Data Science 1

STAT/CS 287
Jim Bagrow, UVM Dept of Math and Statistics

LECTURE 10

Here's an incredibly useful and powerful idea!

Suppose we have a sample of numeric data (a list of numbers)

$$x_1, x_2, \ldots, x_N$$

and we want to know something about its probability distribution P(x)

We could compute the histogram

Here's another idea→

Question

Imagine we take our data and **sort** (or rank) the numbers from smallest to largest.

Meaning we now know that

$$x_1 \le x_2 \le \cdots \le x_N$$
 is true

You can think of this sorting as *computing* a new ordering or indexing of the points (replacing the subscripts)

What else have we computed?

i	x_i	
1	\boldsymbol{x}_1	
2	\boldsymbol{x}_2	
3	\boldsymbol{x}_3	
• •	• •	
N	$\mathcal{X}_{\mathcal{N}}$	

Let me swap these columns

$\boldsymbol{x_i}$	i	
\boldsymbol{x}_1	1	
\mathcal{X}_2	2	
\mathcal{X}_3	3	
• •	• •	
$\mathcal{X}_{\mathcal{N}}$	N	

How does i relate to x_i ?

x_i	i	<i>i</i> -1	
\boldsymbol{x}_1	1	0	
\mathcal{X}_2	2	1	
\mathcal{X}_3	3	2	
•	• •	• •	
x_N	N	N-1	

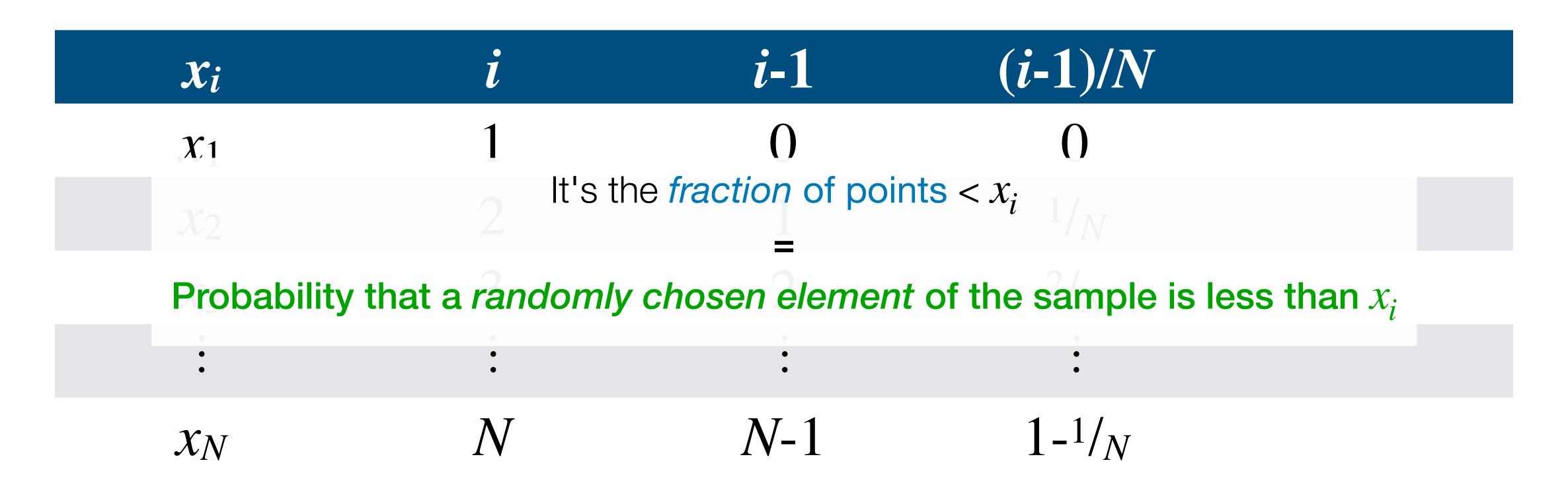
What about i-1?

It's the number of points $< x_i$

x_i	i	<i>i</i> -1	(i-1)/N	
\boldsymbol{x}_1	1	0	0	
\boldsymbol{x}_2	2	1	1/N	
x_3	3	2	2/N	
• •	• •	• •	• •	
$\mathcal{X}_{\mathcal{N}}$	N	N-1	1-1/N	

What about (i-1)/N?

It's the *fraction* of points $< x_i$



$\boldsymbol{x_i}$	i	<i>i</i> -1	(i-1)/N	$\approx P(X < x_i)$
\boldsymbol{x}_1	1	0	0	0
\boldsymbol{x}_2	2	1	1/N	1/N
\mathcal{X}_3	3	2	2/N	2/N
•	• •	•	• •	• •
$\mathcal{X}_{\mathcal{N}}$	N	N-1	1-1/N	1-1/N

It's the *fraction* of points $< x_i$

Probability that a randomly chosen element of the sample is less than x_i

P(X < x) is the cumulative distribution function (CDF)

For a random variable X with associated probability distribution P(x) If X is continuous

$$P(X < x) = \int_{-\infty}^{x} P(x) dx$$

If X is discrete

$$P(X \le x) = \sum_{x_i \le x} P(X = x_i) = \sum_{x_i \le x} P(x_i)$$

Complementary cumulative distribution function (CCDF): $P(X \ge x) = 1 - P(X < x)$

Back to our table

x_i	(i-1)/N	$\approx P(X < x_i)$
\boldsymbol{x}_1	0	0
\mathcal{X}_2	1/N	1/N
x_3	2/N	2/N
• •	• •	• •
x_N	1-1/N	1-1/N

What we have is an empirical estimate of the CDF (Empirical CDF = ECDF)

$$P(X < x) \approx \frac{\text{(number of datapoints} < x)}{\text{(number of datapoints)}}$$
$$= \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}_{x_i < x}$$

The *fraction* of points $< x_i$

Recall our question

Imagine we take our data and **sort** (or rank) the numbers from smallest to largest.

Meaning we now know that

$$x_1 \le x_2 \le \cdots \le x_N$$
 is true

You can think of this sorting as *computing* a new ordering or indexing of the points (replacing the subscripts)

What else have we computed?

Recall our question

Imagine we take our data and **sort** (or rank) the numbers from smallest to largest.

You can think of this sorting as *computing* a new ordering or indexing of the points (replacing the subscripts)

We have computed the cumulative distribution

Meaning we now know that

$$x_1 \le x_2 \le \cdots \le x_N$$
 is true

What else have we computed?

sorting = integrating

Another question

x_i	(i-1)/N	$\approx P(X < x_i)$
\boldsymbol{x}_1	0	0
\mathcal{X}_2	$1/_N$	$1/_N$
χ_3	2/N	$2/_N$
• •	• •	• •
$\mathcal{X}_{\mathcal{N}}$	1-1/N	$1-1/_{N}$

What happens if we plot these columns?

P(X < x) as a function of x...

(x-axis) — (y-axis

Back to the notebook

