



Basics of Differential Privacy

Zhenjie Zhang

Advanced Digital Sciences Center, Singapore
(Thanks to Xiaokui Xiao for contributing slides)

Formulation of Privacy

- What information can be published?
 - Average height of US people 
 - Height of an individual 
- Intuition:
 - If something is insensitive to the change of any individual tuple, then it should not be considered private
- Example:
 - Assume that we arbitrarily change the height of an individual in the US
 - The average height of US people would remain roughly the same
 - i.e., The average height reveals little information about the exact height of any particular individual

ϵ -Differential Privacy

- Definition:

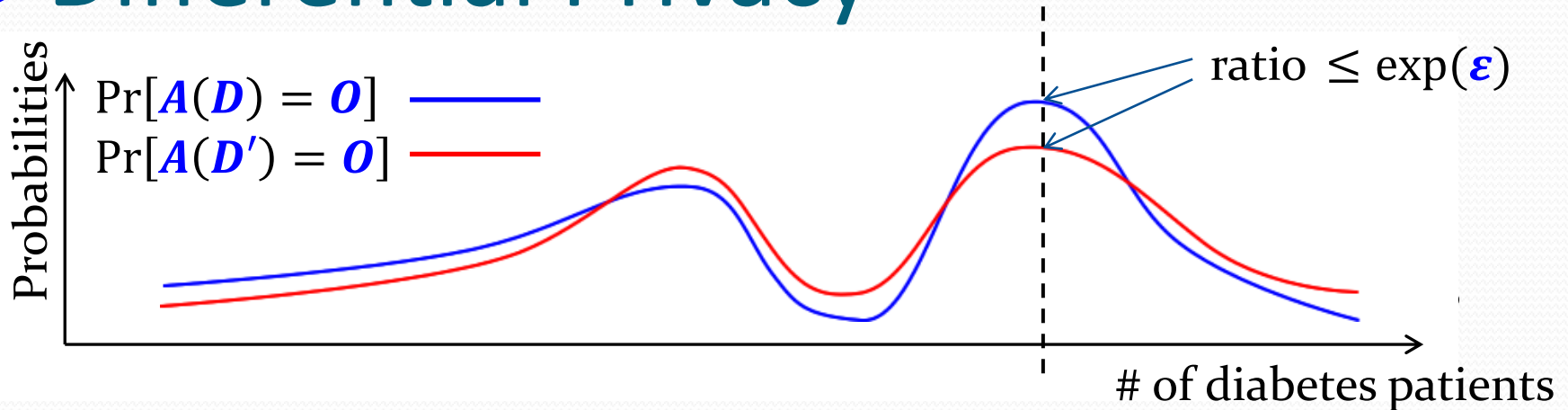
- Neighboring datasets: Two datasets D and D' , such that D' can be obtained by changing one single tuple in D
- A randomized algorithm A satisfies ϵ -differential privacy, iff for any two neighboring datasets D and D' and for any output O of A ,

$$\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$$

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	23	Y
Doug	M	30	N

ϵ -Differential Privacy



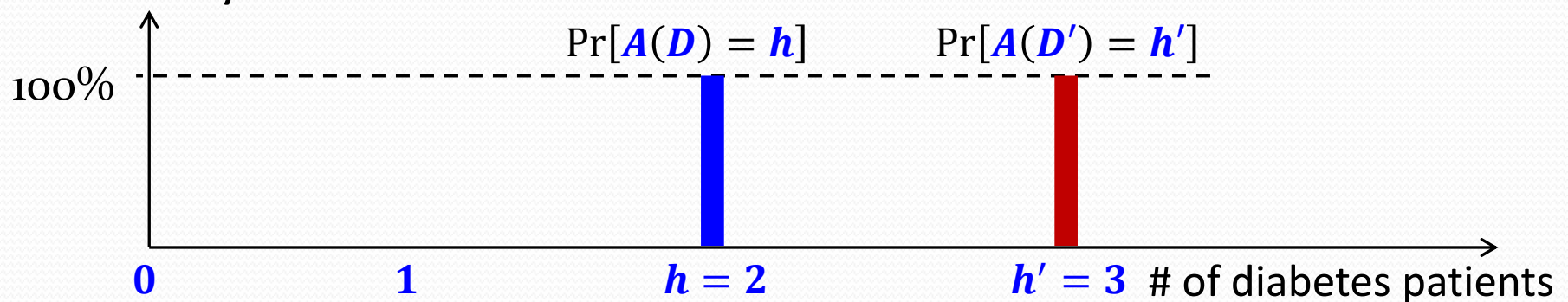
- Definition:
 - Neighboring datasets: Two datasets D and D' , such that D' can be obtained by changing one single tuple in D
 - A **randomized** algorithm A satisfies **ϵ -differential privacy**, iff for **any** two neighboring datasets D and D' and for any output O of A ,
$$\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$$
 - The value of ϵ decides the degree of privacy protection

Achieving ϵ -Differential Privacy

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	23	Y
Doug	M	30	N

- It won't work if we release the number directly:
 - D : the original dataset
 - D' : modify an arbitrary patient in D
 - $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$ does not hold for any ϵ

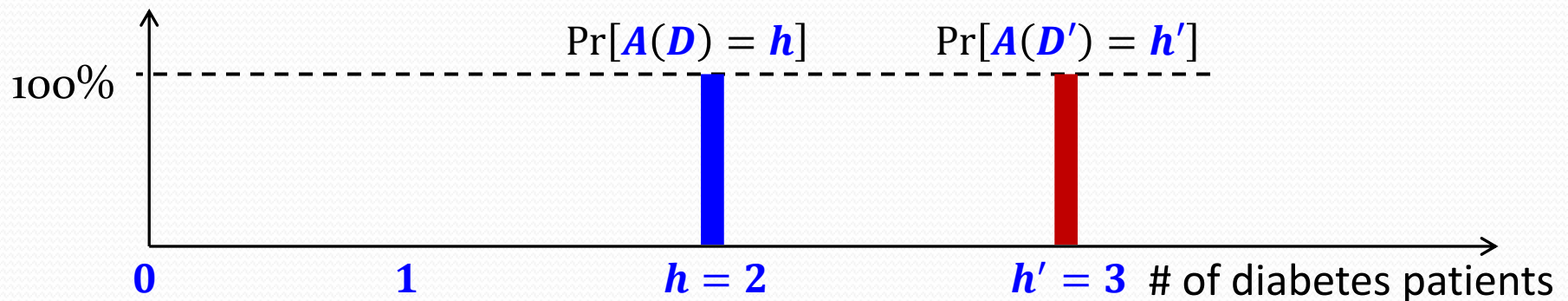


Achieving ϵ -Differential Privacy

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	23	Y
Doug	M	30	N

- Idea:
 - Perturb the number of diabetes patients to obtain a smooth distribution

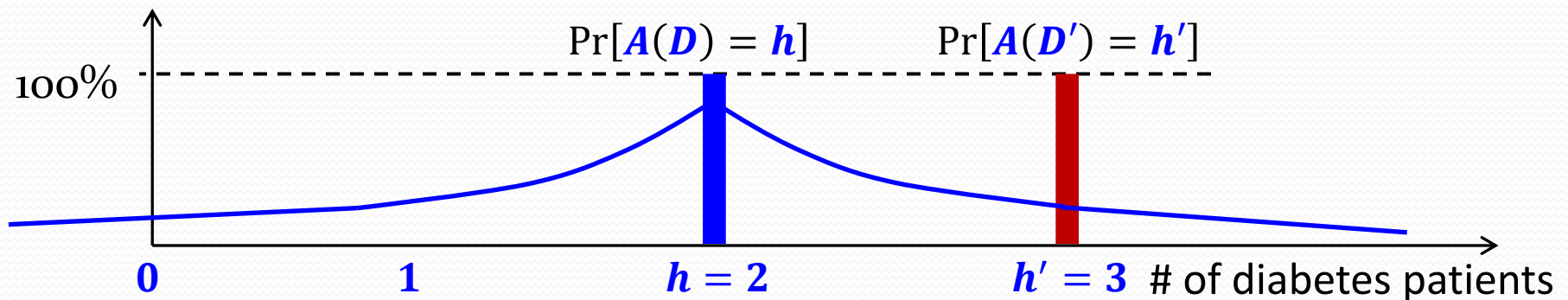


Achieving ϵ -Differential Privacy

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	23	Y
Doug	M	30	N

- Idea:
 - Perturb the number of diabetes patients to obtain a smooth distribution

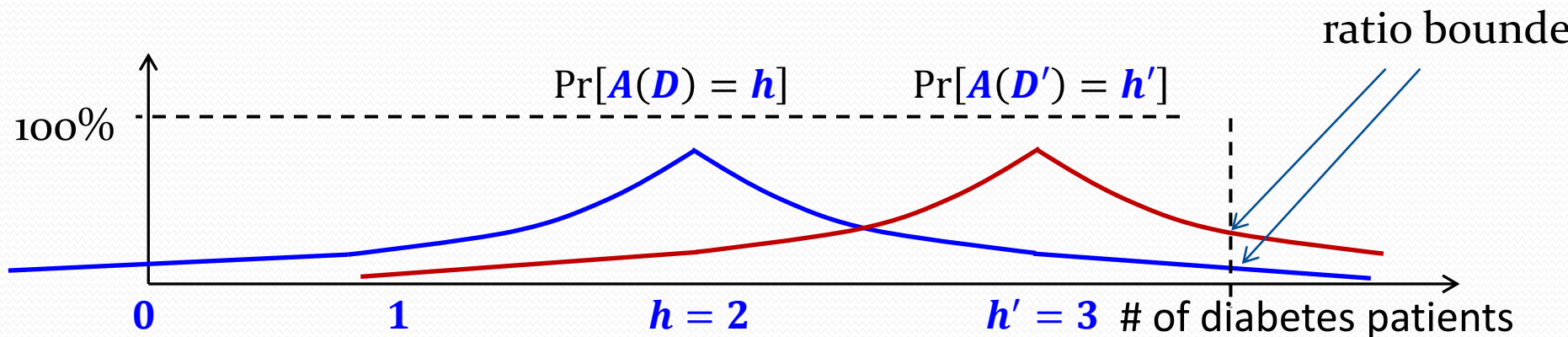


Achieving ϵ -Differential Privacy

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N

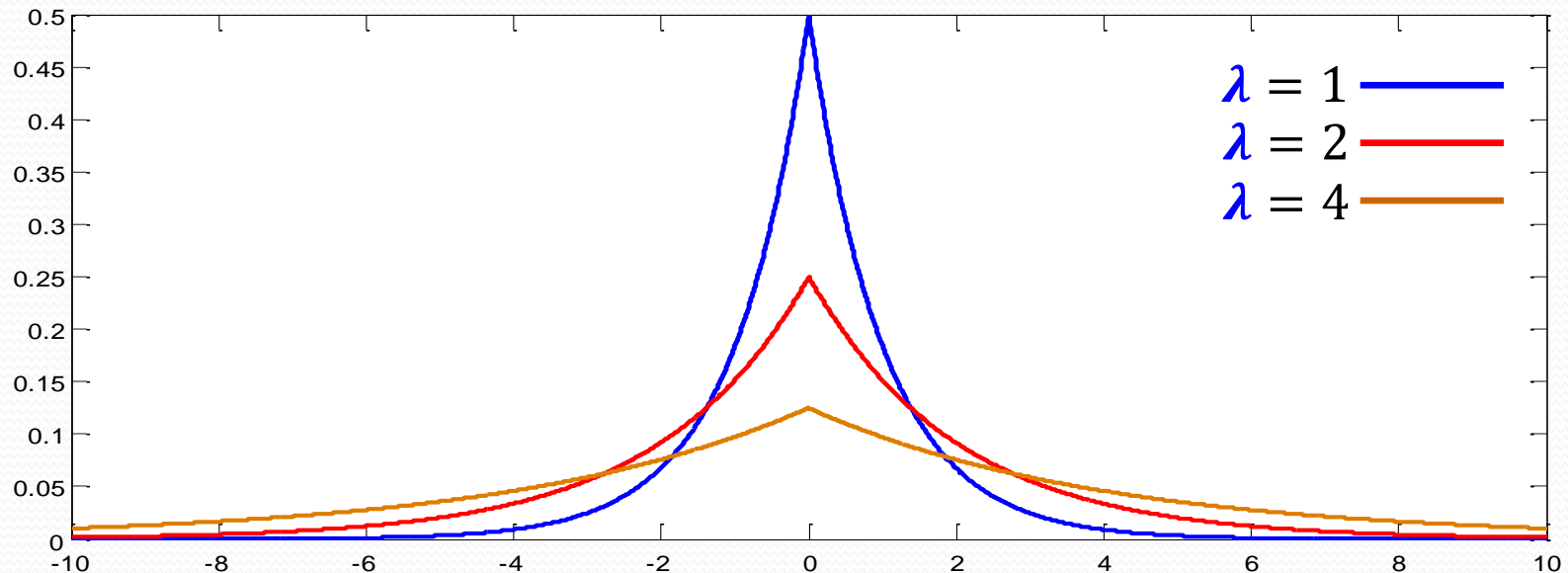
Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	23	Y
Doug	M	30	N

- Idea:
 - Perturb the number of diabetes patients to obtain a smooth distribution



Laplace Distribution

- $pdf(\mathbf{x}) = \exp\left(-\frac{|\mathbf{x}|}{\lambda}\right)/2\lambda$;
- increase/decrease \mathbf{x} by 1
- $\rightarrow pdf(\mathbf{x})$ changes by a factor of $\exp\left(-\frac{1}{\lambda}\right)$
- λ is referred as the *scale*



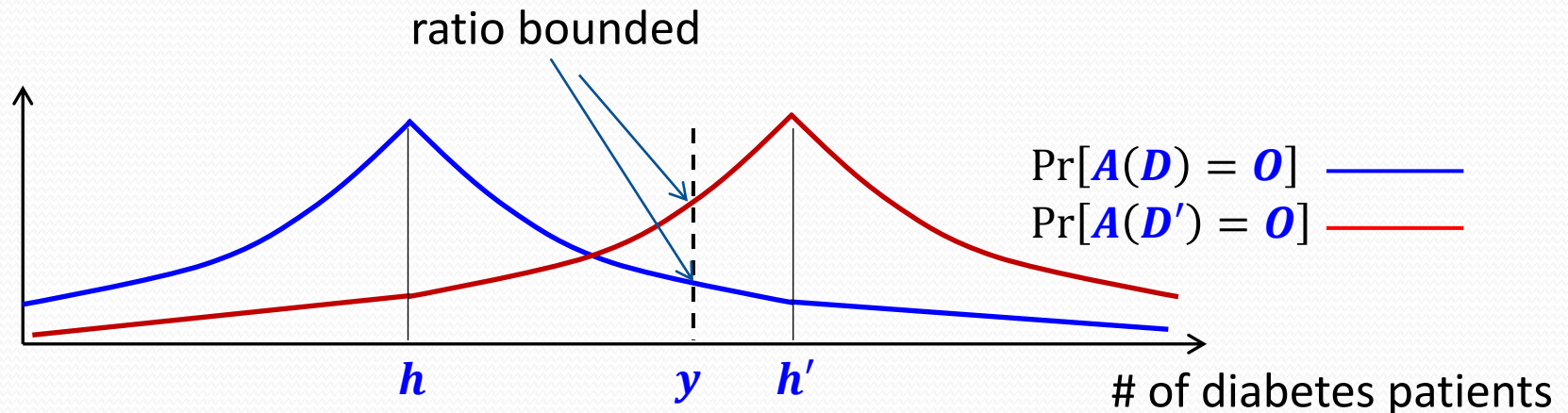
Differential Privacy via Laplace Noise

- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

- D : the original dataset; # of diabetes patients = h
- D' : modify a patient in D ; # of diabetes patients = h'



Differential Privacy via Laplace Noise

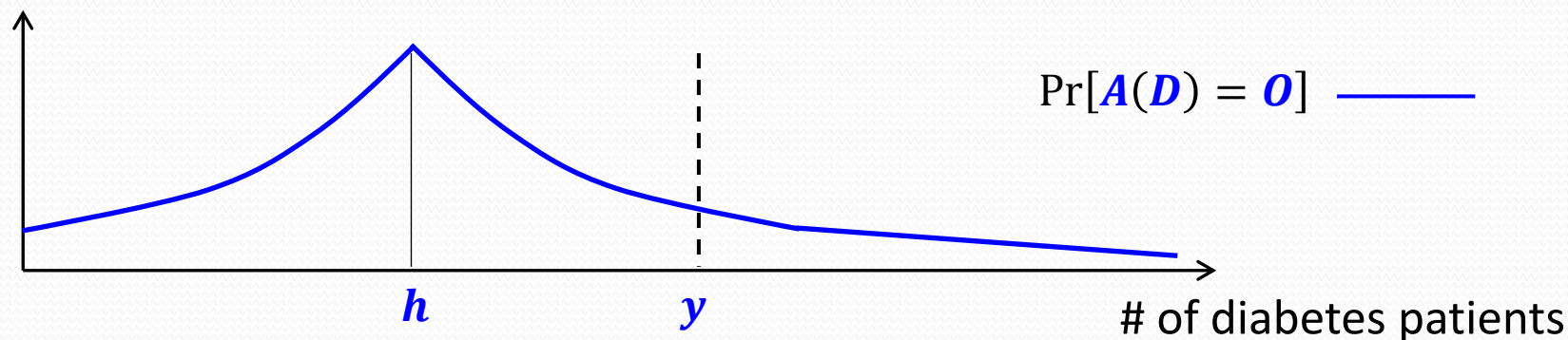
- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

- D : the original dataset; # of diabetes patients = h
- D' : modify a patient in D ; # of diabetes patients = h'

$$\Pr[A(D) = y] = pdf(y - h) = \exp(-|y - h|/\lambda) / 2\lambda$$



Differential Privacy via Laplace Noise

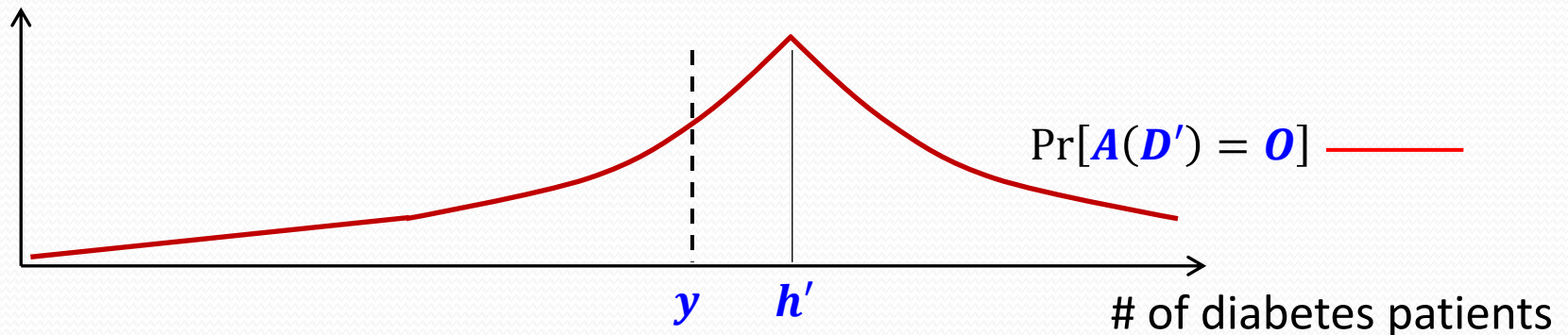
- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

- D : the original dataset; # of diabetes patients = h
- D' : modify the height of an individual in D ; # of diabetes patients = h'

$$\Pr[A(D') = y] = pdf(y - h') = \exp(-|y - h'|/\lambda) / 2\lambda$$



Differential Privacy via Laplace Noise

- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

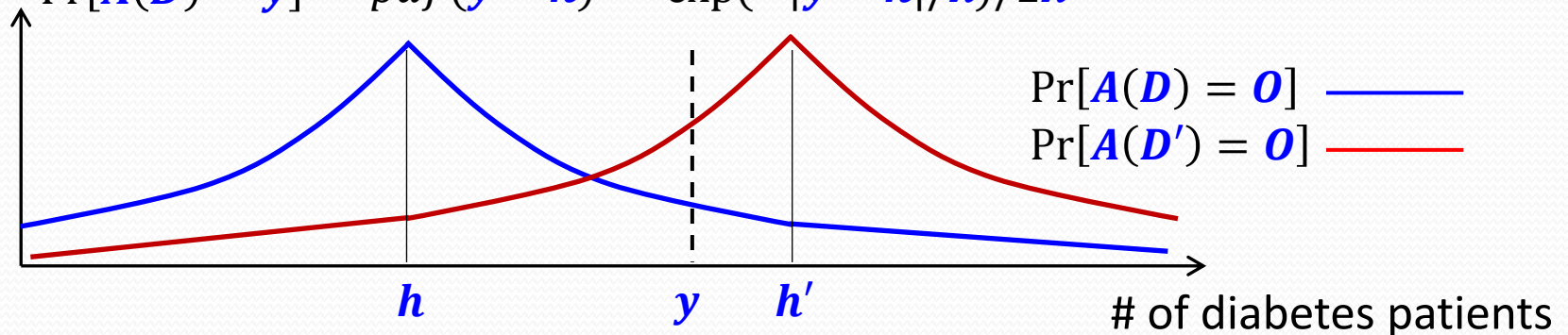
$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

- D : the original dataset; # of diabetes patients = h
- D' : modify the height of an individual in D ; # of diabetes patients = h'

$$\Pr[A(D') = y] = pdf(y - h') = \exp(-|y - h'|/\lambda) / 2\lambda$$

$$\Pr[A(D) = y] = pdf(y - h) = \exp(-|y - h|/\lambda) / 2\lambda$$



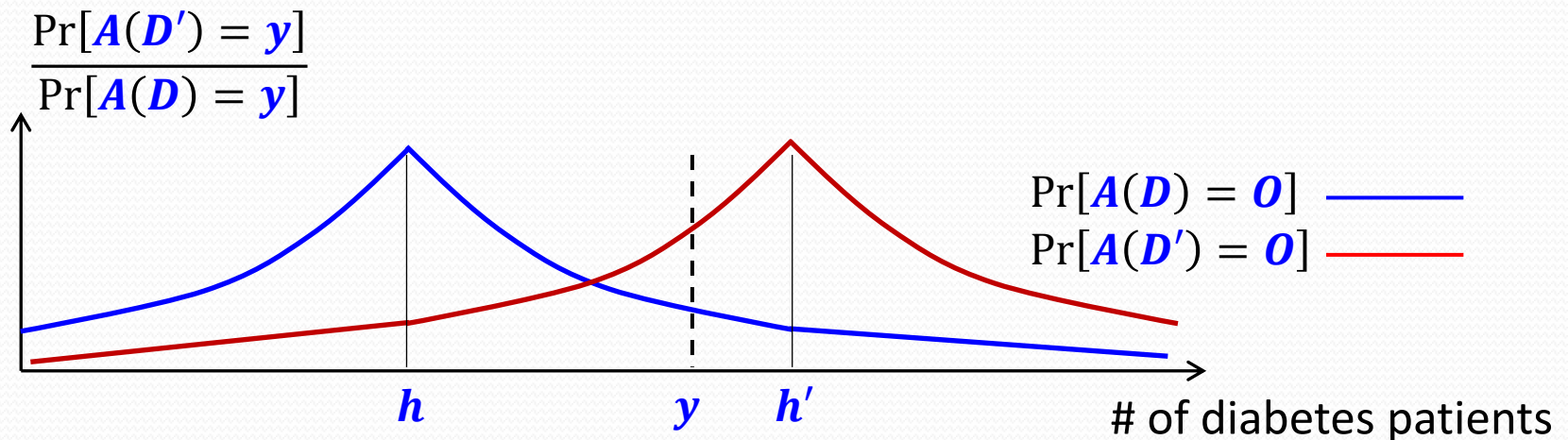
Differential Privacy via Laplace Noise

- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

- D : the original dataset; # of diabetes patients = h
- D' : modify the height of an individual in D ; # of diabetes patients = h'



Differential Privacy via Laplace Noise

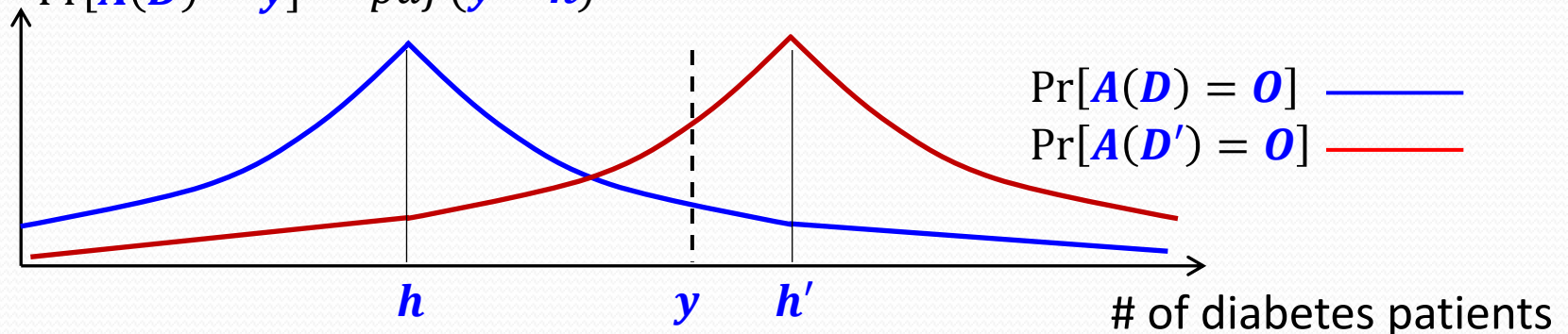
- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

- D : the original dataset; # of diabetes patients = h
- D' : modify the height of an individual in D ; # of diabetes patients = h'

$$\frac{\Pr[A(D') = y]}{\Pr[A(D) = y]} = \frac{pdf(y - h')}{pdf(y - h)}$$



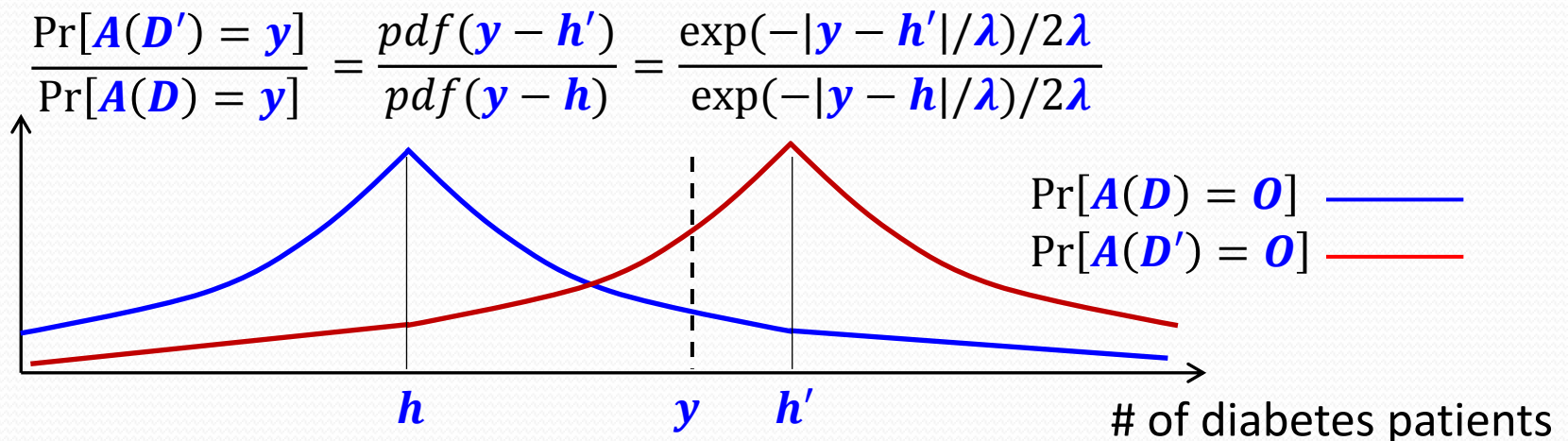
Differential Privacy via Laplace Noise

- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

- D : the original dataset; # of diabetes patients = h
- D' : modify the height of an individual in D ; # of diabetes patients = h'



Differential Privacy via Laplace Noise

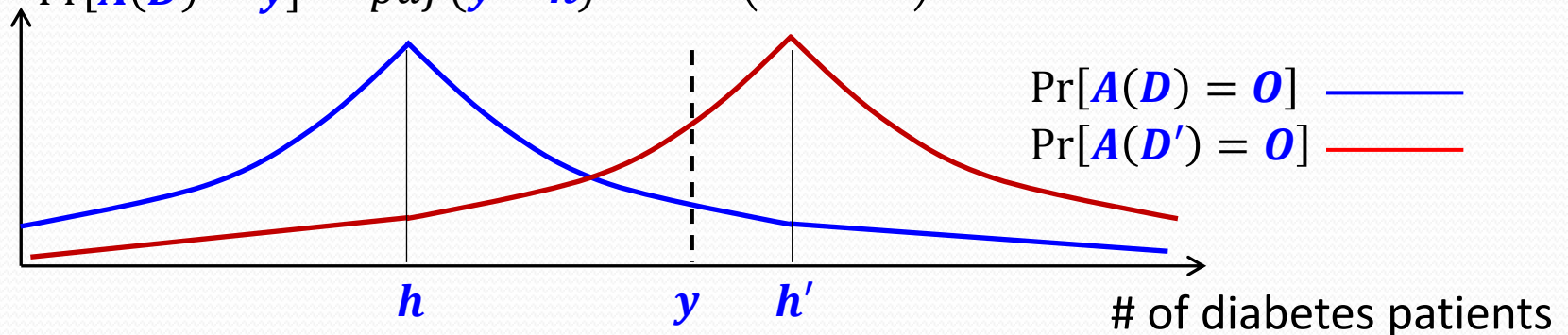
- Dataset: A set of patients
- Objective: Release # of diabetes patients with ϵ -differential privacy
 $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
- Method: Release the number + Laplace noise

$$pdf(x) = \exp\left(-\frac{|x|}{\lambda}\right) / 2\lambda$$

- Rationale:

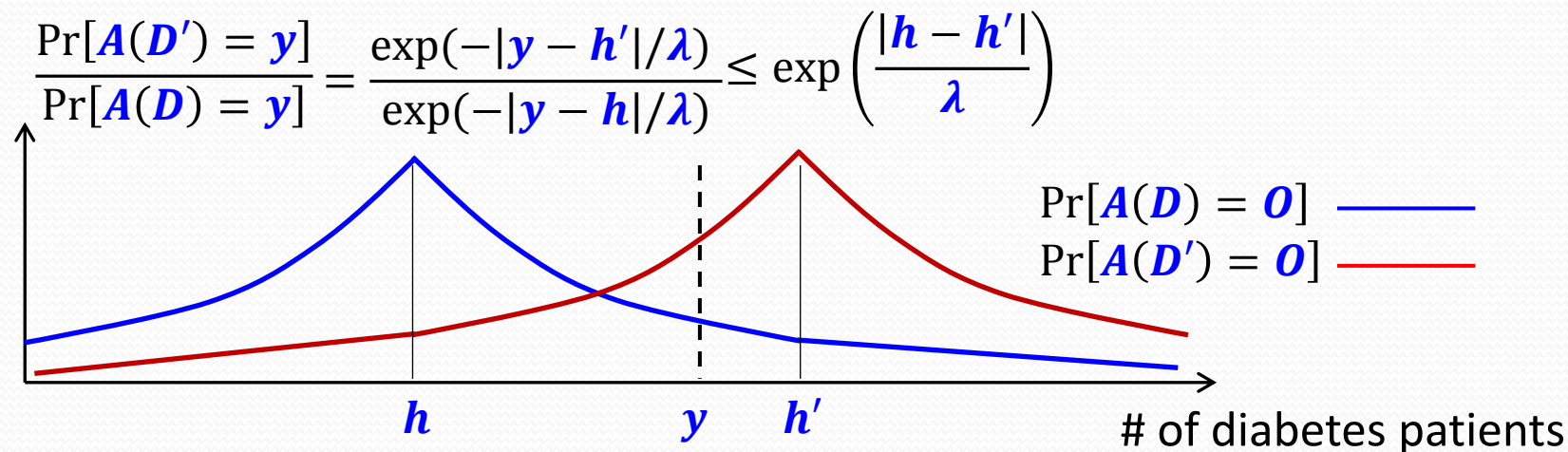
- D : the original dataset; # of diabetes patients = h
- D' : modify the height of an individual in D ; # of diabetes patients = h'

$$\frac{\Pr[A(D') = y]}{\Pr[A(D) = y]} = \frac{pdf(y - h')}{pdf(y - h)} \leq \exp\left(\frac{|h - h'|}{\lambda}\right)$$



Differential Privacy via Laplace Noise

- We aim to ensure ϵ -differential privacy
- How large should λ be?
 - Change of a patient's data would change the number of diabetes patients by at most 1, i.e.,
- Conclusion: Setting $\lambda \geq \frac{|h - h'|}{\epsilon}$ would ensure ϵ -differential privacy



General Mechanism with Laplace Noise

- In general, if the query result v is a real number
 - Add Laplace noise into v
- To decide the scale λ of Laplace noise
 - Look at the maximum change that can occur in v (when we change one tuple in the dataset)
 - Set λ to be proportional to the maximum change

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	23	Y
Doug	M	30	N

General via Laplace Noise

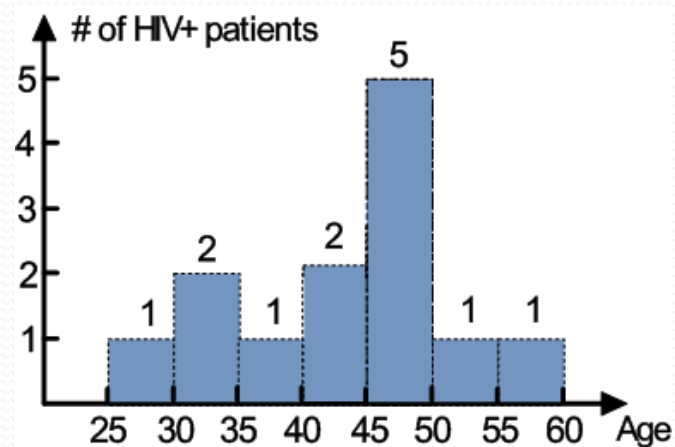
- What if we have multiple queries?
 - Add Laplace noise to each value
- How do we decide the noise scale?
 - Look at the *total change* that can occur in the values when we modify one tuple in the data
 - Total change: sum of the absolute change in each value (i.e., differences in L1 norm)
 - Set the scale of the noise to be proportional to the maximum total change
- The maximum total change is referred to as the *sensitivity* of the values
- Theorem [Dwork et al. 2006]: Adding Laplace noise of scale λ to each value ensures ϵ -differential privacy, if
$$\lambda \geq (\text{the sensitivity of the values}) / \epsilon$$

Sensitivity of Queries

- Histogram

- Sensitivity of the bin counts: 2
- Reason: When we modify a tuple in the dataset, at most two bin counts would change; furthermore, each bin count would change by at most 1
- Scale of Laplace noise required:

Name	Age	HIV+
Frank	42	Y
Bob	31	Y
Mary	28	Y
Dave	43	N
...



- For more complex queries, the derivation of sensitivity can be much more complicated
 - Example: Parameters of a logistic model

Exponential Mechanism

- What if the query result is on discrete space?
 - Example: Which one is a more important factor to diabetic, age or gender?
- Given k items, each item is associated with a score $S(I, D)$, how to pick the one with maximal score under differential privacy?
- Adding Laplace noise is a feasible solution

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N

$$S(\text{Gender}, D) = \text{Corr}(\text{Gender}, \text{Diabetes})$$

$$S(\text{Age}, D) = \text{Corr}(\text{Age}, \text{Diabetes})$$

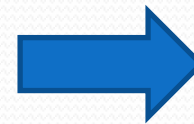
Exponential Mechanism

- Using exponential mechanism, we can directly manipulate the probability of item pickup.
- For each item I_j , the probability is proportional to $\exp(S(I, D)/\lambda)$

$$S(\text{Gender}, D) = \text{Corr}(\text{Gender}, \text{Diabetes}) = 0.5$$

$$S(\text{Age}, D) = \text{Corr}(\text{Age}, \text{Diabetes}) = 0.3$$

Name	Gender	Age	Diabetes
Alice	F	28	Y
Bob	M	19	Y
Chris	M	25	N
Doug	M	30	N



$$\text{Pr}(\text{Gender}) = 0.71$$

$$\text{Pr}(\text{Age}) = 0.39$$

Exponential Mechanism

- Advantage: Improve skewedness on the probabilities
- Limitation: Needs to iterate all possible answers in the solution space. It is thus not applicable when the solution space is too large.
- Example: Pick up the best order of k items with maximal score. The number of possible orders is $k!$.

Variants of Differential Privacy

- Alternative definition of neighboring dataset:
 - Two datasets D and D' , such that D' is obtained by adding/deleting one tuple in D
- $\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$
 - Even if a tuple is added to or removed from the dataset, the output distribution of the algorithm is roughly the same
 - i.e., the output of the algorithm does not reveal the presence of a tuple
- Refer to this version as “unbounded” differential privacy, and the previous version as “bounded” differential privacy

Variants of Differential Privacy

- Bounded: D' is obtained by changing the values of one tuple in D
- Unbounded: D' is obtained by adding/removing one tuple in D
- Observation 1
 - Change of a tuple can be regarded as removing a tuple from the dataset and then inserting a new one
 - Indication: Unbounded ϵ -differential privacy implies bounded (2ϵ) -differential privacy
 - Proof: $\Pr[A(D_1) = O] \leq \exp(\epsilon) \cdot \Pr[A(D_2) = O]$
 $\leq \exp(\epsilon) \cdot \exp(\epsilon) \cdot \Pr[A(D_3) = O]$

Variants of Differential Privacy

- Bounded: D' is obtained by changing the values of one tuple in D
- Unbounded: D' is obtained by adding/removing one tuple in D
- Observation 2
 - Bounded differential privacy allows us to directly publish the number of tuples in the dataset

$$\Pr[A(D) = O] \leq \exp(\epsilon) \cdot \Pr[A(D') = O]$$

- Unbounded differential privacy does not allow this

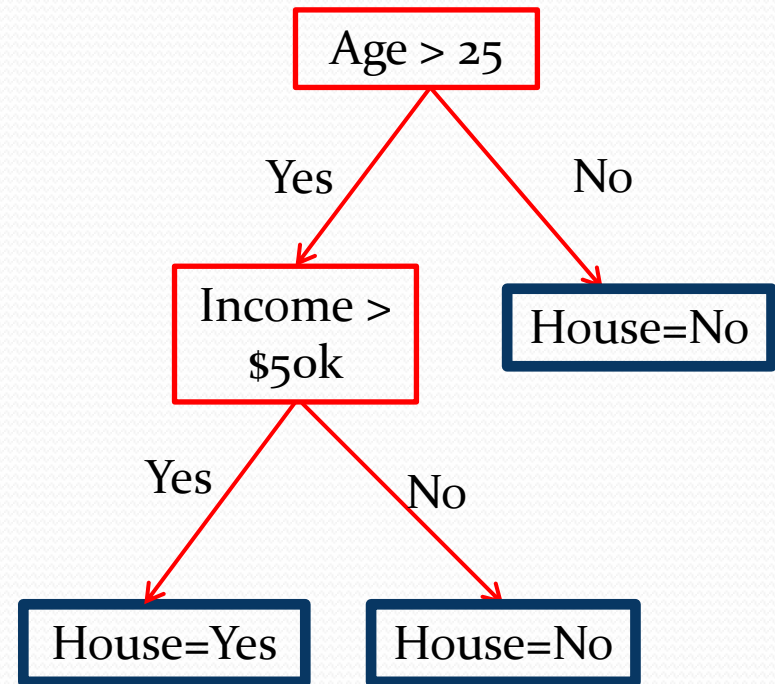
Limitations of Differential Privacy

- Differential privacy tends to be less effective when there exist correlations among the tuples
- Example (from [Kifer and Machanavajjhala 2011]):
 - Bob's family includes 10 people, and all of them are in a database
 - There is a highly contagious disease, such that if one family member contracts the disease, then the whole family will be infected
 - Differential privacy would underestimate the risk of disclosure
- Summary: Amount of noise needed depends on the correlations among the tuples, which is not captured by differential privacy

Decision Tree Classification

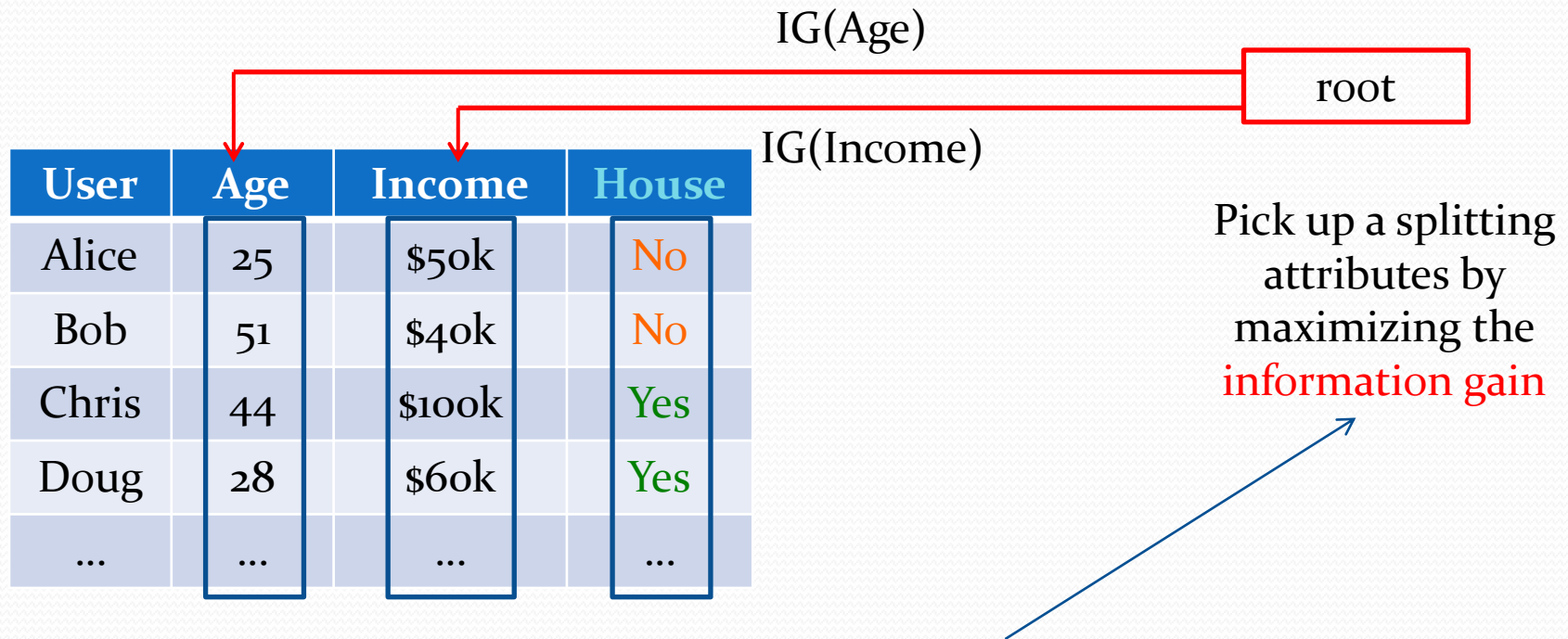
- Problem Definition

User	Age	Income	House
Alice	25	\$50k	No
Bob	51	\$40k	No
Chris	44	\$100k	Yes
Doug	28	\$60k	Yes
...



Decision Tree Classification

- Attribute Selection [*Friedman, 2010*]



$$IG(T, a) = H(T) - \sum_{v \in \text{vals}(a)} \frac{|\{\mathbf{x} \in T | x_a = v\}|}{|T|} \cdot H(\{\mathbf{x} \in T | x_a = v\})$$

Decision Tree Classification

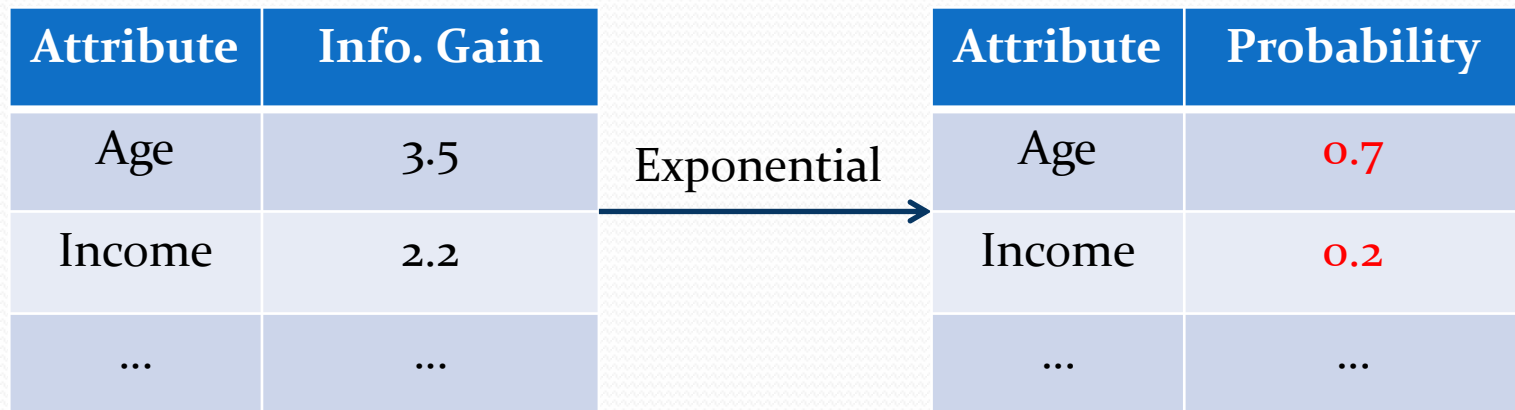
- How to enforce differential privacy in the selection?
 - Laplace Mechanism
 - Exponential Mechanism

Attribute	Info. Gain	Laplace →	Attribute	Info. Gain
Age	3.5		Age	2.9
Income	2.2		Income	2.7
...

Budget consumption:
 $\varepsilon \times m$

Decision Tree Classification

- How to enforce differential privacy in the selection?
 - Laplace Mechanism
 - Exponential Mechanism



Budget consumption: ϵ

Conclusion

- Differential Privacy is a new and robust criterion of privacy detection
- There are simple algorithms enforcing differential privacy
- For a specific query engine, we need to carefully pick up the appropriate place to insert noise.