

# Применение статистического аппарата

При Анализе Данных

Роман  
Гараев

# Задача #1. Данные о Кредитном Лимите и Сравнение Различных Групп Относительно Возврата/Невозврата

— — —

Применить статистический аппарат проверки гипотез, составления доверительных интервалов и визуализации в целях изучения данных и построения выводов.

LIMIT\_BAL: размер кредитного лимита (в том числе и на семью клиента) SEX: пол клиента (1 = мужской, 2 = женский )

EDUCATION: образование (0 = доктор, 1 = магистр; 2 = бакалавр; 3 = выпускник школы; 4 = начальное образование; 5 = прочее; 6 = нет данных ).

MARRIAGE: (0 = отказываюсь отвечать; 1 = замужем/женат; 2 = холост; 3 = нет данных).

AGE: возраст в годах

PAY\_0 - PAY\_6 : История прошлых платежей по кредиту. PAY\_6 - платеж в апреле, ..., PAY\_0 - платеж в сентябре. Платеж = (0 = исправный платеж, 1 = задержка в один месяц, 2 = задержка в 2 месяца ...)

BILL\_AMT1 - BILL\_AMT6: задолженность, BILL\_AMT6 - на апрель, ..., BILL\_AMT1 - на сентябрь

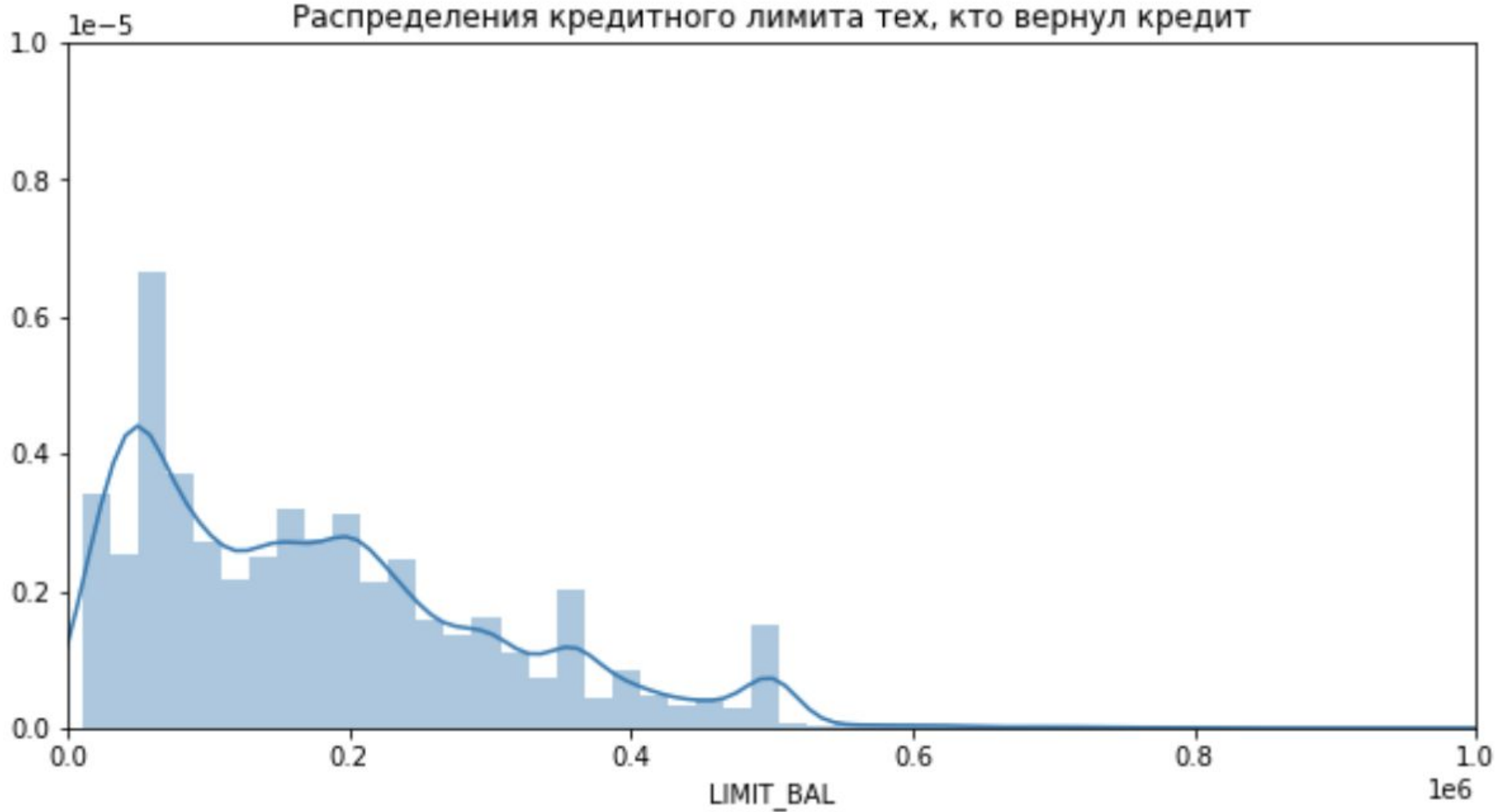
PAY\_AMT1 - PAY\_AMT6: сумма уплаченная в PAY\_AMT6 - апреле, ..., PAY\_AMT1 - сентябре

default - индикатор невозврата денежных средств

# Проверка равенства медианных значений 2 групп



Распределения кредитного лимита тех, кто вернул кредит



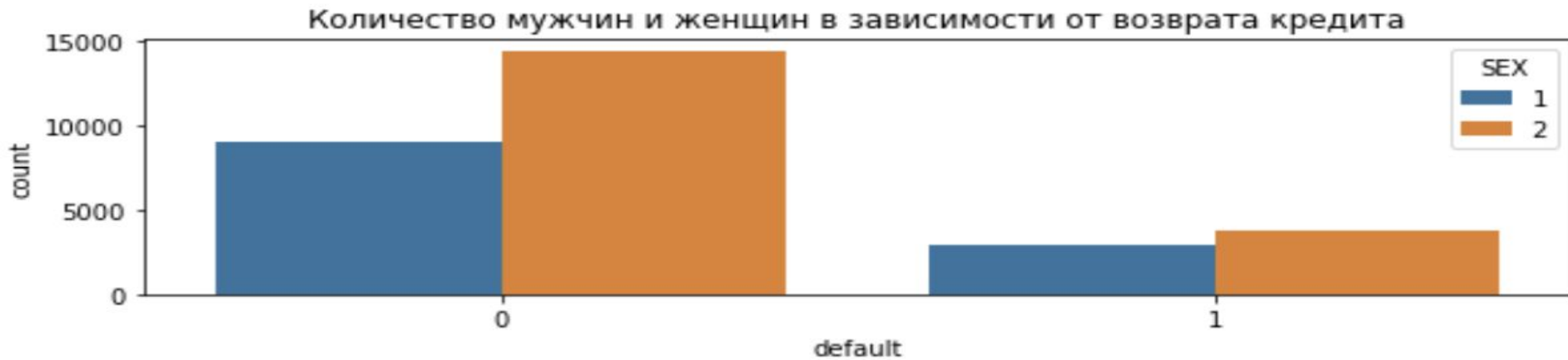
# Использованные техники и Вывод

— — —

Проверка равенства медианных значений с помощью сравнения медианных значений, доверительных интервалов и применения перестановочного критерия.

**Вывод:** Заемщики, которые не возвращают кредит берут большую сумму.

# Сравнение данных о возврата/невозврате относительно пола



# Использованные техники и Вывод

— — —

Построение доверительных интервалов для долей и разности долей

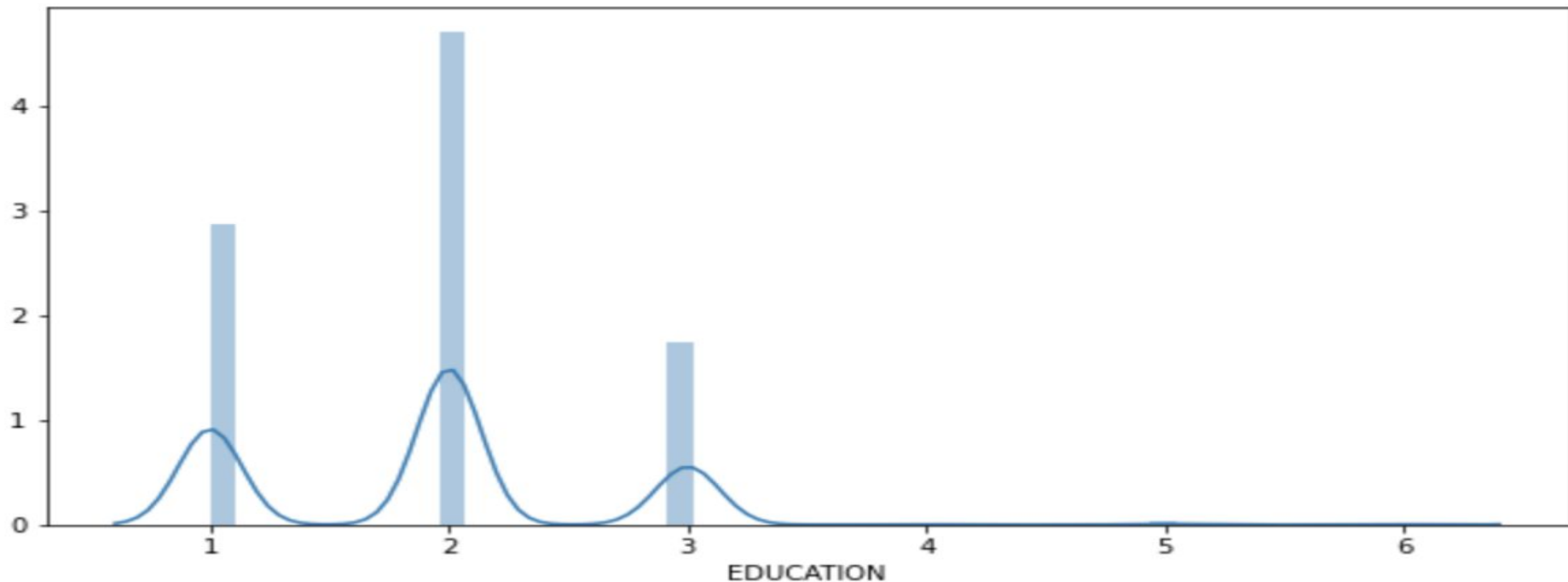
**Вывод:** Доверительный интервал для разности долей не включает ноль, а значит можно говорить о различных долях возврат в зависимости от пола.

**P.S.** Данный вывод справедлив только для данного датасета. На самом деле данный признак никак нельзя использовать для реальной оценки по причине неэтичности.

# Сравнение данных об образовании относительно возврата/невозврата

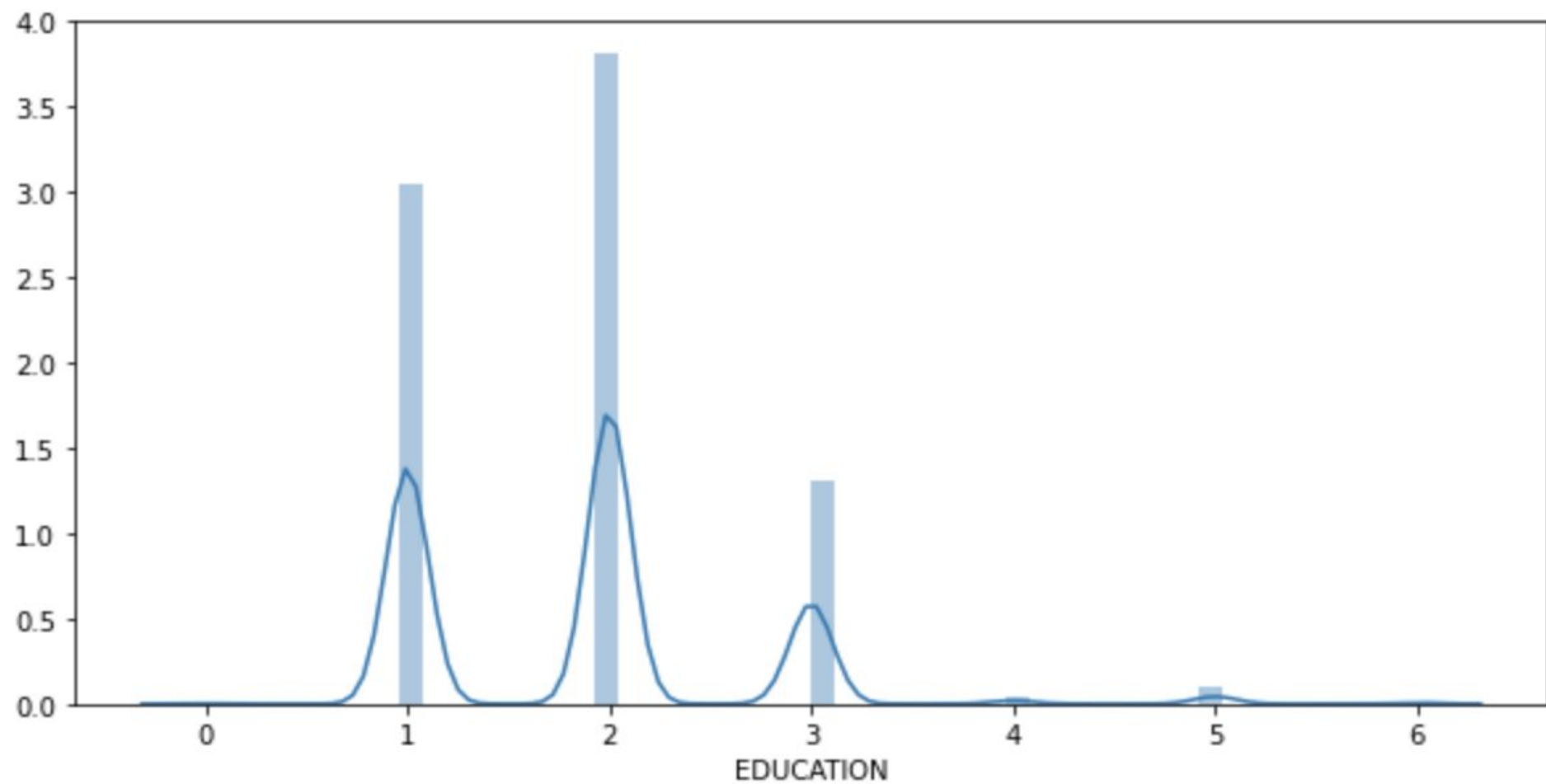
---

Возврат





# Невозврат





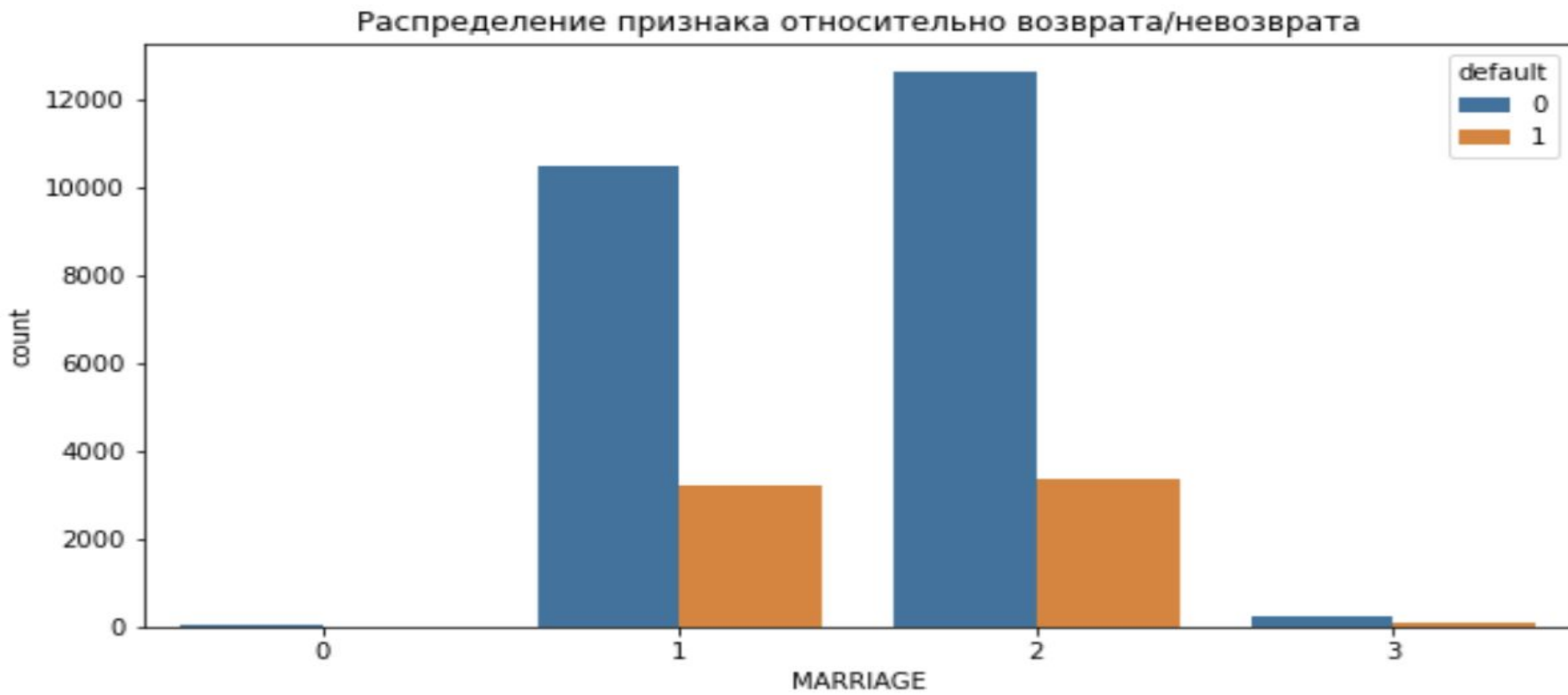
# Использованные техники и Вывод

— — —

Проверка на соответствие распределения нормальному с помощью критерия согласия Пирсона  $\chi^2$ -квадрат

**Вывод:** Распределение признака отличается от нормального.  
Уровень образования влияет на вероятность  
возрата/невозврата

# Сравнение данных о семейном положении относительно возврата/невозврата



# Использованные техники и Вывод

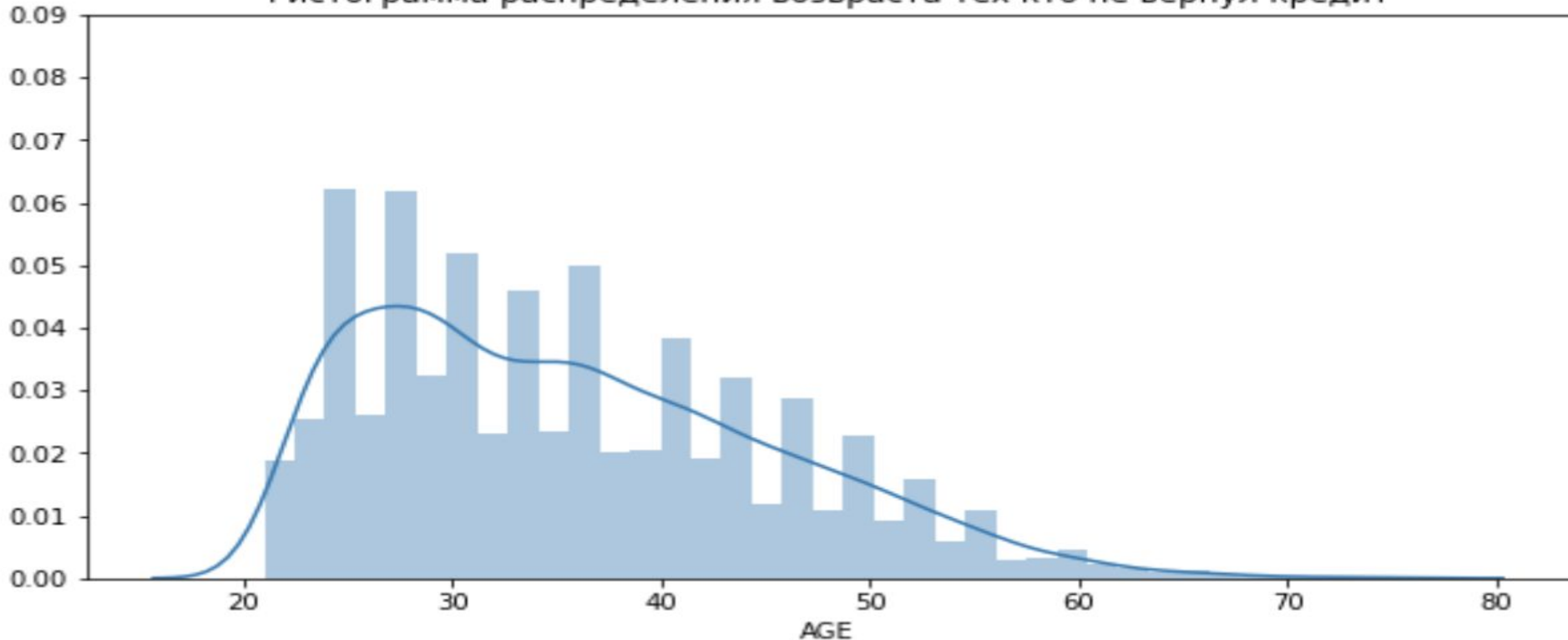
— — —

Коэффициент корреляции “V” Крамера

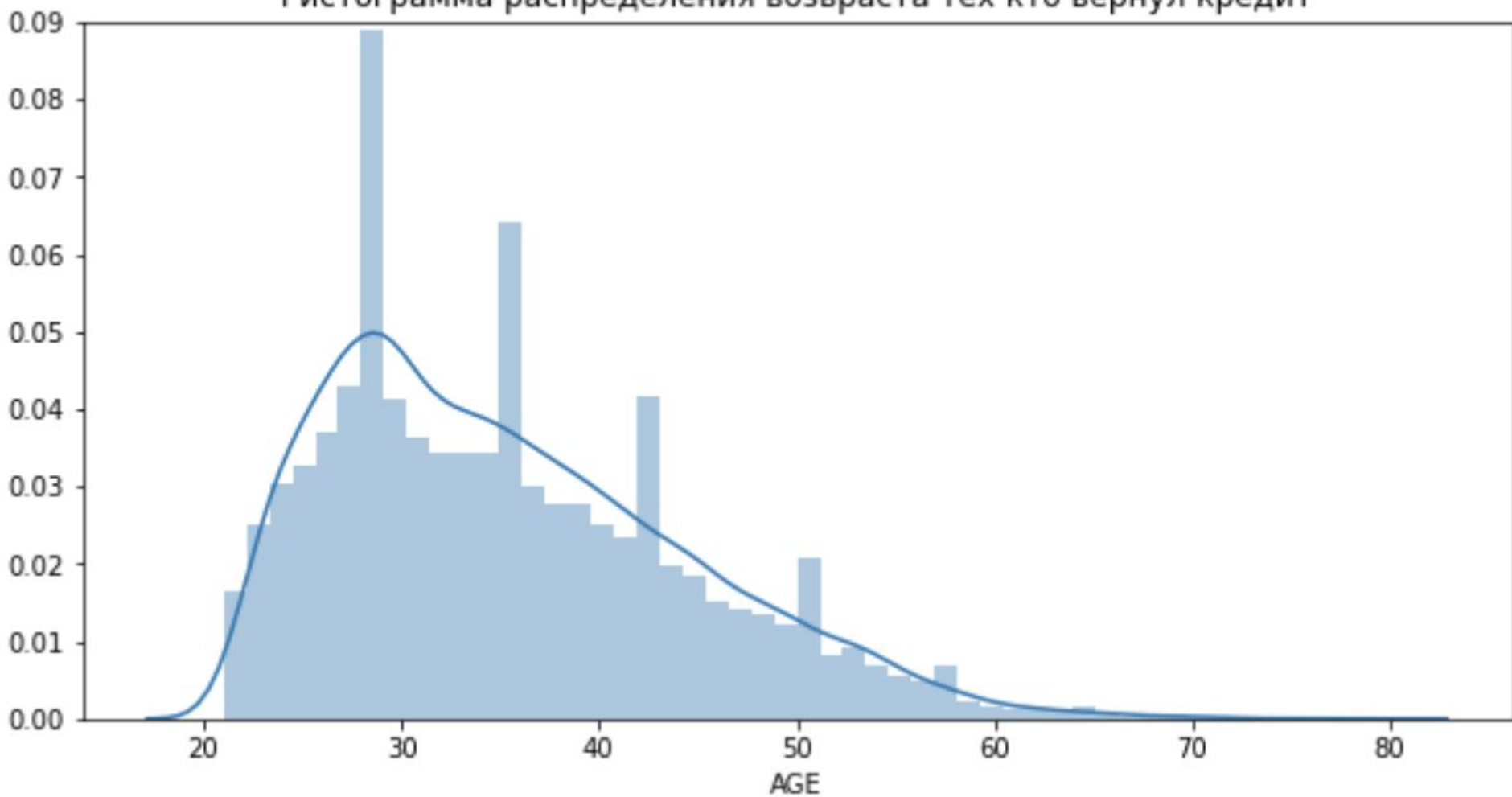
**Вывод:** Судя по коэффициенту Крамера линейной связи между полом и возвратом/невозвратом нет. Это соответствует здравому смыслу, а значит практическая значимость присутствует

# Сравнение данных о возрасте относительно возврата/невозврата

Гистограмма распределения возраста тех кто не вернул кредит



Гистограмма распределения возраста тех кто вернул кредит



# Использованные техники и Вывод

— — —

Интервальная оценка, доверительные интервалы, проверка гипотезы о равенстве распределений при помощи непараметрического перестановочного критерия проверки средних.

**Вывод:** Распределения выборок отличаются, возраст влияет на вероятность возврата/невозврата



# Задача #2. Анализ эффективности удержания

state — штат США

account\_length — длительность использования аккаунта

— — —

area\_code — деление пользователей на псевдорегионы, используемое в телекоме

intl\_plan — подключена ли у пользователя услуга международного общения

vmail\_plan — подключена ли у пользователя услуга голосовых сообщений

vmail\_message — количество голосых сообщений, который пользователь отправил / принял

day\_calls — сколько пользователь совершил дневных звонков

day\_mins — сколько пользователь проговорил минут в течение дня

day\_charge — сколько пользователь заплатил за свою дневную активность

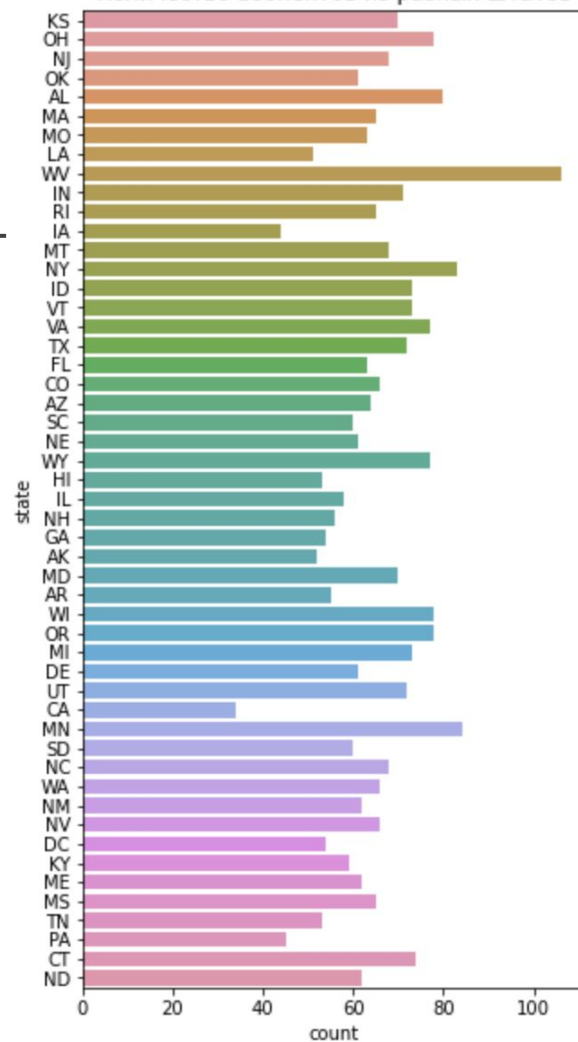
custserv\_calls — сколько раз пользователь позвонил в службу поддержки

treatment — номер стратегии, которая применялись для удержания абонентов (0, 2 = два разных типа воздействия, 1 = контрольная группа)

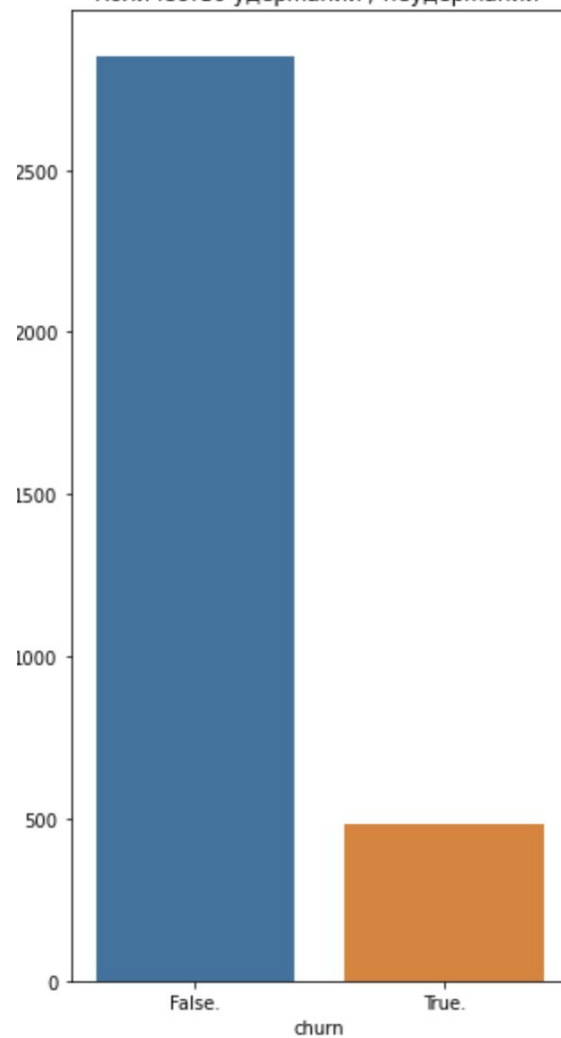
mes\_estim — оценка интенсивности пользования интернет мессенджерами

churn — результат оттока: перестал ли абонент пользоваться услугами оператора

Количество абонентов из разных штатов



Количество удержаний / неудержаний



# Проверка гипотезы о том, что штат абонента не влияет на то, перестанет ли абонент пользоваться услугами оператора

— — —

Проверка проводилась с помощью критерия  $\chi^2$  квадрат для таблицы сопряженности, а в дальнейшем была использована поправка Йетса (коррекция точности аппроксимации на непрерывность)

**Вывод:** до проверки в 34 из 1241 пар штатов последний влиял на результат, однако после поправки 0/1241

# Корреляция между дневными звонками и интенсивностью использования мессенджера

— — —

Для проверки величины линейной взаимосвязи между данными признаками были использованы коэффициенты корреляции Спирмена и Мэтьюса

**Вывод:** Значение коэффициентов корреляции не превышают 0.013 , а значит можно сказать об отсутствии линейной взаимосвязи

# Сравнение методов удержания и контрольной группы

— — —

Для сравнения методов удержания были построены доверительные интервалы Уилсона для каждой пары групп. Далее также 3 пары доверительных интервалов для разности долей. Наконец с помощью  $p$ -value из проверки на разности долей и поправки на множественную проверку гипотез Холма-Бонфферони был сделан окончательный вывод

**Вывод:** На основании данных можно сказать об эффективности метода “2”

# Задача #3 Анализ Результатов АБ теста

— — —

- userID: уникальный идентификатор пользователя
- browser: браузер, который использовал userID
- slot: в каком статусе пользователь участвовал в исследовании (exp = видел измененную страницу, control = видел неизменную страницу)
- n\_clicks: количество кликов, которые пользователь совершил за n\_queries
- n\_queries: количество запросов, который совершил userID, пользуясь браузером browser
- n\_nonclk\_queries: количество запросов пользователя, в которых им не было совершено ни одного клика

# Сравнение среднего количества кликов в контрольной и экспериментальной группах

— — —

Сравнение проводилось с помощью средних значений, доверительных интервалов и их разностей

**Вывод:** Кликов больше было в экспериментальной группе

# Сравнение доли кликов в обеих группах по разным браузерам

— — —

Сравнение проводилось с помощью непараметрического критерия сравнения независимых выборок Манна–Уитни и поправки Холма на множественную проверку гипотез.

**Вывод:** Тестируемое значение можно предложить внедрить для той доли пользователей, для которых улучшение показало себя практически и теоретически значимо