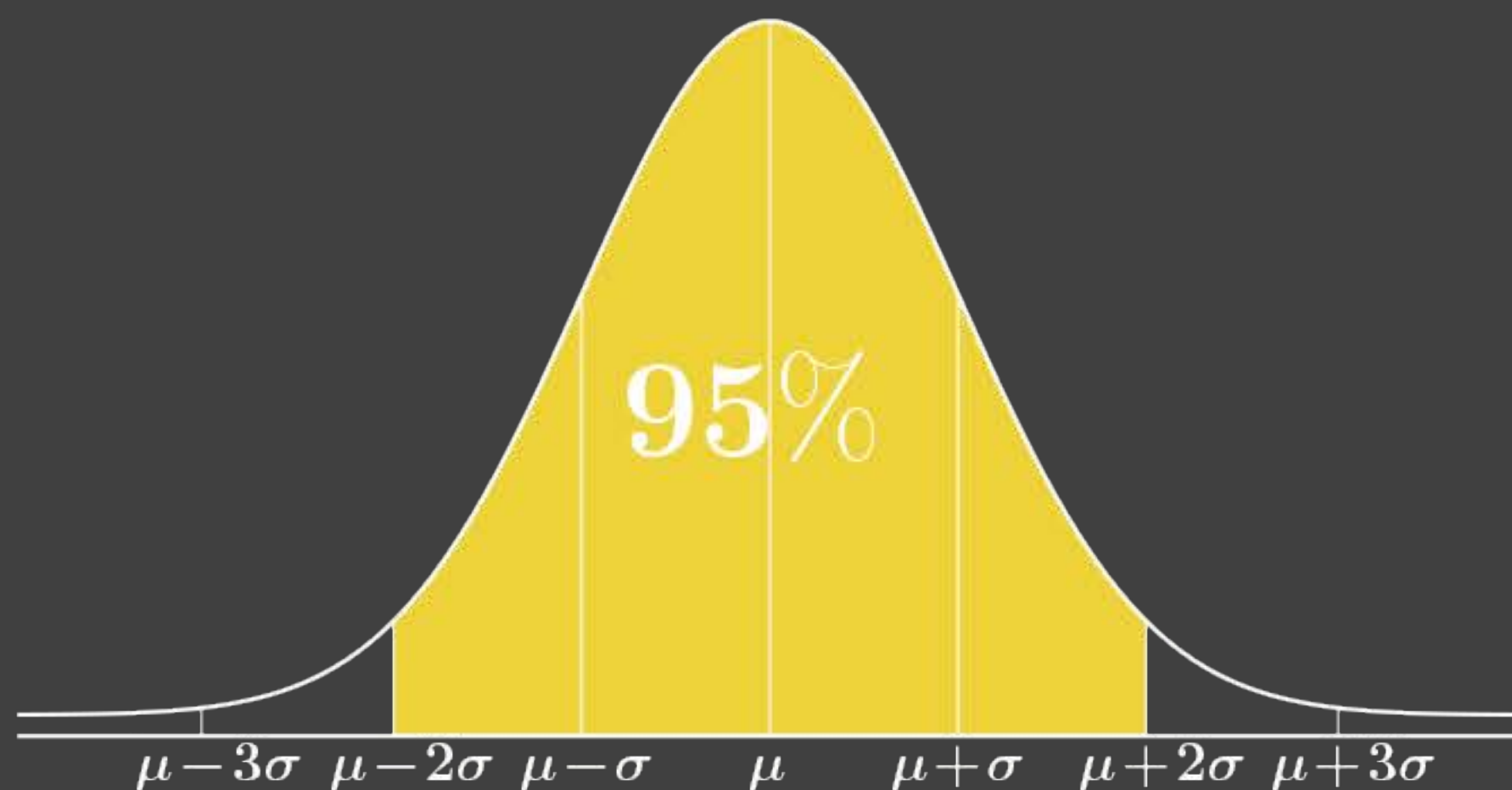


ИНТЕРВАЛЬНЫЕ ОЦЕНКИ

ПРАВИЛО ДВУХ СИГМ

$$X \sim N(\mu, \sigma^2) \Rightarrow P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 0.95$$



Можно получить более точную оценку для любой вероятности

› Квантиль порядка $\alpha \in (0, 1)$ — такое число X_α , что:

$$P(X \leq X_\alpha) \geq \alpha$$

$$P(X \geq X_\alpha) \geq 1 - \alpha$$

› Функция распределения X :

$$F(x) = P(X \leq x) \Rightarrow$$

› Эквивалентное определение квантиля:

$$X_\alpha = F^{-1}(\alpha) = \inf\{x: F(x) \geq \alpha\}$$

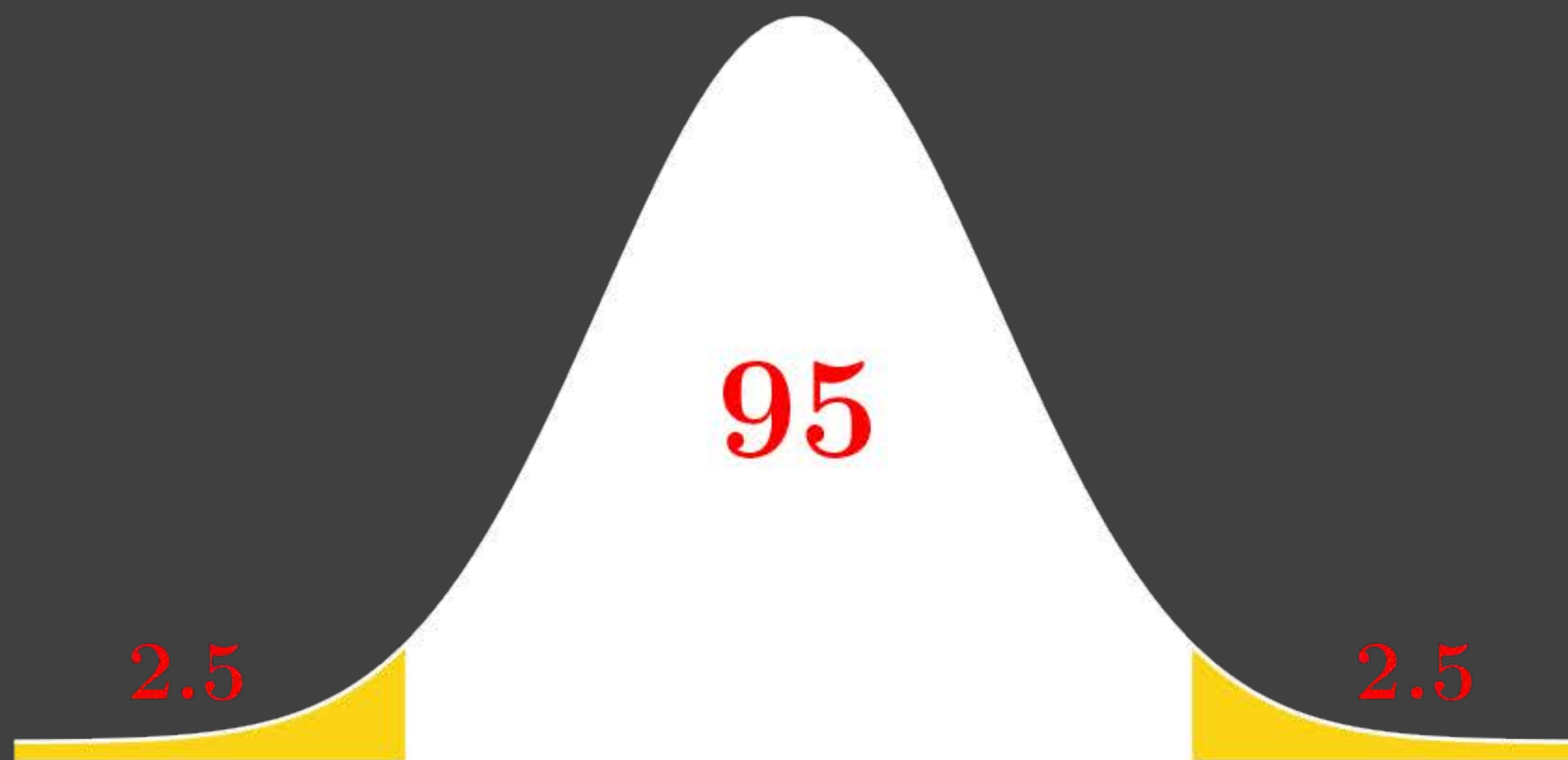
УТОЧНЕНИЕ ПРАВИЛА ДВУХ СИГМ



$$P(? \leq X \leq ?) = 0.95$$

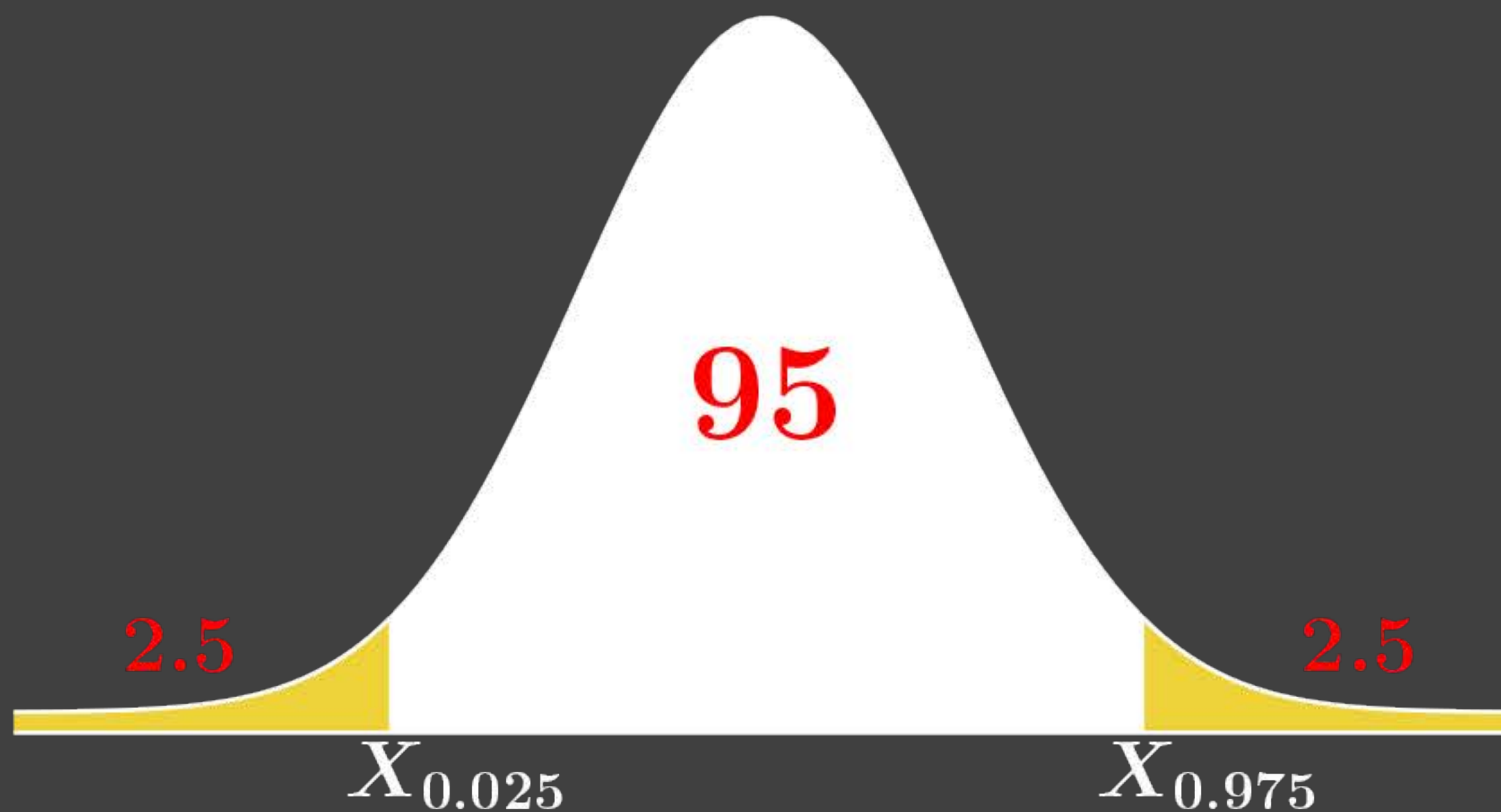
УТОЧНЕНИЕ ПРАВИЛА ДВУХ СИГМ

$$P(? \leq X \leq ?) = 0.95$$



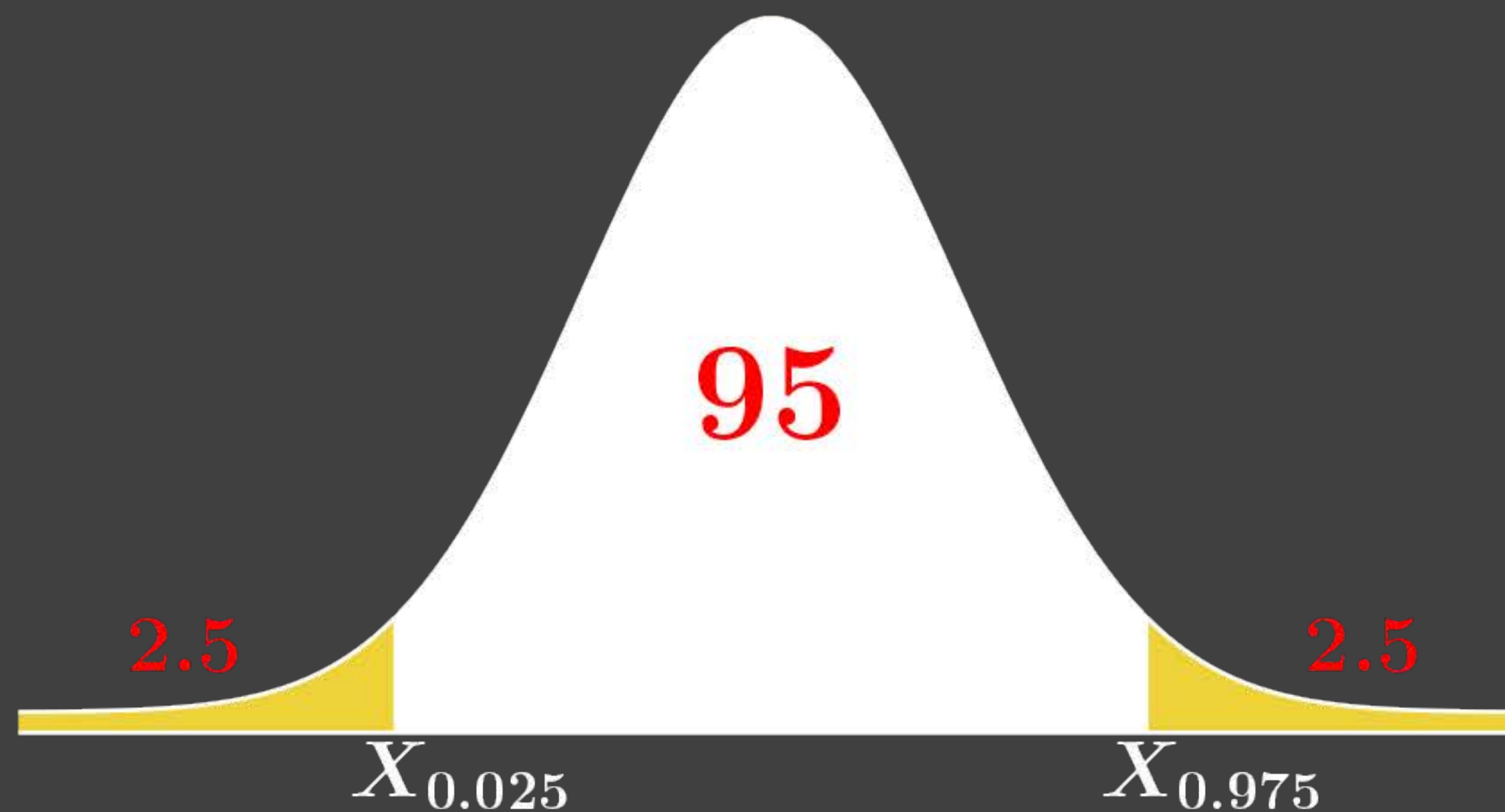
УТОЧНЕНИЕ ПРАВИЛА ДВУХ СИГМ

$$P(? \leq X \leq ?) = 0.95$$



УТОЧНЕНИЕ ПРАВИЛА ДВУХ СИГМ

$$P(X_{0.025} \leq X \leq X_{0.975}) = 0.95$$



ПРЕДСКАЗАТЕЛЬНЫЙ ИНТЕРВАЛ



$$\triangleright X \sim F(x) \Rightarrow P\left(X_{\frac{\alpha}{2}} \leq X \leq X_{1-\frac{\alpha}{2}}\right) = 1 - \alpha$$

$\left[X_{\frac{\alpha}{2}}, X_{1-\frac{\alpha}{2}}\right]$ — предсказательный интервал
порядка $1 - \alpha$

ПРЕДСКАЗАТЕЛЬНЫЙ ИНТЕРВАЛ

► $X \sim F(x) \Rightarrow P\left(X_{\frac{\alpha}{2}} \leq X \leq X_{1-\frac{\alpha}{2}}\right) = 1 - \alpha$

$\left[X_{\frac{\alpha}{2}}, X_{1-\frac{\alpha}{2}}\right]$ — предсказательный интервал
порядка $1 - \alpha$

► $X \sim N(\mu, \sigma^2) \Rightarrow$

$$P\left(\mu - z_{1-\frac{\alpha}{2}}\sigma \leq X \leq \mu + z_{1-\frac{\alpha}{2}}\sigma\right) = 1 - \alpha$$

z_{α} — квантиль стандартного нормального
распределения $N(0, 1)$

$$z_{0.975} \approx 1.95996 \approx 2$$

- › Использование квантилей для построения интервальных оценок
- › Предсказательный интервал
- › Уточнение правила двух сигм
- › Далее: доверительные интервалы

ДОВЕРИТЕЛЬНЫЕ ИНТЕРВАЛЫ

ТОЧЕЧНЫЕ ОЦЕНКИ



› $X \sim F(x, \theta)$, θ — неизвестный параметр

› $\theta = ?$

› $X \sim F(x, \theta)$, θ — неизвестный параметр

› $\theta = ?$

› $X^n = (X_1, \dots, X_n)$

› $\hat{\theta}$ — оценка θ по выборке

- › $X \sim F(x, \theta)$, θ — неизвестный параметр
- › $\theta = ?$
- › $X^n = (X_1, \dots, X_n)$
- › $\hat{\theta}$ — оценка θ по выборке

- › Например, для $\theta = \mathbb{E}X$:

$$\hat{\theta} = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \text{ — хорошая оценка}$$

- › Доверительный интервал для параметра θ — пара таких статистик C_L, C_U , что

$$P(C_L \leq \theta \leq C_U) \geq 1 - \alpha$$

- › Как оценить C_L и C_U по выборке?

› $P(C_L \leq \theta \leq C_U) \geq 1 - \alpha$

› Как оценить C_L и C_U по выборке?

› Если $\hat{\theta}$ — оценка θ и мы знаем её распределение $F_{\hat{\theta}}(x)$, то:

$$P\left(F_{\hat{\theta}}^{-1}\left(\frac{\alpha}{2}\right) \leq \theta \leq F_{\hat{\theta}}^{-1}\left(1 - \frac{\alpha}{2}\right)\right) = 1 - \alpha$$

ДЛЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ

› $X \sim N(\mu, \sigma^2)$, $X^n = (X_1, \dots, X_n)$

› \bar{X}_n — оценка $\mathbb{E}X = \mu$

ДЛЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ

› $X \sim N(\mu, \sigma^2), X^n = (X_1, \dots, X_n)$

› \bar{X}_n — оценка $\mathbb{E}X = \mu$

› $\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right) \Rightarrow$

ДЛЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ

› $\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right) \Rightarrow$

› Предсказательный интервал для \bar{X}_n :

$$P\left(\mu - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \bar{X}_n \leq \mu + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

ДЛЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ

› Предсказательный интервал для \bar{X}_n :

$$P\left(\mu - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \bar{X}_n \leq \mu + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

› Доверительный интервал для μ :

$$P\left(\bar{X}_n - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

ДОВЕРИТЕЛЬНЫЙ И ПРЕДСКАЗАТЕЛЬНЫЙ ИНТЕРВАЛЫ



› $X \sim N(\mu, \sigma^2), X^n = (X_1, \dots, X_n)$

›

ДОВЕРИТЕЛЬНЫЙ И ПРЕДСКАЗАТЕЛЬНЫЙ ИНТЕРВАЛЫ

› $X \sim N(\mu, \sigma^2)$, $X^n = (X_1, \dots, X_n)$

› Предсказательный интервал для X :

$$P\left(\mu - z_{1-\frac{\alpha}{2}}\sigma \leq X \leq \mu + z_{1-\frac{\alpha}{2}}\sigma\right) = 1 - \alpha$$

ДОВЕРИТЕЛЬНЫЙ И ПРЕДСКАЗАТЕЛЬНЫЙ ИНТЕРВАЛЫ

› $X \sim N(\mu, \sigma^2)$, $X^n = (X_1, \dots, X_n)$

› Предсказательный интервал для X :

$$P(\bar{X}_n - z_{1-\frac{\alpha}{2}}\sigma \leq X \leq \bar{X}_n + z_{1-\frac{\alpha}{2}}\sigma) = 1 - \alpha$$

ДОВЕРИТЕЛЬНЫЙ И ПРЕДСКАЗАТЕЛЬНЫЙ ИНТЕРВАЛЫ

› Предсказательный интервал для X :

$$P\left(\bar{X}_n - z_{1-\frac{\alpha}{2}}\sigma \leq X \leq \bar{X}_n + z_{1-\frac{\alpha}{2}}\sigma\right) = 1 - \alpha$$

› Доверительный интервал для μ :

$$P\left(\bar{X}_n - z_{1-\frac{\alpha}{2}}\frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{1-\frac{\alpha}{2}}\frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

ДЛЯ ЛЮБОГО РАСПРЕДЕЛЕНИЯ

› $X \sim F(x), X^n = (X_1, \dots, X_n)$

› \bar{X}_n — оценка $\mathbb{E}X$

› $\bar{X}_n \approx \sim N\left(\mathbb{E}X, \frac{\mathbb{D}X}{n}\right)$ (ЦПТ) \Rightarrow

ДЛЯ ЛЮБОГО РАСПРЕДЕЛЕНИЯ

› $\bar{X}_n \approx \sim N\left(\mathbb{E}X, \frac{\mathbb{D}X}{n}\right)$ (ЦПТ) \Rightarrow

› Доверительный интервал для $\mathbb{E}X$:

$$\mathbb{P}\left(\bar{X}_n - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\mathbb{D}X}{n}} \leq \mathbb{E}X \leq \bar{X}_n + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\mathbb{D}X}{n}}\right) \approx 1 - \alpha$$

- › Построение доверительных интервалов
- › Разница между доверительными и предсказательными интервалами
- › Интервалы для нормального распределения

- › Разница между доверительными и предсказательными интервалами
- › Интервалы для нормального распределения
- › Далее: распределения, производные от нормального

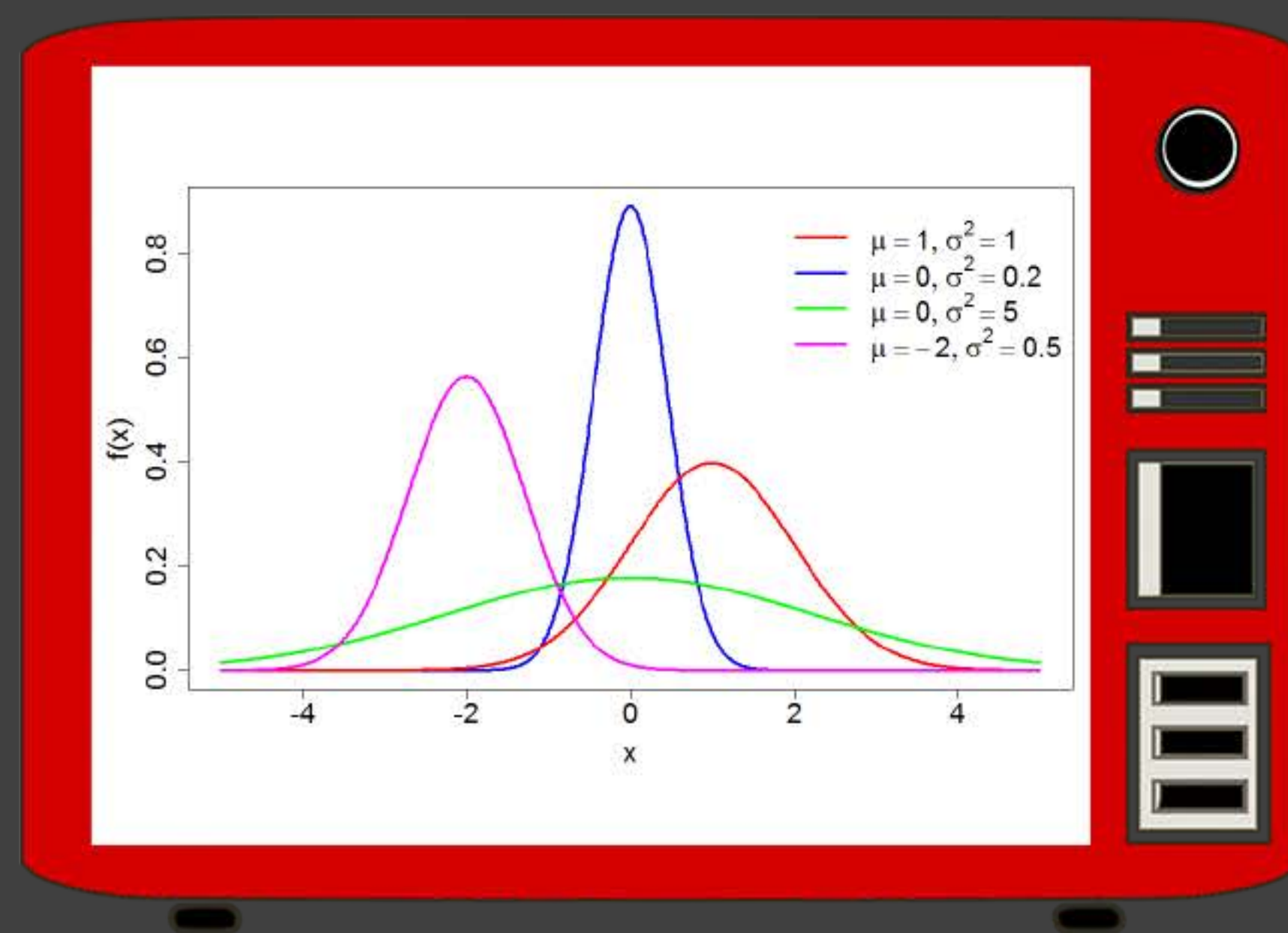
РАСПРЕДЕЛЕНИЯ, ПРОИЗВОДНЫЕ ОТ НОРМАЛЬНОГО

НОРМАЛЬНОЕ РАСПРЕДЕЛЕНИЕ

$$\triangleright X \sim N(\mu, \sigma^2)$$

$$\triangleright f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

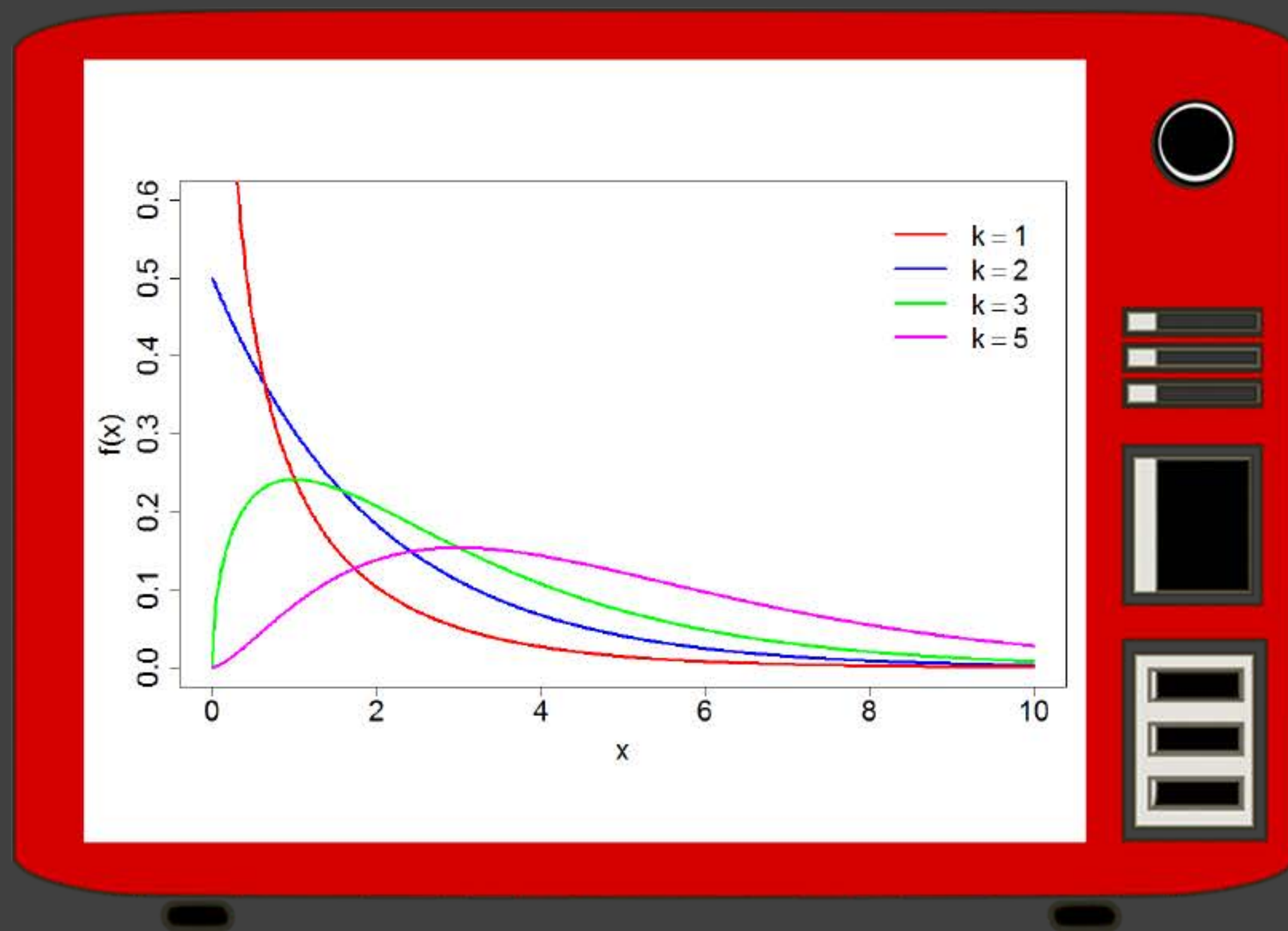
$$\triangleright F(x) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt$$



РАСПРЕДЕЛЕНИЕ ХИ-КВАДРАТ

› $X_1, X_2, \dots, X_k \sim N(0, 1)$ независимы

› $X = \sum_{i=1}^k X_i^2 \sim \chi_k^2$ — распределение хи-квадрат
с k степенями свободы



РАСПРЕДЕЛЕНИЕ СТЬЮДЕНТА

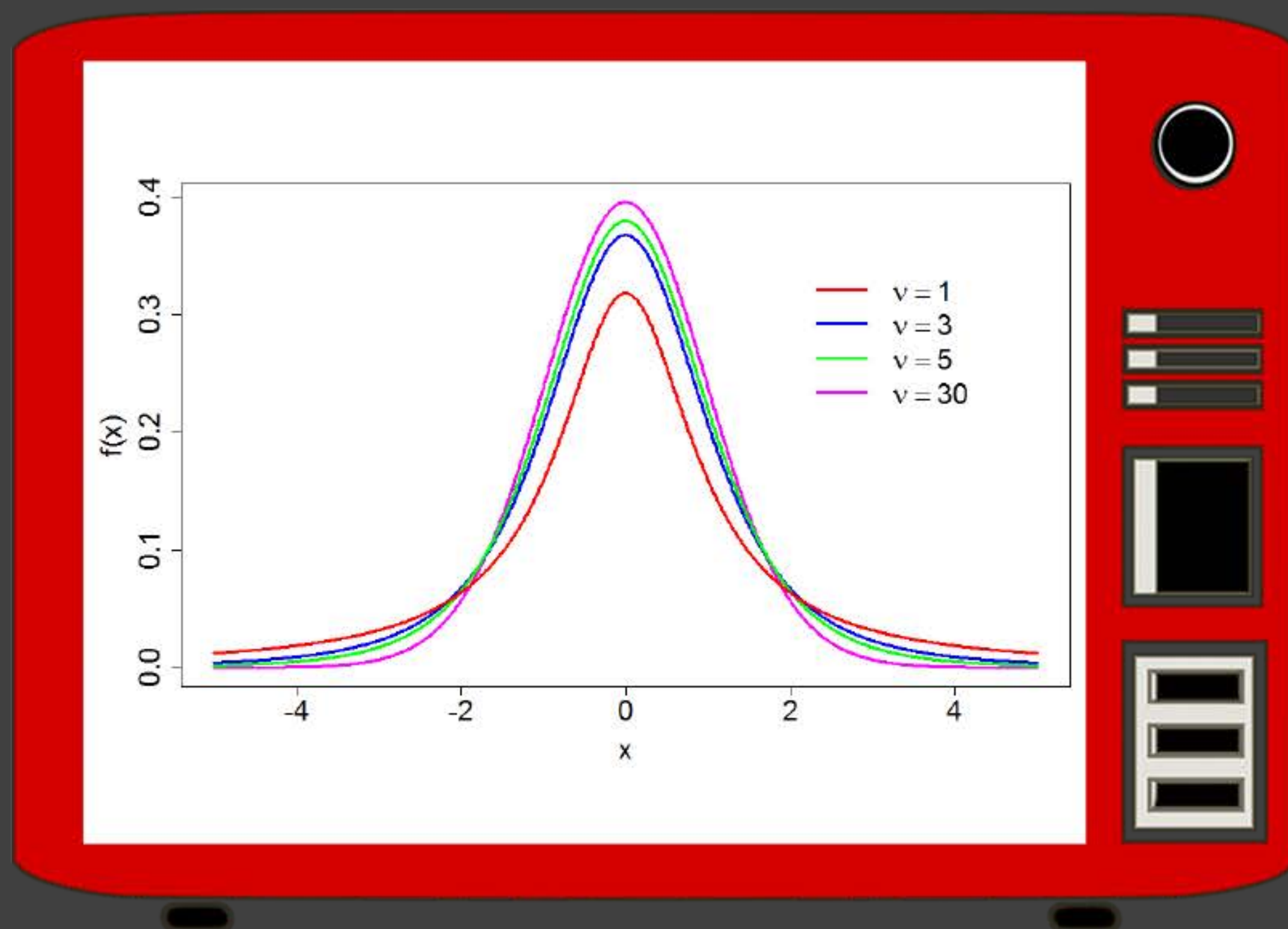
› $X_1 \sim N(0, 1)$, $X_2 \sim \chi^2_\nu$ независимы

› $X = \frac{X_1}{\sqrt{X_2/\nu}} \sim St(\nu)$ — распределение
Стьюдента с ν степенями
свободы

РАСПРЕДЕЛЕНИЕ СТЬЮДЕНТА



› При больших ν очень похоже на $N(0, 1)$

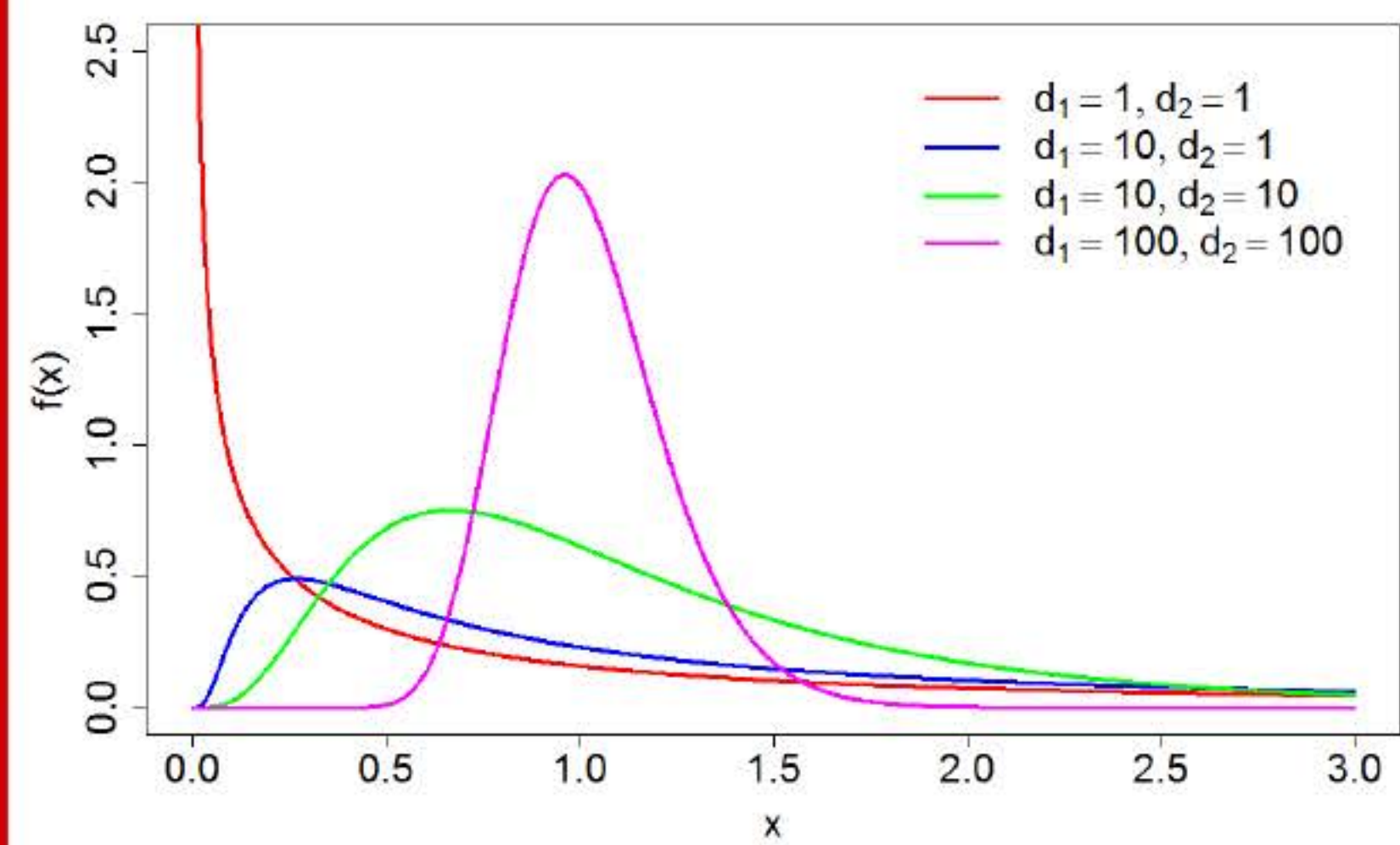


РАСПРЕДЕЛЕНИЕ ФИШЕРА



› $X_1 \sim \chi_{d_1}^2$, $X_2 \sim \chi_{d_2}^2$ независимы

› $X = \frac{X_1/d_1}{X_2/d_2} \sim F(d_1, d_2)$ — распределение Фишера
с d_1, d_2 степенями свободы



ЗАЧЕМ?



› $X \sim N(\mu, \sigma^2), X^n = (X_1, \dots, X_n)$

ЗАЧЕМ?

› $X \sim N(\mu, \sigma^2), X^n = (X_1, \dots, X_n)$

› $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

ЗАЧЕМ?

$$\triangleright X \sim N(\mu, \sigma^2), X^n = (X_1, \dots, X_n)$$

$$\triangleright \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\triangleright S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

ЗАЧЕМ?

$$\triangleright X \sim N(\mu, \sigma^2), X^n = (X_1, \dots, X_n)$$

$$\triangleright \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\triangleright (n-1) \frac{S_n^2}{\sigma^2} = \sum_{i=1}^n (X_i - \bar{X}_n)^2 \sim \chi_{n-1}^2$$

$$\blacktriangleright (n - 1) \frac{S_n^2}{\sigma^2} = \sum_{i=1}^n (X_i - \bar{X}_n)^2 \sim \chi_{n-1}^2$$

$$\blacktriangleright T = \frac{\bar{X}_n - \mu}{S_n / \sqrt{n}} \sim St(n - 1)$$

ЗАЧЕМ?

$$\triangleright X_1 \sim N(\mu_1, \sigma_2^2), X_1^{n_1} = (X_{11}, \dots, X_{1n_1})$$

$$\triangleright X_2 \sim N(\mu_2, \sigma_2^2), X_2^{n_2} = (X_{21}, \dots, X_{2n_2})$$

$$\triangleright \frac{S_1^2}{S_2^2} \sim F(n_1 - 1, n_2 - 1)$$

- › Распределения хи-квадрат, Стьюдента и Фишера
- › Их связь со статистиками нормальных выборок
- › Далее: доверительные интервалы для среднего и дисперсии

ДОВЕРИТЕЛЬНЫЕ ИНТЕРВАЛЫ НА ОСНОВЕ БУТСТРЕПА

ПОСТРОЕНИЕ ДОВЕРИТЕЛЬНЫХ ИНТЕРВАЛОВ



- › Чтобы построить доверительный интервал для статистики $T_n = T(X^n)$, нужно знать её выборочное распределение $F_{T_n}(x)$
- › Как его оценить?

ПОСТРОЕНИЕ ДОВЕРИТЕЛЬНЫХ ИНТЕРВАЛОВ

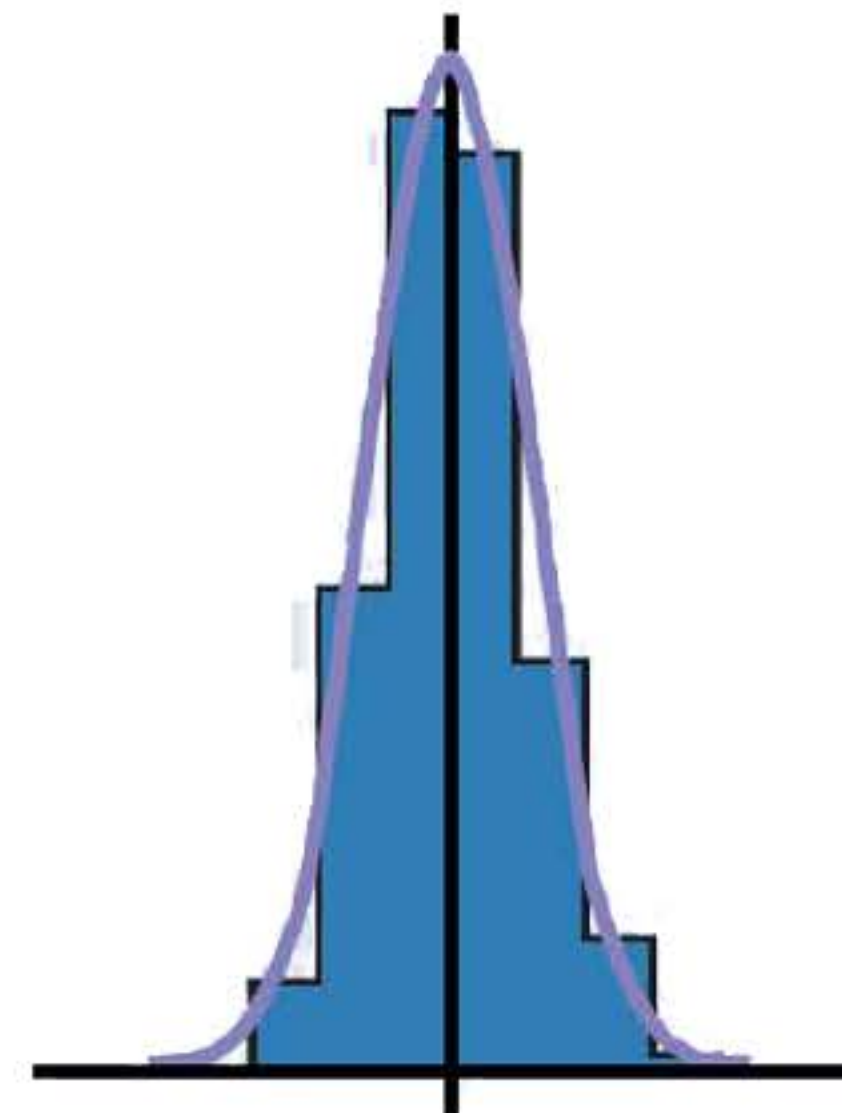
» Наивный метод:

- Извлечь из генеральной совокупности N выборок объёма n и оценить выборочное распределение T_n эмперическим



POPULATION
unknown mean μ

SRS of size n → \bar{x}
SRS of size n → \bar{x}
SRS of size n → \bar{x}
.
.
.



Sampling distribution

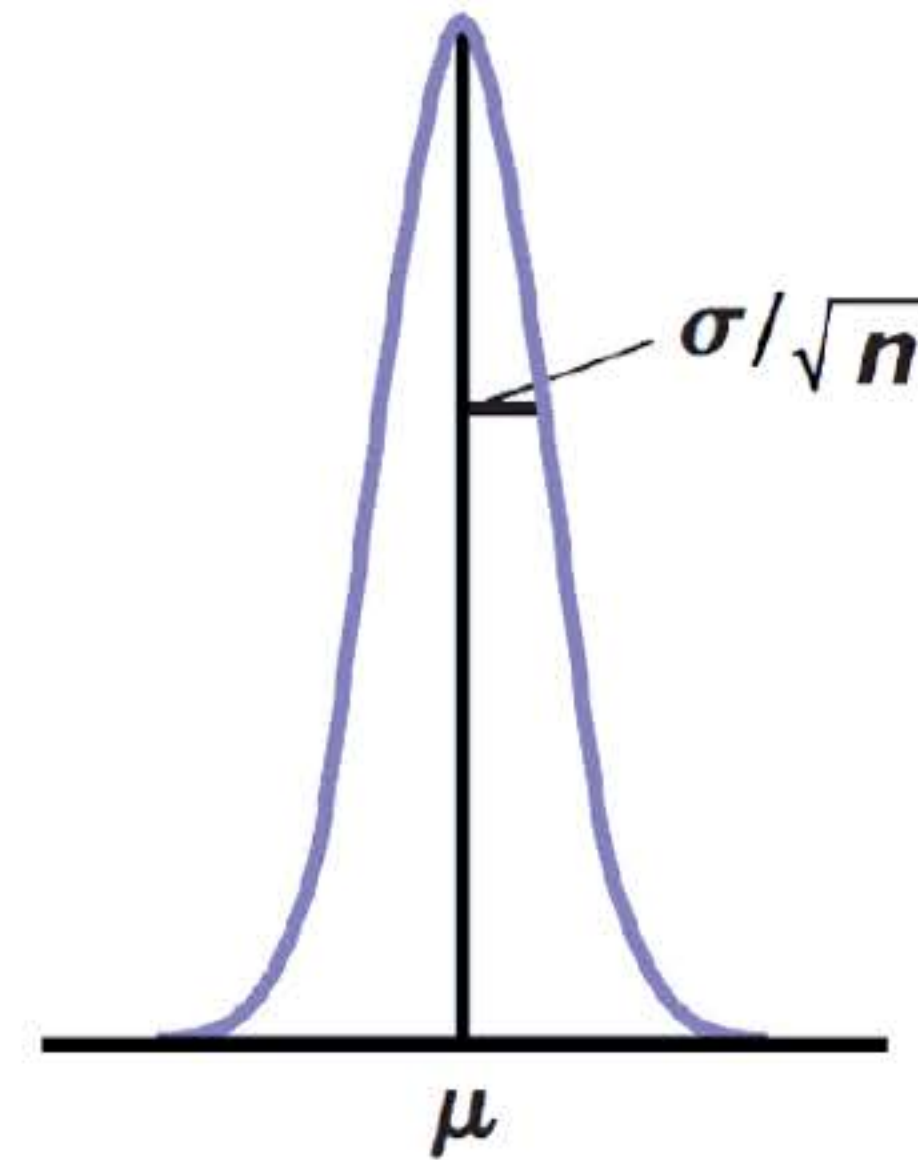
ПОСТРОЕНИЕ ДОВЕРИТЕЛЬНЫХ ИНТЕРВАЛОВ

- Параметрический метод:
 - ▶ Сделать предположение, что X распределена по закону $F_X(x)$, при выполнении которого закон распределения известен



NORMAL POPULATION
unknown mean μ

Theory



Sampling distribution

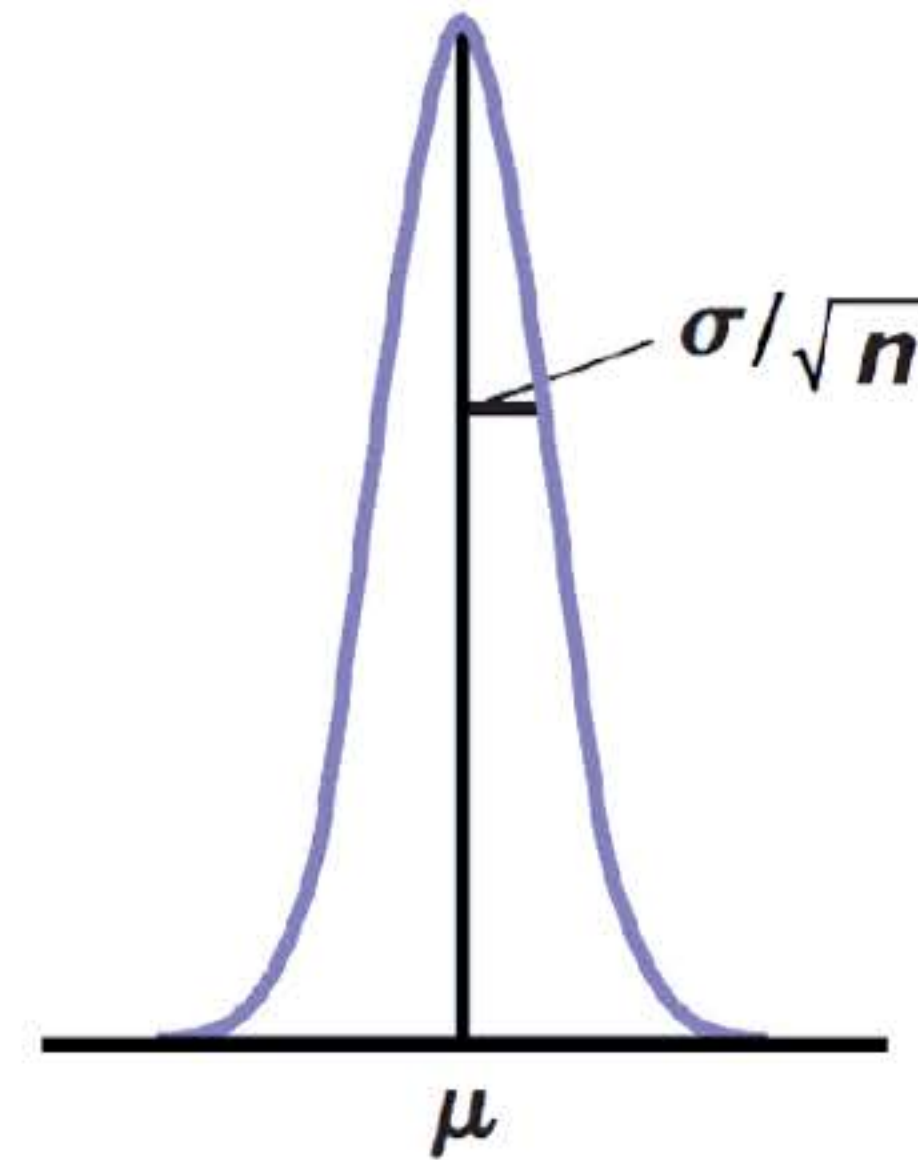
ПОСТРОЕНИЕ ДОВЕРИТЕЛЬНЫХ ИНТЕРВАЛОВ

- Параметрический метод:
 - ▶ Сделать предположение, что X распределена по закону $F_X(x)$, при выполнении которого закон распределения T_n известен



NORMAL POPULATION
unknown mean μ

Theory



Sampling distribution

ПОСТРОЕНИЕ ДОВЕРИТЕЛЬНЫХ ИНТЕРВАЛОВ

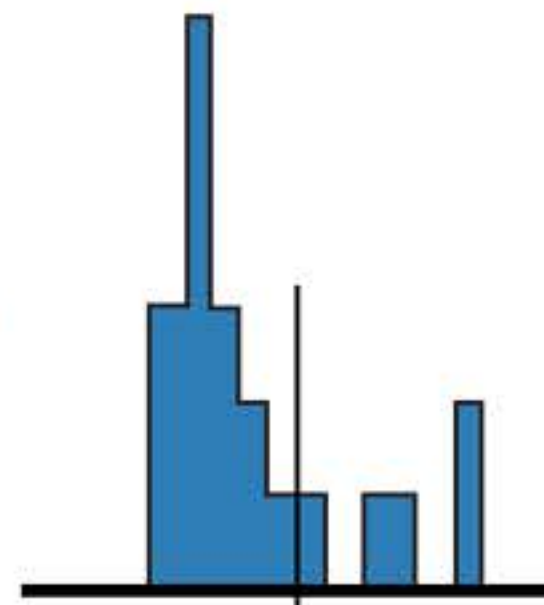
➤ Бутстреп:

- ▶ Сгенерировать N “псевдовыборок” объёма n и оценить выборочное распределение T_n “псевдоэмпирическим”



POPULATION
unknown mean μ

One SRS of size n



Resample of size n

Resample of size n

Resample of size n

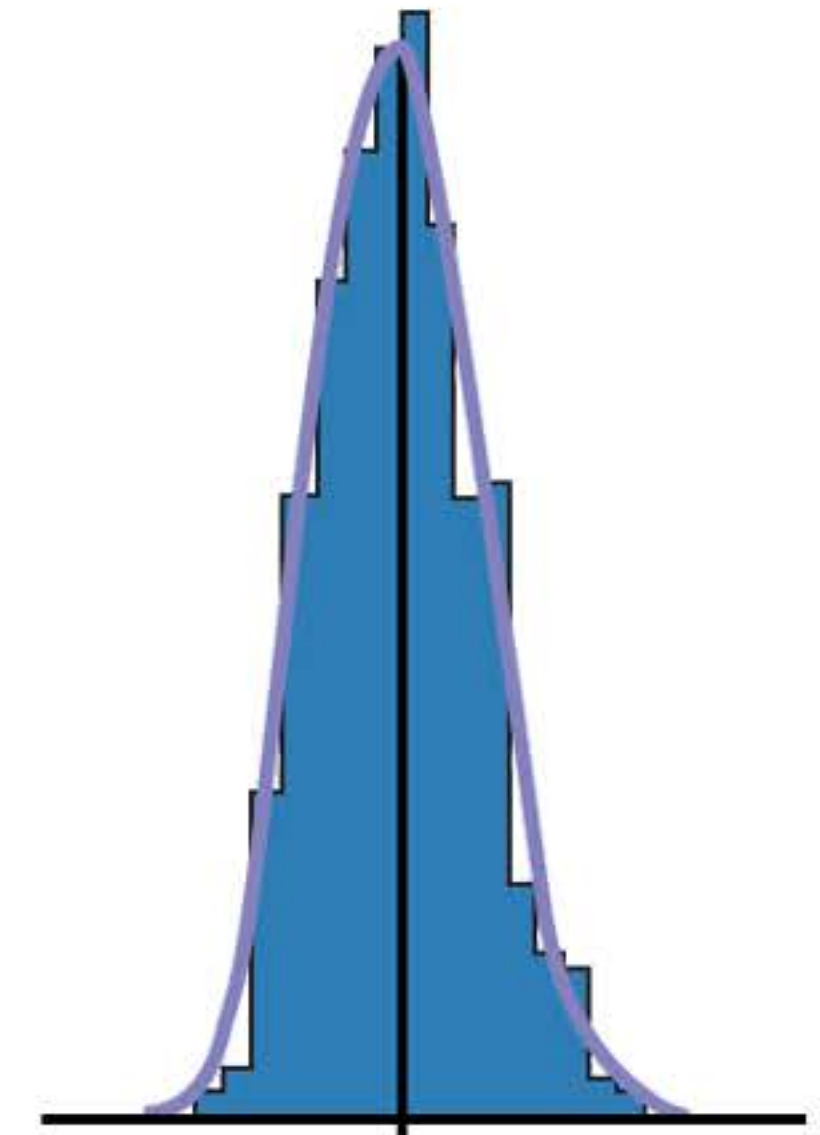
\bar{x}

\bar{x}

\bar{x}

⋮

⋮



Bootstrap distribution

- › Извлечение выборок из генеральной совокупности — сэмплирование из неизвестного распределения $F_X(x)$
- › Лучшая оценка $F_X(x)$, которая у нас есть — $F_{X_n}(x)$
- › Будем сэмплировать из неё. Это то же самое, что делать из X^n выборки с возвращением объёма n