

## Семинар 3: основы статистики!

В первом разделе семинара большая задачка, которую мы с вами решали на паре.

### Упражнение 1

Коллекционер Настя собрала целых 10 наблюдений и записала их в табличку. Теперь Настя хочет стать аналитиком и проанализировать таблицу. Помогите ей.

имя	пол	возраст	вес
Кхал	м	14	80
Санса	ж	16	40
Мелисандра	ж	20	40
Эддард	м	20	80
Сандор	м	14	80
Миссаедея	ж	25	40
Якен	м	30	80
Теон	ж	23	40
Тирион	м	22	80
Станис	м	16	440

- Что такое непрерывная переменная? Что такое категориальная переменная? Какие переменные в табличке относятся к непрерывным? Какие к категориальным? Приведите ещё примеров непрерывных и категориальных переменных!
- Найдите долю мужчин и женщин в выборке. Постройте для пола гистограмму.
- Найдите средний возраст и медианный возраст. Что означают эти числа. В чём они измеряются?
- Найдите дисперсию возраста. В чём измеряется эта величина? Зачем обычно ищут среднее квадратичное отклонение? Найдите его.
- Постройте гистограмму для возраста. Считайте, что ширина одного столбца — 5 лет. Если человек попадает на правую границу отрезка, он попадает в текущий столбец. Изобразите на гистограмме среднее, медиану. Как бы вы нарисовали на гистограмме стандартное отклонение?
- Что такое выброс? Есть ли выбросы в возрасте? Есть ли выбросы в весе? Как выглядит выброс на гистограмме? Найдите средний вес и медианный вес. Чем медиана в данном случае лучше, чем среднее?
- Чувствительна ли дисперсия к выбросам?
- Что такое мода? Почему использовать её для непрерывных переменных не очень хорошая идея? Найдите моду для имени, пола и возраста.

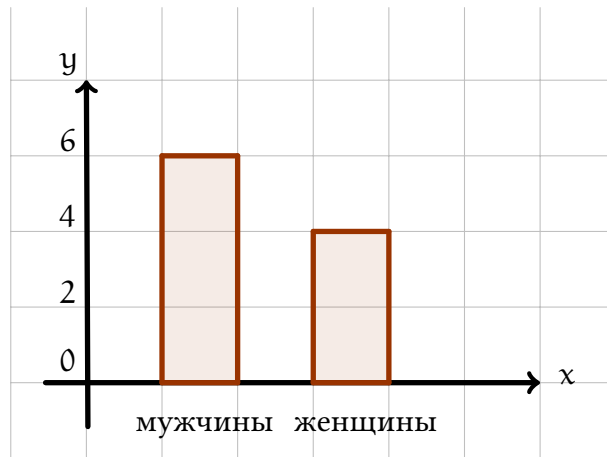
- и) Что такое квантиль? Предложите способ, борьбы с выбросами, основанный на знании того, что такое квантиль.

### Решение:

- а) **Непрерывная переменная** не ограничена каким-то конечным набором значений и может принимать любые числовые значения. Например: цена на квартиру, валютный курс, возраст, число лайков под фоткой и т.п.

**Категориальная переменная** принимает значения из какого-то фиксированного конечного множества. Например: пол, марка машины и т.п.

- б) В выборке 6 мужчин и 4 женщины. Всего 10 человек. Значит доля мужчин  $\frac{6}{10} = 0.6$ , доля женщин  $\frac{4}{10} = 0.4$ . Нарисуем гистограмму. По оси  $x$  будем откладывать возможные значения для нашей переменной, по оси  $y$  насколько часто это значение наблюдается в выборке.



- в) Найдём **средний возраст**. Для этого сложим все числа и поделим их на количество наблюдений

$$\frac{1}{10} \cdot (14 + 16 + 20 + 20 + 14 + 25 + 30 + 23 + 22 + 16) = 20.$$

Средний возраст это 20 лет. Формула для подсчёта среднего выглядела так:

$$\bar{x} = \frac{1}{n}(x_1 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

Привыкайте к формулам. Они будут часто встречаться вам по жизни.

Чтобы найти медиану, нам нужно упорядочить всех людей из выборки по возрасту и посмотреть какое число оказалось в середине:

14 14 16 16 20 20 22 23 25 30

У нас в середине находятся сразу два числа. В случаях, когда такое происходит в качестве медианы берут среднее этих двух чисел. В нашем случае это  $1/2 \cdot (20 + 20) = 20$ . Медиана это число, взятое посередине. Половина выборки оказывается слева от него, а вторая половина справа. Среднее и медиана в нашей задачке измеряются в годах и обозначают типичный возраст, который присущ людям из выборки.

- г) **Дисперсия** — это мера разброса. Она показывает насколько разнообразными могут быть элементы в выборке, насколько сильно они могут отклоняться от своего типичного значения.

Чтобы найти её, нужно посмотреть насколько сильно каждый представитель в выборке отличается от текущего. Величина такого отличия называется отклонением. Предположим, что Алёне 18 лет. Карине 22 года. Тогда отклонением для Алёны от среднего возраста будет  $18 - 20 = -2$  года. Для Карины отклонением будет  $22 - 20 = 2$  года.

Если просуммировать эти отклонения, мы получим  $-2 + 2 = 0$ . То есть в выборке нет никакого разброса. Все люди не отличаются от среднего. Это неправда. Для того, чтобы избежать неправды и жить по правде, отклонения возводят в квадрат, тогда мы получаем, что суммарное отклонение будет  $(-2)^2 + 2^2 = 4 + 4 = 8$ . Посмотрев на такое число мы сразу же поймём, что в выборке есть неоднородность.

Среднее значение квадратов отклонений от среднего и называется дисперсией. Давайте найдём её. Ещё раз выпишем наши наблюдения:

14 14 16 16 20 20 22 23 25 30

Сначала из каждого вычитаем среднее. Это даст нам

-6 -6 -4 -4 0 0 2 3 5 10.

Теперь возводим все отклонения в квадрат

36 36 16 16 0 0 4 9 25 100.

Складываем их. Получается 242. Остаётся разделить это число на 10 (количество наблюдений). Получается, что дисперсия составит 24.2 квадратных года. Из-за того, что мы каждое слагаемое возводили в квадрат, **дисперсия измеряется в квадратных годах**.

Когда мы умножаем одну сторону квадрата, измеренную в метрах, на другую, мы получаем его площадь. Она измеряется в квадратных метрах. Тут похожая ситуация. Мы бы хотели вернуться назад, к обычным годам. Для этого из дисперсии извлекают корень и получают штуку под названием стандартное отклонение. В нашем случае получится  $\sqrt{24.2} \approx 4.9$  года.

Можно найти дисперсию проще и быстрее. Для этого есть специальная формула:

$$\hat{\sigma}^2 = \bar{x}^2 - (\bar{x})^2,$$

то есть дисперсия это среднее квадратов минус квадрат среднего. Эту формулу довольно

просто доказать. На матстате вы её докажете. А пока просто воспользуемся ей. Найдём квадрат среднего:

$$(\bar{x})^2 = 20^2 = 400$$

Теперь среднее квадратов:

$$\bar{x^2} = \frac{1}{10} \cdot (14^2 + 16^2 + 20^2 + 20^2 + 14^2 + 25^2 + 30^2 + 23^2 + 22^2 + 16^2) = 424.2.$$

Остался последний штрих:

$$\hat{\sigma}^2 = 424.2 - 400 = 24.2.$$

Пользуйтесь тем способом, который вам больше нравится. Про дисперсию давайте обсудим ещё пару дополнительных полезных нюансов:

- Мы возводим отклонения в квадрат не только для того, чтобы сделать все числа положительными. Попутно мы подчёркиваем, что чем больше отклоняется возраст от среднего, тем это хуже. Так штраф за отклонение в два года составит 4, а за отклонение в три года, 9. С подобной логикой мы ещё встретимся, когда будем обсуждать различные метрики, используемые в машинном обучении.
- Часто при подсчёте дисперсии вместо формулы

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2,$$

которую использовали мы, используют

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

**Вторая формула, на самом деле, корректнее, чем первая.** В *pandas* используется именно она. У этого есть глубокая причина, которая называется несмещённостью оценки. В полной мере вы узнаете про это в курсе по математической статистике. Мы вкратце поговорим про несмещённые оценки ближе к концу курса, когда будем говорить про АБ-тесты. Пока держите это в голове, как вопрос, на который у вас нет ответа. Надеюсь, что это будет как следует мучать вас по ночам и стимулировать ботать.

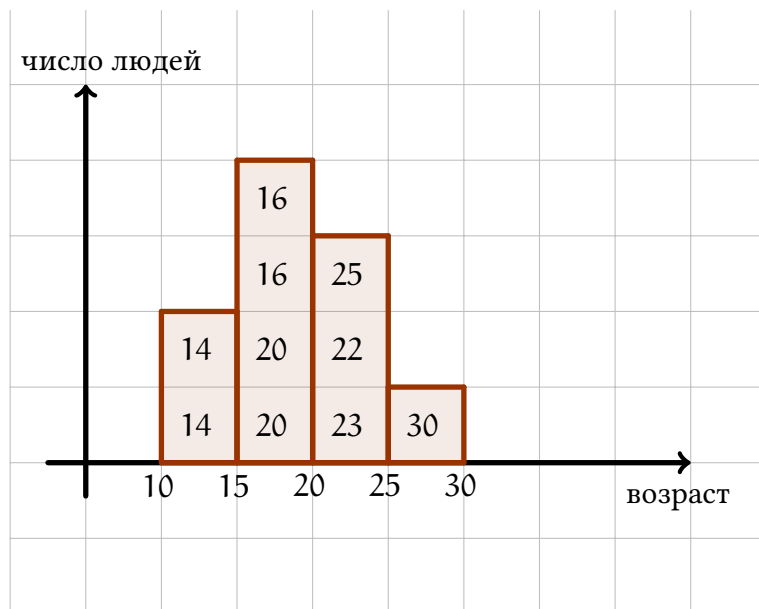
- Если распределение у данных нормальное (что такое нормальное распределение — отдельный и очень важный вопрос), тогда большая часть выборки, а именно 69% кучкуется в диапазоне между  $\bar{x} - \hat{\sigma}$  и  $\bar{x} + \hat{\sigma}$ .

При этом 95% выборки находится между  $\bar{x} - 2 \cdot \hat{\sigma}$  и  $\bar{x} + 2 \cdot \hat{\sigma}$ , а 99.9% выборки находятся между  $\bar{x} - 3 \cdot \hat{\sigma}$  и  $\bar{x} + 3 \cdot \hat{\sigma}$ .

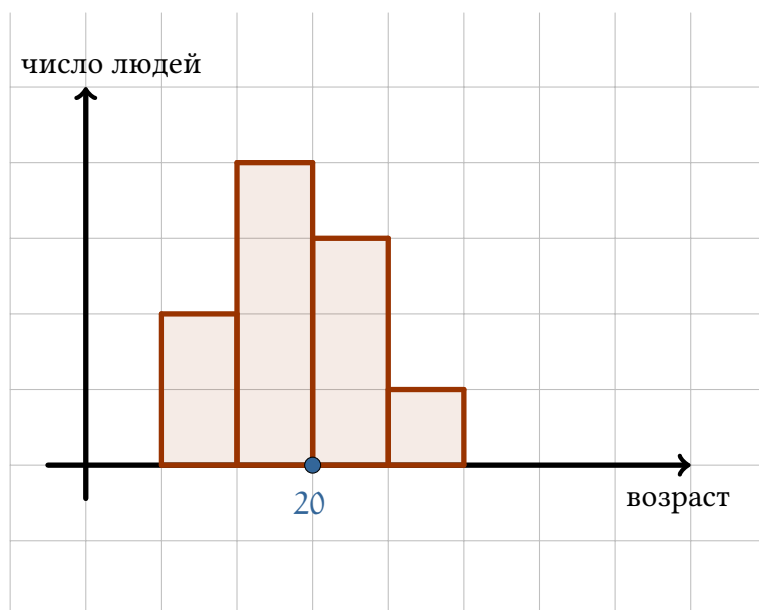
Правила таких кучкований называют правилом одной, двух и трёх сигм. Их часто используют для проведения АБ-тестов. Об этом мы тоже поговорим ближе к концу

курса.

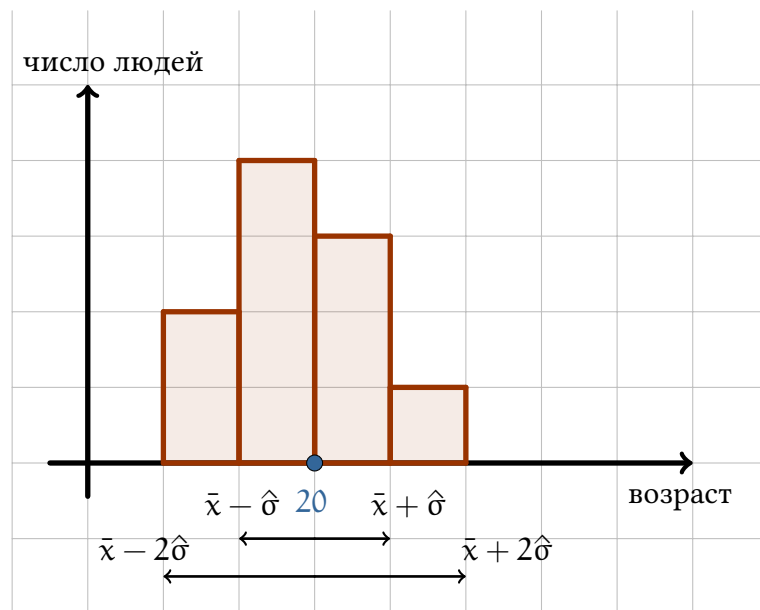
- д) Отмечаем по оси  $x$  каждые 5 лет, как сказано в условии задачи. Для всех людей, попавших в этот отрезок рисуем столбик высоты равной количеству людей, попавших в отрезок. Если человек попадает в правую границу отрезка, он попадает и в столбик. Например, 20 — это правая граница второго отрезка. Все люди, которым 20 лет попадают во второй столбик. Это просто договорённость о том, что делать на границе, в спорной ситуации. Не более того.



Отлично! Гистограмма готова. Каждого человека, которого мы внесли в тот или иной столбец, мы подписали. Давайте отметим на гистограмме медиану и среднее значение. Как это не странно, они оказываются в "центре" распределения.



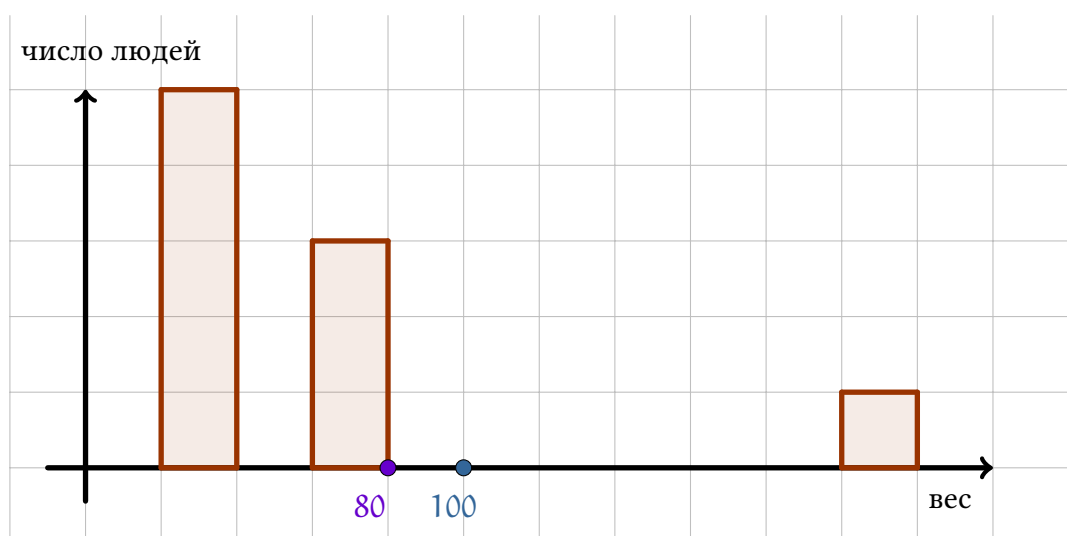
Выше мы обсудили, что стандартное отклонение — величина, которая описывает вариацию выборки вокруг среднего значения и поговорили про правила сигм. Давайте нарисуем от среднего отступы на сигмы вправо и влево.



е) В возрасте всё хорошо. В весе есть выброс, кто-то слишком много ест. Давайте найдём среднее и медиану. Среднее окажется равно  $\frac{1000}{10} = 100$ . Медиана окажется равна 80. Видим, что выброс существенно сдвинул среднее значение веса в большую сторону. Из-за этого оно перестало отражать типичный вес человека из выборки. Наше представление о людях оказалось искажено.

Медиана, в отличие от среднего, оказывается нечувствительна к выбросам. Это происходит из-за способа её поиска. Мы упорядочиваем наблюдения по порядку и смотрим на то, какое в середине. Значение выброса никак не участвует в подсчёте медианы и именно из-за этого не искажает её.

На гистограмме переменным, в которых есть выбросы соответствуют очень длинные хвосты. В нашем случае именно так и произошло:



ж) К несчастью, **дисперсия чувствительна к выбросам**. Когда мы считаем её, мы возводим все разности в квадрат. Разница между средним и выбросом будет большой. Когда мы возведем её в квадрат и прибавим к дисперсии, она очень сильно увеличится.

- з) Мы с вами определили моду как самое часто встречаемое значение признака в выборке. Для пола модной будут мужчины. Для веса модой будет 80. Для возраста модой будет либо 20 либо 14.

Для непрерывных переменных использовать моду в качестве меры типичности довольно глупо. Часто бывает так, что непрерывные признаки довольно близки друг к другу, но немного различаются. **Чаще всего моду используют, чтобы охарактеризовать именно категориальные переменные. Смотрят на пару: мода, её частота.**

На самом деле моду можно определить так, чтобы она была корректна и для непрерывных признаков. Обычно говорят, что мода это самое вероятное значение в выборке. Это позволяет найти её по плотности распределения (грубо говоря, по гистограмме), пытаясь понять какому числу соответствует её самая высокая точка. Подробнее об этом вы узнаете на теории вероятностей.

- и) На вопрос что такое квантиль, нам поможет ответить медиана. Мы сказали с вами, что если отсортировать выборку по возрастанию, то в середине у неё окажется медиана.

14 14 16 16 20 20 22 23 25 30

Получается, что 50% выборки больше медианы, и 50% выборки меньше медианы. Медиана — это 50% квантиль. По аналогии можно придумать другие квантили. Например, ниже красным отмечены 30% и 70% квантили:

14 14 16 16 20 20 22 23 25 30

Ровно 30% выборки  $\leq 16$  и 70% больше 16. И наоборот в случае 22. Среднее и медиана помогают понять какие представители типичны для середины распределения. Квантили помогают понять какие представители типичны для разных кусков распределения. Например, если мы имеем дела со стоимостью недвижимости, мы можем понять какая стоимость квартир типична для элитных районов.

Как мы выяснили выше, **выбросы могут существенным образом исказить наши представления о выборке.** От них нужно выборку очищать. Один из способов: отрубить все наблюдения, которые находятся выше 99% квантиля и все наблюдения, которые находятся ниже 1% квантиля. Все выбросы такой процедурой будут убиты и мы сможем спокойно работать с выборкой. Иногда берут 95% и 5% квантили.

## Ещё задачи!

Тут находится несколько задачек, о которых вам нужно подумать самостоятельно, в домашних условиях, за чашкой чая. Одна из этих задачек точно попадётся вам на самостоятельной работе. Вторая задачка на ней будет совсем новой. Посчитать надо будет что-то похожее на то, что было в первой части pdf-ки.

## Упражнение 2

Ваня любит пить чай. Иногда он пьёт его с сахаром, иногда без. На этой неделе он помечал 1 дни, когда пил чай с сахаром. Получилось 1, 1, 0, 0, 1, 0.

- а) Найдите среднее значение сахарных дней в жизни Вани. Найдите дисперсию сахарных дней.
- б) Правда ли, что среднее число сахарных дней совпало с долей сахарных дней? Почему так вышло? Всегда ли так будет происходить?
- в) Между дисперсией и долей в случаях, когда переменная принимает значения 0 или 1 тоже есть связь. Сможете догадаться как будет выглядеть формула, описывающая эту связь?

### Решение:

- а) Среднее получится  $1/2$ , а дисперсия  $1/4$ .
- б) Да, правда! Долю единиц мы с вами считаем по формуле:

$$p = \frac{\text{количество единиц}}{\text{общее число наблюдений}}.$$

Среднее мы считаем по формуле:

$$\bar{x} = \frac{1 + 1 + 0 + 0 + 1 + 0}{6},$$

но ведь сумма в числителе среднего это и есть количество единиц! Для переменной, которая принимает значения 0 или 1 среднее всегда совпадает с долей.

Кстати говоря, по-другому долю можно проинтерпретировать как вероятность встретить в выборке 1. Запомните этот факт, мы часто будем в питоне считать долю единиц в выборке как `pr.mean(x)`, то есть как среднее.

- в) Самый сложный пункт этой задачки. Сначала скажу вам ответ, а потом расскажу как можно до этого догадаться. **Дисперсия для величины, которая принимает значения 0 и 1 считается по формуле  $p \cdot (1 - p)$ , где  $p$  — доля единиц.**

Дисперсию можно найти по формуле:

$$\hat{\sigma}^2 = \bar{x^2} - (\bar{x})^2.$$

Когда в выборке есть только 0 и 1 всегда  $x_i^2 = x_i$ , потому что  $1^2 = 1$ ,  $0^2 = 0$ . Получается, что

$$\bar{x^2} = \frac{1}{n} \cdot (x_1^2 + x_2^2 + \dots + x_n^2) = \frac{1}{n} \cdot (x_1 + x_2 + \dots + x_n) = \bar{x}.$$

Мы помним, что среднее равно доле, получается, что

$$\hat{\sigma}^2 = p - p^2 = p \cdot (1 - p).$$



Получилась формула, которую мы с вами выписали чуть выше. Этот фокус работает только для долей! То есть только для выборок из нулей и единиц. Эта формула нам понадобится, когда мы будем разбираться с АБ-тестами.

### Упражнение 3

Имеется пять чисел:  $x$ , 9, 5, 4, 7. При каком значении  $x$  медиана будет равна среднему? А можно ли поставить такие цифры в условии задачи, чтобы  $x$  не существовал?

#### Решение:

Расположим числа в порядке возрастания: 4, 5, 7, 9. В зависимости от расположения  $x$  меняется медиана. Так, если мы воткнём  $x$  перед или сразу после 4, медианой будет 5. Если воткнуть  $x$  после 5, то сам  $x$  будет медианой. Если воткнуть  $x$  в конце или перед 9, то медианой окажется 7.

Составим три уравнения:

$$\begin{aligned}\frac{x + 4 + 5 + 7 + 9}{5} &= 5 \Rightarrow x = 0 \\ \frac{4 + 5 + x + 7 + 9}{5} &= x \Rightarrow x = 6.25 \\ \frac{4 + 5 + 7 + 9 + x}{5} &= 7 \Rightarrow x = 10\end{aligned}$$

### Упражнение 4

Измерен рост 25 человек. Средний рост оказался равным 160 см. Медиана оказалась равной 155 см. Машин рост в 163 см был ошибочно внесен как 173 см. Как изменится медиана и среднее после исправления ошибки? А как могут измениться медиана и среднее, если рост Маши равен 153?

#### Решение:

Если рост Маши ошибочно был внесен как 173 см вместо 163, то при исправлении ошибки изменения никак не отразятся на медиане, потому что ошибочно внесенный рост и ее рост больше медианы. Средний рост уменьшится. В случае 153 изменения могут коснуться как среднего, так и медианы.

### Упражнение 5

Деканат утверждает, что если студента N перевести из группы А в группу В, то средний рейтинг каждой группы возрастет. Возможно ли такое?

#### Решение:

Да, возможно. Если средняя оценка N ниже средней оценки группы А, но выше средней в группе В, то после смены студентом группы средняя оценка каждой группы и ее рейтинг возрастут.

Например, группы и состоят из 3 человек, которые имеют оценки 8, 9, 10 и 1, 2, 3 соответственно. Студент N, имеющий оценку 8, желает перейти в группу . Тогда изменения оценки

группы А:  $\frac{(9+10)}{2} - \frac{(9+10+8)}{3} = 0.5$ , то есть рейтинг группы повысится. Изменения для группы В:  $\frac{(1+2+3+8)}{4} - \frac{(1+2+3)}{3} = 1.5$ , а значит рейтинг группы В тоже повысится.

## Упражнение 6

Иногда в качестве меры разброса используют размах. Находят максимальное значение в выборке, минимальное значение выборке, а после вычитают из максимума минимум. Как думаете, такая мера чувствительна к выбросам? Предложите способ сделать её устойчивой к ним.

### Решение:

Да, будет. Для того, чтобы сделать эту меру устойчивой к выбросам, можно считать интерквартильный размах, то есть вычитать из 75% квантиля 25% квантиль. Это поможет срезать все выбросы и учитывать только центральную часть распределения. Если вы не очень поняли что означает эта фраза, вернитесь к гистограмме, поставьте на ней засечки там, где расположены 75% квантиль и 25% квантиль, а потом подумайте где окажутся выбросы.

Упражнение на придумать дисперсию в модулях и медианах либо сказать устойчива ли она к выбросам - ???

Задача на то, кто обычно левее а кто правее (мода медиана среднее - ???