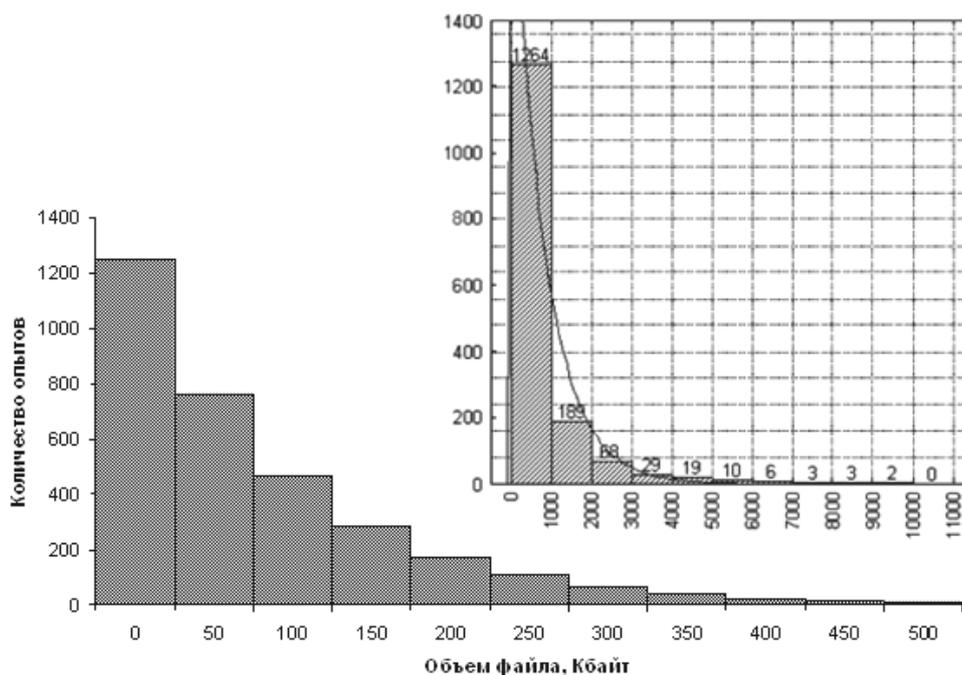


Рис. 2 Гистограмма размера текстового файла
(на примере сайта Lib.ru)

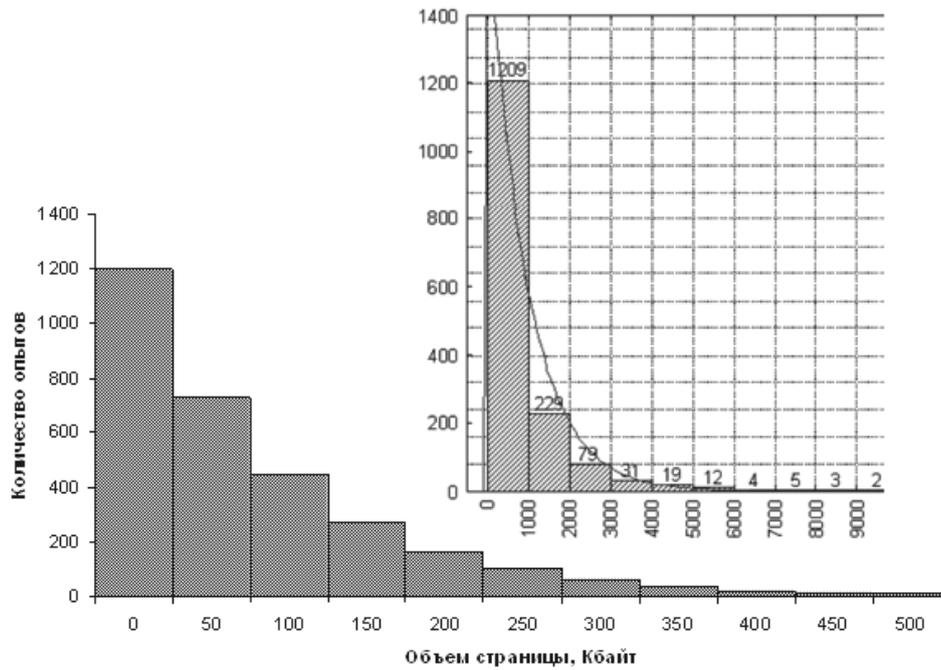


Рис. 3 Гистограмма размера страницы с текстовым контентом
(на примере сайта Lib.ru)

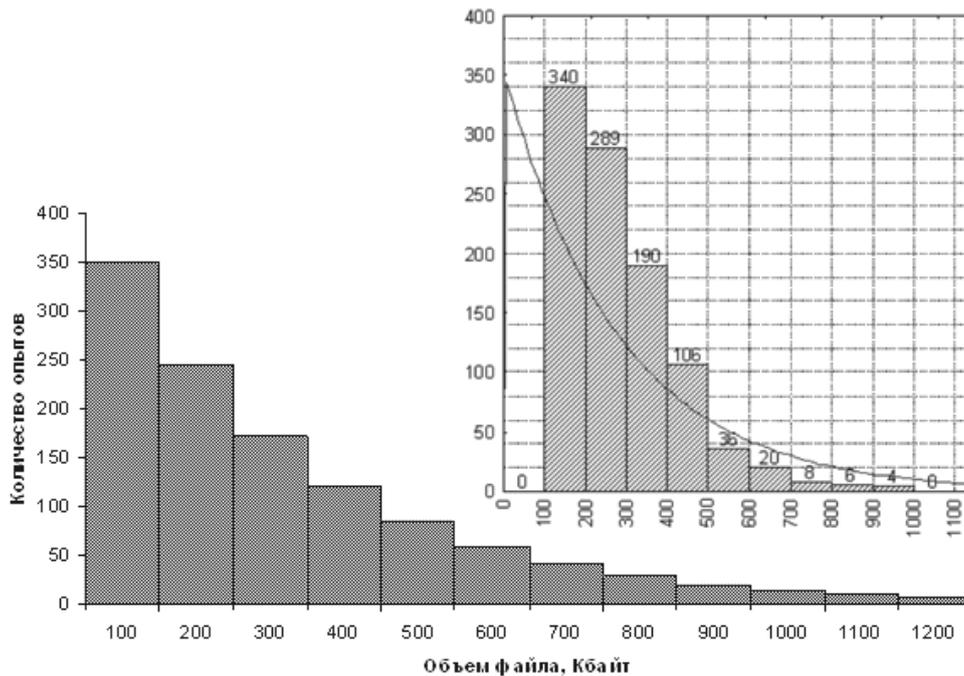


Рис. 4 Гистограмма размера графического файла
(на примере сайта bestwallpapers.net.ru)

Объектом для исследований графического формата был выбран типичный сайт содержащий широкоформатные изображения высокого разрешения, предназначенные для использования в качестве заставок

рабочего стола ОС Windows XP/Vista – bestwallpapers.net.ru. В ходе исследования оценивался объем графического файла и совокупный объем страницы с ресурсом. Результаты представлены на рис. 4, 5.

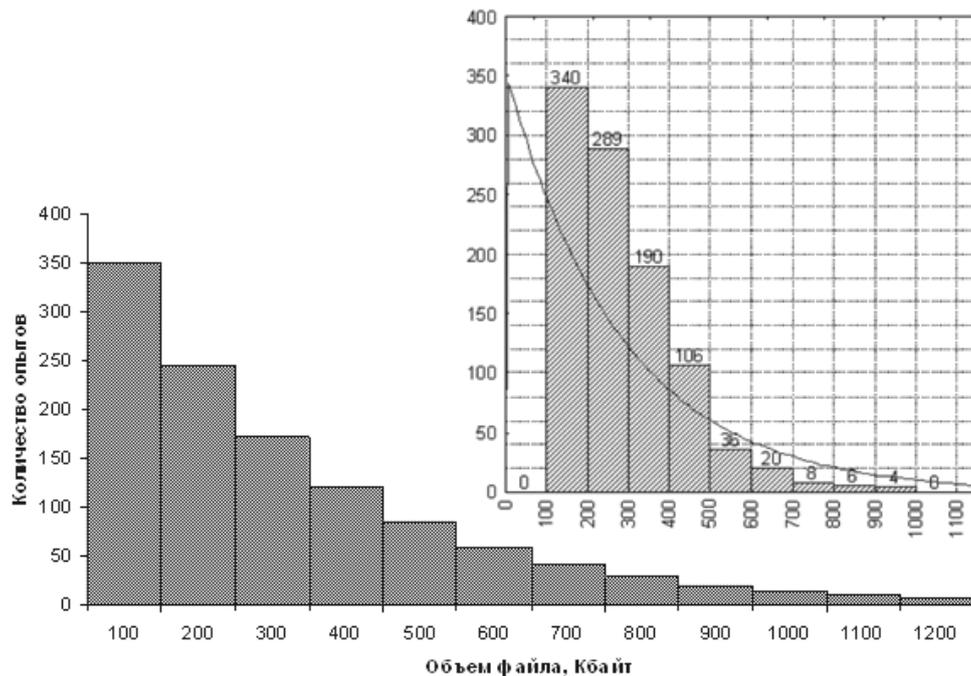


Рис. 5 Гистограмма размера графического файла
(на примере сайта bestwallpapers.net.ru)

Объектом для исследований аудио-ресурсов был выбран сайт 6rb2d.com как яркий представитель соответствующих ресурсов сети Интернет. В ходе исследования объем ресурса и аудиофайла оценивался отдельно. Результаты представлены на рис. 6, 7.

Объектом для исследований видео формата был выбран сайт YouTube.com. В ходе исследования оценивался объем видеофайла и совокупный объем страницы с ресурсом.

Для текстовых ресурсов дополнительно был исследован вариант, когда текстовые файлы представлены в архивированном виде. В этом случае было получено экспоненциальное распределение.

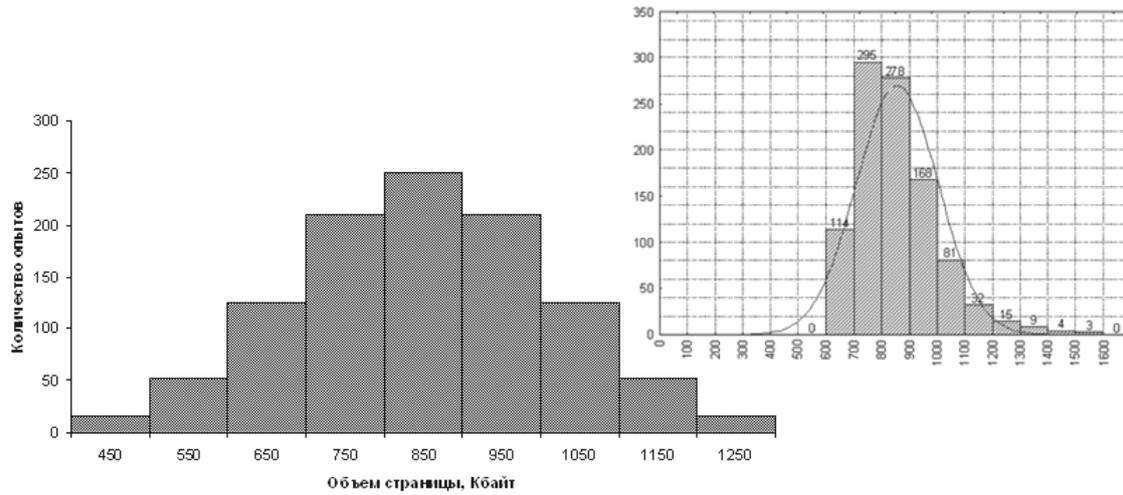


Рис. 6 Гистограмма размера страницы с графическим контентом (на примере сайта bestwallpapers.net.ru)

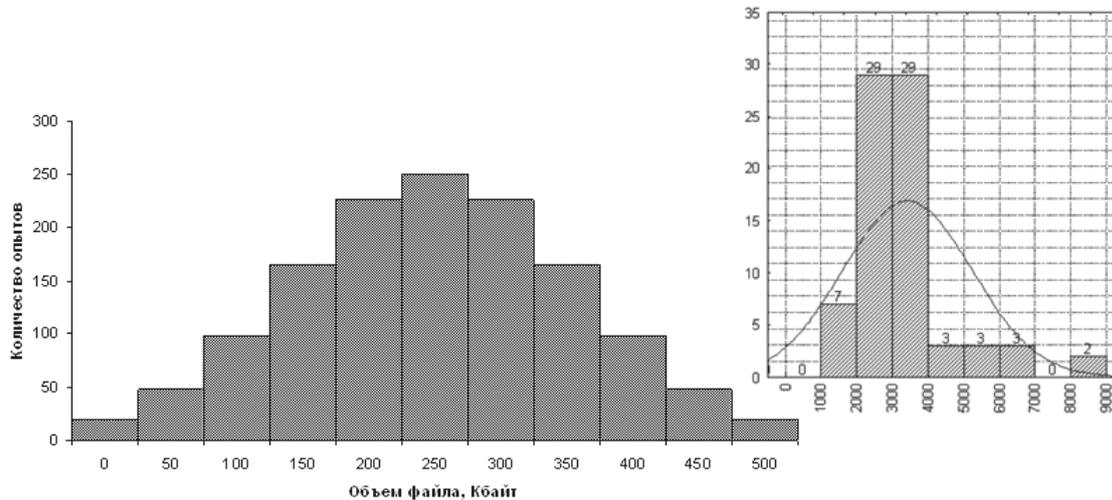


Рис. 7 Гистограмма размера аудиофайла (на примере сайта 6rb2d.com)

На рис. 2–10 приведены также результаты аналогичных исследований (графики в меньшем масштабе), проведенных ранее (в 2008 г.) и опубликованных в работе [1]. Повторное исследование с учетом современных компьютерных систем, с существенно отличающимися параметрами 2008 года, тем не менее, позволяет утверждать в справедливости полученной статистической гипотезы.

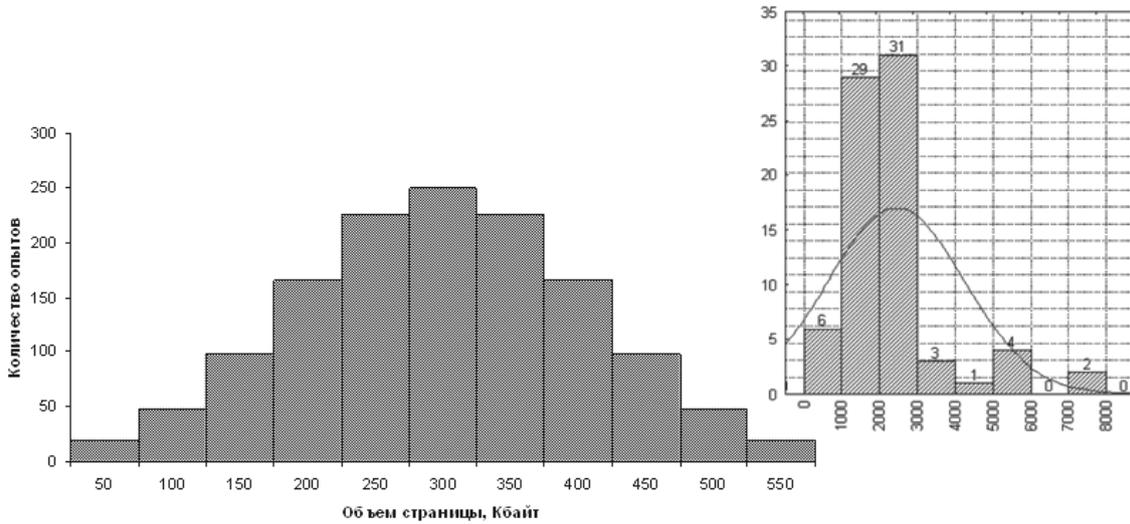


Рис. 8 Гистограмма размера страницы с аудиоконтентом
(на примере сайта 6rb2d.com)

Статистический анализ данных велся с помощью программы STATISTICA 7.0. В работе подбиралось распределение, которое с достаточной степенью точности описывает наблюдаемые данные. Таким образом, проверялась гипотеза, согласно которой распределение X описывается вероятностным законом F .

Наиболее показательным критерием согласия наблюдаемых данных с гипотезой является критерий хи-квадрат (Пирсона). Для применения этого критерия область значений переменной X вначале была разбита на некоторое число интервалов N , содержащих не менее семи чисел.

Затем подсчитывалось число наблюдений, попавших в i -й интервал, что обозначалось как $n(i)$. Полученное значение сравнивалось со средней или ожидаемой при гипотезе частотой, обозначенной как $\bar{n}(i)$.

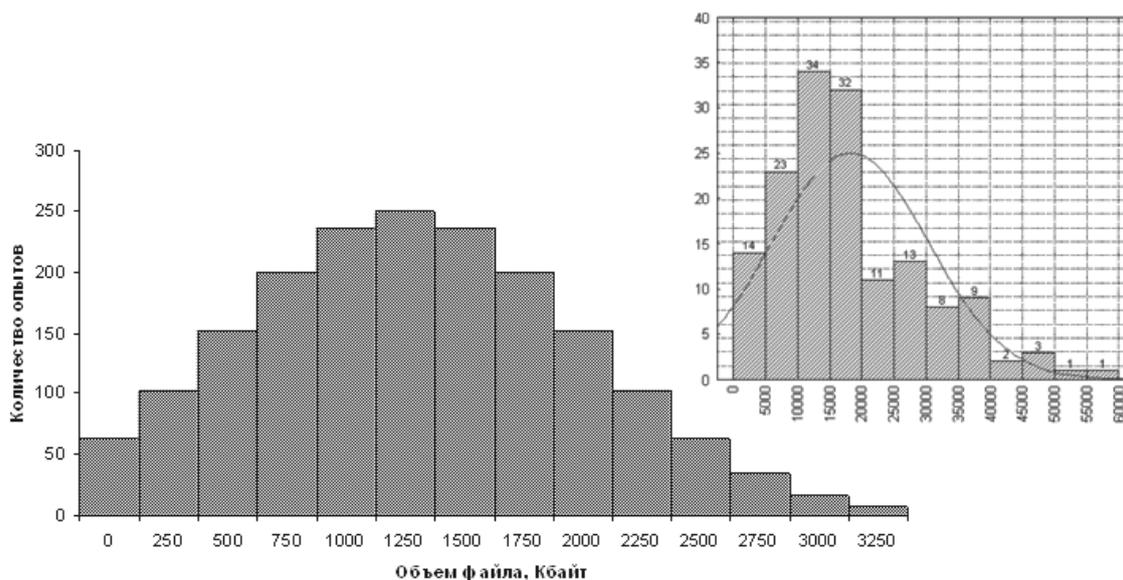


Рис. 9 Гистограмма размера видеофайла (на примере сайта YouTube.com)

Статистика хи-квадрат вычислялась следующим образом:

$$\chi^2 = \sum_{i=1}^N (n(i) - \bar{n}(i))^2 / \bar{n}(i).$$

В этой формуле суммирование распространяется на все интервалы, на которые разбита область значений переменной. При этом сравнивались наблюдаемые и ожидаемые частоты. Статистика принимает значения от нуля до бесконечности. Чем меньше значение статистики хи-квадрат, тем более вероятно, что гипотеза верна, чем больше значение статистики хи-квадрат, тем меньше вероятность того, что гипотеза соответствует данным.

Исходя из вышеизложенного статистика хи-квадрат – это разумная мера согласия (соответствия) данных с гипотезой. Для девяти распределений, изображенных на рис. 2.9-2.17, степени свободы r составили 9, 11, 10, 6, 8, 8, 11, 10, 7 соответственно, уровень значимости $\alpha = 0.09$. В соответствии с r и α из выбирались критические значения распределения хи-квадрат $\chi_{r,1-\alpha}^2$ и сравнивались с рассчитанными значениями χ^2 . Для всех девяти распределений выполняется условие $\chi^2 < \chi_{r,1-\alpha}^2$ и, следовательно, данные согласуются с гипотезой о законе распределения.

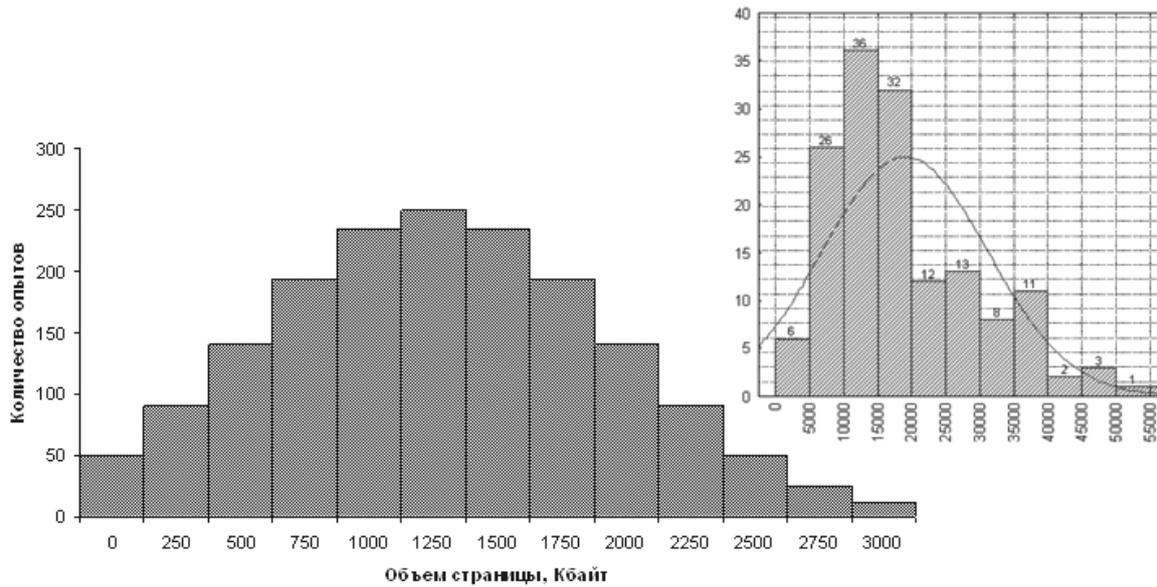


Рис. 10 Гистограмма размера страницы с видеоконтентом
(на примере сайта YouTube.com)

Выводы из данного исследования и перспективы дальнейших исследований в данном направлении. В ходе анализа, на основе вышеизложенных критериев, сделан вывод о том, что распределение большинства исследованных величин подчиняется нормальному закону. Полученные результаты вполне согласуются с тем, что нормальное распределение является краеугольным камнем математической статистики в силу следующих причин:

- схема его возникновения соответствует многим реальным физическим процессам, порождающим результаты обрабатываемых наблюдений;
- при возрастании объема выборки предельное распределение для большинства распределений является нормальным и с успехом может использоваться для аппроксимации последних;
- нормальное распределение обладает рядом благоприятных математико-статистических свойств (легко нормируется и аппроксимируется, обладает свойством аддитивности).

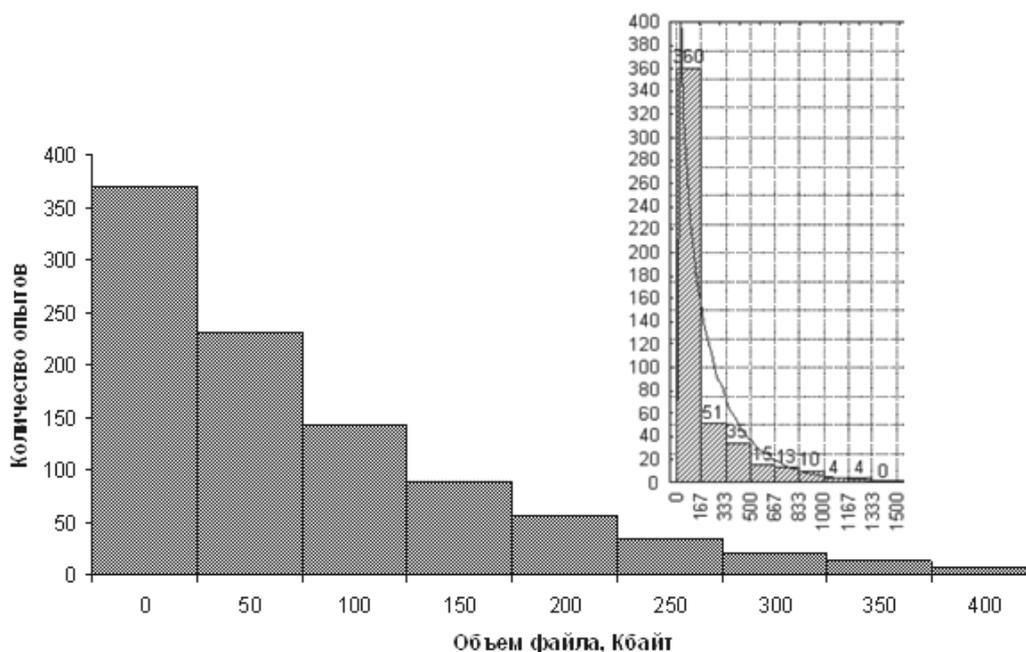


Рис. 11 Гистограмма размера текстового файла, представленного в архивированном виде

Таким образом, можно сделать следующие выводы о массивах полученных экспериментальных данных:

- имеется сильная тенденция группироваться вокруг центра (в нашем случае таким центром является среднее значение каждой выборки);
- положительные и отрицательные отклонения от центра равновероятны;
- частота отклонений быстро падает, когда отклонения от центра становятся большими.

В ряде случаев (текстовые ресурсы) было получено экспоненциальное распределение, часто встречающееся в теории надежности и в теории массового обслуживания. Например, наработка на отказ большой многокомпонентной системы может быть описана экспоненциальным распределением при любом распределении наработки на отказ компонентов системы.

Литература:

1. Аноприенко А.Я. Модели нагрузки в веб-ориентированных компьютерных сетях / А.Я. Аноприенко, Х. Аль-Абабнех // Научные труды Донецкого национального технического университета. Серия «Проблемы моделирования и автоматизации проектирования динамических систем». – 2008. – выпуск 7 (150). – С.258–274.
2. Аноприенко А. Я., Аль Абабнех Хасан. Повышение эффективности Интернет-ориентированной сетевой инфраструктуры: Методы, задачи и инструменты / Научные труды Донецкого национального технического университета. Серия "Проблемы моделирования и автоматизации проектирования динамических систем" (МАП-2007). Выпуск 6 (127): Донецк: ДонНТУ, 2007. С. 228–233.
3. Аноприенко А. Я., Потапенко В. А. WEB-ориентированная среда для интеграции моделирующих, вычислительных и информационных сервисов / Научные труды Донецкого национального технического университета. Выпуск 70. Серия: «Информатика, кибернетика и вычислительная техника» (ИКВТ-2002): - Донецк: ДонНТУ, 2003. С. 61-70.
4. Менаске Д., Алмейда В. Производительность Web-служб. Анализ, оценка и планирование: Пер. с англ. – СПб: ООО «ДиаСофтЮП», 2003. – 480 с.
5. Афанасьев В. В. Теория вероятностей: учеб. пособие для студентов вузов, обучающихся по специальности «Математика» / В. В. Афанасьев. — М.: Гуманитар. изд. центр ВЛАДОС, 2007. — 350 с. — (Учебник для вузов).

References:

1. Anoprienko A.Ya. Modeli nagruzki v veb-orientirovannykh kompyuternykh setyakh / A.Ya. Anoprienko, Kh. Al-Ababnekh //

- Nauchnye trudy Donetskogo natsionalnogo tekhnicheskogo universiteta. Seriya «Problemy modelirovaniya i avtomatizatsii proektirovaniya dinamicheskikh sistem». – 2008. – vypusk 7 (150). – S.258–274.
2. Anoprienko A. Ya., Al Ababnekh Khasan. Povyshenie effektivnosti Internet-orientirovannoy setevoy infrastruktury: Metody, zadachi i instrumenty // Nauchnye trudy Donetskogo natsionalnogo tekhnicheskogo universiteta. Seriya "Problemy modelirovaniya i avtomatizatsii proektirovaniya dinamicheskikh sistem" (MAP-2007). Vypusk 6 (127): Donetsk: DonNTU, 2007. S. 228–233.
 3. Anoprienko A. Ya., Potapenko V. A. WEB-orientirovannaya sreda dlya integratsii modeliruyushchikh, vychislitelnykh i informatsionnykh servisov // Nauchnye trudy Donetskogo natsionalnogo tekhnicheskogo universiteta. Vypusk 70. Seriya: «Informatika, kibernetika i vychislitel'naya tekhnika» (IKVT-2002): - Donetsk: DonNTU, 2003. C. 61-70.
 4. Menaske D., Almeyda V. Proizvoditelnost Web-sluzhb. Analiz, otsenka i planirovanie: Per. s angl. – SPb: OOO «DiaSoftYuP», 2003. – 480 s.
 5. Afanasev V. V. Teoriya veroyatnostey: ucheb. posobie dlya studentov vuzov, obuchayushchikhsya po spetsialnosti «Matematika» / V. V. Afanasev. — M.: Gumanitar. izd. tsentr VLADOS, 2007. — 350 s. — (Uchebnik dlya vuzov).