# PrivBayes: Private Data Release via Bayesian Networks

**Jun Zhang[1]**    **Graham Cormode[2]**    **Cecilia Procopiuc[3]**    **Divesh Srivastava[3]**    **Xiaokui Xiao[1]**
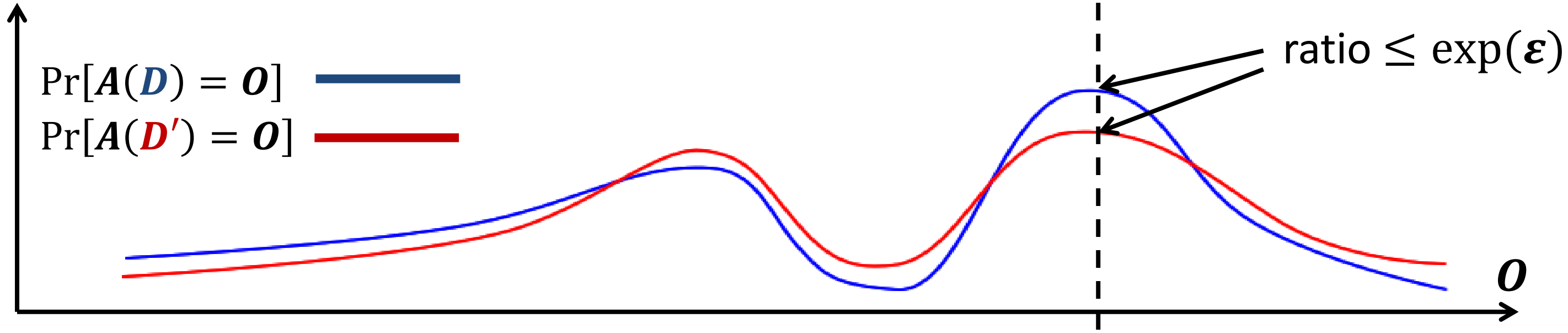
1. Nanyang Technological University
{jzhang027, xkxiao}@ntu.edu.sg

2. University of Warwick
g.cormode@warwick.ac.uk

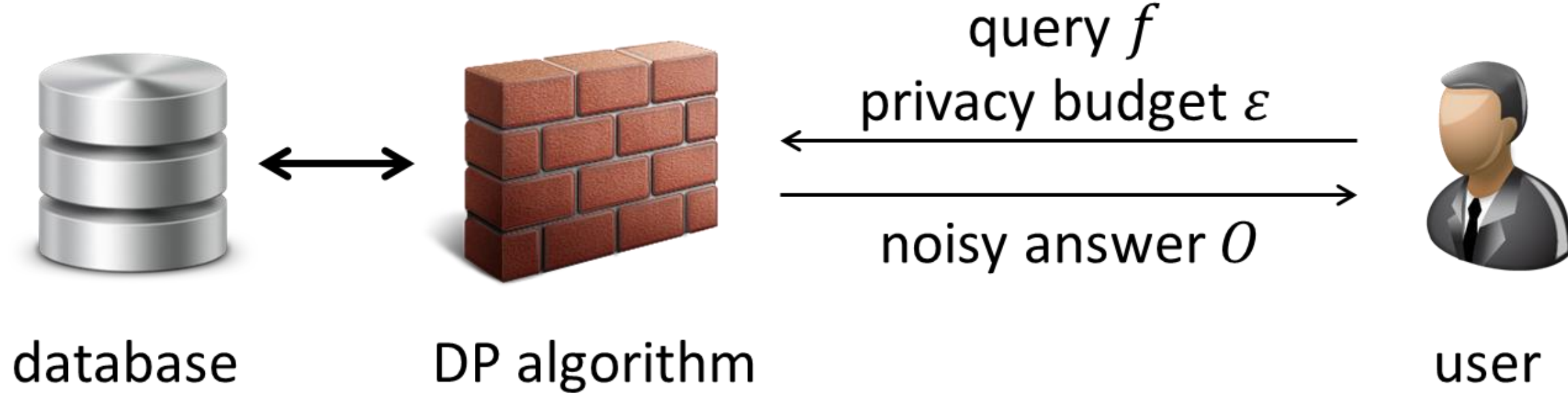3. AT&T Labs -- Research
{magda, divesh}@research.att.com

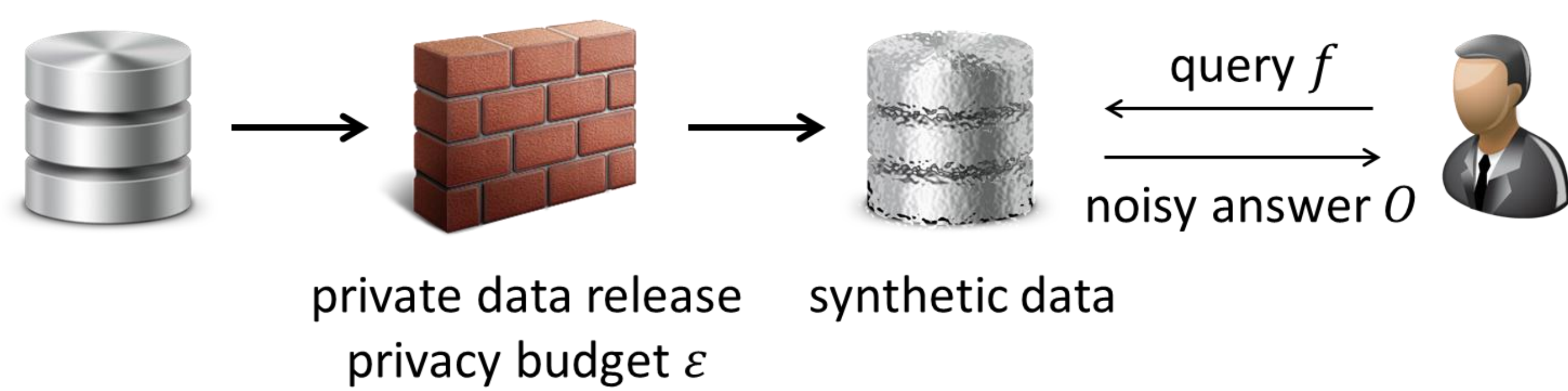## 1. Private Data Release

❑ Differential privacy



Pr[$A(D) = O$]
Pr[$A(D') = O$]

ratio ≤ exp($\varepsilon$)

where $D$ and $D'$ are neighboring databases that differ by **at most one** tuple

$$\exp(-\varepsilon) \le \frac{\Pr[A(D) = O]}{\Pr[A(D') = O]} \le \exp(\varepsilon)$$

❑ Interactive mode



query $f$
privacy budget $\varepsilon$
noisy answer $O$

database        DP algorithm        user

❑ Non-interactive mode (synthetic data release)



query $f$
noisy answer $O$

private data release
privacy budget $\varepsilon$
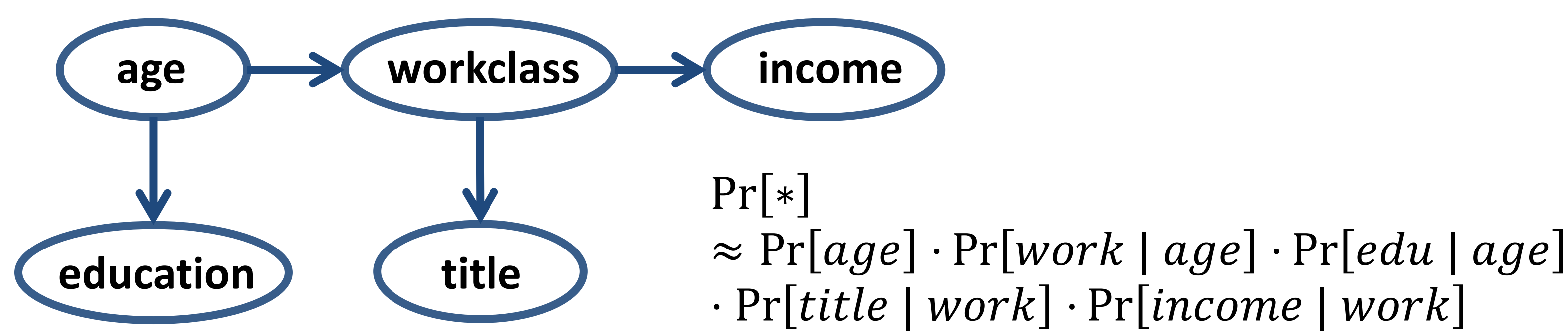
synthetic data

Reusability: only access sensitive data once    Generality: support most queries
However, the tuple distribution has a huge domain (exponential to dimension), which leads to high computational cost and low signal-to-noise ratio.

## 2. Private Bayesian Network

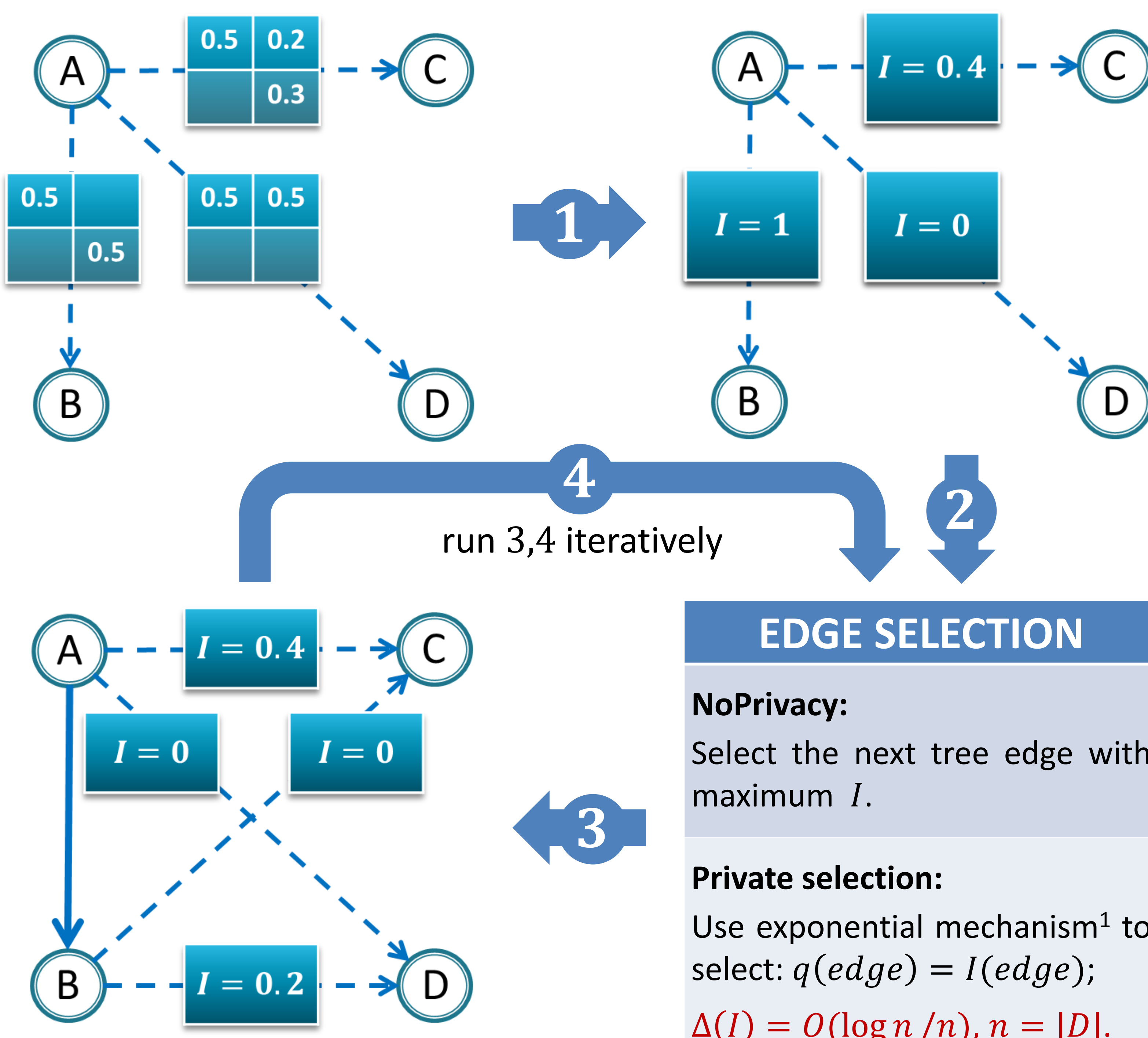❑ Approximate full distribution by low dimensional ones



age → workclass → income
education    title

$$\Pr[*] \approx \Pr[age] \cdot \Pr[work \mid age] \cdot \Pr[edu \mid age] \cdot \Pr[title \mid work] \cdot \Pr[income \mid work]$$

❑ Build a Bayesian network

The quality of a Bayesian network is measured by sum of mutual information $I$ of its edges. $I$ is defined as

$$I(X,Y) = \sum_{y \in Y} \sum_{x \in X} \Pr[x,y] \log\left(\frac{\Pr[x,y]}{\Pr[x]\Pr[y]}\right).$$



run 3,4 iteratively

### EDGE SELECTION

**NoPrivacy:**
Select the next tree edge with maximum $I$.

**Private selection:**
Use exponential mechanism[1] to select: $q(edge) = I(edge)$;
$\Delta(I) = O(\log n / n), n = |D|$.

[1] Frank McSherry and Kunal Talwar. "Mechanism design via differential privacy." FOCS'07

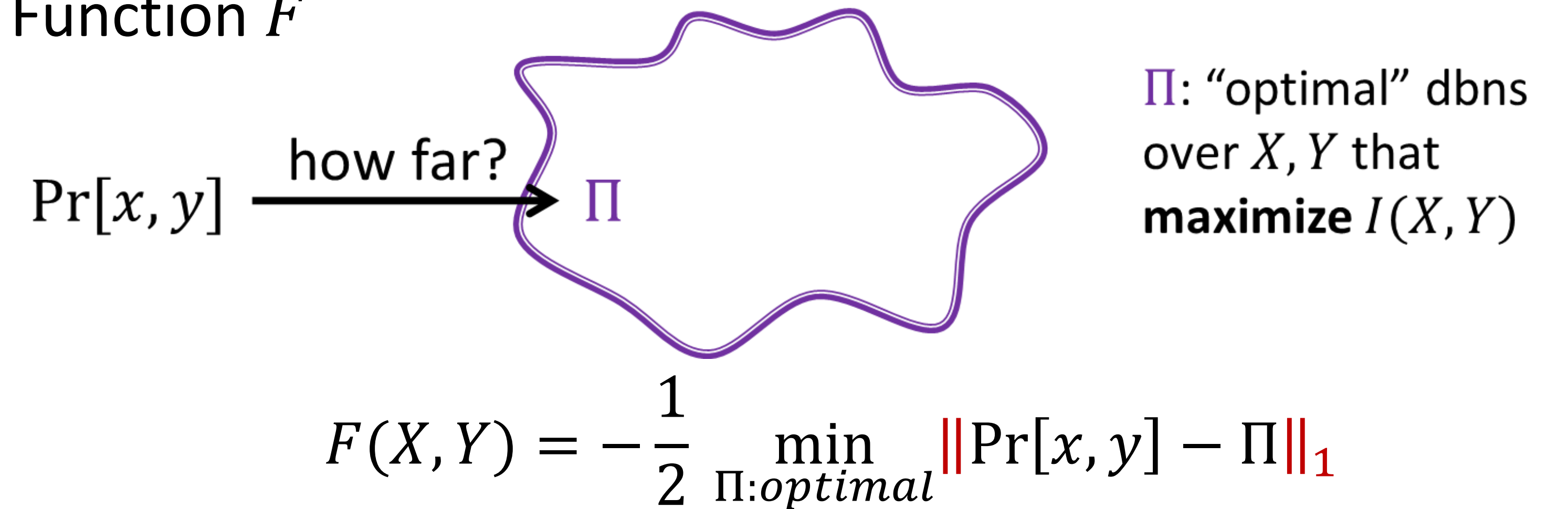## 3. Function $F$:  Linear vs. Logarithmic

❑ Drawback of Function $I$

| range (info) | $O(1)$ | sensitivity (noise) | $O(\log n / n)$ |
|---|---|---|---|

**Problem**: low info-noise ratio

**Solution**: design a new function $F$ that (i) has a higher info-noise ratio; (ii) has a strong positive correlation with $I$
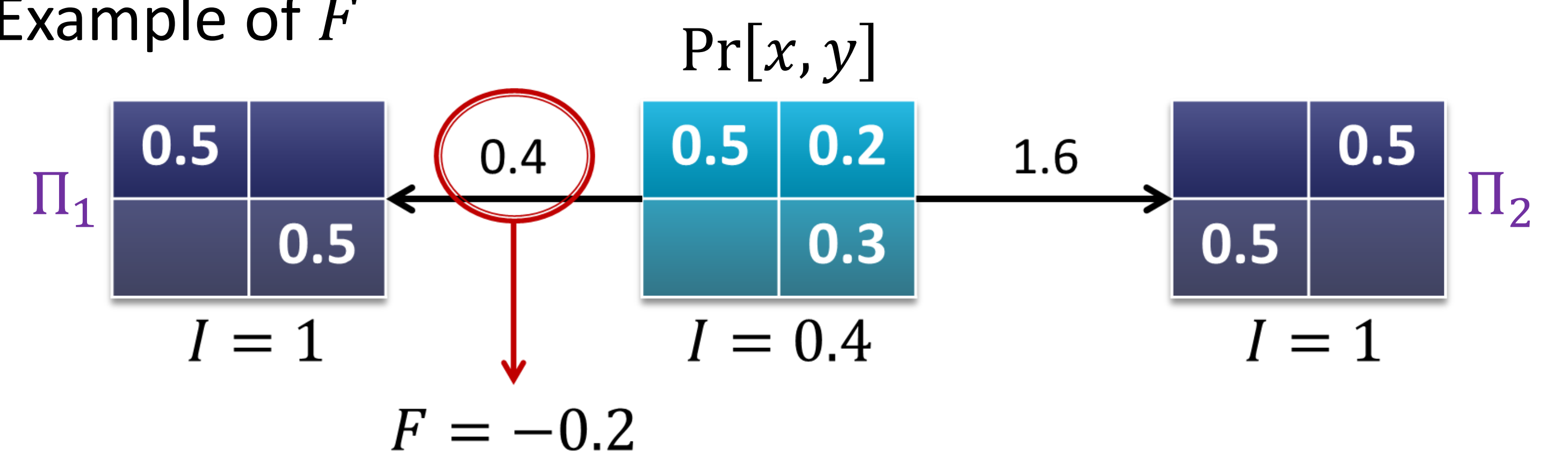
**Idea:** define $F$ to agree with $I$ at maximum values and interpolate linearly in-between
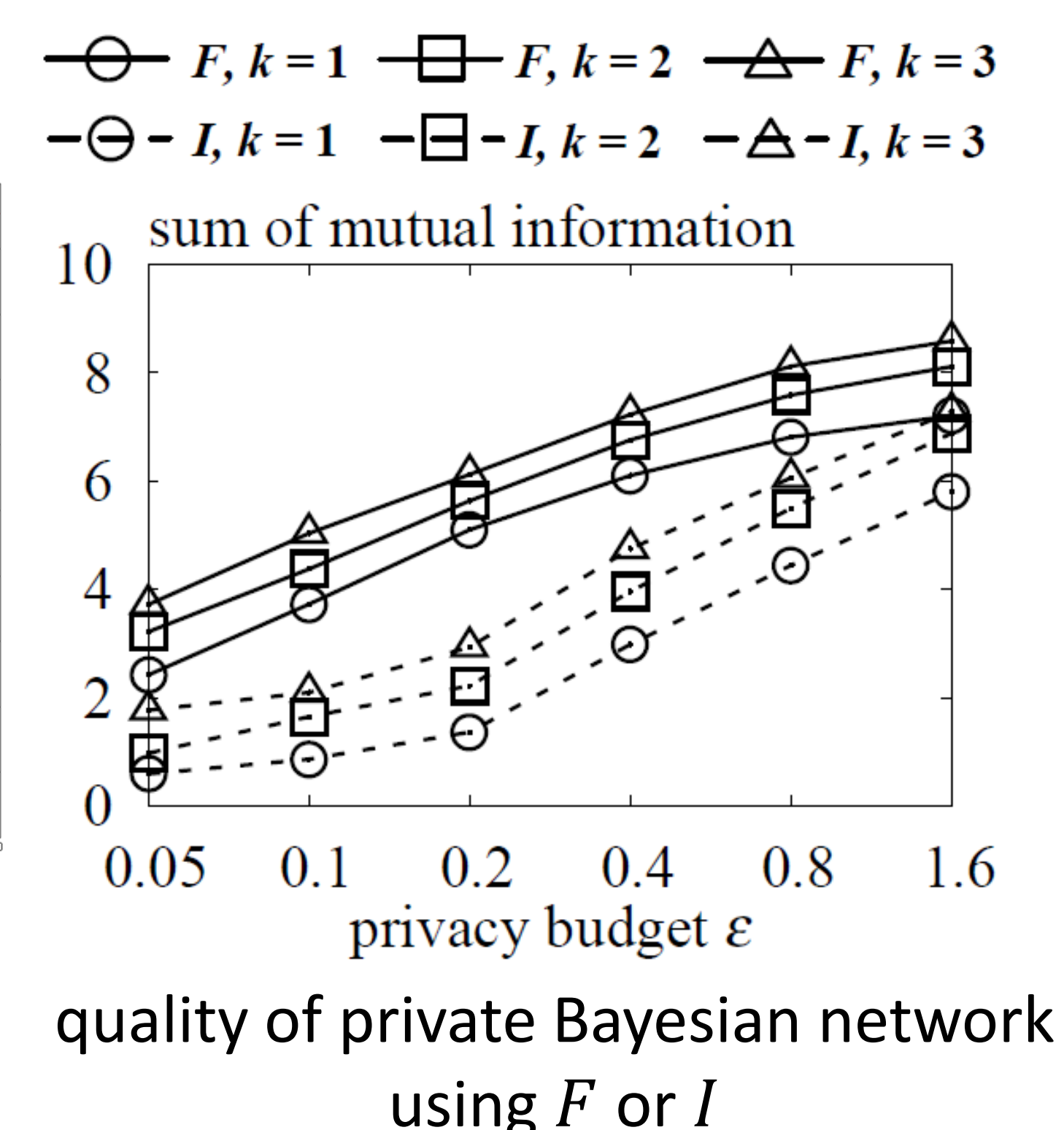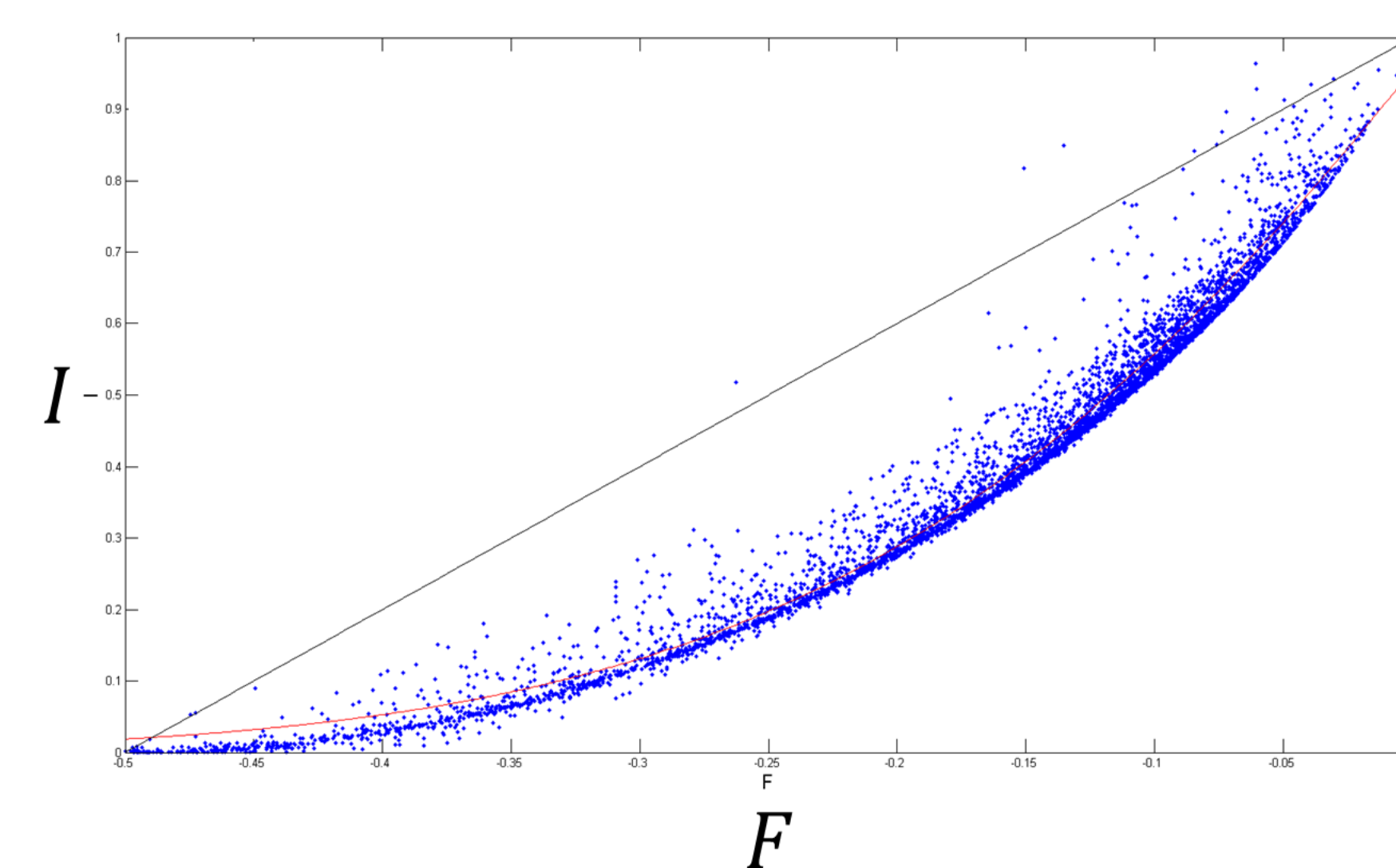
❑ Function $F$



$\Pr[x,y]$ — how far? → $\Pi$

$\Pi$: "optimal" dbns over $X, Y$ that **maximize** $I(X,Y)$

$$F(X,Y) = -\frac{1}{2} \min_{\Pi:optimal} \|\Pr[x,y] - \Pi\|_1$$

| range (info) | $O(1)$ | sensitivity (noise) | $O(1/n)$ |
|---|---|---|---|

❑ Example of $F$



Pr[$x,y$]

$\Pi_1$ | 0.5 / 0.5 |    0.4    | 0.5  0.2 / 0.3 |    1.6    | 0.5 / 0.5 | $\Pi_2$

$I = 1$        $I = 0.4$        $I = 1$

$F = -0.2$

❑ $F$ vs. $I$



$F$ and $I$ of random distributions
correlation coefficient $r = 0.9472$

quality of private Bayesian network using $F$ or $I$

sum of mutual information

privacy budget $\varepsilon$

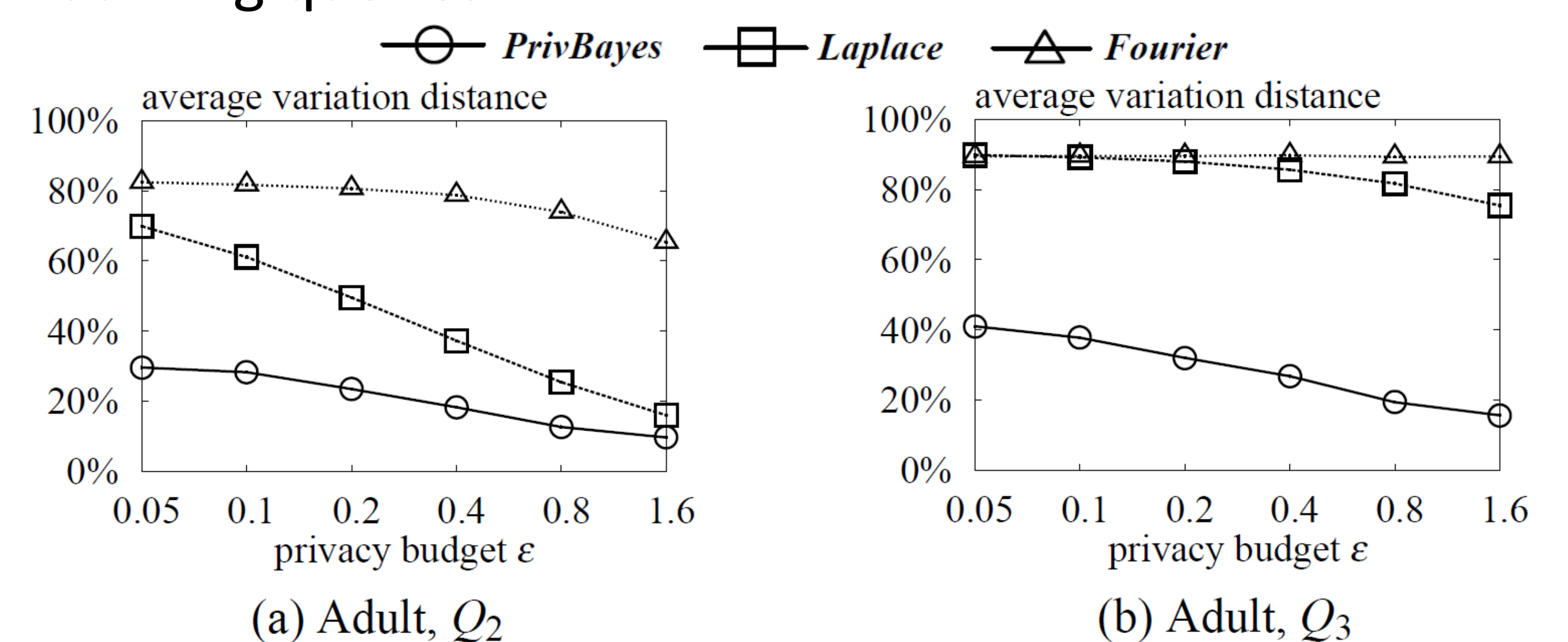— $F, k=1$  — $F, k=2$  — $F, k=3$
— $I, k=1$  — $I, k=2$  — $I, k=3$

## 4. Experiments

We apply PrivBayes to generate one synthetic dataset, to answer a set of counting and SVM training queries simultaneously.

❑ Counting queries



*PrivBayes*    *Laplace*    *Fourier*

average variation distance

privacy budget $\varepsilon$

(a) Adult, $Q_2$

(b) Adult, $Q_3$

❑ Multiple SVMs



*NoPrivacy*    *PrivBayes*    *PrivateERM*
*PrivateERM (Single)*    *PrivGene*    *Majority*

misclassification rate

privacy budget $\varepsilon$

(a) Adult, $Y$ = gender

(b) Adult, $Y$ = education