

ПРИКЛАДНАЯ СТАТИСТИКА



Посиделка 1: схема статистики

Ульянкин Ппилиф *

You may hate the dictator, but something... far worse is gonna fill that void if you depose of him. I've lived a million lifetimes. I've gone through every, every scenario. TVA is the only way.

Nathaniel Richards a.k.a. He who remains

В этой посиделке мы поговорим про то, зачем мы учим тервер. Мы обсудим, как с помощью него можно замоделировать своё невежество и попытаться найти ответы на вопросы, которые нас мучают. Мы обсудим основы проверки гипотез и построим несколько простых критериев, основанных на комбинаторике. По мотивам этих примеров мы нарисуем схему матстата, которой в следующих посиделках будем активно пользоваться.

1 Зачем мы учим тервер

Зачем мы учим теорию вероятностей? Нам на ней говорят, что наша жизнь полна случайностей, а теория вероятностей помогает их описывать. Например, есть казино. В казино есть рулетка. Шарик катится по рулетке и выпадает на какое-то число. Результат абсолютно случаен **...или нет.**

Подкидывание шарика на рулетке — это **субъективная случайность**. Если у вас есть деньги и желание, можно измерить с какой именно скоростью шарик падает на рулетку, скорректировать это на то, как рулетка крутится, как в комнате дует ветер, учесть ещё кучу разных факторов и идеально предсказать, куда упадёт шарик.

*https://github.com/FUlyankin/matstat_lec

Проблема в том, что со всем этим измерительным оборудованием вас в казино не пустят. А ещё у человечества пока нет такого оборудования. Поэтому мы говорим, что падение шарика — случайно, а **вероятность описывает наше невежество**.

Точно также неслучайно то время, через которое автобус приедет на остановку. Если у нас будет информация о том, какие в городе пробки, насколько часто ломается автобус, в какое время безбилетники пытаются залезть в автобус, и водитель никуда не едет, пока они не уйдут, тогда можно попробовать предсказать время приезда. Однако мы всего этого не знаем. Более того, в течение поездки может произойти столько всего непредсказуемого, что нам проще считать время приезда автобуса случайным.

Так можно рассуждать про любое событие, которое нам кажется случайным. Можно даже подумать, что **объективной случайности** в природе просто-напросто не существует. Однако это не так. Можно придумать как минимум два примера объективной случайности. Первый связан с квантовой физикой.

Если бы человек мог в любой момент времени измерить положение и скорость каждой частицы во вселенной, никакой случайности бы не было. Вымышленное разумное существо, которое способно это делать называют **демоном Лапласа**. Лаплас придумал его в 1814 году для демонстрации нашего невежества. Так он обосновывал необходимость статистически описывать реальные процессы¹.

Более того, проблема даже не в человеке, а в природе. Современная физика, в частности **принцип неопределённости Гейзинберга**, говорит нам о том, что невозможно одновременно точно измерить и положение и скорость частицы. Получается, что случайность вшита в нашу природу. По крайней мере, люди сейчас так думают. Возможно, через двадцать лет физики придумают новые способы описывать реальный мир, в котором не будет никаких случайностей.

Второй пример объективной случайности — свобода воли. Если она существует, вряд ли получится предсказать, что именно взбредёт человеку сделать в следующую секунду. Даже если воткнуть ему в голову кучу электродов. Можно рискнуть сказать, что поведению людей свойственна объективная неопределённость.

Возможно, наше будущее предопределено, но человек не может этого осознать, и у него возникает ощущение, что он на что-то влияет. В таком случае свобода воли порождается субъективным незнанием, а не объективной неопределённостью. В фильме Дени Вильнёва «Притяжение (2016)» как раз обыгрывается то, что у нас нет органа, который чувствовал бы будущее, и из-за этого эволюционного невежества возникает иллюзия свободы воли².

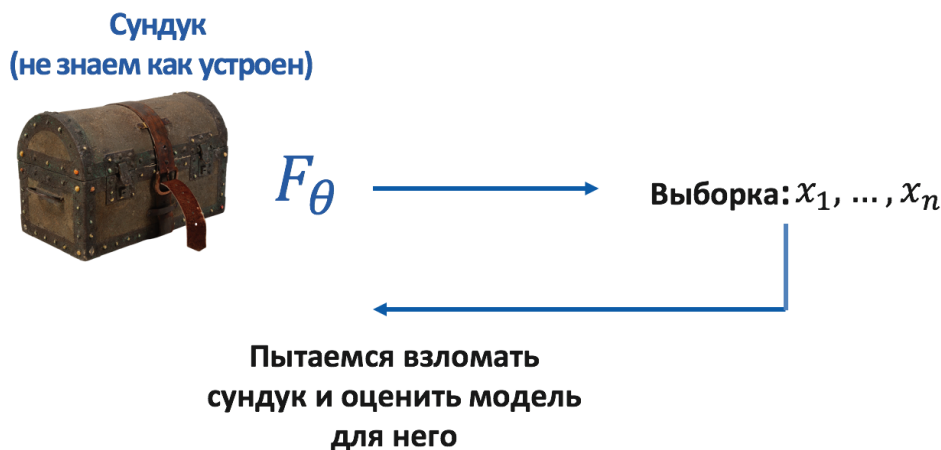
Теория вероятностей пытается предложить нам инструменты, которые помогают описать наше невежество с помощью субъективных случайностей. На ней изучают разные модели, которые можно использовать на практике для того, чтобы принимать решения. У каждой модели

¹Именно этот деман нарисован в шапке посиделки. Другая версия состоит в том, что там нарисован сэр Томас Байес. Возможно, что после смерти он стал демоном Лапласа.

²А ещё наша свобода воли может быть искусственно ограничена священным таймлайном. Тот, кто остаётся (he, who remains) сидит в своём замке на конце времени и с помощью могущественной TVA следит за тем, чтобы не возникало событий-некусов. Он делает это ради нашего блага, иначе Канг Завоеватель всех поработит. Если это так, то свобода воли довольно детерминирована. Диктатор ограничил её.

есть свои предпосылки. Мы собираем данные, выбираем модель и смотрим подходит ли она нам.

Весь наш мир — это сундук. Мы не знаем, как он устроен. Он выплёвывает на нас данные. Посмотрев на эти данные, мы можем предположить, как устроены внутренности сундука. **Наше предположение — это модель.** Модели помогают нам принимать решения и предсказывать, что сундук выплюнет на нас дальше.



Теория вероятностей занимается тем, что предлагает простейшие модели, которые могли бы описать внутренности сундука. Математическая статистика предлагает кучу способов смотреть на выборки, которые сундук выплюнул. Давайте посмотрим на пару примеров того, как можно формализовать модель, которая бы описывала сундук, и принять решение на основе данных.

2 Тест Фишера

На работу нужно взять 8 человек. Работодатель говорит, что берёт людей на работу абсолютно беспристрастно. Ему неважно какого пола кандидат. Прособеседовались 20 человек. Из них 8 были мужчинами, 12 женщинами. На работу взяли 7 мужчин и 1 женщину. Есть ли на рынке труда дискриминация?

Мы можем попытаться найти ответ на этот вопрос с помощью данных. Предположения, которые можно проверить с помощью данных, называют **гипотезами**. Нарисуем табличку. Внизу откладываются суммы по столбцам, справа по строкам. Всего собеседовалось 20 человек. Эту цифру запишем на диагоналке.

	мальчик	девочка	
взяли	7	1	8
не взяли	1	11	12
	8	12	20

Допустим, что работодатель не врёт и отбор правда был беспристрастен. Это наш статус-кво. По-другому статус-кво в статистике называют **нулевой гипотезой**. Чтобы отказаться от нулевой гипотезы, нужны весомые доказательства. Какова вероятность получить именно такую выборку, если работодатель не врёт, на рынке нет дискриминации и нулевая гипотеза выполнена?

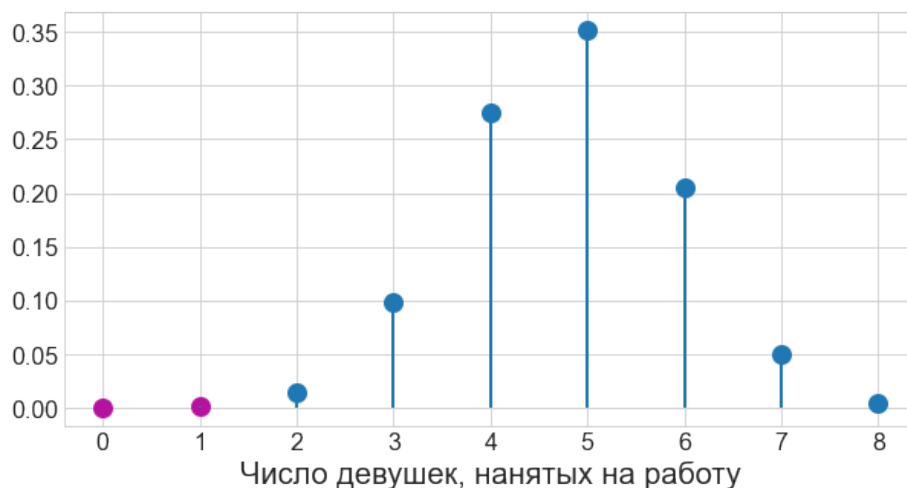
Всего вариантов выбрать 8 человек из 20 существует C_{20}^8 способов. Мы взяли 7 мальчиков и 1 девочку. Сделать это есть $C_8^7 \cdot C_{12}^1$ способов. Получается, вероятность того, что на 8 вакансий будет взята только одна девушка составит

$$\frac{C_8^7 \cdot C_{12}^1}{C_{20}^8} \approx 0.00076.$$

Вероятность того, что всё будет ещё хуже и девушек на работу вообще не возьмут составит

$$\frac{C_8^8 \cdot C_{12}^0}{C_{20}^8} \approx 7 \cdot 10^{-6}.$$

По аналогии можно посчитать все вероятности для всех возможных комбинаций и нарисовать их на картинке. По оси x отложим число девушек, взятых на работу, по оси y вероятность того, что такая ситуация возможна при беспристрастном отборе, то есть при верности нулевой гипотезы. В сумме все вероятности будут давать единицу.



Вероятность того, что при беспристрастном отборе мы увидим нашу ситуацию, либо ситуацию ещё хуже оказывается очень низкой. На картинке исходы, которые соответствуют этому, закрашены лиловым. Вероятность нашего, либо более плохого исхода, при верности нулевой гипотезы, обычно называют **p-значением (p-value)**. Для одной взятой на работу девушки р-значение составляет 0.00077.

Вероятность оказалась низкой. Это не означает дискриминации на рынке труда, но заставляет нас задуматься. Если нулевая гипотеза верна, и дискриминации нет, то именно такой исход эксперимента, как у нас, должен получаться очень-очень редко.

Ситуация, в которой на работу берут двух девушек или ещё меньше, при верности нулевой гипотезы, происходит с вероятностью 0.015. Это уже заметно больше, и при беспристрастном отборе происходит чаще. Если бы результаты нашего эксперимента дали бы три или четыре девушки, взятые на работу, мы бы оказались близко к центру распределения, и никаких вопросов к беспристрастности отбора у нас бы не было.

К сожалению, в нашей ситуации на работу взяли только одну девушку. Мы анализируем случайную величину — число девушек взятых на работу. Одна девушка — это **наблюдаемое значение** нашей случайной величины. Отталкиваясь от него мы должны принять решение, верим ли мы в нулевую гипотезу о беспристрастности. Вопрос в том, как его принять.

На практике выбирают какое-то место, начиная с которого перестают верить в нулевую гипотезу. Это место называют **критическим значением**. Например, мы можем сказать, что если на работу взяли 2 девушки и больше, мы верим в беспристрастность. Такая засечка и будет критическим значением, с которым сравнивают значение, полученное в эксперименте. Всё, что левее критического значения, мы считаем аномальным, и отказываемся от нулевой гипотезы.

Чтобы подобрать критическое значение, перед проведением эксперимента фиксируют **уровень значимости (ошибку первого рода)**. Уровень значимости — это наша готовность ошибаться и зря отказываться от статуса-кво. Если мы возьмём $\alpha = 0.01$, тогда при постоянном повторении эксперимента, мы в 1% случаев ошибочно откажемся от нулевой гипотезы, потому что при беспристрастном отборе такие редкие события иногда происходят.

Каждому значению α соответствует критическое значение случайной величины. Если мы возьмём $\alpha = 0.01$, оно будет между 2 и 1. Это рубеж для принятия решения. Если на работу взяли две девушки, мы верим в беспристрастность. Если одну, то уже не верим.

Можно рассуждать о том же самом в терминах р-значения. Если оно большое, значит фиолетовая площадь на графике больше уровня значимости. Это означает, что наблюдаемое значение попало правее критического, то есть мы оказались в хорошей зоне и гипотезы не отвергается. Мы продолжаем верить в статус-кво.

В нашем эксперименте $p\text{value} = 0.00077$. Р-значение оказалось меньше уровня значимости $\alpha = 0.01$. Значит наблюдаемое значение левее критического и от нулевой гипотезы нам надо отказаться. Возможно, что это ошибочное решение. Однако весь смысл этой процедуры заключается в том, что мы хотим мало ошибаться в долгосрочном периоде, при постоянном проведении эксперимента. С 1% ошибок мы согласны смириться.

Чем меньше р-значение, тем аномальнее наши наблюдения. Чем меньше значение уровня значимости, тем в более аномальные наблюдения мы верим. Чем меньше уровень значимости, тем сильнее мы боимся отказаться от нулевой гипотезы, и тем сильнее в хвост распределения углубляется засечка, на основании которой мы принимаем решение.

Кроме ошибки первого рода, уровня значимости, есть ещё и ошибка второго рода. **Ошибка первого рода** — это когда мы отказались от нулевой гипотезы (статуса-кво), а она была верна. Вероятность такой ошибки обычно обозначают буквой α .

Ошибка второго рода — это когда мы остались верны статусу-кво, а он на самом деле неверен. Её вероятность обычно обозначают буквой β .

	H_0 верна	H_0 не верна
выбрали H_0	ok	ошибка 2 рода β
отказались от H_0	ошибка 1 рода α	ok

Ошибка первого и второго рода связаны друг с другом. Когда растёт одна, вторая падает. Жёстко контролировать мы можем только одну из этих ошибок. Вторая минимизируется по остаточному принципу. Для простых критериев можно вывести формулы, которые показывают как эти ошибки связаны друг с другом. Обычно, чем больше собрано данных, тем меньше обе ошибки. В будущем мы научимся такие формулы выводить.

На практике обычно не хотят отказываться от статуса-кво. Все мысли формулируют именно так, чтобы было сложно от него отказаться. Например, если тестируют какое-то лекарство, намного страшнее выпустить на рынок плохое лекарство, чем не выпустить хорошее. Из-за этого в качестве гипотезы H_0 рассматривают ситуацию, когда новое лекарство бесполезно. Нам нужны весомые статистические доказательства, чтобы отказаться от такого статуса-кво.

Перекося в пользу гипотезы H_0 обычно называют **презумпцией нулевой гипотезы**. Это как презумпция невиновности в суде. Нельзя называть человека убийцей, пока его вина не доказана. Нельзя отказываться от нулевой гипотезы, пока мы не собрали данные и не увидели в них противоречия. Данные — наш судья.

Итак, наша нулевая гипотеза — на рынке нет дискриминации девушек. Альтернативная гипотеза — дискриминация есть. Если выбрать $\alpha = 0.01$, мы будем вынуждены отвергнуть нулевую гипотезу. Данные говорят против неё. Такую процедуру, которая помогает понять взаимосвязаны ли между собой два дискретных признака называют **тестом Фишера**. В примере выше мы искали взаимосвязь между полом и наймом на работу.

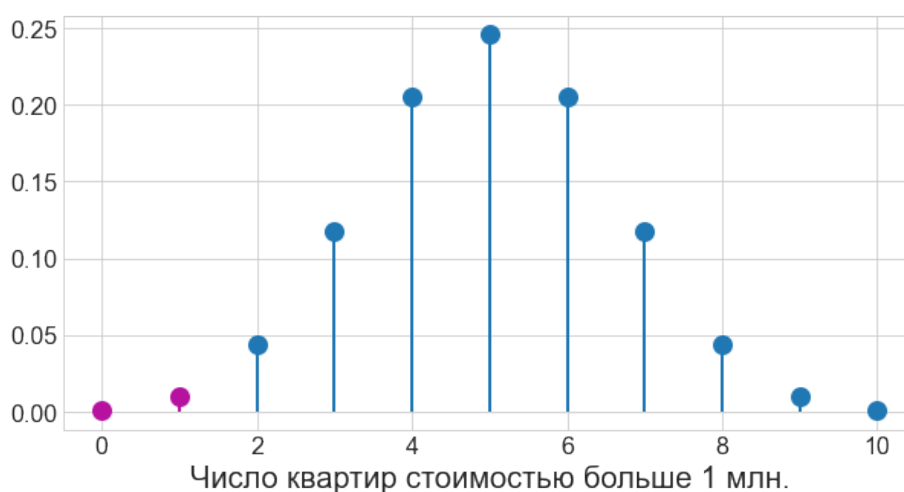
Выводы, которые мы сделали могут оказаться неверными. Мы действовали в рамках конкретной модели. Когда мы строили тест, мы предполагали, что есть только два фактора — пол и найм. Вполне возможно, что на самом деле есть какой-то третий признак, который повлиял на эксперимент, а мы его нигде не учли. Например, девушки могли сами отказаться от работы, хотя работодатель был готов их взять.

Когда проводят эксперимент, пытаются исключить все такие третьи факторы. Если это оказывается невозможным, либо если эксперимент провести невозможно, приходится привлекать более мощные статистические методы. Модели, описывающие сундук могут быть довольно сложными. При этом не факт, что они тоже хорошо будут описывать реальность.

3 Критерий знаков

Мэр сказал в новостях, что больше половины квартир в его городе стоит дешевле 1 млн. рублей. Как проверить, не врёт ли он? Мы можем собрать случайную выборку из квартир. Целых n штук. Будем ставить \oplus , если квартира окажется дешевле 1 млн. Если квартира будет дороже, будем ставить \ominus .

Мэр не врёт, если плюсики много. Мы изначально верим ему. Это наш статус-кво, наша нулевая гипотеза. Чтобы отказаться от неё, мы должны увидеть в данных что-то ужасное. Предположим, что мы посмотрели на 10 квартир. У нас есть 9 плюсиков и 1 минусик. Какова вероятность получить такую выборку либо выборку ещё хуже, при верности нулевой гипотезы?



Наш эксперимент представляет из себя неправильную монетку. Если монетка выпала орлом, квартира дешёвая. Если решкой, дорогая. Мы предполагаем, что вероятность орла $p \geq 0.5$. То есть больше половины квартир будут дешёвыми. Альтернатива заключается в том, что $p < 0.5$. Кратко можно записать это как:

$$H_0 : p \geq 0.5$$

$$H_a : p < 0.5.$$

Обычно в таких сложных нулевых гипотезах рассматривают самый шаткий статус-кво. В нашей ситуации таким будет утверждение, что ровно половина квартир оказалась дешевле 1 млн.

$$H_0 : p = 0.5$$

$$H_a : p < 0.5.$$

Зафиксируем **уровень значимости** $\alpha = 0.01$. Мы можем выбрать его любым. Чем меньше мы его возьмём, тем сильнее мы боимся зря отказаться от статуса-кво по ошибке. Ошибки второго рода мы здесь боимся не очень сильно. Лучше ошибочно не назвать мэра лжецом, чем

ошибочно назвать. Посчитаем **p-значение**, то есть вероятность получить либо наш исход, либо исход ещё хуже при верности нулевой гипотезы

$$pvalue = \frac{10}{2^{10}} + \frac{1}{2^{10}} \approx 0.0107 > \alpha = 0.01.$$

Получается, что гипотеза о том, что мэр не врёт, не отвергается. Многие скажут, что это противоречит здравому смыслу. И будут правы. На самом деле на такой маленькой выборке, при таком маленьком уровне значимости, ошибка второго рода будет зашкаливать. Чтобы сделать её меньше, не меняя ошибку первого рода, можно собрать более большую выборку. Тогда решение будет более устойчивым. В будущих посиделках мы выведем формулу для ошибки второго рода для данной задачи.

Такая процедура с расстановкой плюсиков и минусиков называется **критерием знаков**. Её довольно часто используют, когда измерения оказались не очень точными, но при этом нам хочется понять, в какую именно сторону направлены изменения. Например, так довольно часто происходит в биологии.

Если ваши измерения были сделаны точно, то воспользовавшись критерием знаков, вы потеряете довольно много полезной информации, так как превратите все цифры в единицы и нули. В таких ситуациях лучше использовать другие тесты. Например, в будущем мы будем говорить про тесты для средних, основанные на центральной предельной теореме.

4 Ранговый критерий (критерий Манна-Уитни)

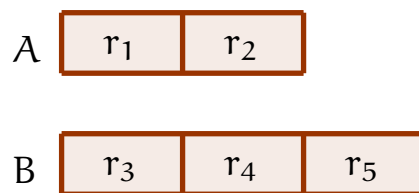
Нам срочно нужно протестировать новое лекарство от гипертонии. Для этого мы воспользуемся **двойным слепым тестированием**. В такой ситуации врач не знает, кому что даёт. Пациент тоже не знает, что именно ему дали. Если врач знает какие таблетки лекарство, а какие плацебо, он пытается давать лекарство тем, кто на его взгляд болен сильнее. Это может исказить результат. Поэтому приходится держать врачей в неведении.

Вопрос заключается в том, есть ли от лекарства какая-то польза. В качестве статуса-кво выберем гипотезу, что никакой пользы нет. Будем давать лекарство людям из группы А. Группе В дадим плацебо. У всех людей измерим давление.

Разбираться в том, есть ли разница между группами будем с помощью следующей процедуры: свалим все наблюдения в одну кучу и отсортируем их по возрастанию. Каждому наблюдению **припишем свой ранг**. Если все наблюдения разные — это просто порядковый номер. Если встречаются одинаковые наблюдения, то надо сложить их порядковые номера и поделить на число повторов.

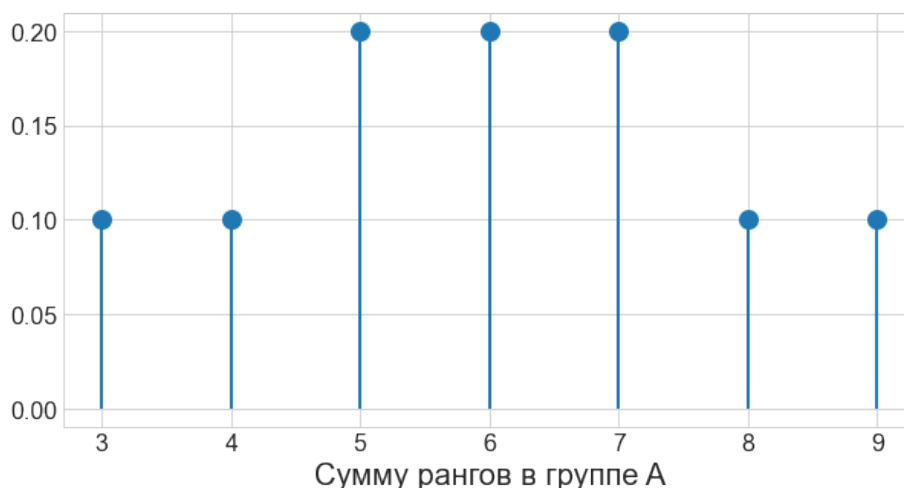
Например, для выборки 5, −4, 5, 3, 22 каждому наблюдению будет приписан порядковый номер 3, 1, 4, 2, 5. Третий и четвёртый элементы — это пятёрки. Если мы хотим превратить порядковые номера в ранги, нужно сложить эти номера и поделить на два. То есть рангами будут числа 3.5, 1, 3.5, 2, 5.

Пусть в группе А ранги оказались равны r_1 и r_2 . В группе В r_3, r_4, r_5 .



Если получается так, что сумма $r_1 + r_2$ оказалась маленькой, это сигнализирует нам о том, что лекарство помогает, и давление понижается. Как понять, как будет распределена такая случайная величина?

Предположим, что все наши наблюдения разные. Тогда ранги принимают значения 1, 2, 3, 4, 5. Самое маленькое значение для $r_1 + r_2$ будет равно $1 + 2 = 3$. Самое большое $5 + 4 = 9$. Посчитаем вероятность того, что $r_1 + r_2 = 6$. Всего два числа из пяти можно выбрать $C_5^2 = 10$ способами. В сумме 6 можно получить двумя способами. Получается, что $P(r_1 + r_2 = 6) = 0.2$. По аналогии можно посчитать все другие вероятности.



Осталось понять, насколько наша ситуация плохая и принять решение в соответствии с тем уровнем значимости, который мы выбрали до начала эксперимента.

Обратите внимание, что в этой задачке мы впервые явно выписали случайную величину, на основе которой мы принимаем решение. Это $T = r_1 + r_2$. В первых двух пунктах можно также выписать такую случайную величину. Для критерия знаков — это число плюсов. Для критерия Фишера — это число нанятых девушек.

Каждый раз мы принимаем решение на основе того, как эта случайная величина себя ведёт при верности нулевой гипотезы. Если она показывает аномальные значения, мы не верим в нулевую гипотезу. Такие случайные величины называют **статистиками**.

Принимая решение мы можем рассуждать, как в терминах уровня значимости и р-значения, так и в терминах **наблюдаемого значения статистики и критического значения**

статистики. Наблюдаемое значение рассчитывается по выборке, критическое получается исходя из уровня значимости. Если наблюдаемое значение слишком далеко от центра распределения, данные говорят против статуса-кво. Критическое значение — это порог, начиная с которого мы прекращаем верить в нулевую гипотезу. Для данной ситуации нам надо найти такое значение T , левее которого лежит α процентов распределения. Это значение и будет критическим.

Многие статистические процедуры можно построить отталкиваясь от таких статистик, случайных величин — помощников. Самое главное — знать их распределения при статусе-кво, при верности нулевой гипотезы. Есть разные теоремы, которые нам эти распределения подсказывают. Во всех трёх рассмотренных критериях нам не понадобилось каких-то особенных теорем. Мы смогли выписать все распределения в явном виде благодаря простой комбинаторике.

Критерий, который мы описали выше, обычно называют **критерием рангов**. В нём мы точно также отказываемся от изначальных измерений, предполагая что они не очень точные. Но в отличие от критерия знаков, мы пытаемся сохранить чуть больше информации, используя вместо плюсов и минусов, порядковые номера измерений.

5 Схема статистики

Любое исследование в анализе данных происходит примерно по одной и той же схеме. Сначала у нас есть **вопросы**. Мы хотим найти на них ответы. Нас могут интересовать дискриминация, изменение цен, работоспособность лекарств и многое-многое другое.

Чтобы грамотно ответить на каждый из этих вопросов, нужно собрать данные. Для этого мы проводим **эксперимент**. Эксперименты нужно планировать аккуратно. Продумыванию условий эксперимента в статистике уделяется особое внимание. Выборка, которая получается в результате его проведения, должна быть **репрезентативной**. То есть она должна отражать те свойства, которые действительно присущи всей **генеральной совокупности**, всем тем объектам, которые нас интересуют в исследовании. В ходе эксперимента мы собираем **данные**.

Сундук выплюнул на нас данные. Пора выяснить что именно внутри сундука. Нужно предположить как именно устроен мир вокруг нас, то есть надо придумать какую-то **модель**. Любая модель работает только в той ситуации, когда выполнены все её предпосылки. Про них нельзя забывать. В будущих посиделках мы постоянно будем обращать внимание на предпосылки и разбираться с тем, как их можно проверить.

$$\begin{array}{ccc} \text{Вопрос} & \Rightarrow & \text{Эксперимент} \\ x_1, \dots, x_n & \sim & F_\theta \\ \text{Данные} & & \text{Модель} \end{array}$$

С помощью собранных данных мы можем **оценить неизвестные параметры θ** нашей модели F_θ . Есть довольно много методов для оценки. В будущем мы будем обсуждать каждый

из них более подробно. Оценки неизвестных параметров обычно отмечают колпачком. Так $\hat{\theta}$ — это оценка неизвестного параметра θ .

Когда у нас есть готовая оценка, нам хочется понимать насколько точной она получилась. Для этого надо знать как она распределена, то есть нужно найти её плотность распределения $f_{\hat{\theta}}(t)$. Тут к нам на помощь будут приходить разные **теоремы — союзники**. Это может быть центральная предельная теорема или какие-то специфические распределения вроде распределения Стьюдента.

На основе распределения оценки $\hat{\theta}$ мы можем проверять гипотезы и строить доверительные интервалы, то есть пытаться понять насколько наша оценка получилась точной. Например, в трёх критериях выше, мы пользовались дискретными распределениями, которые мы выводили с помощью простейшей комбинаторики. Она была нашим союзником. Никакие параметры в этих распределениях нам оценивать не надо было, достаточно было просто посчитать вероятности.

Использованные нами критерии — очень простые. Они игнорируют довольно много информации из выборки. Они слишком упрощают её. Чем больше у модели параметров, тем лучше она может описать данные и тем качественнее может получиться ответ на мучающий нас вопрос. Если конечно мы не **переобучимся**. Переобучением называют ситуацию, когда модель вместо того, чтобы обобщить данные, запомнила их. Чем больше у модели параметров, тем легче ей переобучиться.

Хотим:

- несмещённость
- состоятельность
- эффективность



- прогнозы
- насколько точны прогнозы
- ответы на вопросы (гипотезы)

Союзники:

- метод моментов
- метод максимального правдоподобия

Союзники:

- ЦПТ
- Дельта-метод
- χ^2_n , t_n , $F_{n,k}$
- Теорема Фишера

Ещё нам часто будет хотеться, чтобы оценки неизвестных параметров обладали хорошими свойствами: **несмещённостью, состоятельностью и эффективностью**. Методы оценивания придумывают так, чтобы добиваться этих свойств. В следующих посиделках мы углубимся в схему матстата и подробно обсудим каждую её часть.

Отдельно хочется обратить внимание на то, что идти по схеме матстата можно очень по-разному. Можно отталкиваться от средних и строить всю науку через закон больших чисел и центральную предельную теорему. Можно использовать метод максимального правдоподобия и его свойства. Можно предположить, что параметр θ на самом деле какая-то случайная

величина и пойти по дороге байесовской статистики. Можно пойти по дороге информационных подходов. Мы подробно поговорим про каждый из этих подходов.