



**ОБРАЗОВАТЕЛЬНЫЙ  
ЦЕНТР** МГТУ им. Н. Э. Баумана

# Data Science

## Описательная статистика

Корпоративное обучение на базе  
Образовательного центра МГТУ им. Н. Э. Баумана  
под управлением МИЦ «Композиты России»  
Докладчик: Панфилов И.А. канд. техн. наук, доцент

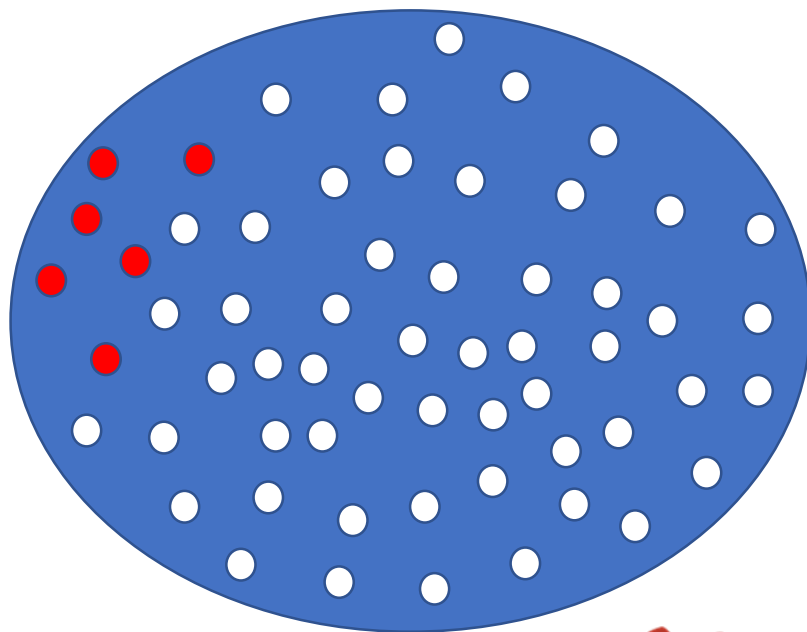
# Генеральная совокупность объектов и выборка

Генеральная совокупность:  
ВСЕ автомобили в городе

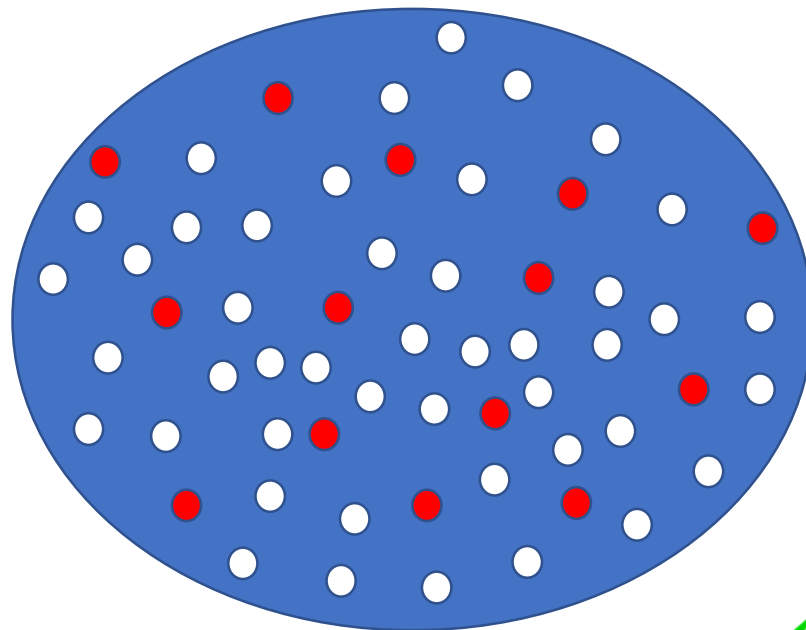
Выборка:  
5% от ВСЕХ автомобилей

# Репрезентативная выборка

простая случайная выборка (simple random sample)

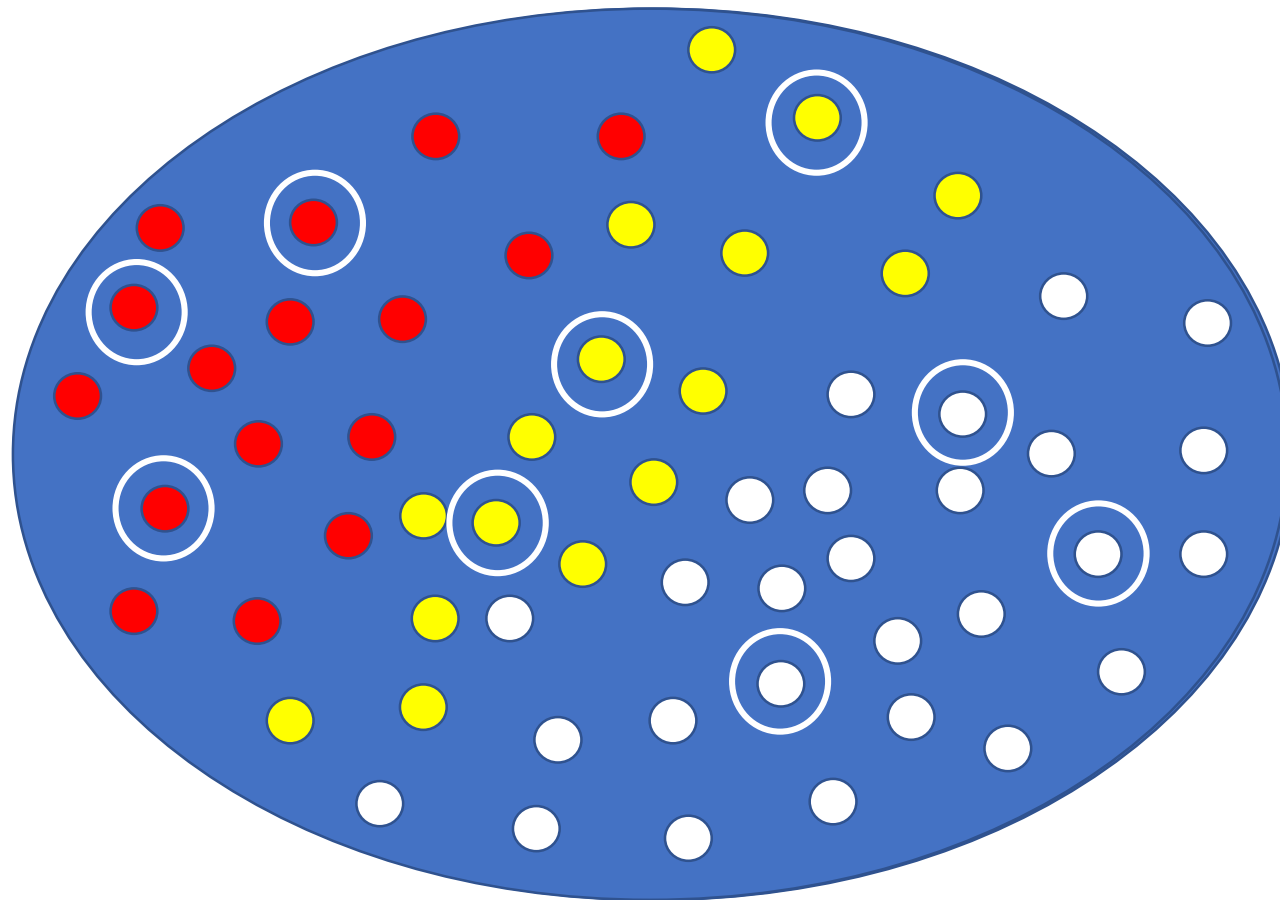


✗



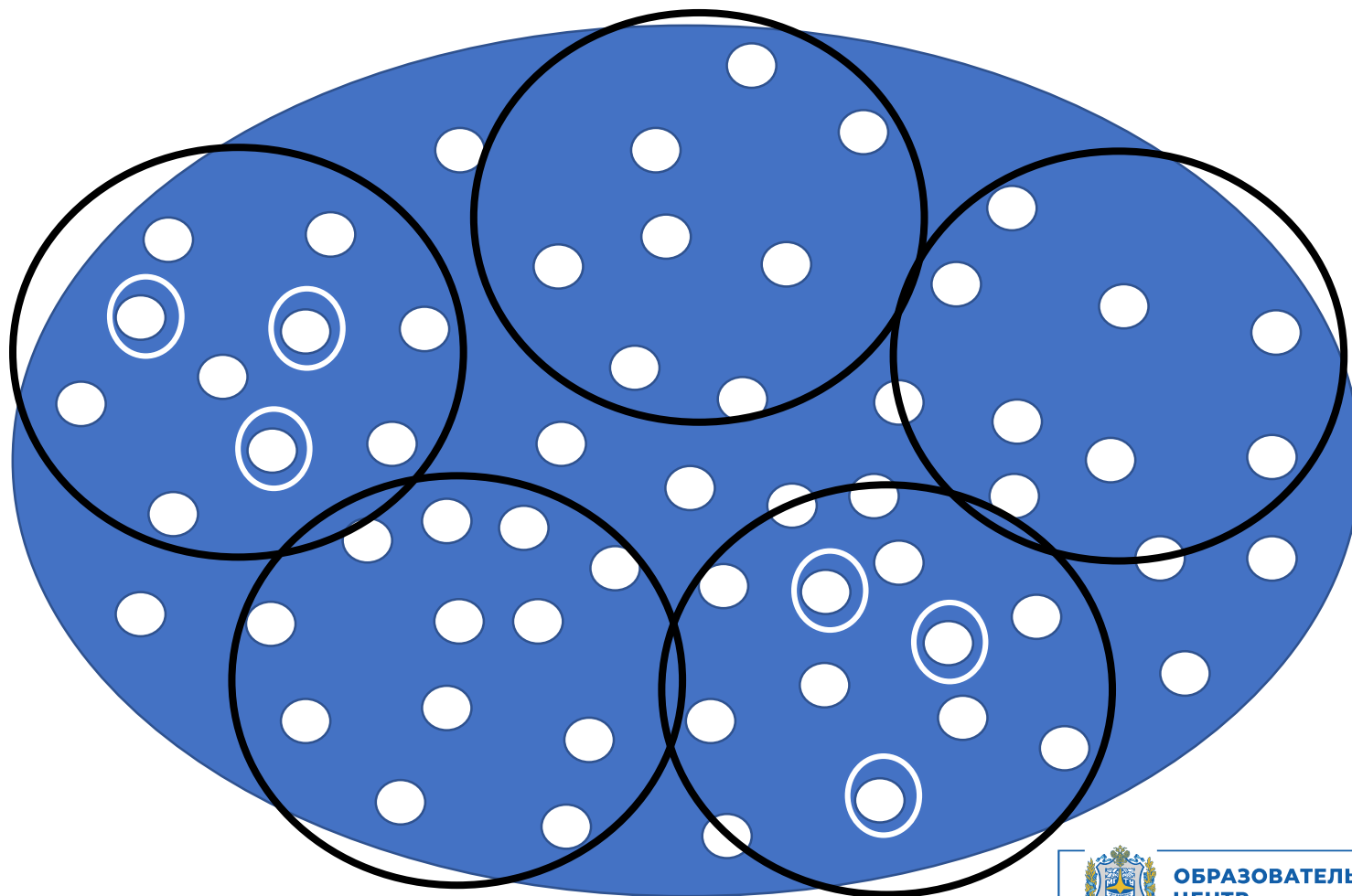
# Стратифицированная выборка

(stratified sample)

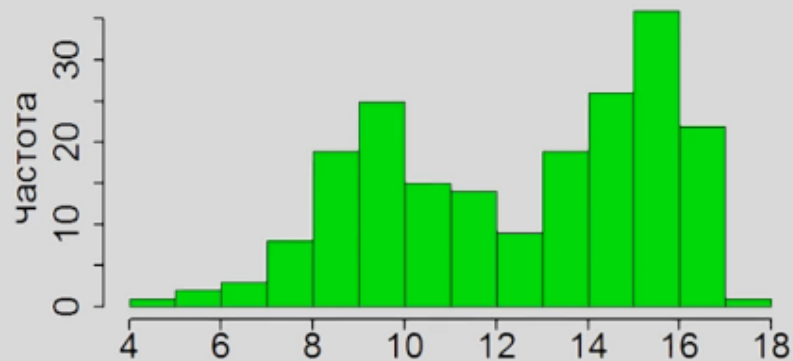
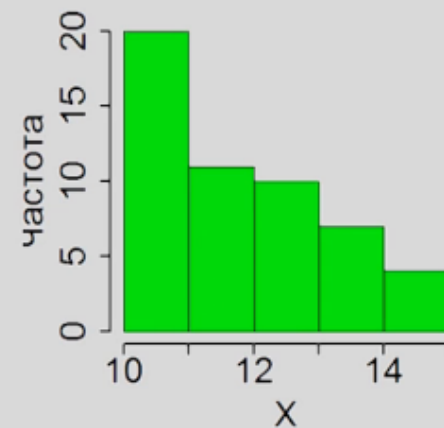
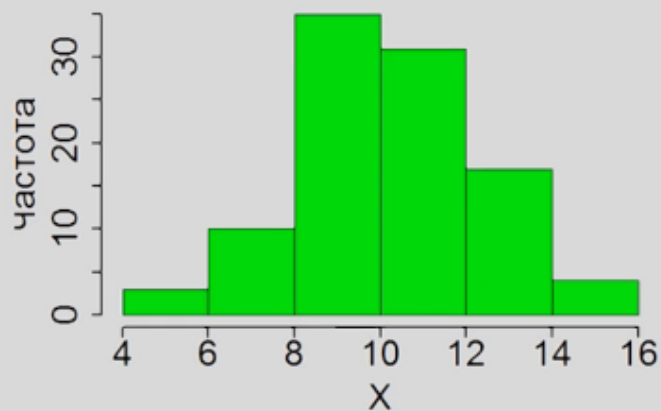


# Групповая выборка

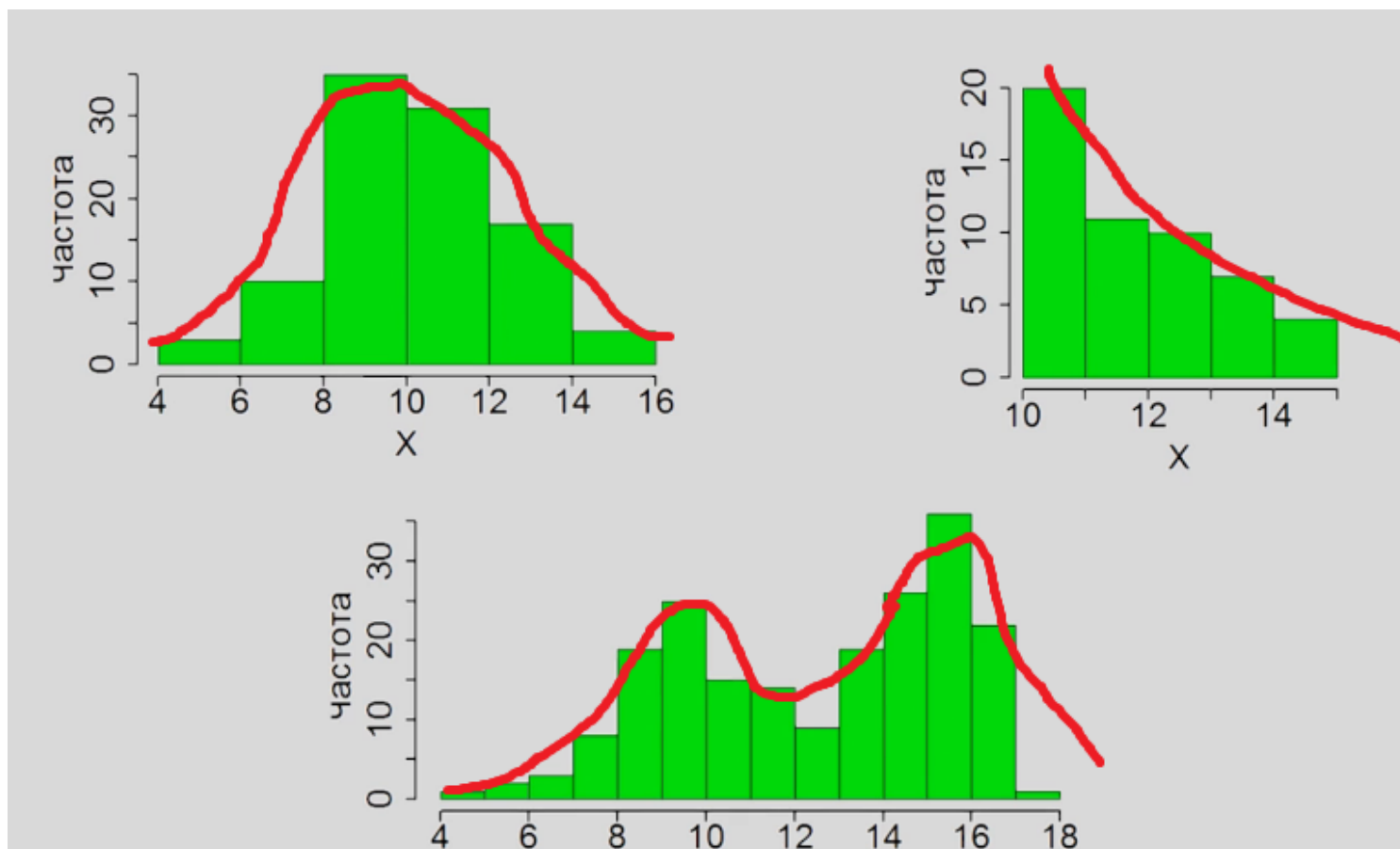
cluster sample



# Гистограмма частот



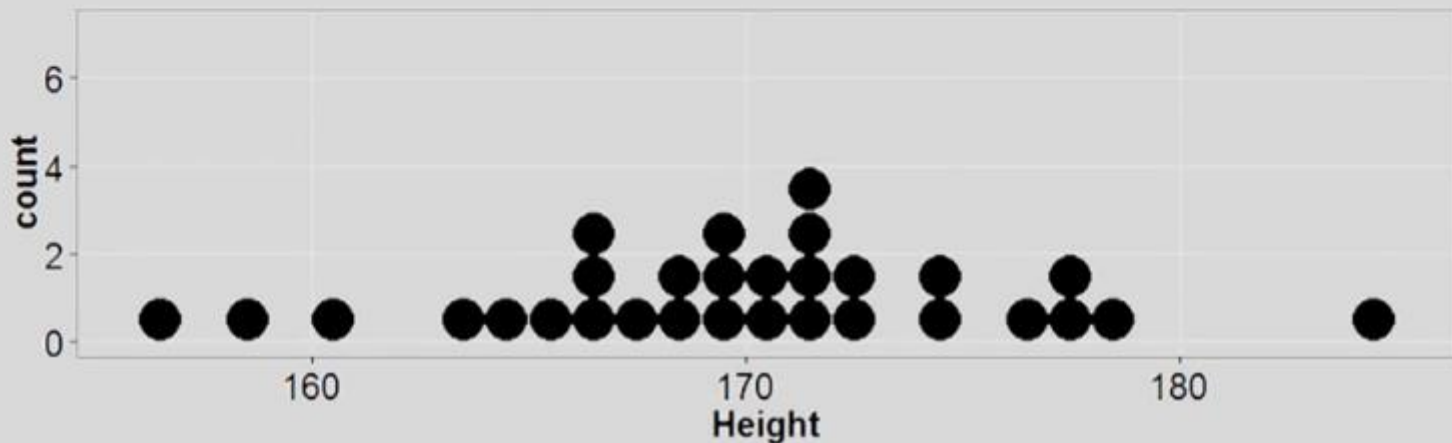
# Характер распределения



# Меры центральной тенденции

**Мода** (Mode) – значение измеряемого признака, которое встречается максимально часто.

185 175 170 169 171 172 175 157 170 172 167 173 168 167 166  
167 169 172 177 178 165 161 179 159 164 178 172 170 173 171





# Медиана

Медиана (median) – значение признака, которое делит упорядоченное множество данных пополам.

157 159 161 164 165 166 167 167 167

157 159 161 164 165 166 167 167 167 168 169 169 170 170 170

171 171 172 172 172 172 173 173 175 175 177 178 178 179 185

$$M_c = \frac{170 + 171}{2} = 170,5$$

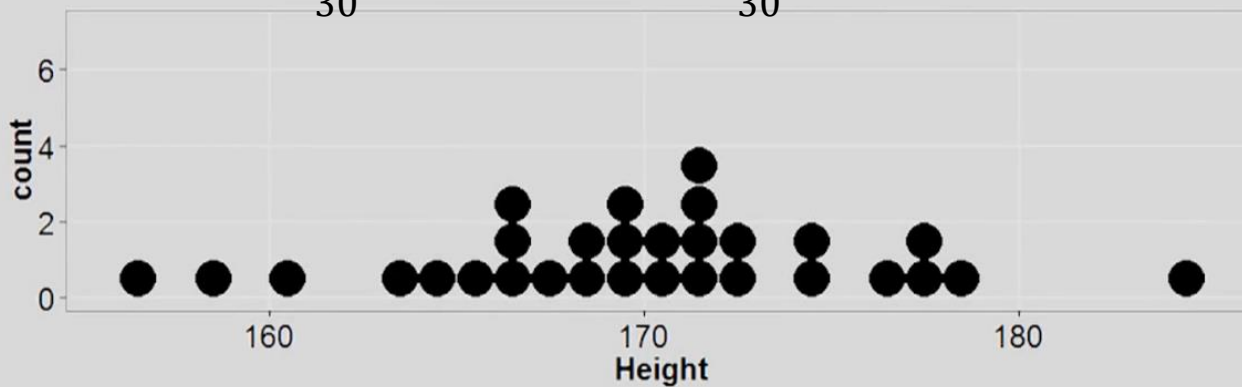


# Среднее значение

**Среднее значение** (mean, среднее арифметическое)  
сумма всех значений измеренного признака, деленная  
на количество измеренных значений.

185 175 170 169 171 172 175 157 170 172 167 173 168 167 166  
167 169 172 177 178 165 161 179 159 164 178 172 170 173 171

$$\bar{x} = \frac{x_1 + x_1 + \dots + x_1}{30} = \frac{185 + 175 + 170 + \dots + 171}{30} = 170,4$$

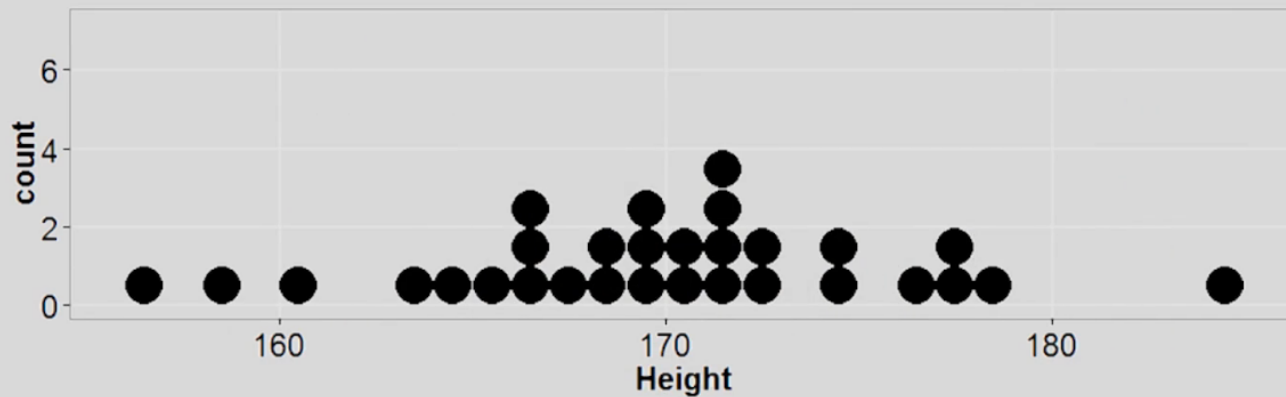


# Мера изменчивости

## размах (range)

**Размах (Range)** - разность максимального и минимального значения.

$$R = X_{\max} - X_{\min}$$

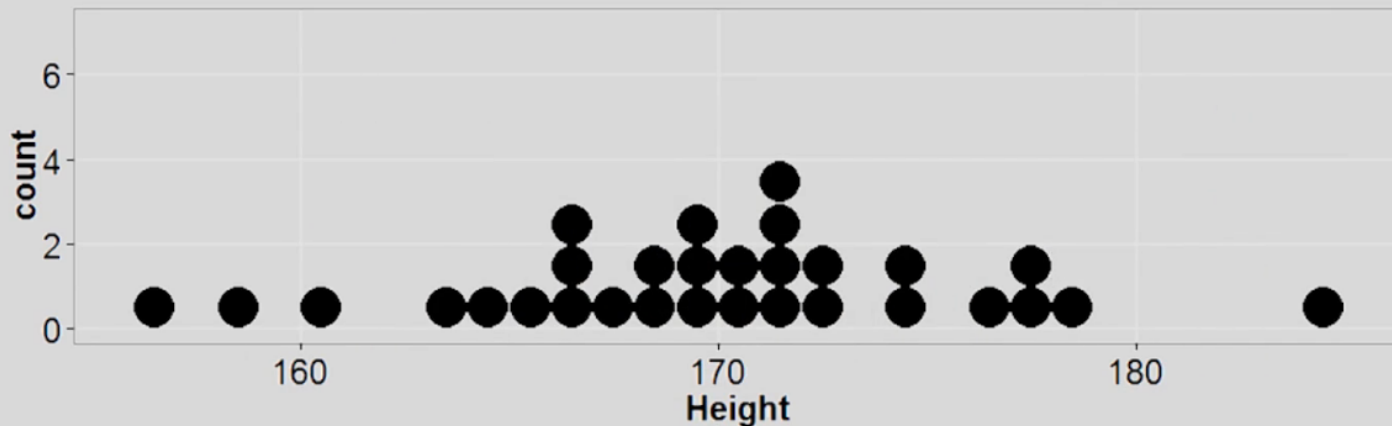


# Мера изменчивости

дисперсия (variance),  
среднеквадратическое отклонение (standard deviation)

**Дисперсия** (variance) – средний квадрат отклонений индивидуальных значений признака от их средней величины.

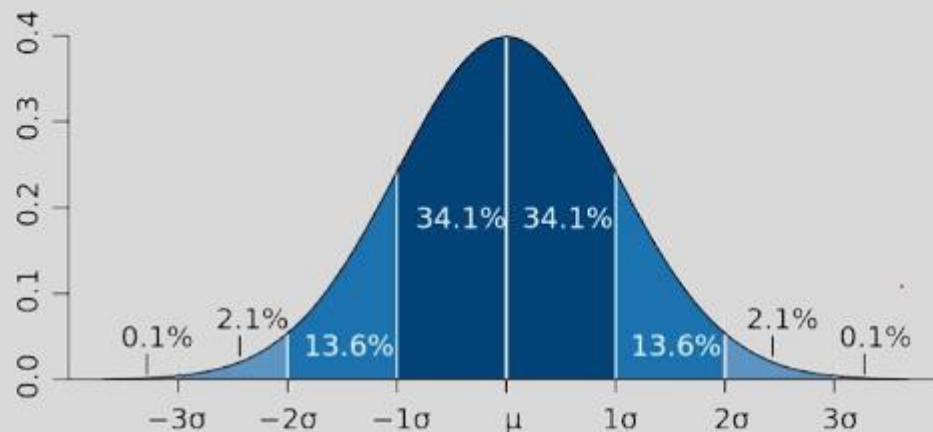
$$D = \frac{\sum (x_i - \bar{x})^2}{n} \quad \delta = \sqrt{D} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$



# Нормальное распределение

- Унимодально

- Симметрично

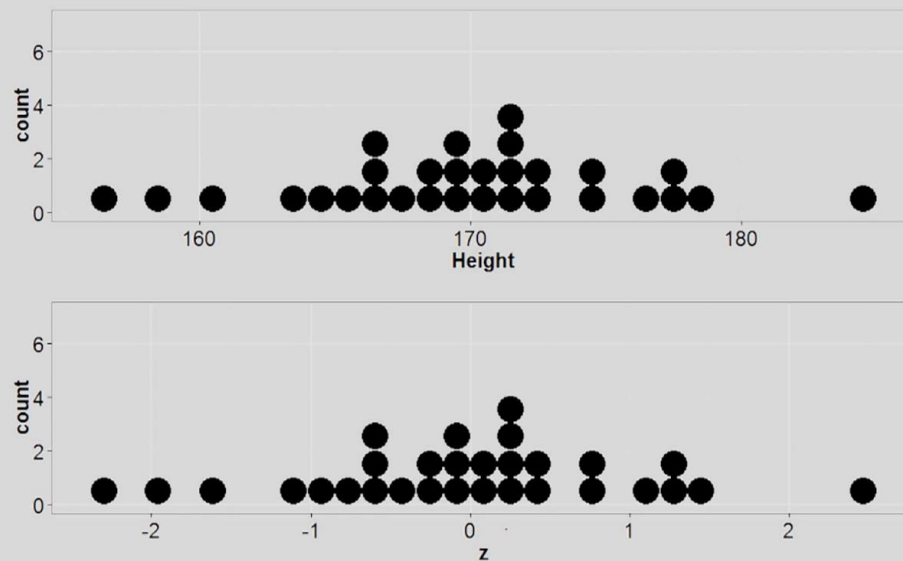


- Отклонения наблюдений от среднего подчиняются определенному вероятностному закону.

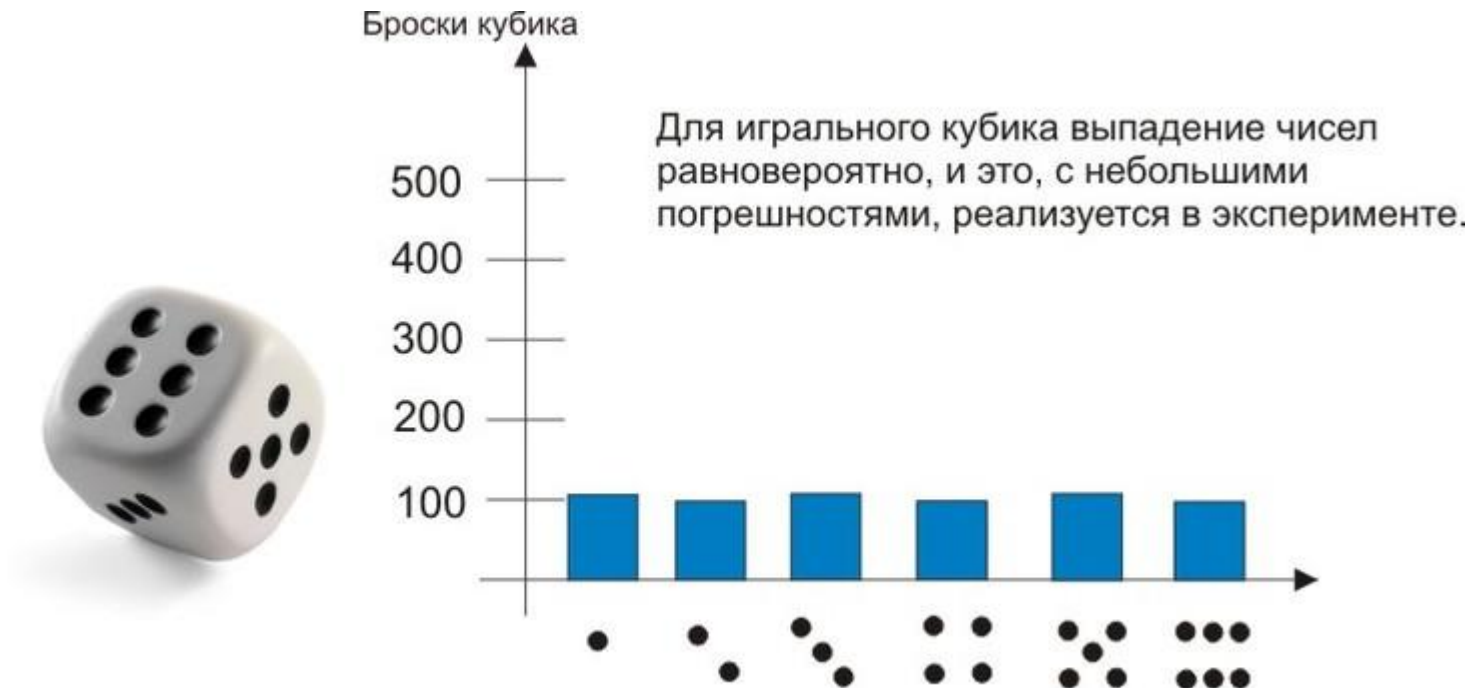
# Нормирование

Стандартизация или *z-преобразование* – преобразование полученных данных в стандартную Z-шкалу (Z-scores) со средним  $M_z = 0$  и  $D_z = 1$

$$Z_i = \frac{x_i - \bar{x}}{\delta_x}$$



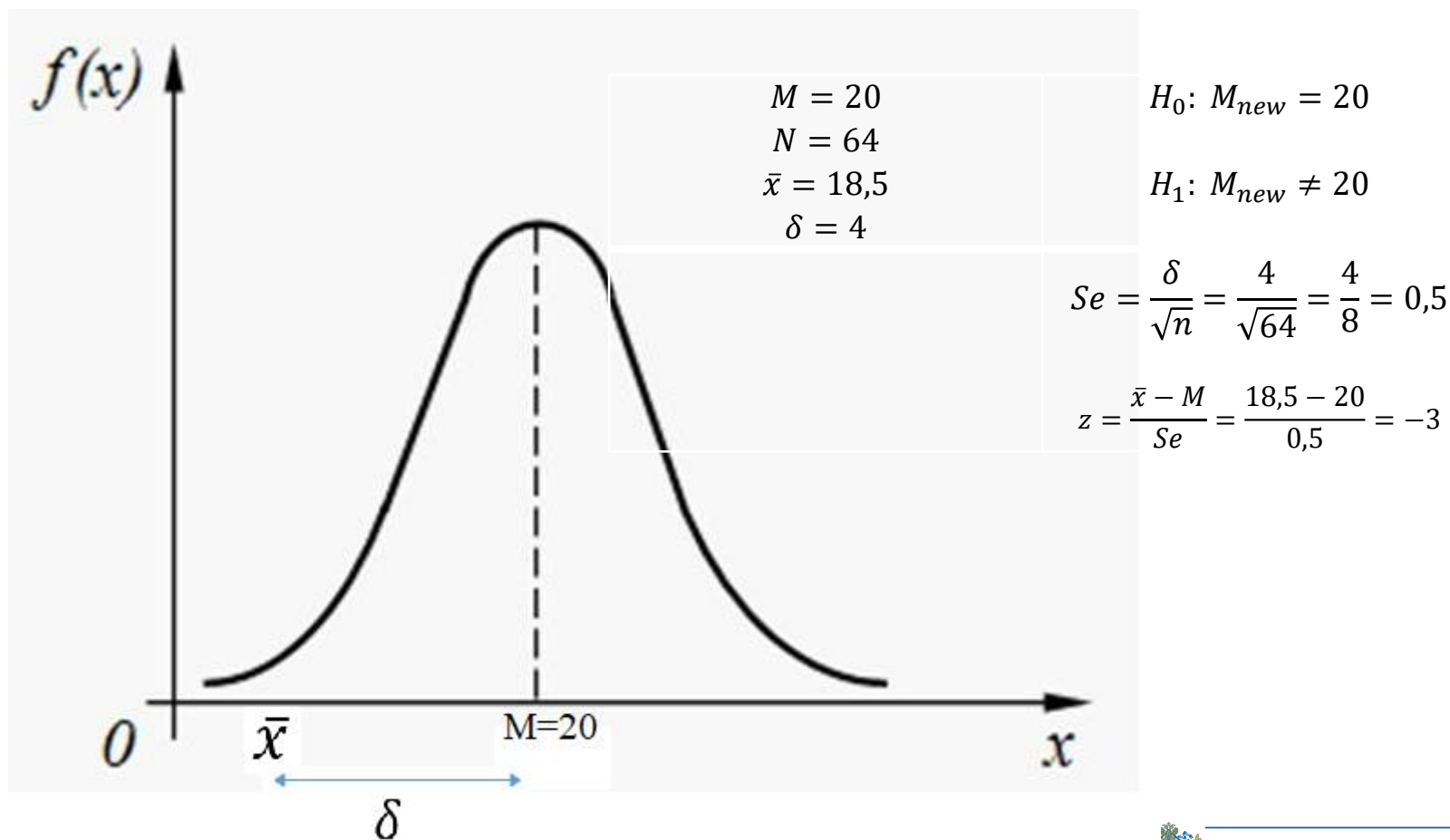
# Равномерное распределение







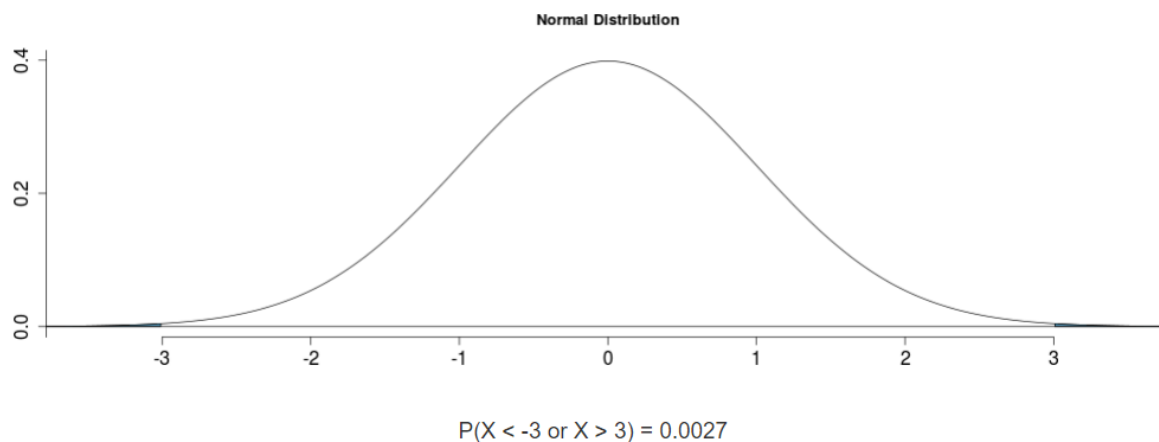
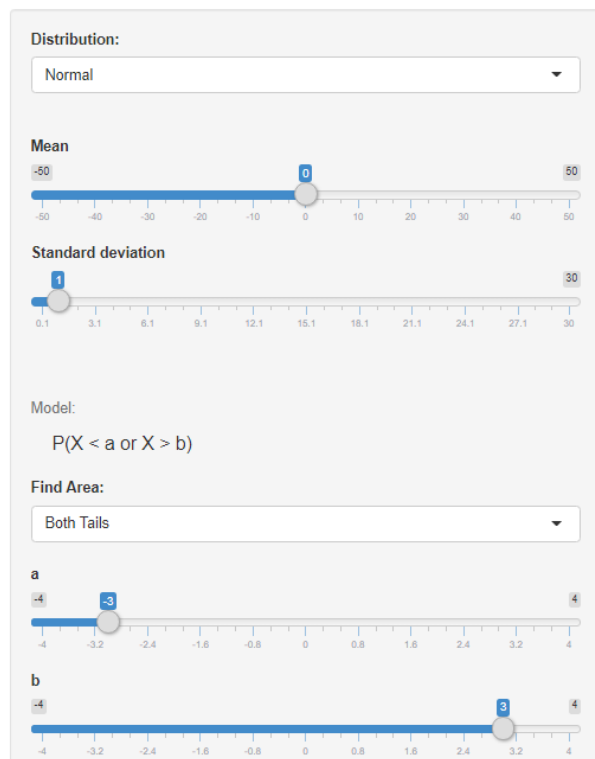
# Статистическая проверка гипотез



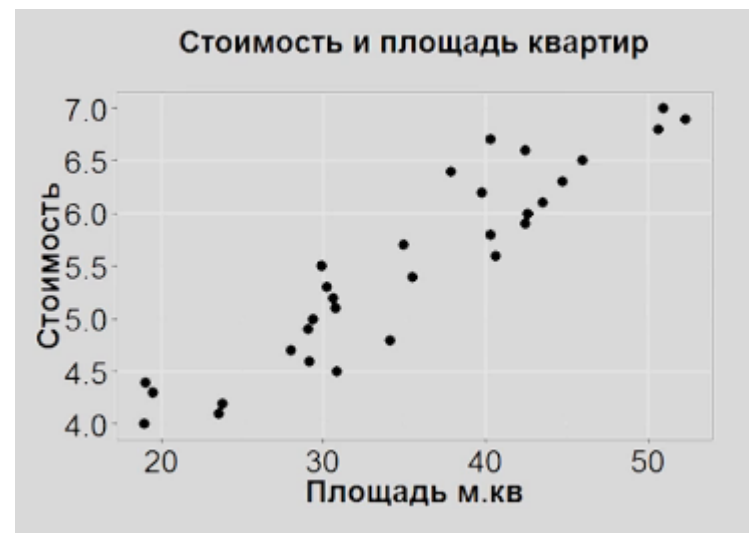
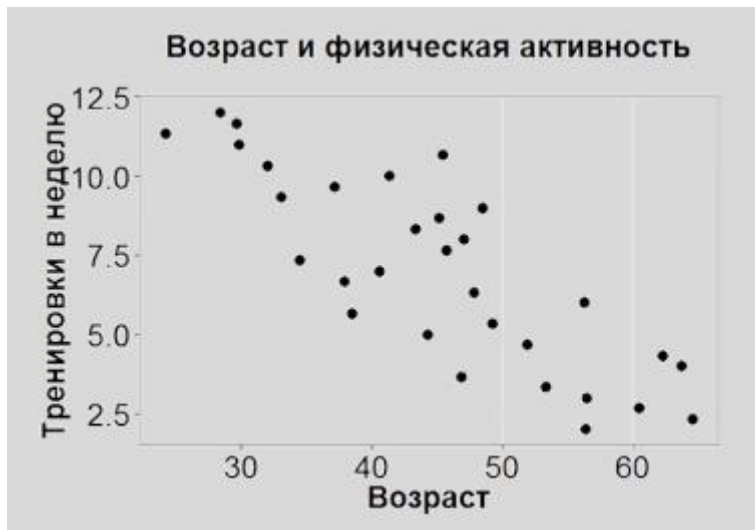
# Статистическая проверка гипотез

[https://gallery.shinyapps.io/dist\\_calc/](https://gallery.shinyapps.io/dist_calc/)

## Distribution Calculator

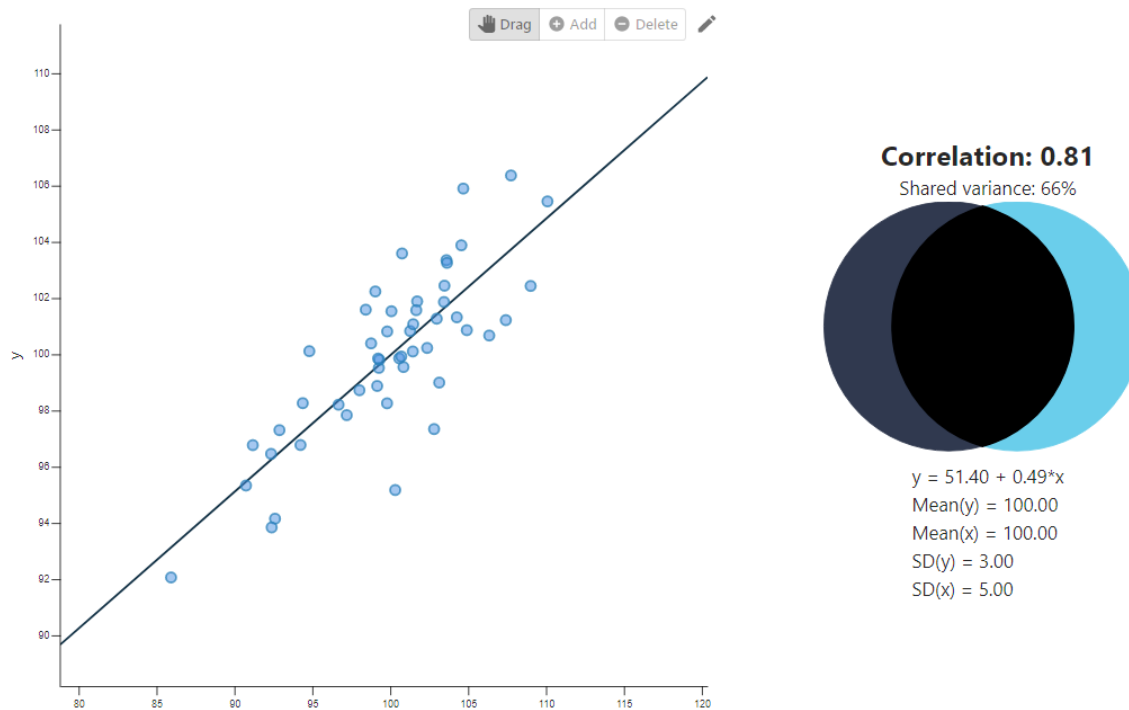


# Коэффициент корреляции



# Коэффициент корреляции

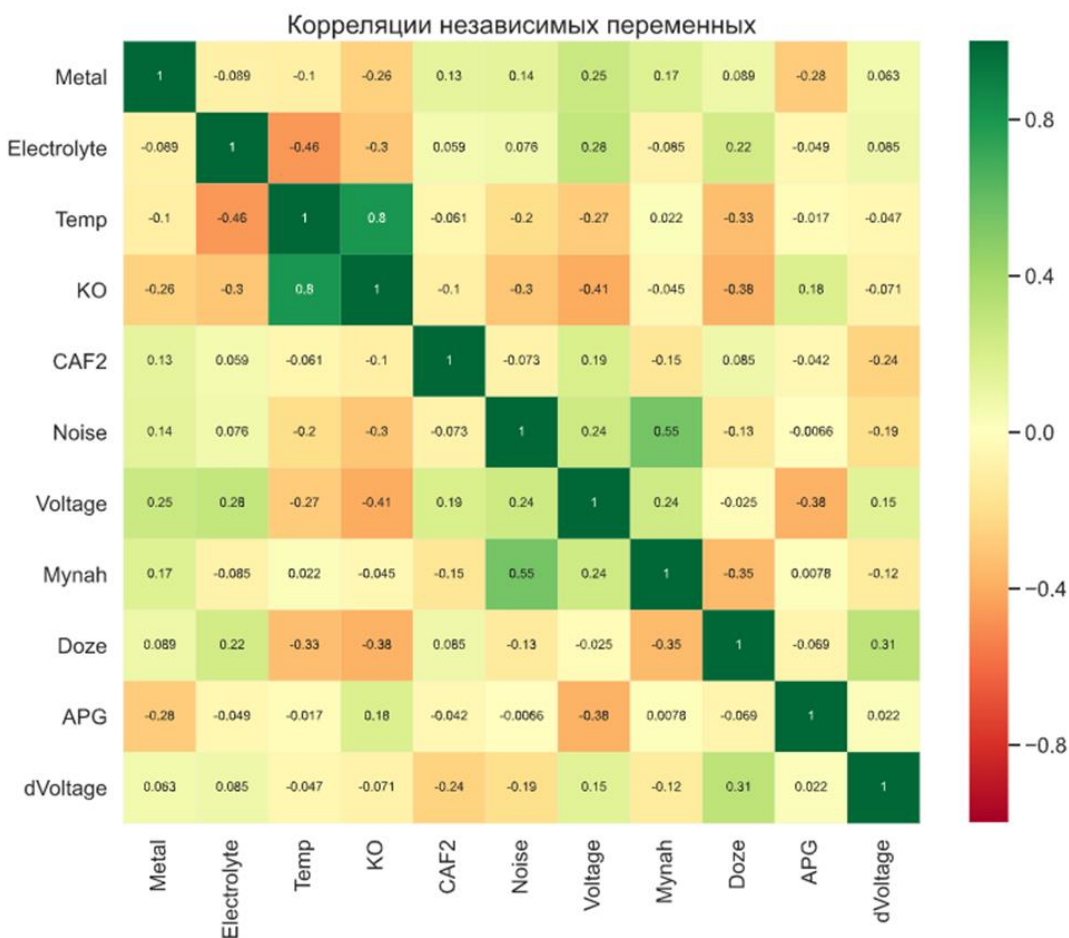
<https://rpsychologist.com/correlation/>



$$r_{xy} = \frac{\sum (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\sum (x_i - \bar{X})^2 \sum (y_i - \bar{Y})^2}}$$

# Коэффициент корреляции

## Тепловая карта





**edu.bmstu.ru**

**+7 (495) 120-30-75**

**E-mail: edu@bmstu.ru**

**Москва, ул. 2-я Бауманская,  
дом 5, стр. 1**