

## ***Expected value over countably infinite sets***

$S$  = a countably infinite subset of  $\mathbb{R}$

$$S = \{s_1, s_2, \dots\}$$

$X$  = a random variable which gets values in  $S$

$$E(X) \doteq \sum_{i=1}^{\infty} s_i \Pr(X(\omega) = s_i)$$

Consider the distribution

*distribution is over  
the natural numbers = positive integers*

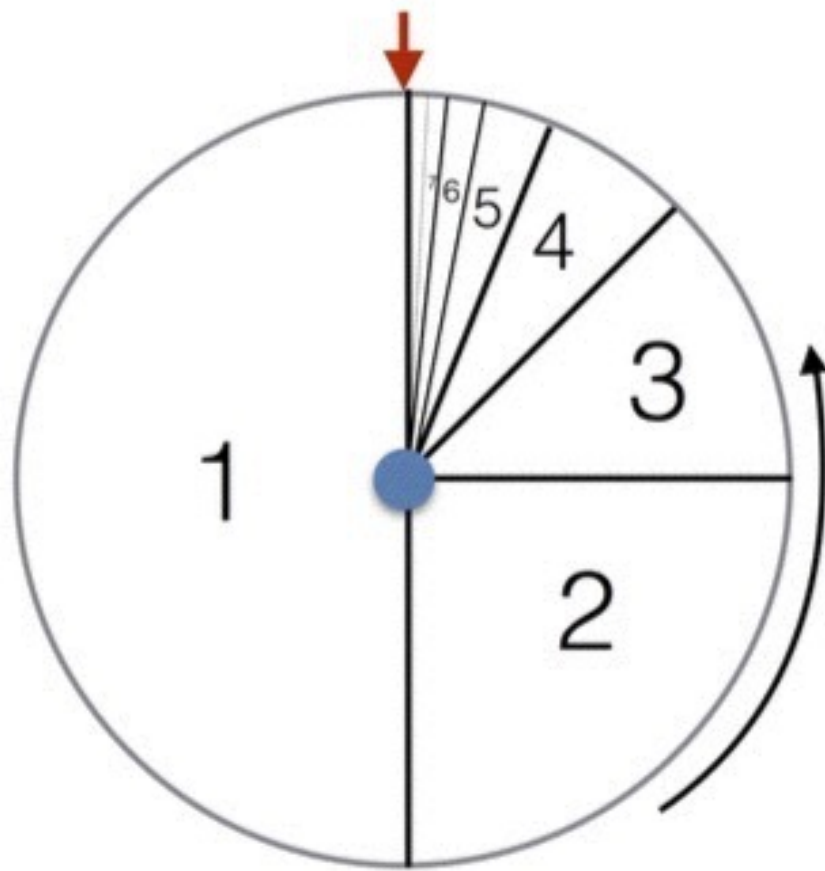
$$P(X = i) = \frac{1}{zi^3}; \quad z = \sum_{i=1}^{\infty} \frac{1}{i^3} < \infty$$

$$E(X) = \sum_{i=1}^{\infty} \frac{i}{zi^3} = \sum_{i=1}^{\infty} \frac{1}{zi^2} < \infty \quad \text{Expectation is finite}$$

Consider next the distribution

$$P(X = i) = \frac{1}{zi^2}; \quad z = \sum_{i=1}^{\infty} \frac{1}{i^2} < \infty$$

$$E(X) = \sum_{i=1}^{\infty} \frac{i}{zi^2} = \sum_{i=1}^{\infty} \frac{1}{zi} = \infty \quad \text{Distribution is well defined but Expectation is *infinite*}$$



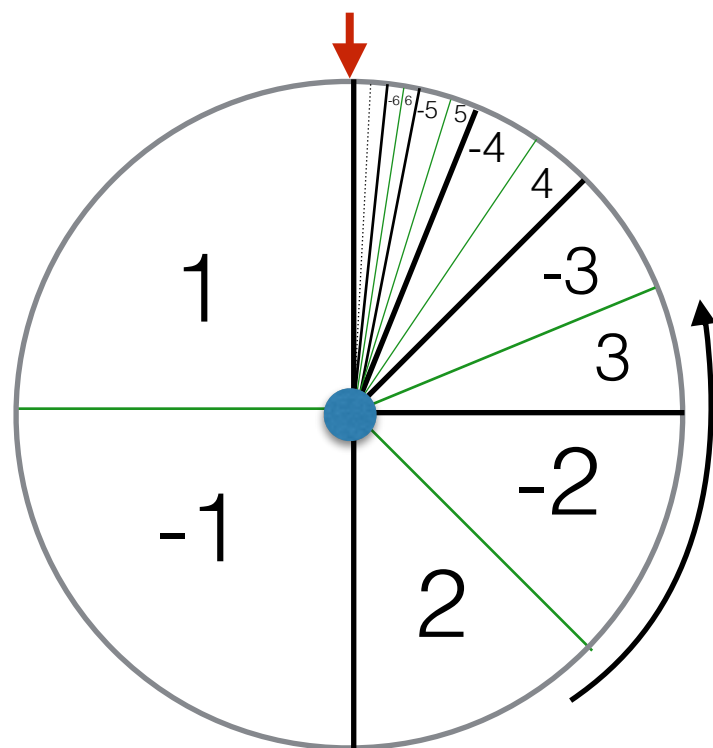
(b) A wheel with  
Infinitely many outcomes

$$P(X = i) = \frac{6}{\pi^2 i^2};$$

$$\sum_{i=1}^{\infty} P(X = i) = 1$$

$$\sum_{i=1}^{\infty} iP(X = i) = \infty$$

Participation in this game is worth any price  
(on the long term)



A wheel with  
Infinitely many outcomes  
both positive and negative

Consider a game with both wins and losses  
 $i \in \{0, -1, +1, -2, +2, \dots\}$

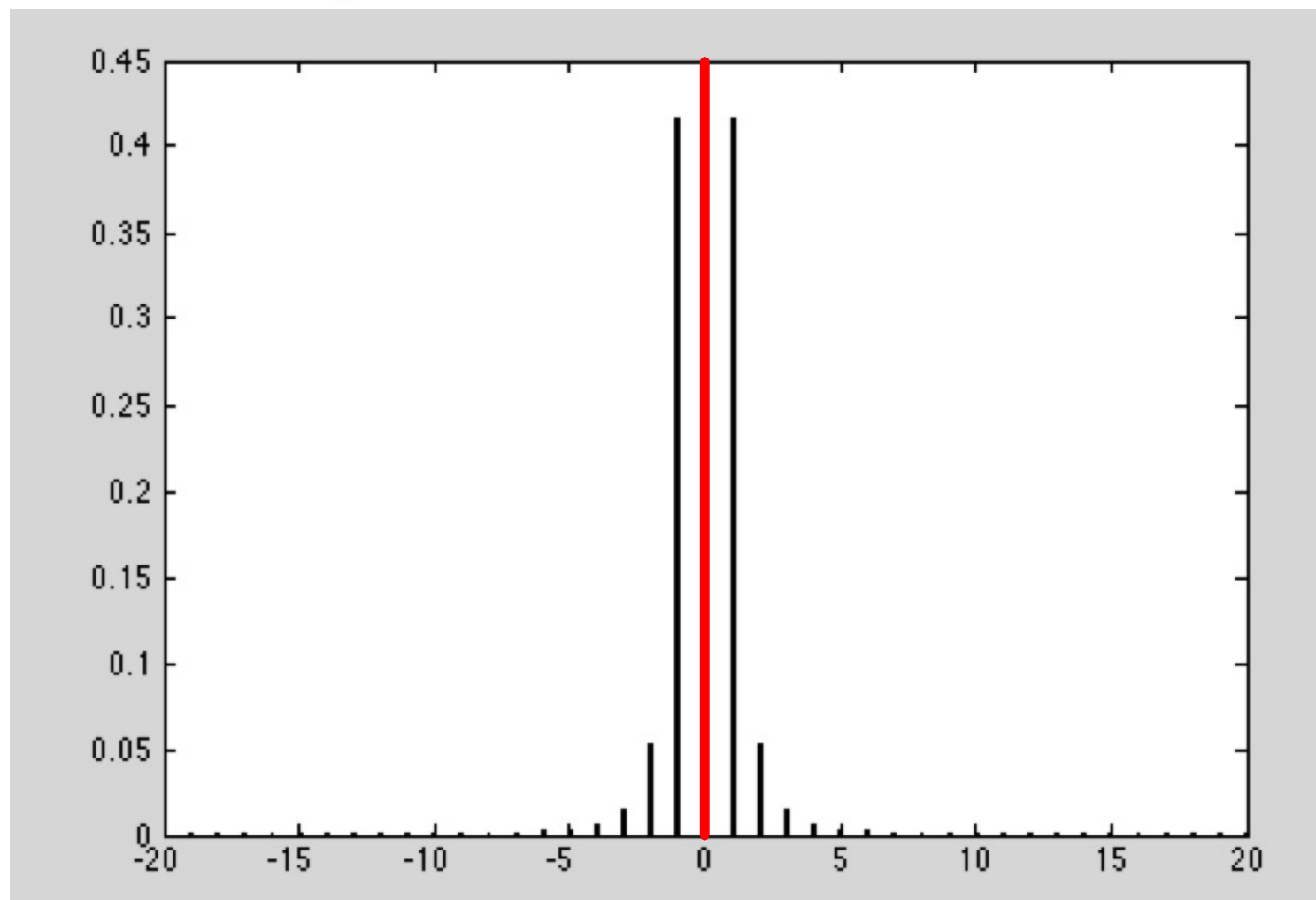
$$P(X = i) = \begin{cases} \frac{1}{Z} \frac{1}{i^{1.5}} & \text{if } i \neq 0 \\ 0 & \text{if } i = 0 \end{cases}, \quad Z = 2 \sum_{i=1}^{\infty} \frac{1}{i^{1.5}}$$

$$\sum_{i=-\infty}^{\infty} P(X = i) = 1$$

$$\sum_{i=-\infty}^{\infty} iP(X = i) \text{ is undefined}$$

## ***Expectation over pos and neg integers: the good case***

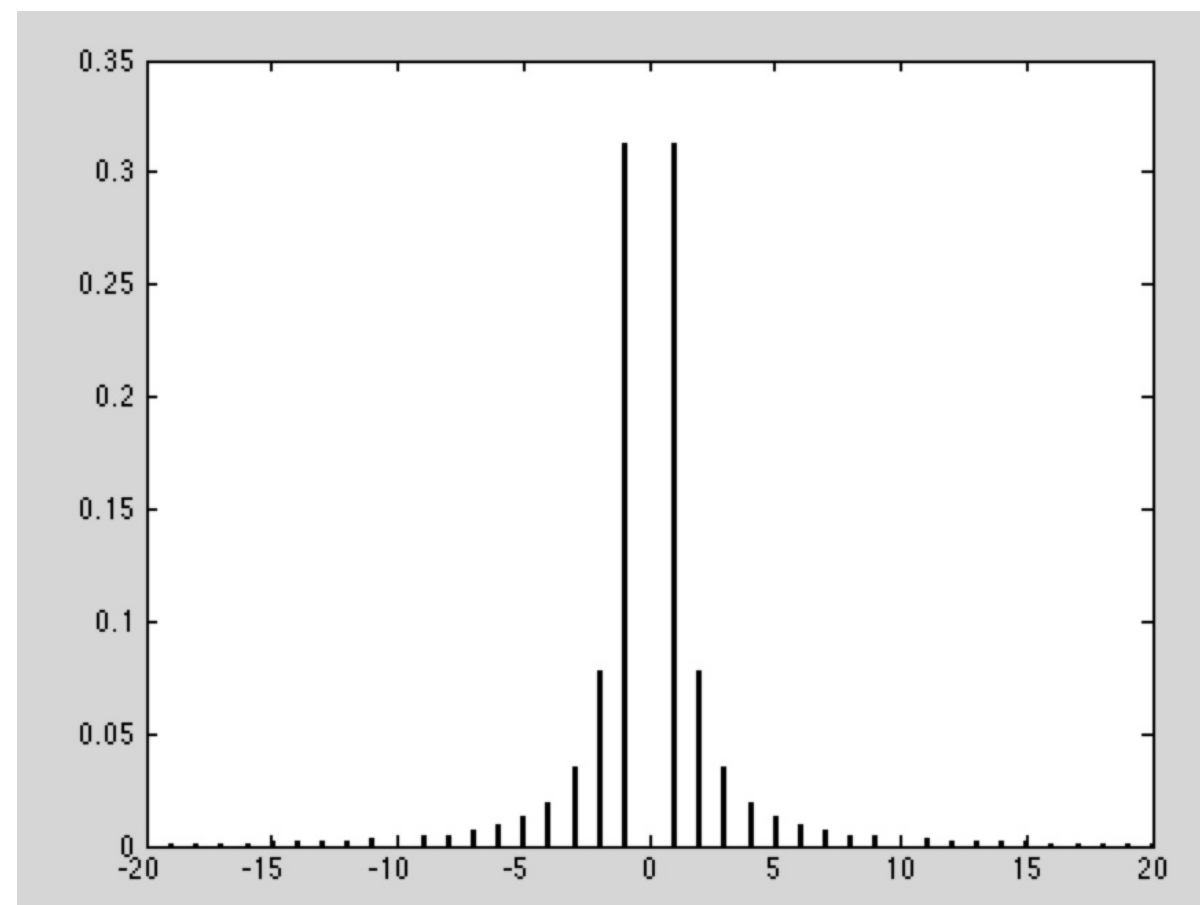
$$P(X = i) = \begin{cases} 0 & \text{if } i = 0 \\ \frac{1}{Z|i|^3} & \text{if } i \neq 0 \end{cases} ; \quad Z = 2 \sum_{i=1}^{\infty} \frac{1}{|i|^3} < \infty$$



$$E(X) = \sum_{i=1}^{\infty} iP(X=i) + \sum_{i=-1}^{-\infty} iP(X=i) = \frac{1}{Z} \left( \sum_{i=1}^{\infty} \frac{1}{i^2} - \sum_{i=1}^{\infty} \frac{1}{i^2} \right) = \frac{c-c}{Z} = 0$$

***A symmetric distribution on pos and neg integers,  
the bad case***

$$P(X = i) = \begin{cases} 0 & \text{if } i = 0 \\ \frac{1}{Zi^2} & \text{if } i \neq 0 \end{cases} ; \quad Z = 2 \sum_{i=1}^{\infty} \frac{1}{i^2} < \infty$$



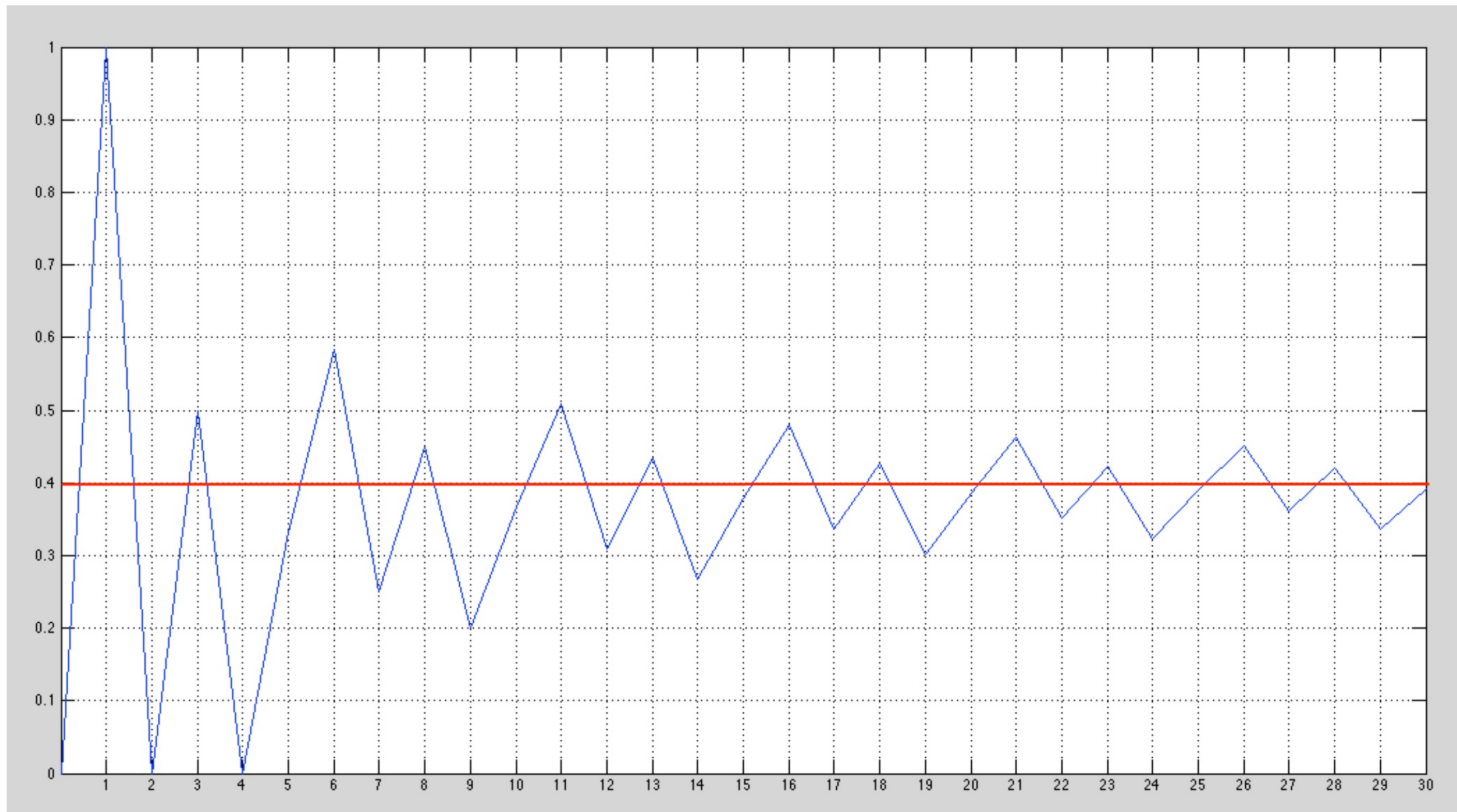
$$E(X) = \sum_{i=1}^{\infty} iP(X = i) + \sum_{i=-1}^{-\infty} iP(X = i) = \frac{1}{Z} \left( \sum_{i=1}^{\infty} \frac{1}{i} - \sum_{i=1}^{\infty} \frac{1}{i} \right) = \frac{\infty - \infty}{Z} = \text{undefined}$$

***Undefined limit means you can get the limit of your choice by changing the order of summation.***

***You have at your disposal two infinitely large sums with shrinkingly small pieces:  
 $1/1, 1/2, 1/3, 1/4, \dots$      $-1/1, -1/2, -1/3, -1/4, \dots$***

***Suppose you want the limit to be 0.4, by alternating between positives and negatives  
you can get arbitrarily close to 0.4 (or to any other number)***

$$\begin{aligned} &1/1 - 1/1 + 1/2 - 1/2 + 1/3 + 1/4 - 1/3 + 1/5 - 1/4 + 1/6 + 1/7 - 1/5 + 1/8 - 1/6 + 1/9 + 1/10 - 1/7 \\ &+ 1/11 - 1/8 + 1/12 + 1/13 - 1/9 + 1/14 - 1/10 + 1/15 + 1/16 - 1/11 + 1/17 - 1/12 + 1/18 = 0.3919 \end{aligned}$$





Let  $X$  be a random variable whose probability distribution is defined as:

$$P(X = i) = \begin{cases} \frac{1}{|i|^\alpha} & \text{if } i \neq 0 \\ 0 & \text{if } i = 0. \end{cases}$$

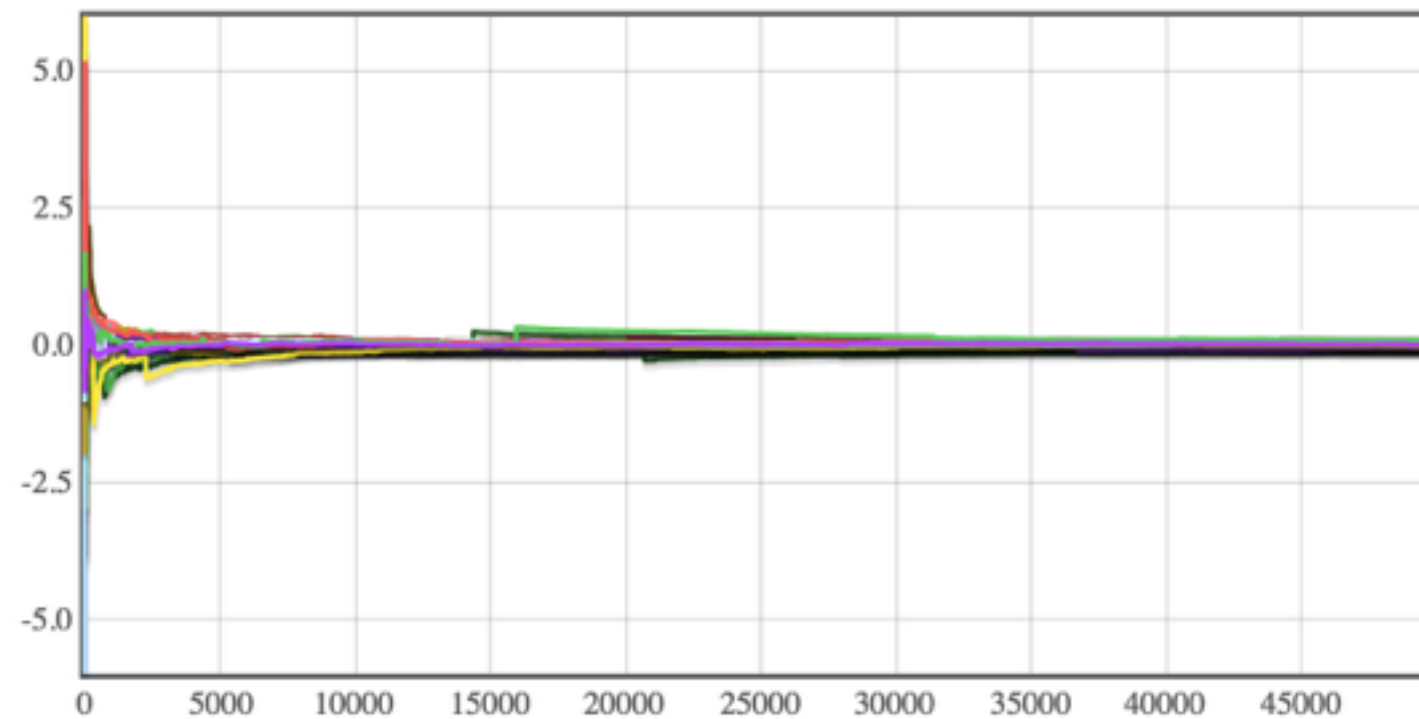
**Simulation parameters:**

$\alpha$ : 2.5

Number of trajectories: 50

Number of data points: 50000

Run





Let  $X$  be a random variable whose probability distribution is defined as:

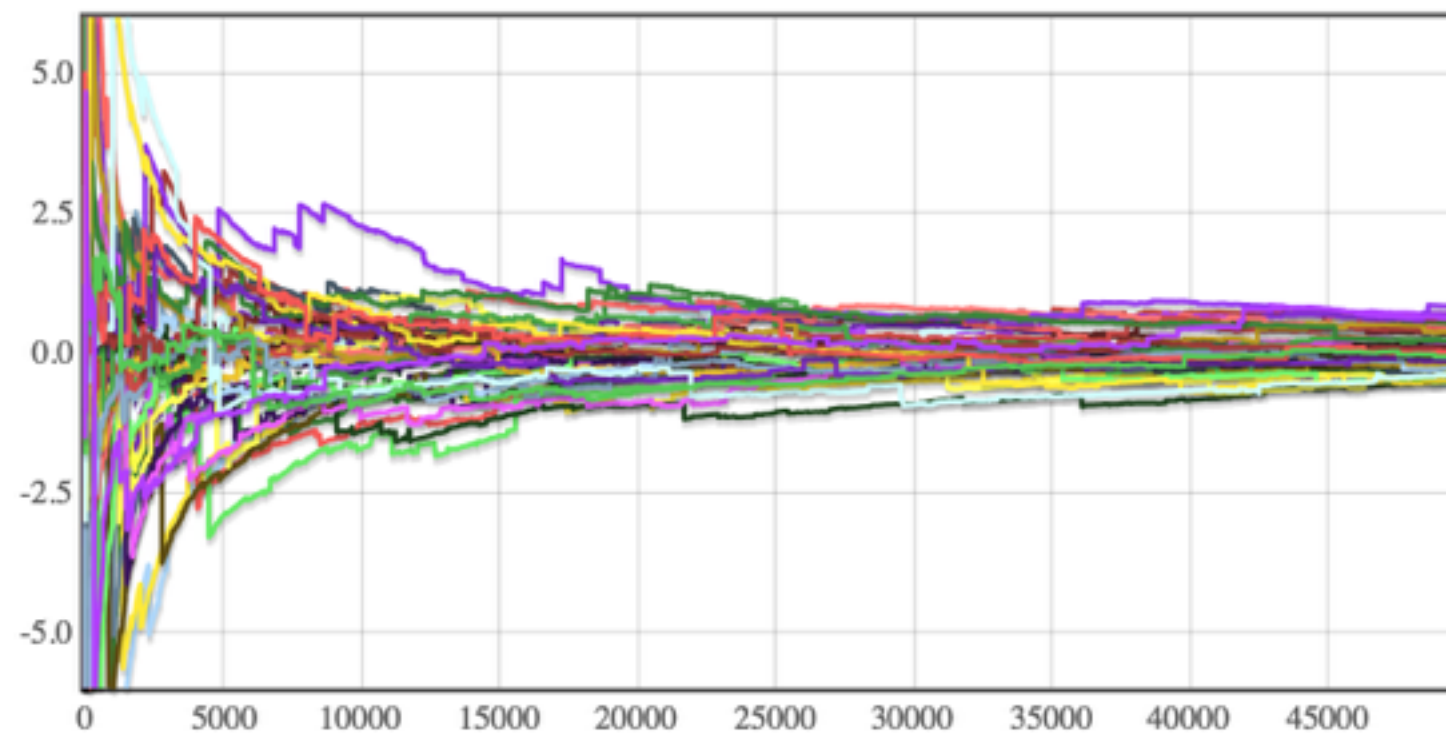
$$P(X = i) = \begin{cases} \frac{1}{|i|^\alpha} & \text{if } i \neq 0 \\ 0 & \text{if } i = 0. \end{cases}$$

**Simulation parameters:**

$\alpha$ :

Number of trajectories:

Number of data points:



Let  $X$  be a random variable whose probability distribution is defined as:

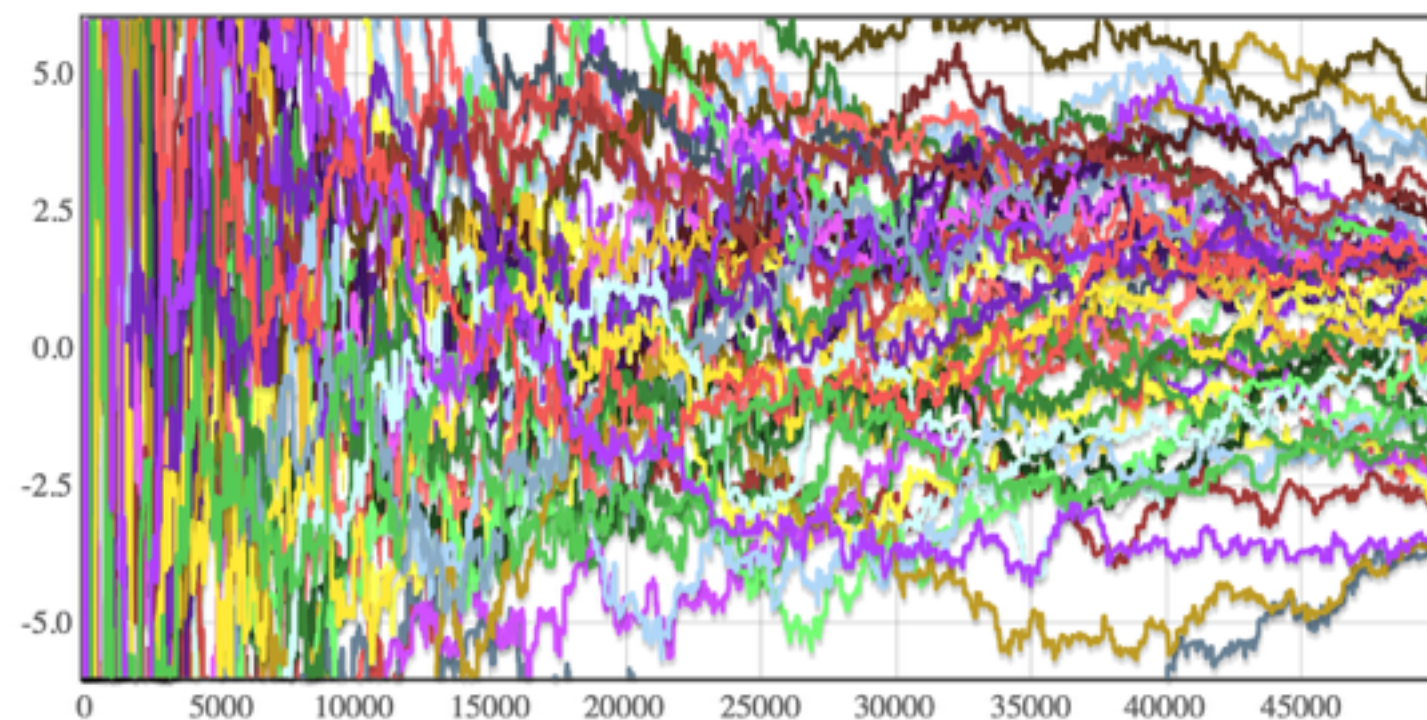
$$P(X = i) = \begin{cases} \frac{1}{|i|^\alpha} & \text{if } i \neq 0 \\ 0 & \text{if } i = 0. \end{cases}$$

**Simulation parameters:**

$\alpha$ :

Number of trajectories:

Number of data points:



# KQS cards

I want to know how to improve my teaching  
Your feedback is important to me.

On the index card provided, please give a one-sentence to each  
of the following questions

Keep: One thing I should **keep** doing?

Quit: One thing I should **quit** doing?

Start: One thing I should **start** doing?

A random algorithm  
for computing percentiles

# A problem with the average

$$\textit{Average}(X_1, \dots, X_n) \doteq \frac{1}{n} \sum_{i=1}^n X_i$$

The average is the most common estimator of the "center" of a distribution. It takes linear time to compute.

However, the average is "sensitive to outliers" :

Suppose that you have a company in which 1000 employees earn 1\$/day and one employee earns 1000\$/day. The average daily pay is  $2000/1001 \sim 2\$/\text{day}$ , but that is double what most people earn.

# Using the Median instead of the average

To compute the median sort all  $n$  elements from smallest to largest and take the value of the element that is the middle of the list (position  $n/2$ ) (take the average of the two middle elements if the list length is even).

In the earlier example, the median will be 1\$ regardless of how big is the largest salary - outliers are ignored.

# The median is a special case of percentiles

To compute the **P-percentile** sort all  $n$  elements from smallest to largest and take the value of the element that is in the  $\text{floor}(Pn)$  position.

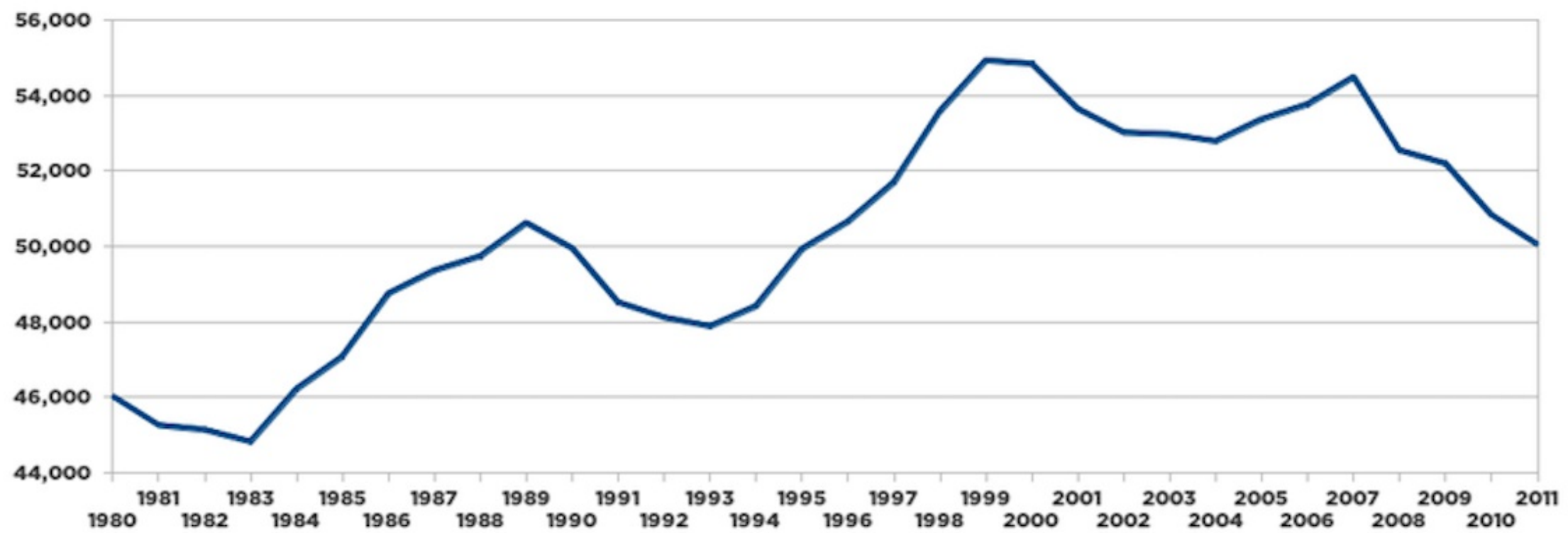
The Median is the  $1/2$ -percentile

Often used when describing distribution of income:

The top



**US Median Household Income, Inflation Adjusted**



**Table 1: Income, net worth, and financial worth in the U.S. by percentile, in 2010 dollars**

<b>Wealth or income class</b>	<b>Mean household income</b>	<b>Mean household net worth</b>	<b>Mean household financial (non-home) wealth</b>
Top 1 percent	\$1,318,200	\$16,439,400	\$15,171,600
Top 20 percent	\$226,200	\$2,061,600	\$1,719,800
60th-80th percentile	\$72,000	\$216,900	\$100,700
40th-60th percentile	\$41,700	\$61,000	\$12,200
Bottom 40 percent	\$17,300	-\$10,600	-\$14,800

From Wolff (2012); only mean figures are available, not medians. Note that income and wealth are separate measures; so, for example, the top 1% of income-earners is not exactly the same group of people as the top 1% of wealth-holders, although there is considerable overlap.

A linear time algorithm for computing percentiles

We can calculate P-percentile by sorting and then picking the element in location  $Pn$ . But this requires time  $O(n \log n)$  (unless distribution of data is known).

We will now describe a randomized algorithm whose expected running time is  $O(n)$

Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Solution 1, sort and locate ---- takes worst case time  $O(n \log n)$ :

<b>3</b>	<b>7</b>	<b>7</b>	<b>10</b>	<b>15</b>	<b>16</b>	<b>20</b>	<b>30</b>	<b>33</b>	<b>70</b>
----------	----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------

Expected time is also  $\Omega(n \log n)$ :

Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Solution 2, randomized algorithm ---- takes ***expected*** time  $O(n)$ :

Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Choose a random element as pivot

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Choose a random element as pivot

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Partition list into 3 lists:  $<7$  ,  $=7$  ,  $>7$



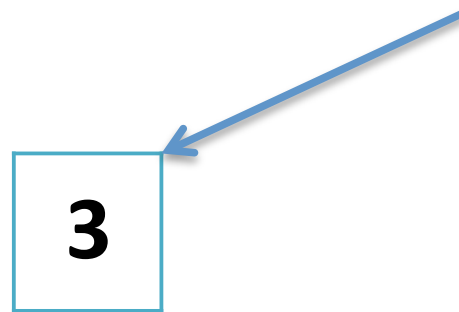
Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Choose a random element as pivot

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Partition list into 3 lists:  $<7$  ,  $=7$  ,  $>7$



$S_L$

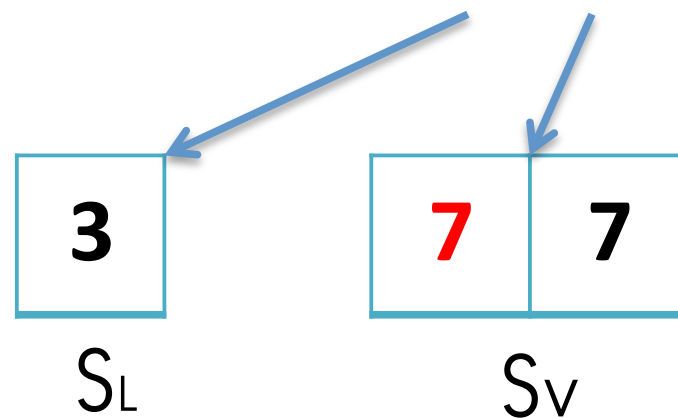
Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Choose a random element as pivot

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Partition list into 3 lists:  $<7$  ,  $=7$  ,  $>7$



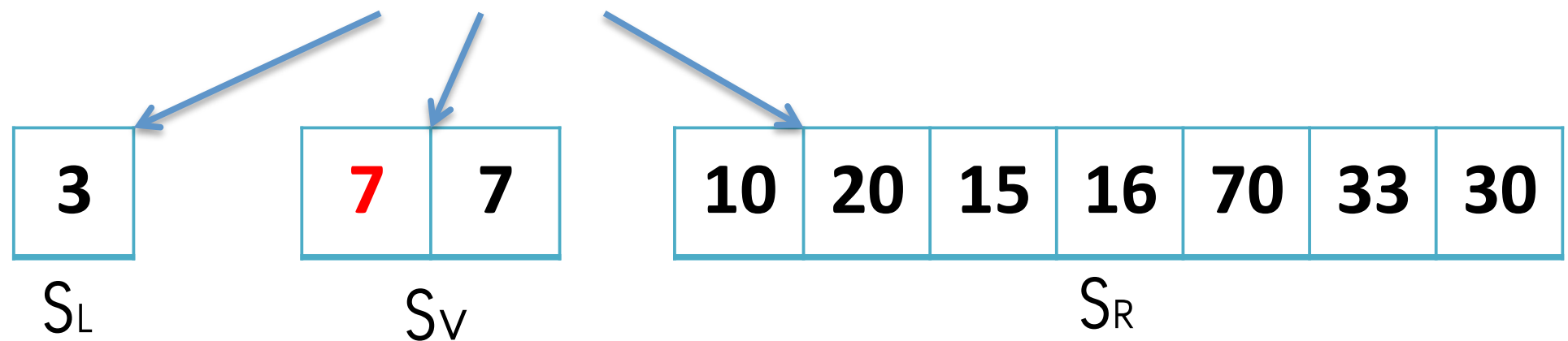
Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Choose a random element as pivot

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Partition list into 3 lists:  $<7$  ,  $=7$  ,  $>7$



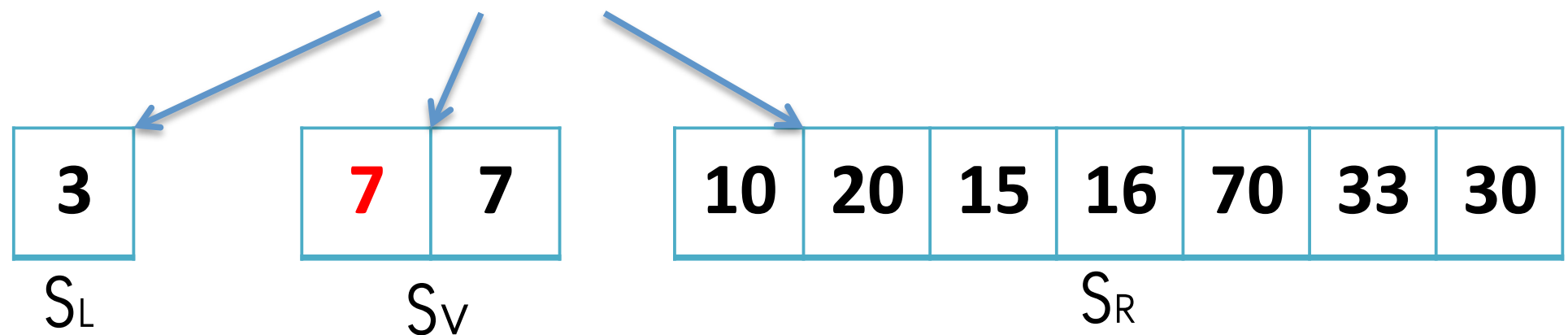
Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Choose a random element as pivot

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Partition list into 3 lists:  $<7$  ,  $=7$  ,  $>7$



Find the 2<sup>nd</sup> smallest element in the following table:

<b>10</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>
-----------	-----------	-----------	-----------	-----------	-----------	-----------

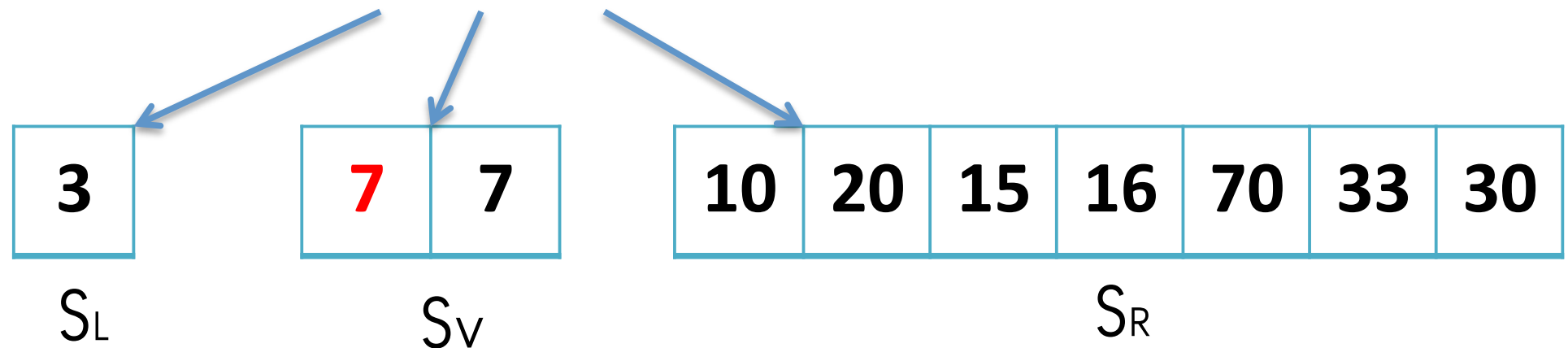
Find the 5<sup>th</sup> smallest element in the following table:

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Choose a random element as pivot

<b>7</b>	<b>10</b>	<b>3</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>	<b>7</b>
----------	-----------	----------	-----------	-----------	-----------	-----------	-----------	-----------	----------

Partition list into 3 lists:  $<7$  ,  $=7$  ,  $>7$



Find the 2<sup>nd</sup> smallest element in the following table:

<b>10</b>	<b>20</b>	<b>15</b>	<b>16</b>	<b>70</b>	<b>33</b>	<b>30</b>
-----------	-----------	-----------	-----------	-----------	-----------	-----------

recurse

# The split operation

- Choosing a split value  $v$  we divide the set  $S$  into three subsets:
  - $S_L = \{x \in S \mid x < v\}$
  - $S_V = \{x \in S \mid x = v\}$
  - $S_R = \{x \in S \mid x > v\}$

# After we split, we know how to continue

- We know the size of the three sets:  $|S_L|$ ,  $|S_v|$ ,  $|S_R|$ .
- If median is in  $S_v$ , then we are done.
- Otherwise we continue with either  $S_L$  or  $S_R$



# Randomized algorithm

- Let  $S$  be a set of  $N$  different numbers. Suppose we pick a random element of  $S$  to be the pivot  $v$ . What is the probability that the size of  $S_L$  is equal to 10 ?
  - $1/10$
  - $9/10$
  - $1/N$

# Randomized algorithm

- Let  $S$  be a set of  $N$  different numbers. Suppose we pick a random element of  $S$  to be the pivot  $v$ . What is the probability that the size of  $S_L$  is between  $\lceil N/4 \rceil$  and  $\lceil 3N/4 \rceil$ ?
  - A. About  $\frac{1}{4}$
  - B. About  $\frac{1}{2}$
  - C. About  $\frac{3}{4}$
  - D. About  $1/N$

# Lucky splits

- We say that the split of a set  $S$  of size  $N$  is lucky if
- $\frac{1}{4}N \leq |S_L| \leq (3/4)N$
- Which implies also that  $\frac{1}{4}N \leq |S_u| \leq (3/4)N$
- If the split is lucky then the size of the set we operate on decreases by a factor of  $(3/4)$
- In order to reduce the set to all-equal elements we need at most  $k$  lucky splits:

$$\left(\frac{3}{4}\right)^k N \leq 1 \Rightarrow k \log \frac{3}{4} + \log N \leq 0 \Rightarrow k \geq \frac{\log N}{\log(3/4)}$$

# Expected time to first lucky split

- What is the expected number of random splits until we get a lucky split?
  - A. 1
  - B. 2
  - C.  $1/2$

## Expected Running time

$n$  = The number of elements in the input array.

$T(n)$  = The expected running time of the algorithm

$$T(n) \leq n + \frac{1}{2}T(n) + \frac{1}{2}T\left(\frac{3}{4}n\right)$$

Multiply both sides by 2 and rearrange:

$$2T(n) \leq 2n + T(n) + T\left(\frac{3}{4}n\right); \quad T(n) \leq 2n + T\left(\frac{3}{4}n\right)$$

$$T(n) \leq 2n + T\left(\frac{3}{4}n\right)$$

$$T(n) \leq 2n + \frac{3}{4}2n + T\left(\left(\frac{3}{4}\right)^2 n\right)$$

$$T(n) \leq 2n + \left(\frac{3}{4}\right)2n + \left(\frac{3}{4}\right)^2 2n + T\left(\left(\frac{3}{4}\right)^2 n\right)$$

$$T(n) \leq 2n \left(1 + \left(\frac{3}{4}\right) + \left(\frac{3}{4}\right)^2 + \left(\frac{3}{4}\right)^3 + \dots\right) + T(\leq 1)$$

$$T(n) \leq 2n \frac{1}{1 - (3/4)} = 8n$$

This is an upper bound - the actual constant is smaller.

Before you leave,  
please bring me your filled-in  
KQS cards



# Two styles of random algorithms

- **Las Vegas**: always gives the correct answer, time to terminate varies from run to run, known bound on **expected** running time.
- **Monte Carlo**: Takes a pre-defined time to terminate, gives the correct answer with non-zero probability.

Tip for remembering:

Monte **C**arlo - outcome not always **C**orrect

Las **V**egas - run-time **V**aries

# From Las Vegas to Monte-Carlo

A Las Vegas algorithm always outputs the correct answer but its run-time varies, the expected run time is bounded.

We want to transform the algorithm into a Monte-Carlo algorithm that Always takes the same same time to complete but has a small probability  $p$  of generating an incorrect answer.

Method: Time out the LV algorithm at time  $s$  and output “fail”.

$T$  = the random variable equal to the run-time of the algorithm.

$T$  is non-negative and  $E(T)$  is known

Markov inequality:  $P(T \geq s) \leq \frac{E(T)}{s}$

Prescription: choose time out to be  $s \geq \frac{E(T)}{p}$

$P(T \leq t)$

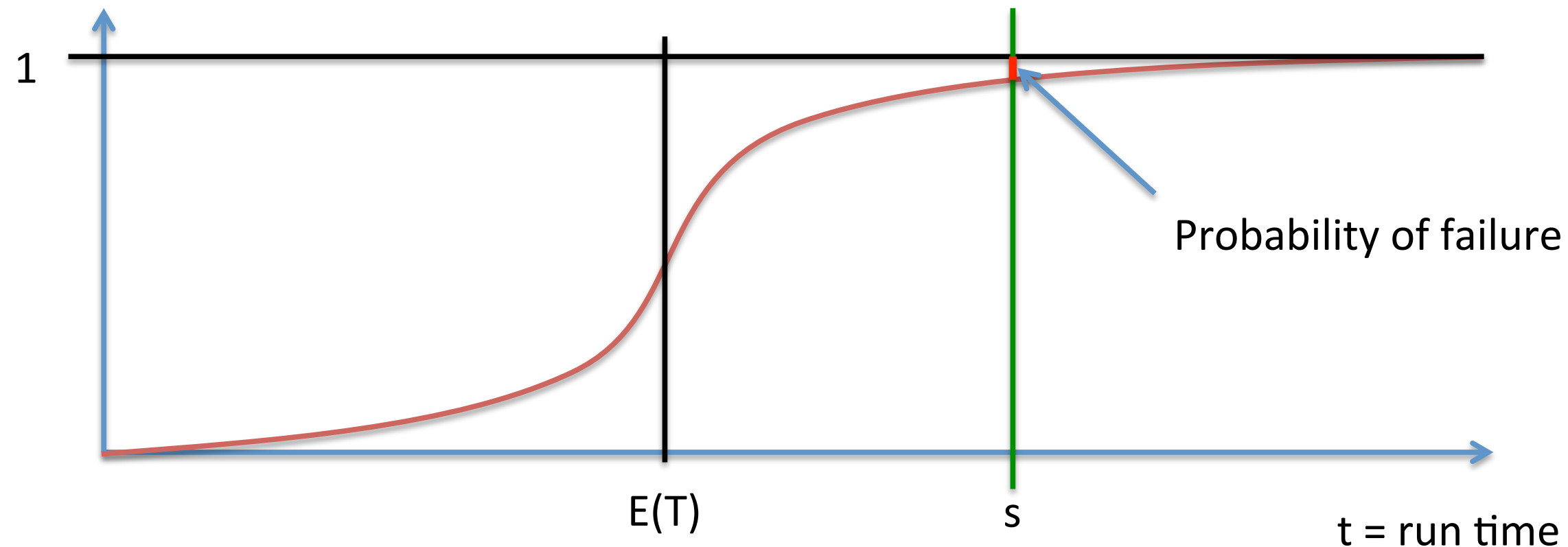
1

$E(T)$

$s$

$t = \text{run time}$

Probability of failure



# From Monte-Carlo to Las Vegas

A Monte-Carlo algorithm always takes the same same time to complete but has probability  $p$  of not generating the correct answer.

We want to transform it into a Las Vegas algorithm, i.e. one that always outputs the correct answer has a varying run time with finite expected value.

Method: run the MC algorithm repeatedly until the answer that it outputs is correct.

Same analysis as the expected number of flips of a biased coin until the first “heads”.

$p$  = the probability that MC fails.

$I$  = the number of times MC is run until the first correct answer.

$$\Pr(I = i) = p^{i-1}(1 - p)$$

$$E(I) = (1 - p) + 2p(1 - p) + 3p^2(1 - p) + \dots =$$

$$= (1 - p) \sum_{i=1}^{\infty} ip^{i-1} = \frac{1-p}{p} \sum_{i=1}^{\infty} ip^i = \frac{1-p}{p} \frac{p}{(1-p)^2} = \frac{1}{1-p}$$

Example: What is the expected number of iteration If the MC algorithm Succeeds with probability  $\frac{1}{4}$  ?