



# **DS-GA 3001.009**

## **Responsible Data Science Lab 7**

**Center for Data Science**  
**Haoyue Ping | Tandon School of Engineering**

## **Introduction**

## **Download DataSynthesizer and setup the running environment**

## **DataSynthesizer usage**

- **Random mode**
- **Independent attribute mode**
- **Correlated attribute mode**

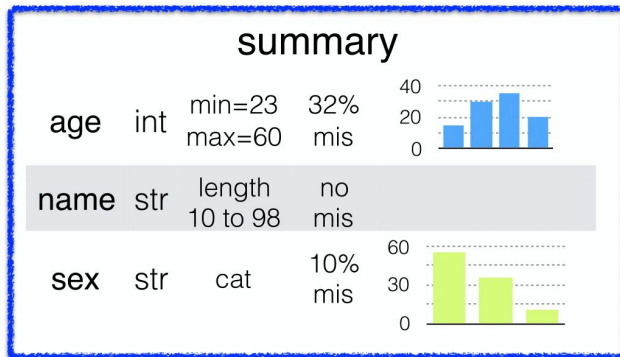
## **Some useful statistical measures**

# Introduction

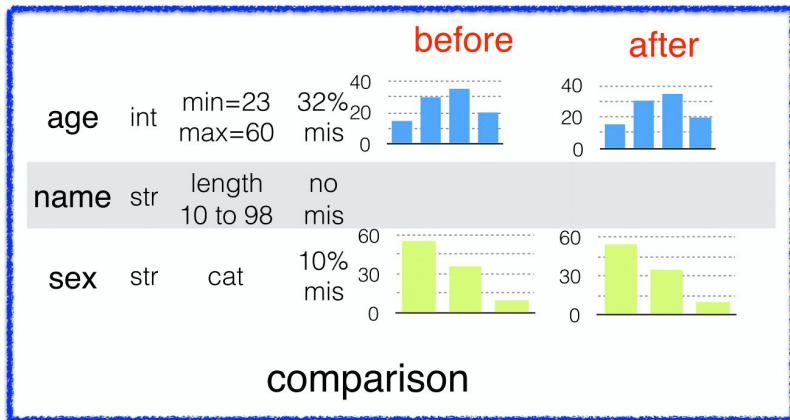
input

id	sex	name	MarriageDate	BirthDate	age	sex	MarriageDate	BirthDate	age
1	1	1	1/1/17	1/1/17	0	1	1	1/1/17	0
2	2	2	1/1/17	1/1/17	0	2	2	1/1/17	0
3	3	3	1/1/17	1/1/17	0	3	3	1/1/17	0
4	4	4	1/1/17	1/1/17	0	4	4	1/1/17	0
5	5	5	1/1/17	1/1/17	0	5	5	1/1/17	0
6	6	6	1/1/17	1/1/17	0	6	6	1/1/17	0
7	7	7	1/1/17	1/1/17	0	7	7	1/1/17	0
8	8	8	1/1/17	1/1/17	0	8	8	1/1/17	0
9	9	9	1/1/17	1/1/17	0	9	9	1/1/17	0
10	10	10	1/1/17	1/1/17	0	10	10	1/1/17	0
11	11	11	1/1/17	1/1/17	0	11	11	1/1/17	0
12	12	12	1/1/17	1/1/17	0	12	12	1/1/17	0
13	13	13	1/1/17	1/1/17	0	13	13	1/1/17	0
14	14	14	1/1/17	1/1/17	0	14	14	1/1/17	0
15	15	15	1/1/17	1/1/17	0	15	15	1/1/17	0
16	16	16	1/1/17	1/1/17	0	16	16	1/1/17	0
17	17	17	1/1/17	1/1/17	0	17	17	1/1/17	0
18	18	18	1/1/17	1/1/17	0	18	18	1/1/17	0
19	19	19	1/1/17	1/1/17	0	19	19	1/1/17	0
20	20	20	1/1/17	1/1/17	0	20	20	1/1/17	0
21	21	21	1/1/17	1/1/17	0	21	21	1/1/17	0
22	22	22	1/1/17	1/1/17	0	22	22	1/1/17	0
23	23	23	1/1/17	1/1/17	0	23	23	1/1/17	0
24	24	24	1/1/17	1/1/17	0	24	24	1/1/17	0
25	25	25	1/1/17	1/1/17	0	25	25	1/1/17	0
26	26	26	1/1/17	1/1/17	0	26	26	1/1/17	0
27	27	27	1/1/17	1/1/17	0	27	27	1/1/17	0
28	28	28	1/1/17	1/1/17	0	28	28	1/1/17	0

Data  
Describer



Data  
Generator



Model  
Inspector



output

id	sex	name	MarriageDate	BirthDate	age	sex	MarriageDate	BirthDate	age
1	1	1	1/1/17	1/1/17	0	1	1	1/1/17	0
2	2	2	1/1/17	1/1/17	0	2	2	1/1/17	0
3	3	3	1/1/17	1/1/17	0	3	3	1/1/17	0
4	4	4	1/1/17	1/1/17	0	4	4	1/1/17	0
5	5	5	1/1/17	1/1/17	0	5	5	1/1/17	0
6	6	6	1/1/17	1/1/17	0	6	6	1/1/17	0
7	7	7	1/1/17	1/1/17	0	7	7	1/1/17	0
8	8	8	1/1/17	1/1/17	0	8	8	1/1/17	0
9	9	9	1/1/17	1/1/17	0	9	9	1/1/17	0
10	10	10	1/1/17	1/1/17	0	10	10	1/1/17	0
11	11	11	1/1/17	1/1/17	0	11	11	1/1/17	0
12	12	12	1/1/17	1/1/17	0	12	12	1/1/17	0
13	13	13	1/1/17	1/1/17	0	13	13	1/1/17	0
14	14	14	1/1/17	1/1/17	0	14	14	1/1/17	0
15	15	15	1/1/17	1/1/17	0	15	15	1/1/17	0
16	16	16	1/1/17	1/1/17	0	16	16	1/1/17	0
17	17	17	1/1/17	1/1/17	0	17	17	1/1/17	0
18	18	18	1/1/17	1/1/17	0	18	18	1/1/17	0
19	19	19	1/1/17	1/1/17	0	19	19	1/1/17	0
20	20	20	1/1/17	1/1/17	0	20	20	1/1/17	0
21	21	21	1/1/17	1/1/17	0	21	21	1/1/17	0
22	22	22	1/1/17	1/1/17	0	22	22	1/1/17	0
23	23	23	1/1/17	1/1/17	0	23	23	1/1/17	0
24	24	24	1/1/17	1/1/17	0	24	24	1/1/17	0
25	25	25	1/1/17	1/1/17	0	25	25	1/1/17	0
26	26	26	1/1/17	1/1/17	0	26	26	1/1/17	0
27	27	27	1/1/17	1/1/17	0	27	27	1/1/17	0
28	28	28	1/1/17	1/1/17	0	28	28	1/1/17	0

## GitHub repo

<https://github.com/DataResponsibly/DataSynthesizer>

- **Download it**
- **Add `./DataSynthesizer/` into `sys.path`**

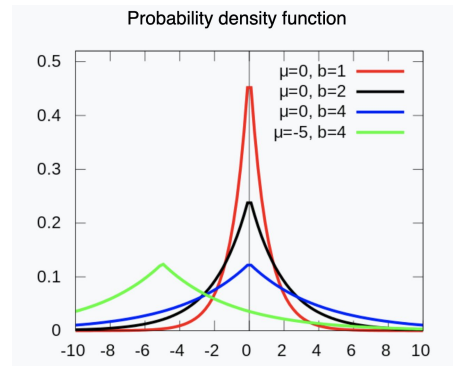
- **Generate type-consistent data**
- **Learn the domains of attributes**
  - Data type
  - Categorical vs non-categorical
    - Threshold = 20 by default
    - True for rating, gender
    - False for score, name
  - Numerical vs non-numerical
    - Integer, Float, Datetime are numerical
    - Datetimes → timestamps if non-categorical
  - Active domain
    - if is\_categorical:
      - Attribute values in dataset
    - else if is\_numerical:
      - Range(min, max)

Data Type	Example
Integer	ID, age
Float	Score, rating
String	Name, gender
Datetime	Birthday, event time

# Independent attribute mode

**Assume the attributes (or columns) are independent.**

- **Run random mode first to get the attribute domains**
- **Model attribute distributions**
  - Bar charts for categorical attributes
  - Histograms for numerical attributes
- **Inject Laplace noise into the bar charts / histograms.**
  - Sensitivity =  $2/n$
  - $d = \text{\#attributes}$ , then privacy budget is  $\epsilon/d$  for each attribute.
  - Inject  $\text{Lap}(2d/n\epsilon)$



## Parameters

- **epsilon**: the privacy budget
- **k**: #parents in Bayesian network (BN)

## Run GreedyBayes to construct a BN

- Connect attributes with high mutual information
- Randomize the attribute connections
- Cost  **$\epsilon/2$** , half of the privacy budget

## Populate conditional probability tables (CPTs)

- Inject Laplace noise into CPTs
- Cost  **$\epsilon/2$** , half of the privacy budget

# Randomize BN structure

---

**Algorithm 1** GreedyBayes( $D, A, k$ )
 

---

**Require:** Dataset  $D$ , set of attributes  $A$ , maximum number of parents  $k$

- 1: Initialize  $\mathcal{N} = \emptyset$  and  $V = \emptyset$ .
  - 2: Randomly select an attribute  $X_1$  from  $A$ .
  - 3: Add  $(X_1, \emptyset)$  to  $\mathcal{N}$ ; add  $X_1$  to  $V$ .
  - 4: **for**  $i = 2, \dots, |A|$  **do**
  - 5:     Initialize  $\Omega = \emptyset$
  - 6:      $p = \min(k, |V|)$
  - 7:     **for** each  $X \in A \setminus V$  and each  $\Pi \in \binom{V}{p}$  **do**
  - 8:         Add  $(X, \Pi)$  to  $\Omega$
  - 9:     **end for**
  - 10:     ~~Compute mutual information based on  $D$  for all pairs in  $\Omega$ .~~
  - 11:     Select  $(X_i, \Pi_i)$  from  $\Omega$  with maximal mutual information.
  - 12:     Add  $(X_i, \Pi_i)$  to  $\mathcal{N}$ .
  - 13: **end for**
  - 14: **return**  $\mathcal{N}$
- 

Select the (child, parents) among all combinations in  $\Omega$  **with a probability proportional to**  $\exp(I(X, \Pi)/2\Delta)$

Where  $I(\cdot)$  is mutual information.

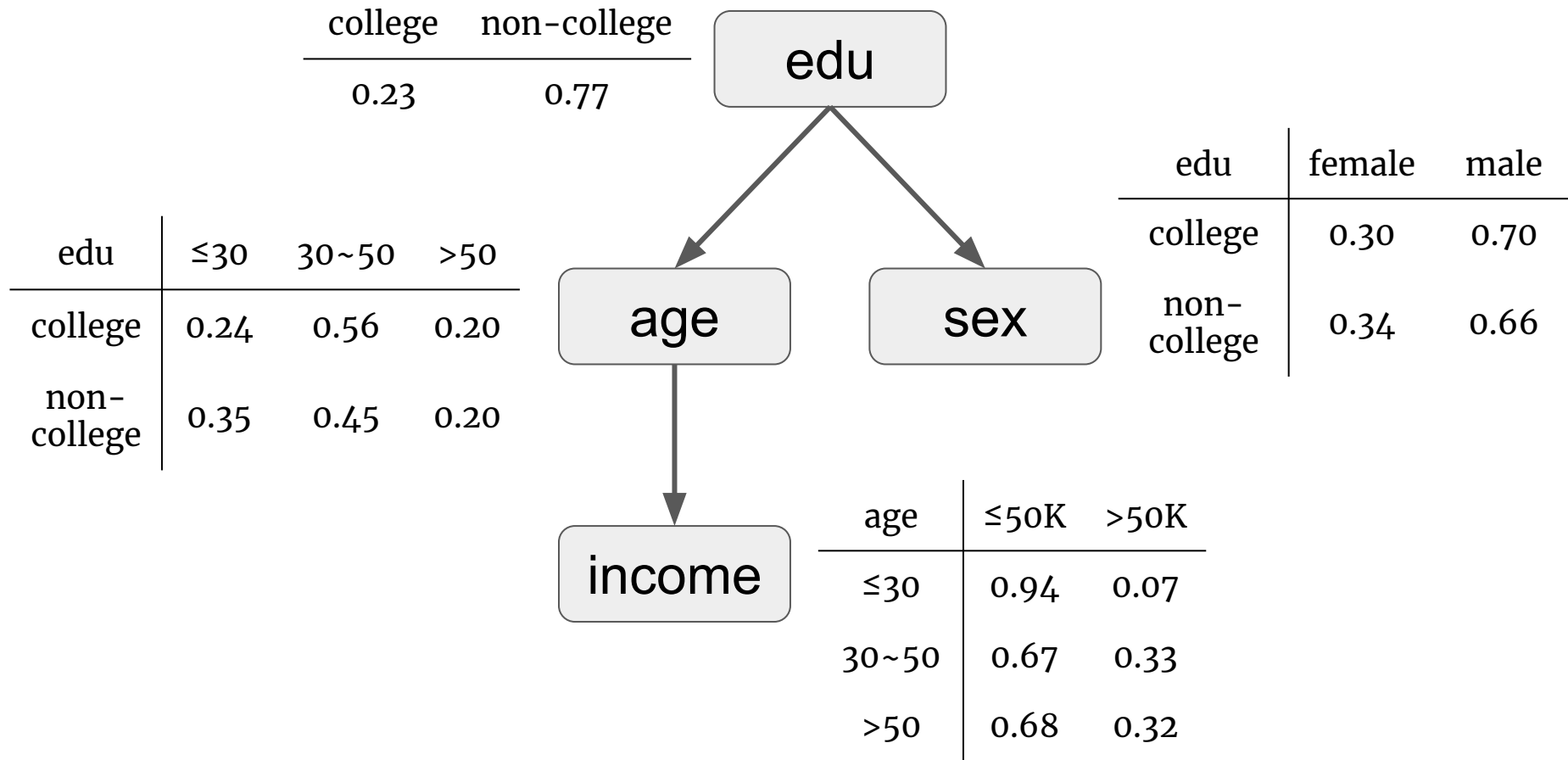
$$\Delta = (d - 1)S(I)/\varepsilon$$

$$S(I(X, \Pi)) = \begin{cases} \frac{1}{n} \log(n) + \frac{n-1}{n} \log\left(\frac{n}{n-1}\right), & \text{if } X \text{ or } \Pi \text{ is binary;} \\ \frac{2}{n} \log\left(\frac{n+1}{2}\right) + \frac{n-1}{n} \log\left(\frac{n+1}{n-1}\right), & \text{otherwise,} \end{cases}$$

$n$  is the number of tuples in  $D$ .



# Randomize BN structure



# Step 0: add root

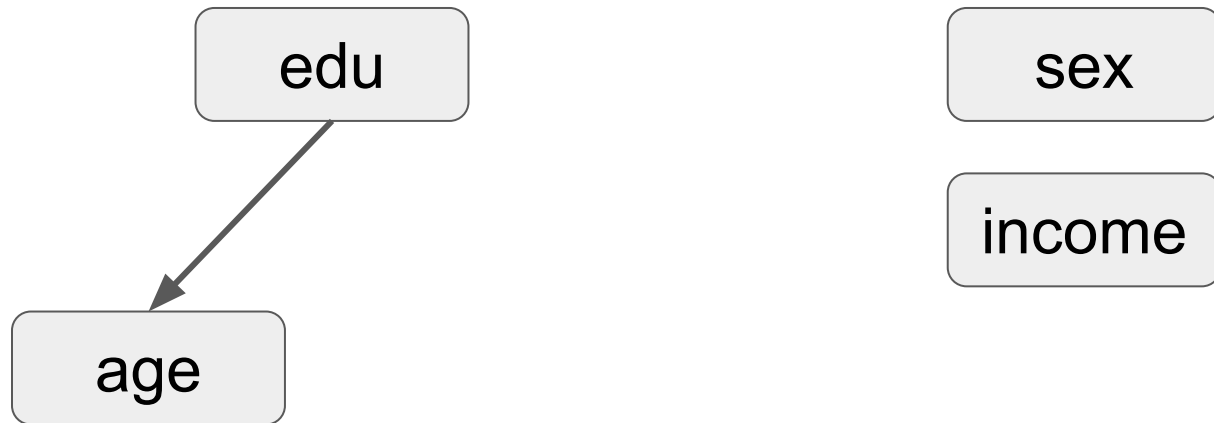
edu

age

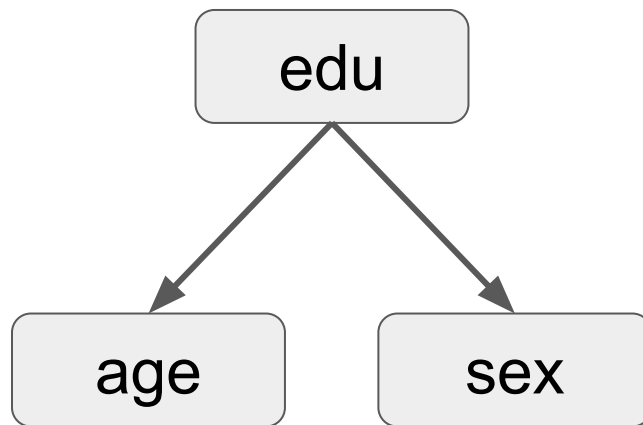
sex

income

# Step 1: add the 1st child

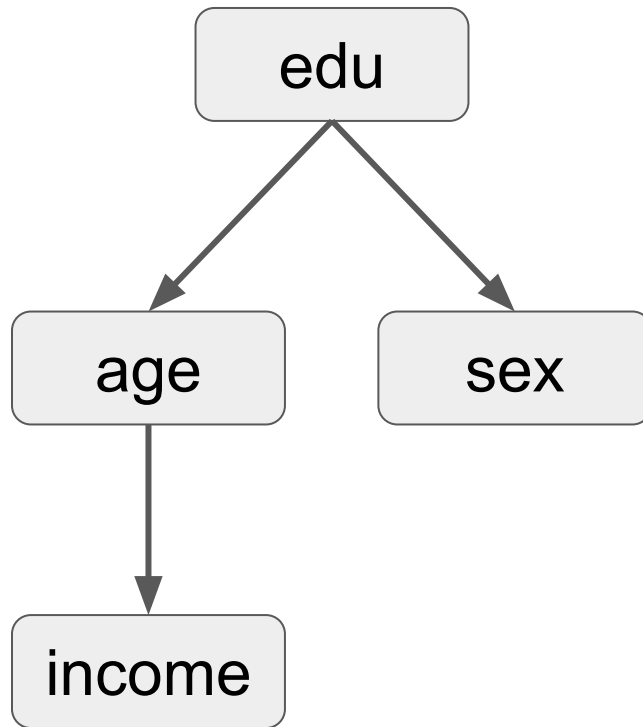


## Step 2: add the 2nd child



income

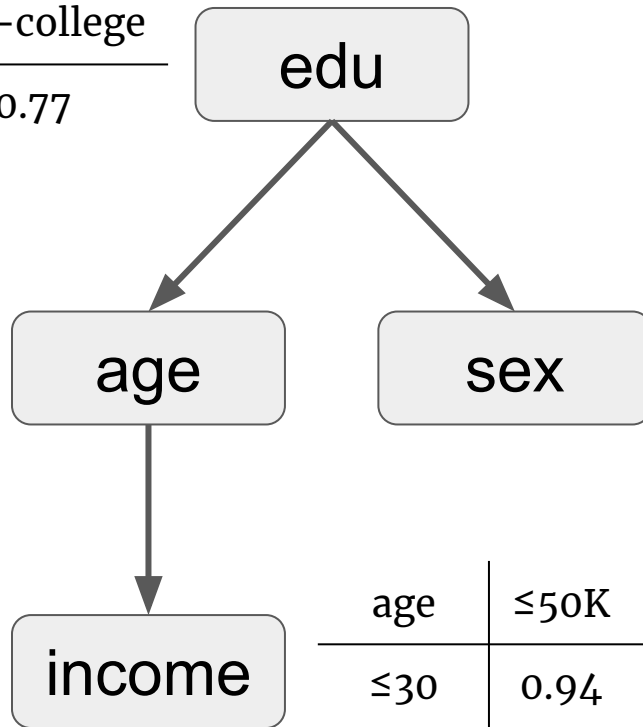
# Step 3: add the 3rd child



# CPTs with $\text{Lap}(4(d-k)/(n \cdot \epsilon))^*$ noise

college	non-college
0.23	0.77

edu	$\leq 30$	30~50	$> 50$
college	0.24	0.56	0.20
non-college	0.35	0.45	0.20



edu	female	male
college	0.30	0.70
non-college	0.34	0.66

age	$\leq 50K$	$> 50K$
$\leq 30$	0.94	0.07
30~50	0.67	0.33
$> 50$	0.68	0.32

\* $\text{Lap}(4(d-k)/(n \cdot \epsilon))$

- $d$ =#attributes in BN
- $k$ =#parents in BN
- $n$ =sensitive dataset size

## **Mutual information**

- **How much information can be obtained from one random variable about another random variable?**

## **Two-sample Kolmogorov–Smirnov test**

- **How different are two continuous distributions?**

## **KL-divergence**

- **How different are two categorical distribution?**

# Mutual information\*

- **The “amount of information” obtained from one random variable about another random variable.**

$$I(X; Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x, y) \log \left( \frac{p(x, y)}{p(x) p(y)} \right)$$

- **MI(X, Y) = 0 if random variables X and Y are independent**

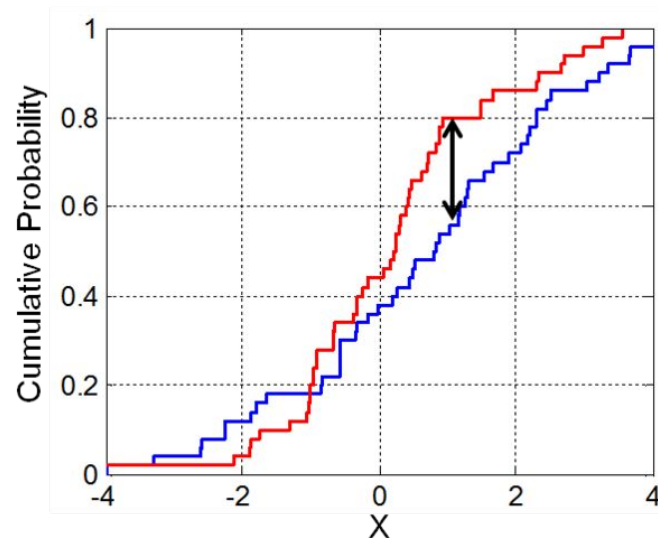
$$\log \left( \frac{p(x, y)}{p(x) p(y)} \right) = \log 1 = 0$$



# Two-sample Kolmogorov–Smirnov test\*

- **Test whether two underlying one-dimensional probability distributions differ.**

$$D_{n,m} = \sup_x |F_{1,n}(x) - F_{2,m}(x)|$$



# KL-divergence\*

- How different are two categorical distribution P and Q?

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

- $D_{\text{KL}}(P \parallel Q) = 0$  if P and Q are identical.
- The KL-divergence is defined **only if** for all x,  $Q(x)=0$  implies  $P(x)=0$

**Thank you!**