

ДАННЫЕ СТАТИСТИКА АНАЛИТИКА

Виды аналитики

- Описательная  
"Что случилось?"  
Графики, сводные таблицы
- Диагностическая  
"Почему это случилось?"  
Анализ данных
- Прогнозная  
"Что может случиться?"  
Статистика и машинное обучение
- Предписательная  
"Что делать?"  
Статистика и машинное обучение.  
В отличие от прогнозной носит декларативный характер.

Этапы стат. анализа

- Описательная статистика - описание полученного в ходе исследования массива данных
- Аналитическая статистика - анализ данных и проверка различных статистических гипотез

Задачи и статметоды для их решения

- Гипотеза - это догадка, предположение, недоказанное утверждение.
  - опровергается
  - ... подтверждается
- Открытая проблема - это неопровернутая и недоказанная гипотеза
- Большой уровень значимости => Менше уверенности в H0. Больше уверенности в H1 => Большой риск не отвергнуть ложную гипотезу. Риск ошибки второго рода.
- Меньший уровень значимости => Больше уверенности в H0. Менше уверенности в H1 => Большой риск принять ложную гипотезу. Риск ошибки первого рода.
- Уровень значимости - это величина для оценки гипотезы. Вероятность отклонить гипотезу, если она на самом деле истинна. Популярные значения: 5%, 1% и 0.1%
- Основная/Нулевая гипотеза (H0) - это гипотеза, которой мы придерживаемся, пока наблюдения не заставят признать обратное.
- Альтернативная гипотеза (H1)
- Ошибки
  - Первого рода - ситуация когда H0 опровергается, хотя она на самом деле верна
  - Второго рода - ситуация когда H0 принимается, хотя он на самом деле неверна

Виды распределения

- Нормальные (гауссове) распределения
  - Стандартное нормальное распределение (z-распределение) - это нормальное распределение со средним = 0, и стандартным отклонением = 1.
  - Процедура нормализация (или стандартизации). Вычест из данных среднее значение и разделить на стандартное отклонение.
  - z-Оценка (z-score) Результат стандартизации отдельной точки данных.
  - 68% данных находится в пределах одного стандартного отклонения от среднего и 95% — в двух стандартных отклонениях.
- Длиннохвостые распределения
  - Это распределения с длинными хвостами - длинная ушка часть частотного распределения, где относительно предельные значения встречаются с низкой частотой.
  - Такие распределения могут быть асимметричными (skew) - когда один хвост длиннее другого.
  - Ключевая идея. Принятие нормального распределения может привести к недооценке предельных событий ("черных лебедей").
  - Нассим Талеб. Теорию черного лебеда - аномальные события, такие как обвал фондового рынка, будут возникать с намного большей вероятностью, чем их предсказание нормальным распределением.
- t-Распределение Стьюдента - это распределение нормальной формы, но немного толще и длиннее в хвостах
- Биномиальное распределение
- Распределение Пуассона

Анализ данных. Инструменты оценки

- Качество
  - Отсутствие дублей и противоречий
  - Отсутствие пропущенных значение
  - Адекватность
- Цели:
  - приведение различных данных в самых разных единицах измерения и диапазонах значений к единому виду, который позволит сравнивать их между собой или использовать для расчета схожести объектов
  - ускорение получения желаемого результата за счет того, что машине приходится обрабатывать меньший диапазон данных
- Аналитически любая нормализация сводится к формуле  $X_{norm} = (X_i - X_{min}) / (X_{max} - X_{min})$ , где  $X_i$  — текущее значение,  $X_{max}$  — величина смещения значений,  $X_{min}$  — величина интервала, который будет преобразован к "единице". По сути всё сводится к тому, что исходный набор значений сперва смещается, а потом масштабируется.
- Нормализация. Это преобразование данных к неким безразмерным единицам.
- Методы нормализации числовых данных
  - sklearn.preprocessing.normalize() Получим данные от 0 до 1
  - sklearn.preprocessing.MinMaxScaler() устанавливает наименьшее наблюдаемое значение равным 0, а наибольшее наблюдаемое значение — 1
  - Использование максимального абсолютного масштабирования. Все значения в столбце делим на макс. число по модулю.  $df[column] = df[column] / df[column].abs().max()$  Таким образом все данные нормализуются между 0 и 1.
  - z-стандартизация. Вычест из данных среднее значение и поделить на стандартное отклонение.  $(X - X_{mean}) / X_{std}$  Получим новый со средним значением равным 0 и стандартным отклонением равным 1. Изменяет форму распределения данных (приводится к нормальному распределению).

Подготовка данных

- Обработка пропусков
  - Удаление строк/столбцов с большим количеством пропусков
  - Восстановление данных
    - Подстановка среднего, медианы или моды
    - Взять значение из наиболее близкой строки (использу меры близости объектов)
    - С помощью коэф. корреляции
- Обработка дубликатов
- Обнаружение выбросов (аномальных значений) - Удаление
- Кодирование данных
  - Многие алгоритмы ML ожидают числовые входные данные, поэтому категориальные данные нужно представить в виде чисел.

Виды данных

- Структурированные (Данные в табличной форме)
  - Количественные (quantitative)
    - Непрерывные значения (continuous) - Рост, вес, продажи, время
    - Дискретные данные (discrete). Могут принимать только целочисленные значения. - например количество возникновений события
  - Категориальные/номинальные/качественные (categorical/nominal/qualitative)
    - Местонахождение: Америка, Европа, Азия.
    - Классификация сотрудников: руководитель, менеджер и пр.
    - Только определенный набор значений, в частности набор возможных категорий. Не имеет количественного отношения друг к другу.
    - Двоичные/Бинарные/дихотомические. Особый случай категориальных данных всего с двумя категориями значений (0/1, истина/ложь).
    - Порядковые. Для ранжирования данных по отношению друг к другу
    - Размеры футболок: S, M, L, XL
    - Оценки: плохо, хорошо, отлично
- Неструктурированные - Данные в произвольной форме (картинка, статья и пр.)

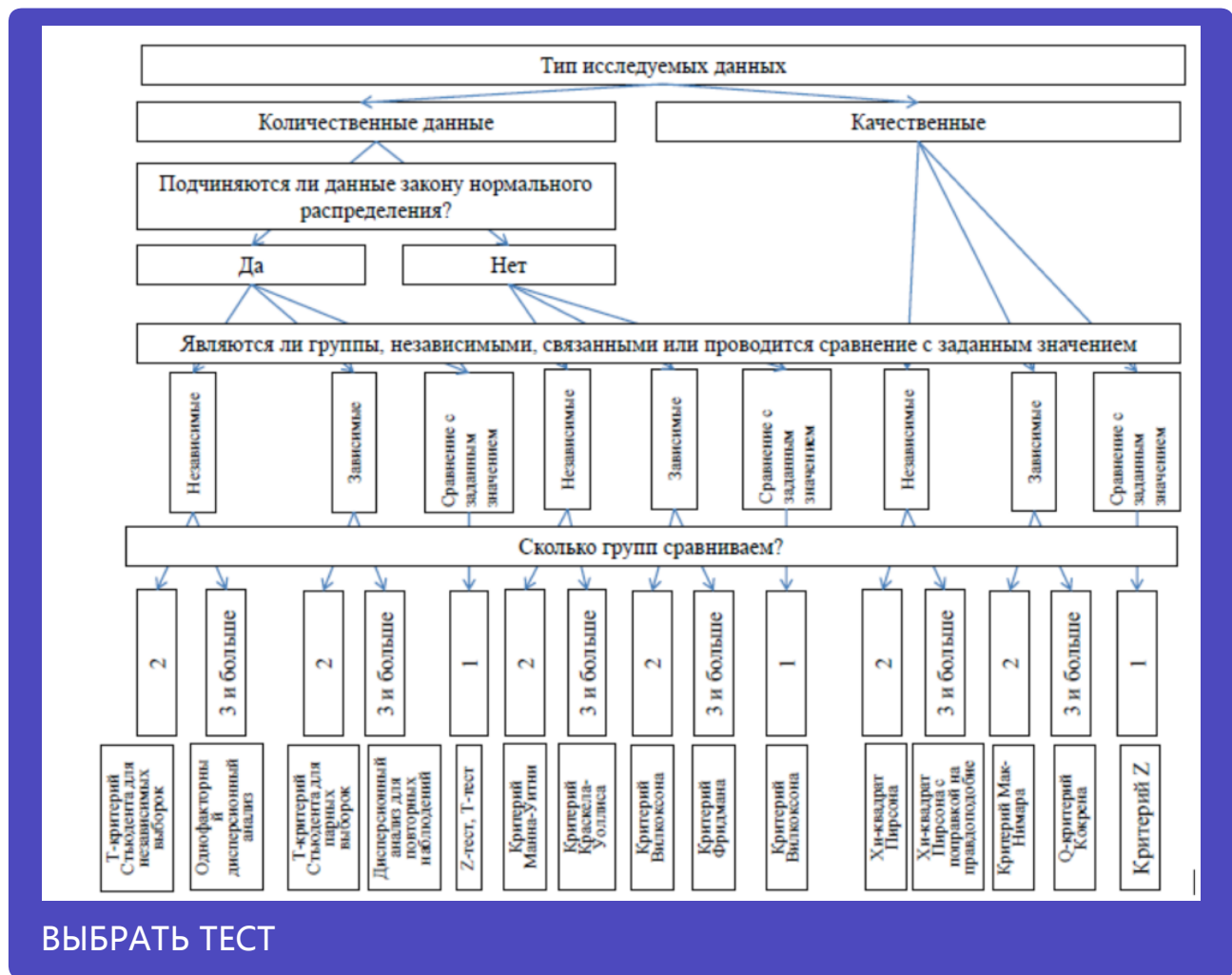
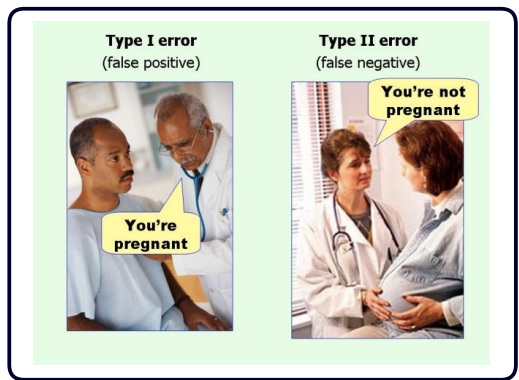
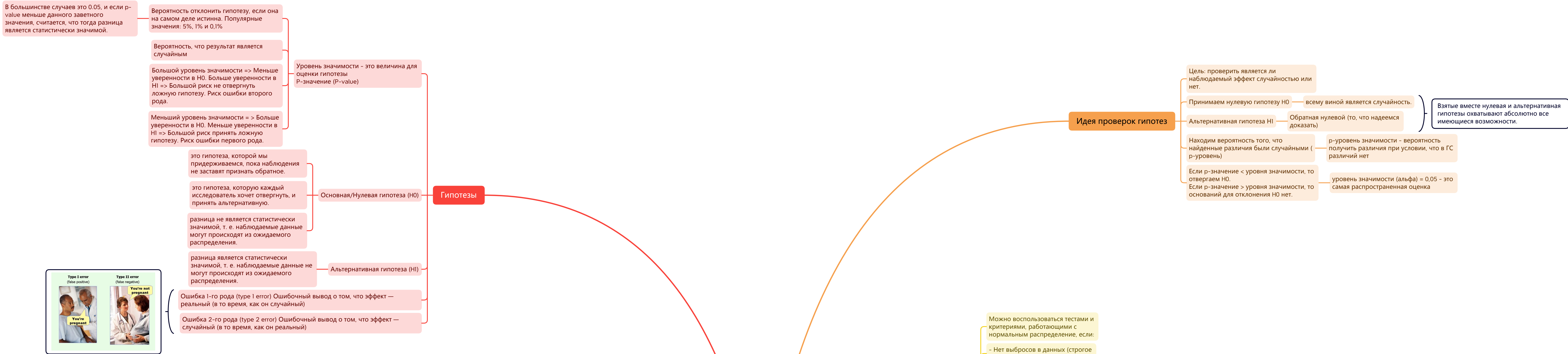
Структура данных

- Data frame (по сути электронная таблица)
  - Столбец (признак, feature, атрибут, вход, предиктор, переменная)
  - Столбец Исход (outcome, зависимая переменная, отклик, цель, выход)
  - Строки (записи, records, экземпляр, наблюдение, шаблон, паттерн)
- Временной (динамический) ряд (последовательные данные измерений одной и той же переменной)
- Графовые (или сетевые). Пример, граф соцсети может представлять связи между людьми в сети.

Выборка и генеральная совокупность

- Выборка - это подмножество данных из более крупного набора данных - популяции (или генеральной совокупности)
- Смещение возникает, когда данные измерений, или наблюдений систематически ошибочны
- Смещенная выборка (sample bias) Выборка, которая представляет популяцию в искаженном виде.
- Зависимые и независимые выборки
  - Если можно установить томоморфную пару (то есть, когда одному случаю из выборки X соответствует один и только один случай из выборки Y и наоборот) для каждого случая в двух выборках (и это основание взаимосвязи является важным для измеряемого на выборках признака), такие выборки называются зависимыми. Примеры зависимых выборок: пары близнецов, два измерения какого-либо признака до и после экспериментального воздействия, мужья и жены и т. п. Зависимые выборки всегда имеют одинаковый объем
  - В случае, если такая взаимосвязь между выборками отсутствует, то эти выборки считаются независимыми, например: мужчины и женщины, психологи и математики. Объемы независимых выборок могут отличаться.
- Отборы
  - Случайный отбор (random sampling) - процесс, в котором каждый доступный член популяции, подвергаемой отбору, имеет равную возможность попасть в выборку при каждой выемке
  - Стратифицированный отбор (stratified sampling) Разделение популяции на страты и случайный отбор элементов из каждой страты
  - С возвратом - после каждой выемки наблюдения кладутся назад в популяцию для возможно повторного отбора в будущем.
  - Без возврата - однажды выбранные наблюдения недоступны для будущих выемок.
- Возвраты при отборе
- Распределение данных (data distribution) Частотное распределение индивидуальных значений в наборе данных
- Выборочная статистика (sample statistic) - показатель статистики, вычисленный для для выборки из популяции.
- Выборочное распределение (sampling distribution) - распределение выборочной статистики на большом числе выборок, вынимаемой из одной и той же популяции.
- Центральная предельная теорема (ЦПТ, (central limit theorem) - тенденция выборочного распределения принимать нормальную форму по мере увеличения размера выборки, даже если исходная популяция не является нормально распределенной.
- По другому: средние значения, вынутые из многочисленных выборок, будут принимать нормальную форму, при условии, что размер выборки достаточно крупный.
- Стандартная ошибка (standard error: se) Стандартное отклонение выборочной статистики на многочисленных выборках.
- se показывает вариабельность выборочного метрического показателя. Не путать со стандартным отклонением, которое показывает вариабельность отдельных точек данных.
- se = s / sqrt(n), где s - стандартное отклонение значений выборки n - размер выборки
- По мере увеличения размера выборки se уменьшается. Для сокращения стандартной ошибки в 2 раза, размер выборки должен быть увеличен в 4 раза.
- Еще один способ понять потенциальную ошибку в оценке выборки
- ДИ показывает насколько параметры из выборки могут отличаться от реально существующих данных в ГС. Насколько мы ошибаемся при формировании той или иной выборки, мы закладываем в так называемую ошибку репрезентативности, в ошибку средней и вокруг нее собственно и строим доверительный интервал.
- Доверительный интервал
- ДИ задается самостоятельно исследователем. Чаще всего он равен 95 %.
- Один из простых и эффективных способов оценки выборочного распределения статистики состоит в том, чтобы вынимать дополнительные выборки с возвратом из самой выборки и повторно вычислять статистику для каждой повторной выборки.
- Бустрап
- Весь статистический анализ строится на основе предположений о свойствах ГС по некоторой выборки из этой ГС. Так как мы не можем взять всю ГС и оценить ее параметр (напр. средний рост), то мы берем случайную выборку из ГС и оцениваем параметр выборки и делаем предположения о том, как параметр может быть устроен в ГС.
- Статистический анализ - использование различных методов для того, чтобы определить свойства генеральной совокупности по выборке.





## Тестирование гипотез

### Идея проверок гипотез

### "Хитрости"

- Можно воспользоваться тестами и критериями, работающими с нормальным распределение, если:
- Нет выбросов в данных (строгое правило)
  - Нет явной асимметрии плотности распределения (иногда можно исправить с помощью логарифмирования или преобразования Бокса-Кокса)
  - Нет сильного отклонения от колоколообразной формы. Для бимодального распределения - попробовать разделить на выборки.

### Алгоритм

1. Формулируем нулевую гипотезу — Что изучаем и какие есть предположения относительно результата?
2. Определяем тип данных — Какие данные имеются?
3. Соотносим гипотезу и тип данных с критериями проверки — Какой критерий выбрать?
4. Подставляем данные в формулы — Какие результаты получили?
5. Отклоняем или не отклоняем нулевую гипотезу — Какие выводы можно сделать?

### A/B тестирование

- Это эксперимент с двумя группами для определения того, какой из двух сравниваемых объектов лучше.
- В некоторых случаях - Единовременность (для исключения влияния случайных факторов)
- Однородность выборок
- Достаточный объем выборки (исходя из стандартной ошибки)
- Независимость выборок
- Разбиение на группы
- Контрольная
- Тестовая
- Проверочная статистика — метрический показатель, который используется для сравнения группы А с группой В (как правило, двоичная переменная: нажатие/отсутствие нажатия, купить/воздержаться от покупки, мошенничество/не мошенничество и т. д.)
- A/B-тест, как правило, конструируется с учетом гипотезы. Например, гипотеза может заключаться в том, что цена В приносит более высокую прибыль.
- Статистическая проверка гипотез является дальнейшим анализом A/B-теста с целью установить, является ли случайная возможность выбора разумным объяснением наблюдаемой разницы между группами А и В.

### ВЫБРАТЬ ТЕСТ:





