

Statistical Data Analysis

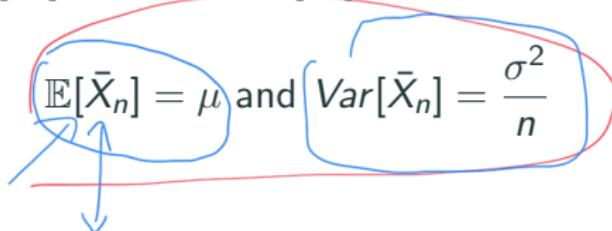
Dr. Jana de Wiljes

3. November 2021

Universität Potsdam

Random variables

Proposition: Let X_1, \dots, X_n be independent and identical random variables with $\mathbb{E}[X_i] = \mu$ and $\text{Var}[X_i] = \sigma^2$. Then

$$\mathbb{E}[\bar{X}_n] = \mu \text{ and } \text{Var}[\bar{X}_n] = \frac{\sigma^2}{n} \quad (1)$$


Remark : if the expected value of the estimator is equal to the true parameter (to be estimated) then we call this estimator **unbiased**

Proof

Law of large numbers

Proposition: Let X_1, \dots, X_n be independent and identical random variables with $\mathbb{E}[X_i] = \mu$. Then

$$\bar{X}_n \rightarrow \mu \text{ for } n \rightarrow \infty \text{ (almost sure)} \quad (2)$$

means $P(\lim_{n \rightarrow \infty} \bar{X}_n = \mu) = 1$

almost everywhere

Proof : Law of large numbers

Empirical variance

Definition: The empirical variance is defined by

$$s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \quad (3)$$

Note: we will also use an analog notation for the random variables:

$$\boxed{s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2} \quad (4)$$

Empirical variance

Proposition: Let X_1, \dots, X_n be independent and identical random variables. Then

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i^2) - n\bar{X}_n^2 \quad (5)$$

The handwritten annotations include a red oval enclosing the formula, a blue circle around the fraction $\frac{1}{n-1}$, a red arrow pointing from the red circle to the red circled term $n\bar{X}_n^2$, and a blue question mark $n?$ written below the oval.

Proof

$$\begin{aligned} \text{Proof: } s_n^2 &= \frac{1}{n-1} \left(\sum_{i=1}^n (x_i - \bar{x}_n)^2 \right) \\ &= \frac{1}{n-1} \left(\left(\sum_{i=1}^n x_i^2 \right) - 2 \cancel{\bar{x}_n x_n} + \cancel{(\bar{x}_n^2)} \right) \\ &= \frac{1}{n-1} \left(\left(\sum_{i=1}^n x_i^2 \right) - 2 \cancel{\bar{x}_n n \bar{x}_n} + \cancel{n \bar{x}_n^2} \right) \\ &= \frac{1}{n-1} \left(\left(\sum_{i=1}^n x_i^2 \right) - n \bar{x}_n^2 \right) \end{aligned}$$

□

$$n \bar{x}_n = \left(\sum_{i=1}^n x_i \right)$$

$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$

Proof

Empirical variance

Proposition: Let X_1, \dots, X_n be independent and identical random variables with $\mathbb{E}[X_i] = \mu$ and $\text{Var}[X_i] = \sigma^2$. Then

$$\mathbb{E}[S_n^2] = \sigma^2 \quad (6)$$

$$\begin{aligned} \mathbb{E}[X_i^2] &= \text{Var}(X_i) + \mathbb{E}[(X_i)]^2 = \sigma^2 + \mu^2 \\ \mathbb{E}[\bar{X}_n^2] &= \text{Var}(\bar{X}_n) + \mathbb{E}[(\bar{X}_n)]^2 = \frac{\sigma^2}{n} + \mu^2 \end{aligned}$$

Proof

$$\text{Proof: } \mathbb{E}[S_n^2] = \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right]$$

$$= \frac{1}{n-1} \sum_{i=1}^n (\mathbb{E}[x_i^2]) - n\mathbb{E}[\bar{x}_n^2]$$

other formulation of
 S_n^2 + Linearity of \mathbb{E}

$$= \frac{1}{n-1} \left(\underbrace{\sum_{i=1}^n \sigma^2 + \mu^2}_{n(\sigma^2 + \mu^2)} - n \left(\frac{\sigma^2}{n} + \underline{\mu^2} \right) \right)$$

$$= \frac{1}{n-1} \left(n(\sigma^2 + \mu^2) - \cancel{n\sigma^2} - \cancel{n\mu^2} \right)$$

$$= \frac{n-1}{n-1} \sigma^2$$

$$= \sigma^2$$

□

$$\frac{1}{n}$$

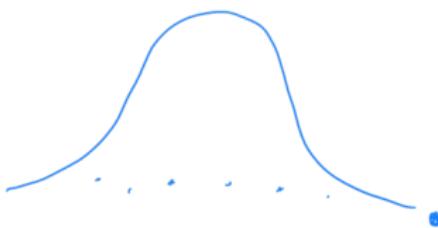
$$\boxed{\frac{1}{n-1}}$$

Proof

Note : that $\tilde{S}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$

$$\leadsto E[\tilde{S}_n] = E\left[\frac{n-1}{n} S_n^2\right] = \frac{n-1}{n} E[S_n^2] = \boxed{\frac{n-1}{n} \sigma^2 < \sigma^2}$$

↑
this estimator biased \leadsto underestimating



Empirical standard deviation

Def: The empirical standard deviation is defined by

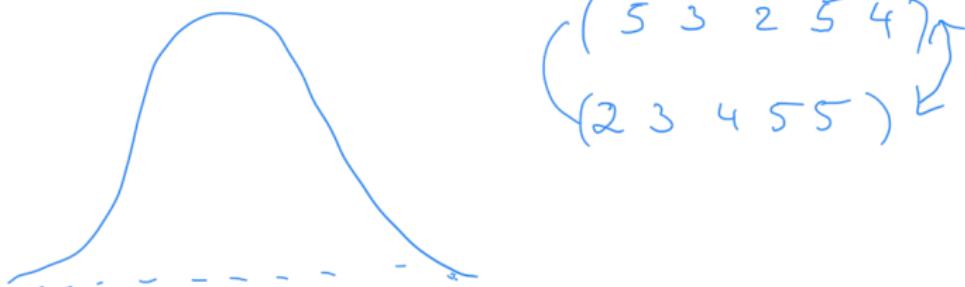
$$s_n = \sqrt{s_n^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2} \quad (7)$$

Order statistic

Def: Let $(x_1, \dots, x_n) \in \mathbb{R}^n$ be a sample set. One can order the elements in an increasing manner:

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)} \quad (8)$$

Then $x_{(i)}$ is referred to as the i-th order statistic of the sample set.



Sample median

Def: The sample median of a set of samples if given by

$$\text{Med}_n = \text{Med}_n(x_1, \dots, x_n) = \begin{cases} x_{\left(\frac{n+1}{2}\right)} & \text{in case } n \text{ is uneven} \\ \frac{1}{2} \left(x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)} \right) & \text{in case } n \text{ is even} \end{cases}$$

Then $x_{(i)}$ is referred to as the i-th order statistic of the sample set.

Note: $\text{Med}_n \neq \bar{x}_n$
in general

$$\rightarrow (1, 2, 2, 1, 1, 1, 2) \rightsquigarrow \bar{x}_n = 1.5$$

$$\rightarrow (1, 2, 2, 1, 1, 1, 200) \rightsquigarrow \bar{x}_n = 20.25$$

$$(1, 1, 1, 1, 2, 2, 2) \rightsquigarrow \frac{1+2}{2} = 1.5$$

$$(1, 1, 1, 1, 2, 2, 200) \rightsquigarrow \frac{1+2}{2} = 1.5$$

robust with respect
to this change of
sample set

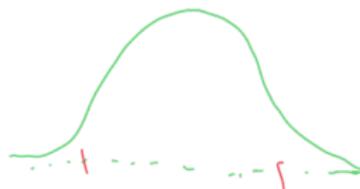
Example

Truncated mean

Def: The truncated mean samples $(x_1, \dots, x_n) \in \mathbb{R}^n$ is defined by

$$\frac{1}{n-2k} \sum_{i=k+1}^{n-k} x_{(i)}$$

$$k = \lfloor 0.05 \cdot n \rfloor$$



Empirische α -Quantile

Def: Let $(x_1, \dots, x_n) \in \mathbb{R}^n$ be a set of samples and $\alpha \in (0, 1)$. The empirical α Quantil is defined by

$$q_\alpha = \begin{cases} x_{\lfloor n\alpha \rfloor + 1} & \text{falls } n\alpha \notin \mathbb{N} \\ \frac{1}{2}(x_{\lfloor n\alpha \rfloor} + x_{\lfloor n\alpha \rfloor + 1}) & \text{falls } n\alpha \in \mathbb{N} \end{cases}$$

Distribution of the order statistic

Proposition: Let X_1, X_2, \dots, X_n be independent and identical distributed random variables, that are absolute continuous with a density f and cumulative distribution function F . Let

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)} \quad (9)$$

be the order statistics. Then the density of the random variable $X_{(i)}$ is

$$f_{X_{(i)}}(t) = \frac{n!}{(i-1)!(n-i)!} f(t) F(t)^{(i-1)} (1 - F(t))^{n-i} \quad (10)$$

Beta distribution

Def: For a and b larger than zero and

$$f(x) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}.$$

where the normalization is given by

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} = \int_0^1 u^{a-1} (1-u)^{b-1} du$$

with $\Gamma(n) = (n-1)!$ being the gamma function.

Excuse to Bandits

Multi-armed bandits

Choose from K options
to receive a
high reward and
reduce loss after T rounds



Examples:

- Which advertising campaign generates the largest revenue
- Which restaurant to pick ?
- Which netflix series to streamen?
- Which vaccination should be further developed ?

Multi-armed bandits

A stochastic K-Armed Bandit is defined via the tuple $\langle \mathcal{A}, \mathcal{Y}, P, r \rangle$ where

- \mathcal{A} is the set of actions (arms) and $|\mathcal{A}| = K$
- \mathcal{Y} is the set of possible outcomes
- $P(\cdot|a) \in \mathcal{P}(\mathcal{Y})$ is the outcome probability, conditioned on action $a \in \mathcal{A}$ being taken,
- $r(\mathcal{Y}) \in \mathcal{R}$ represents the reward obtained when outcome $Y \in \mathcal{Y}$ is observed

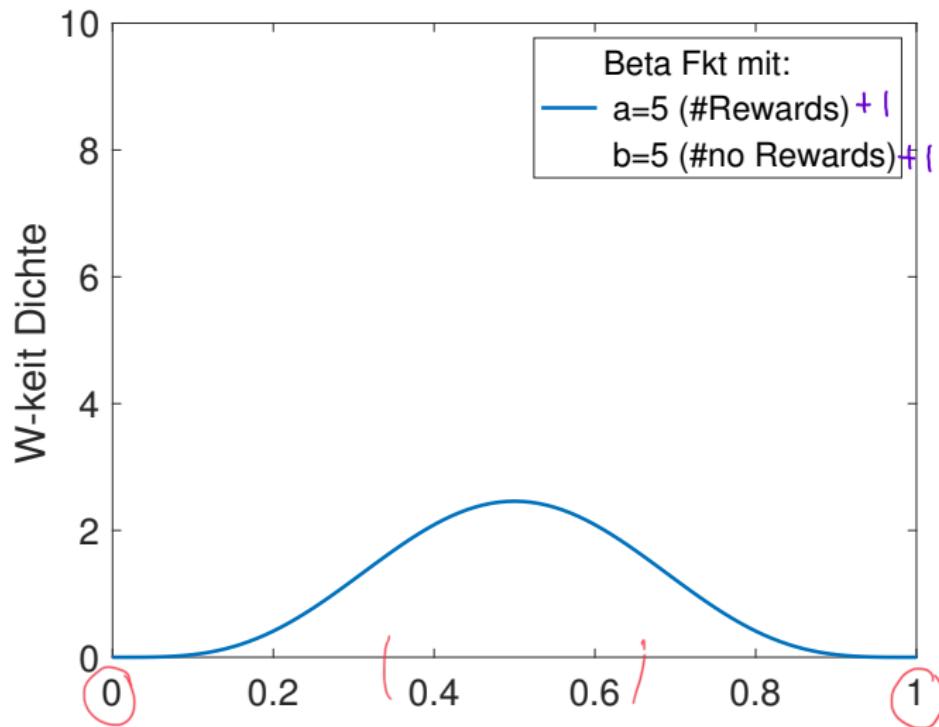
Regret

Def: Let $a^* \in \arg \max_{a \in \mathcal{A}} \mathbb{E}_{y \sim P(\cdot|a)}[r(y)]$ denote the optimal arm. The T-period regret of the sequence of actions a_1, \dots, a_T is the random variable

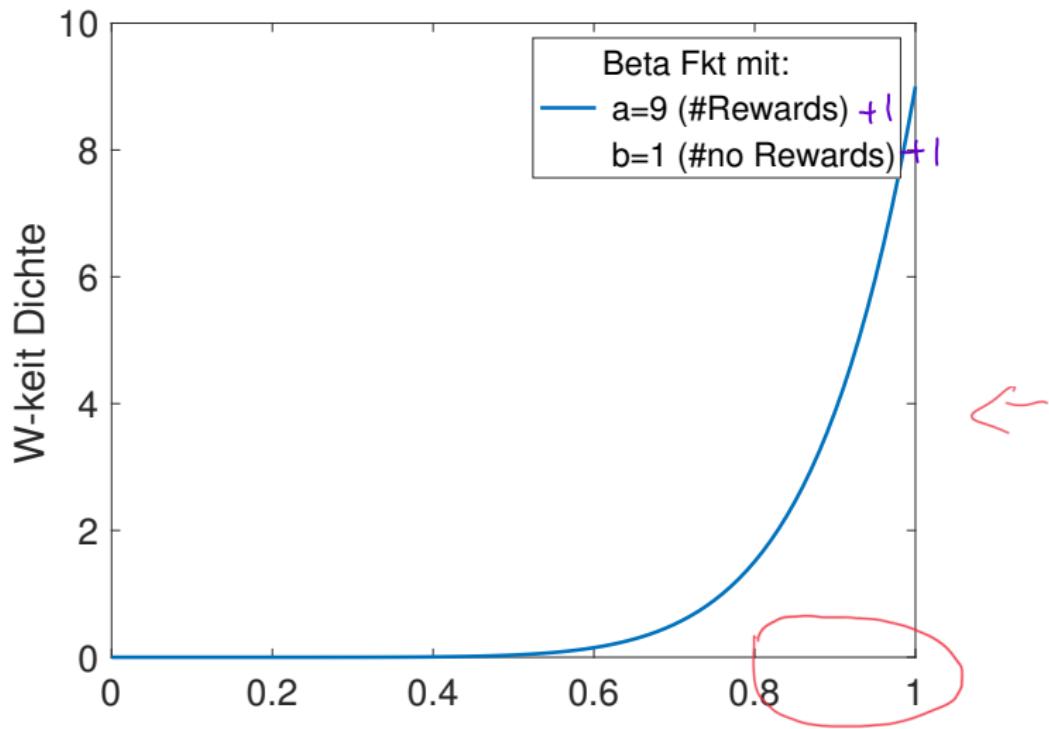
$$\text{Regret}(T) = \sum_{t=1}^T [r(Y_t(a^*)) - r(Y_t(a_t))] \quad (11)$$

Thompson Sampling

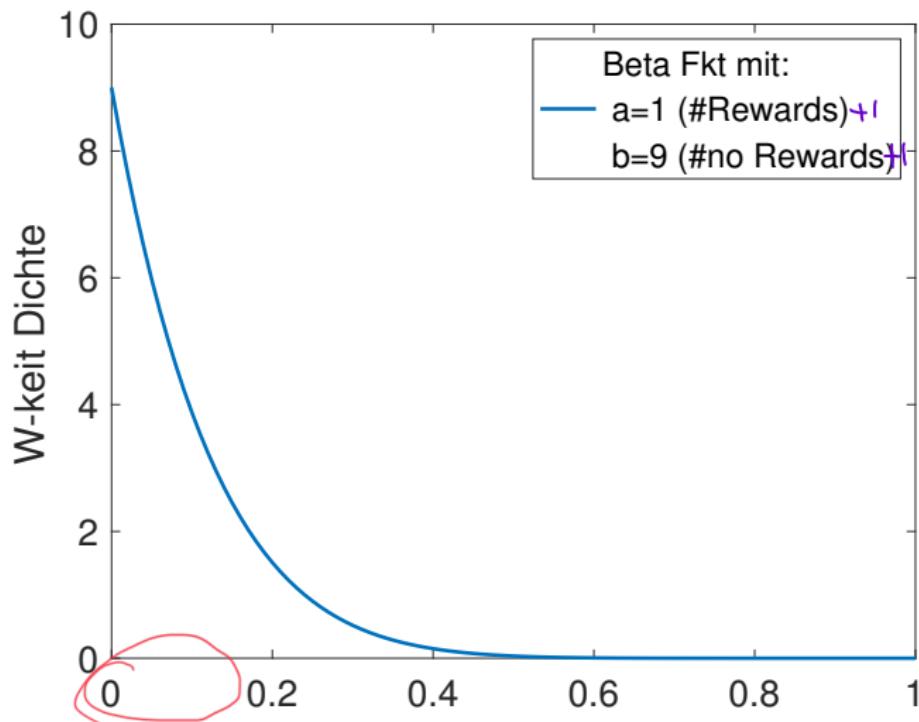
Beta distribution



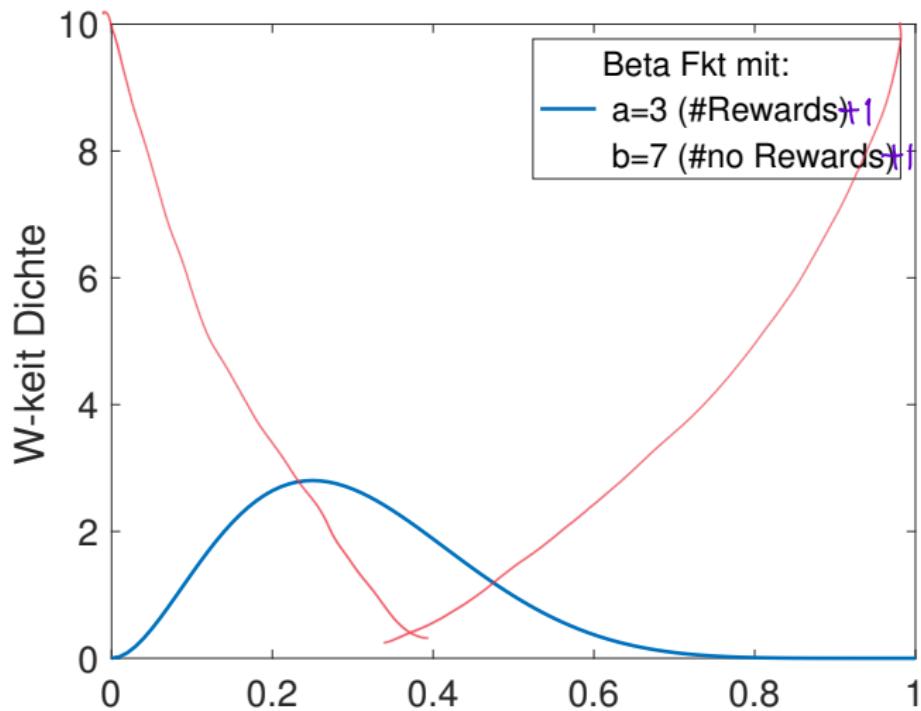
Beta distribution



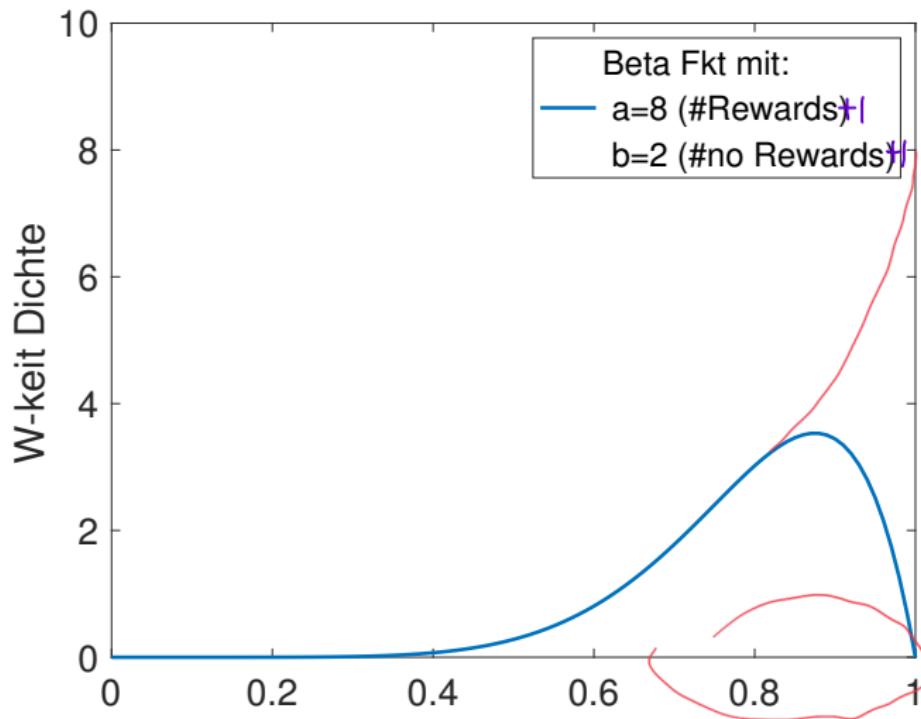
Beta distribution



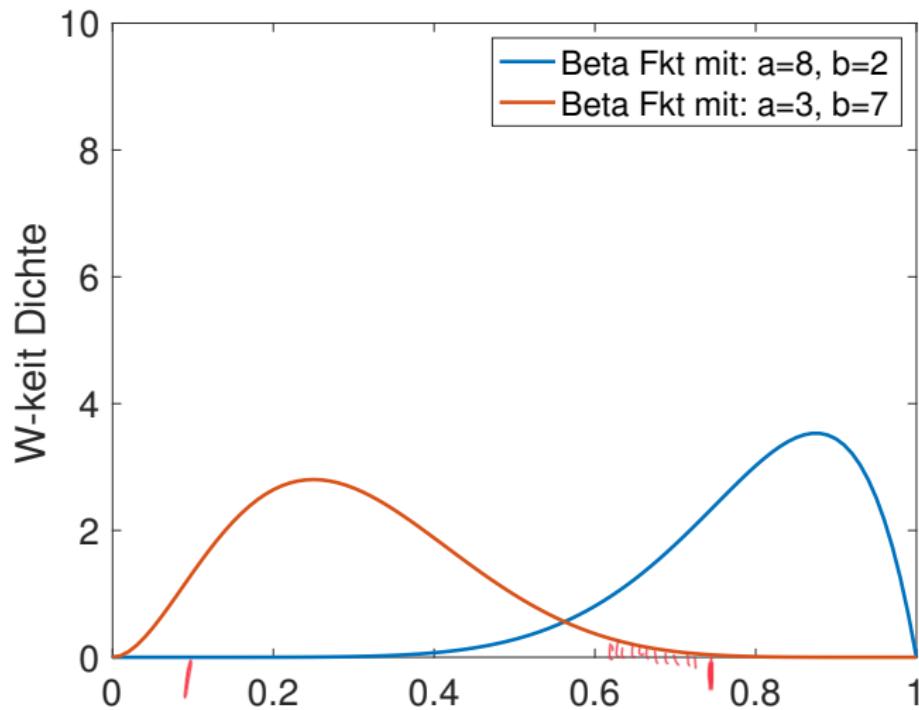
Beta distribution



Beta distribution



Beta distribution



Thompson Sampling

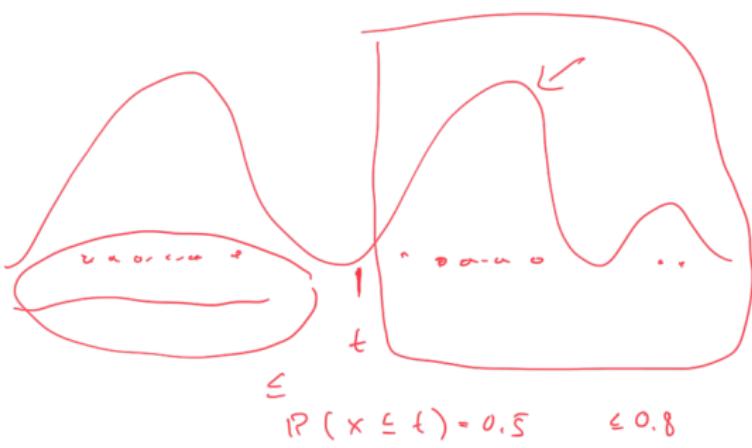
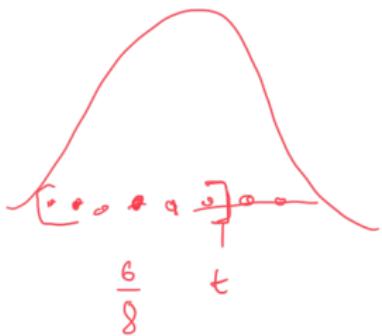
- **Problem setting:** Choose from K options to receive a high reward
- **Algorithm:** Iterated over the following steps:
 1. In each round save information on the choice of action and if a reward was received
 2. Draw from the beta distribution: defined via for each action by
 - a) how often performing action resulted in a reward
 - b) how often performing action did not result in a reward
 3. Choose the action that has the highest beta function value

Empirical cdf

The empirical cdf of a sample set $(x_1, \dots, x_n) \in \mathbb{R}^n$ is defined through

$$\hat{F}_n(t) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i \leq t} = \frac{1}{n} \#\{i \in \{1, \dots, n\} : x_i \leq t\}, \quad t \in \mathbb{R}$$

(12)



Empirical cdf

