

Parameter Learning in Probabilistic Databases

Bernd Gutmann

joint work with

Angelika Kimmig, Kristian Kersting, Luc De Raedt

Motivation

DiseaseX

What is the
reason for
DiseaseX?

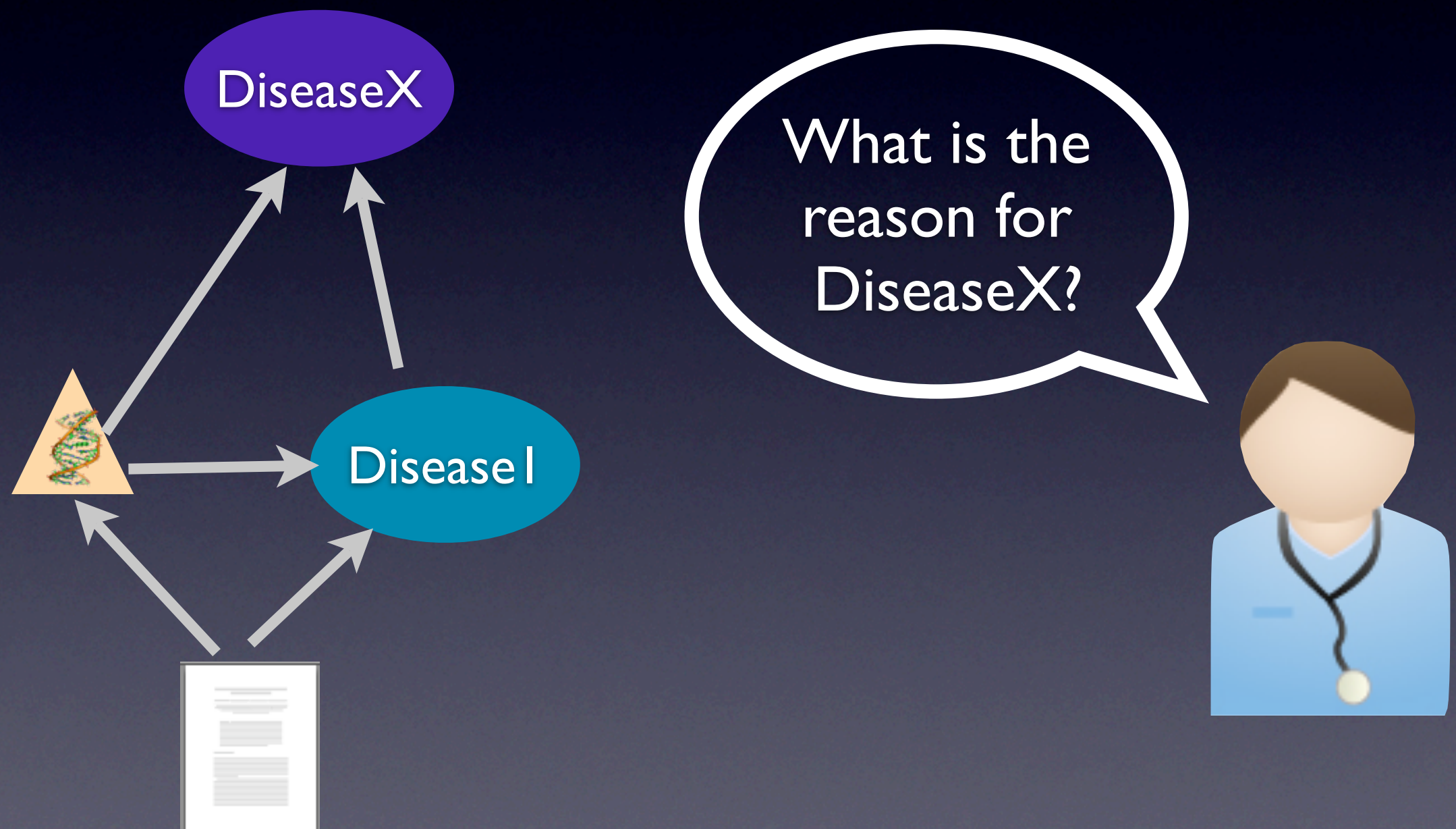
Disease I



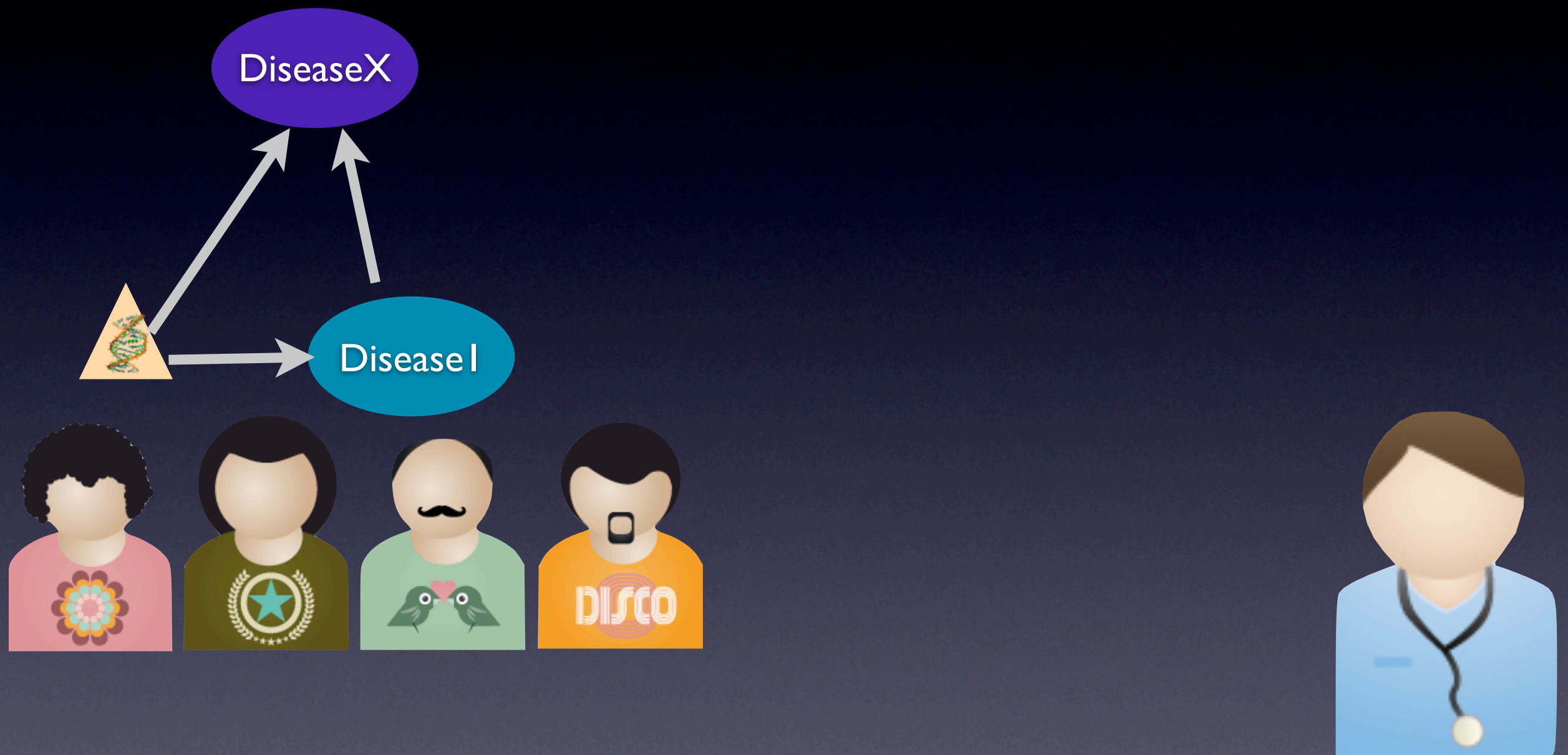
Motivation



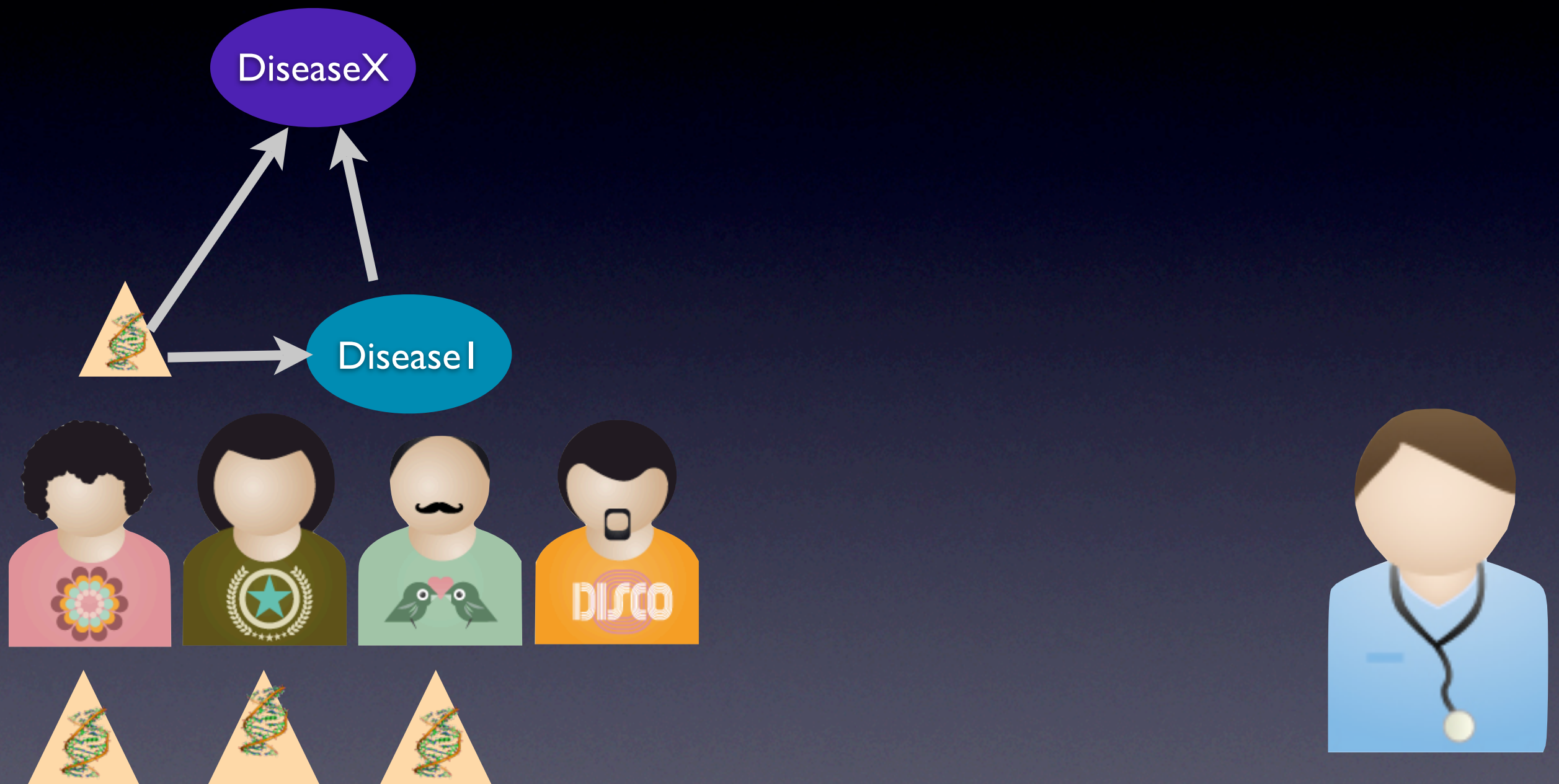
Motivation



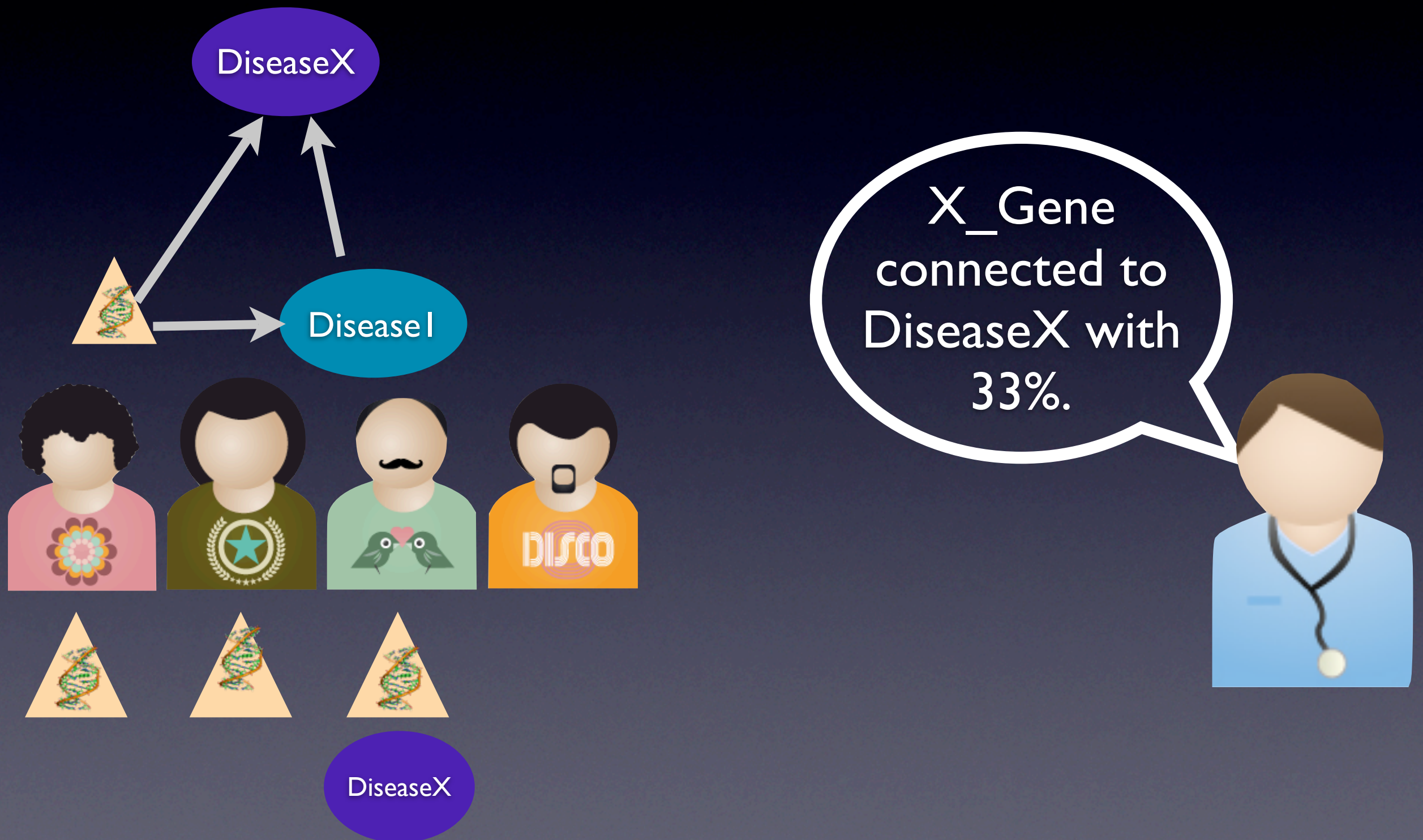
Motivation



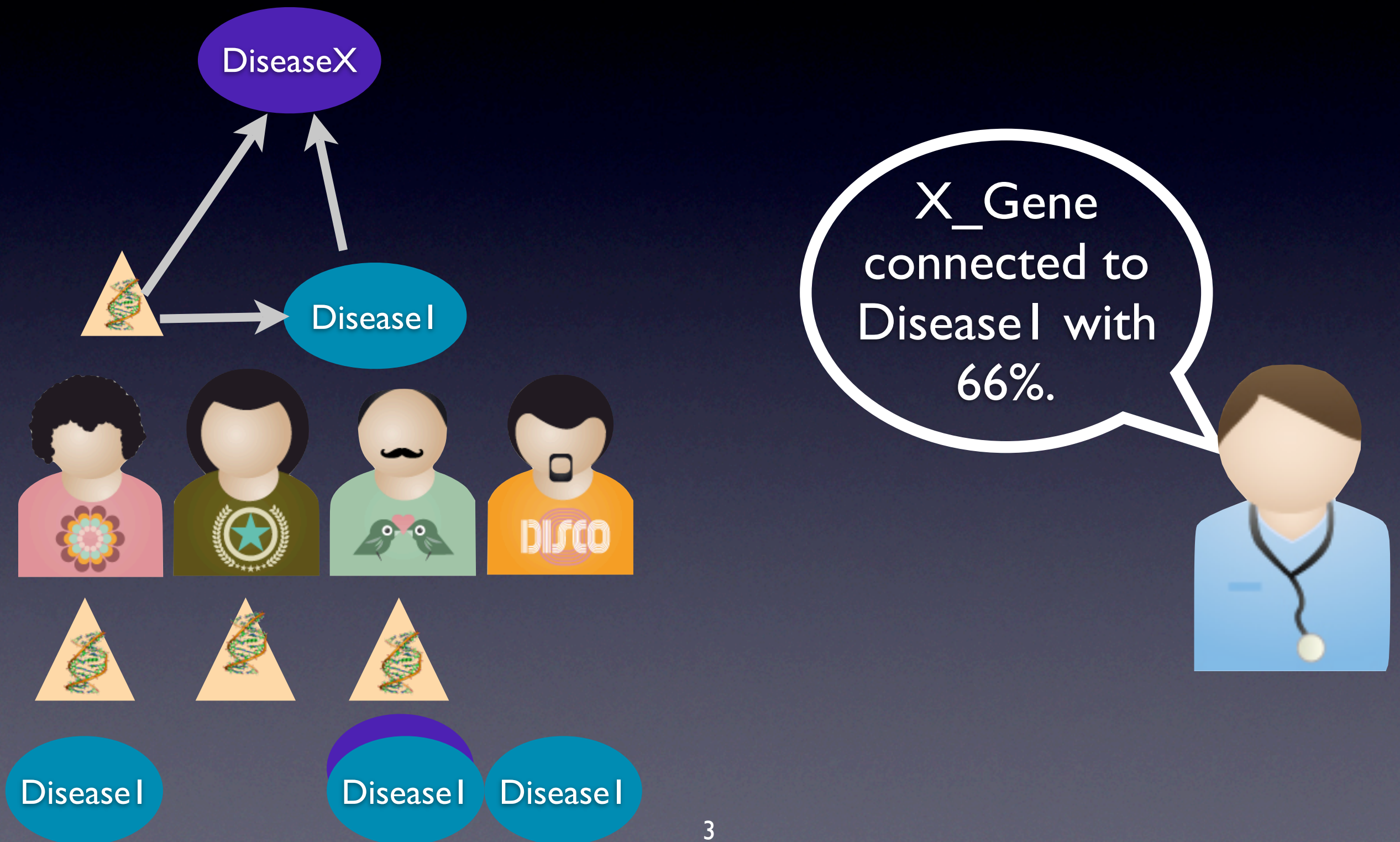
Motivation



Motivation



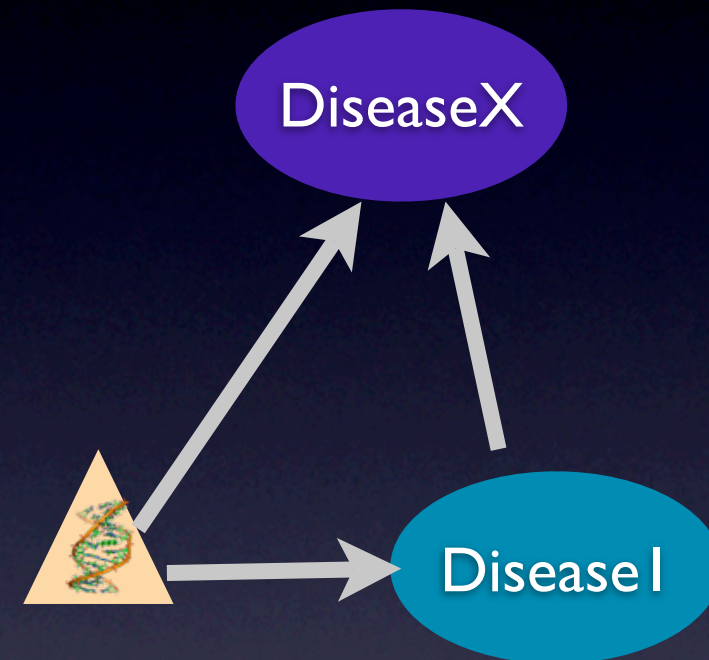
Motivation



The Situation

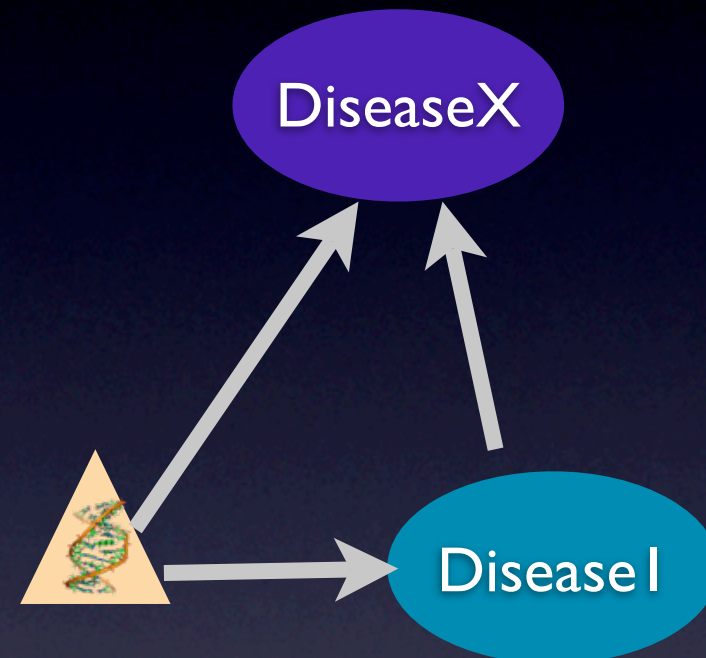
The Situation

- Domain knowledge



The Situation

- Domain knowledge

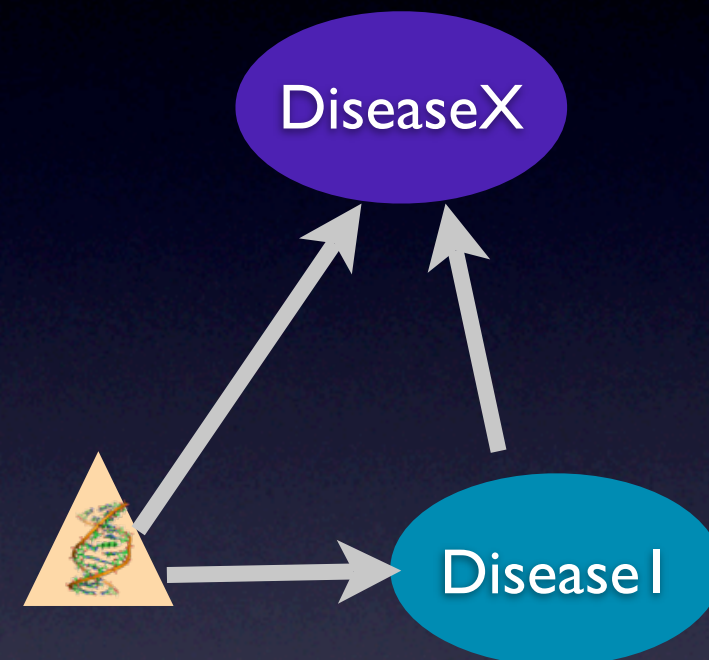


- Examples with probabilities



The Situation

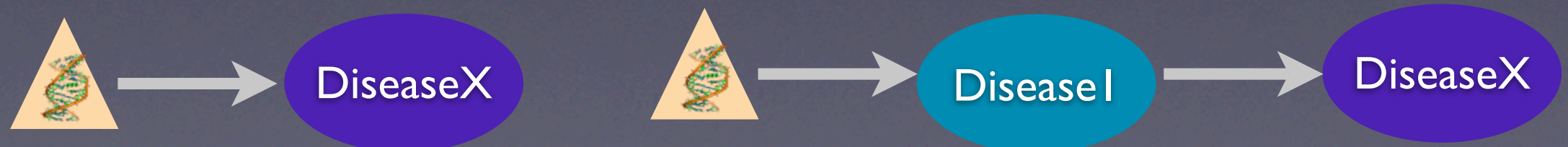
- Domain knowledge



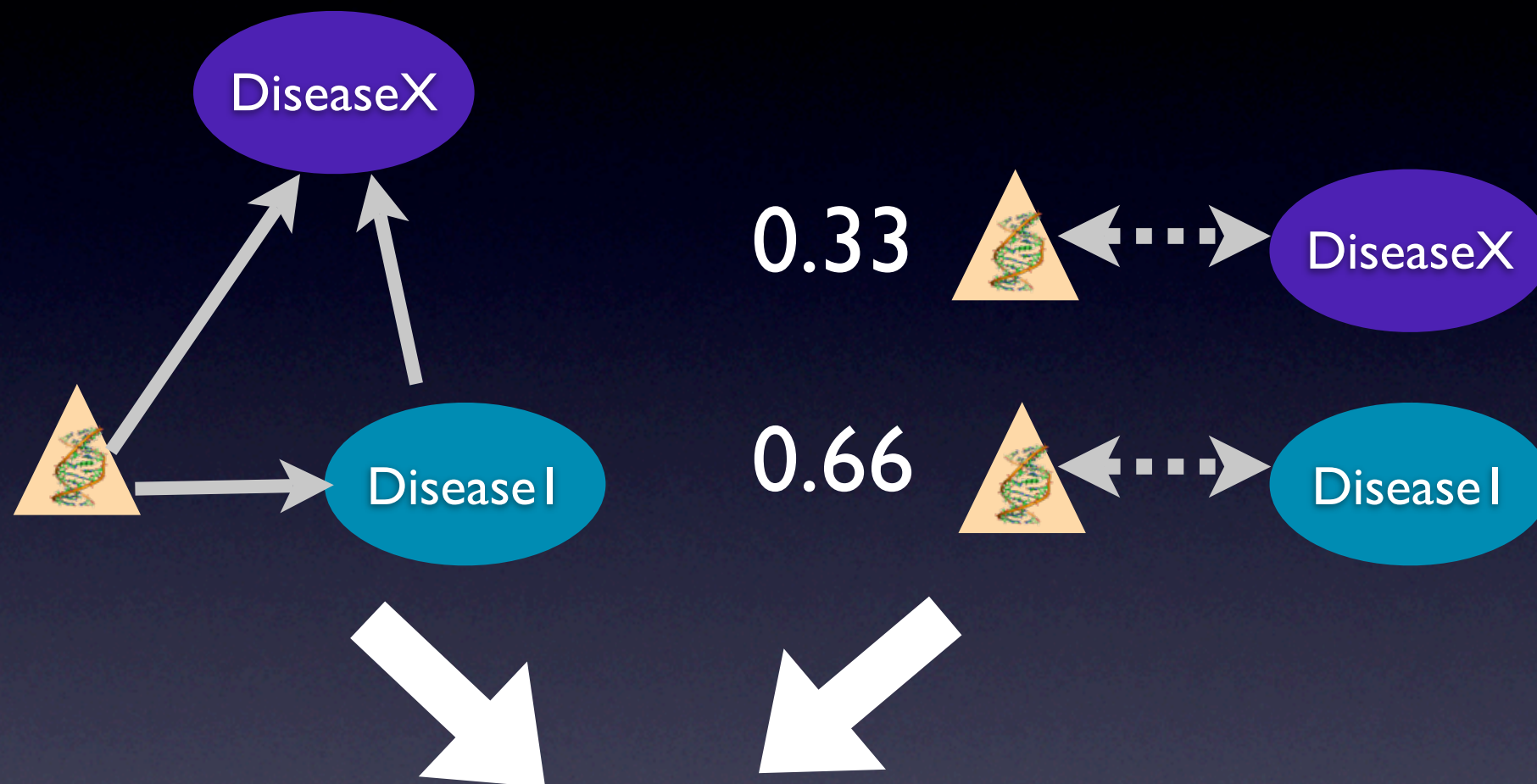
- Examples with probabilities



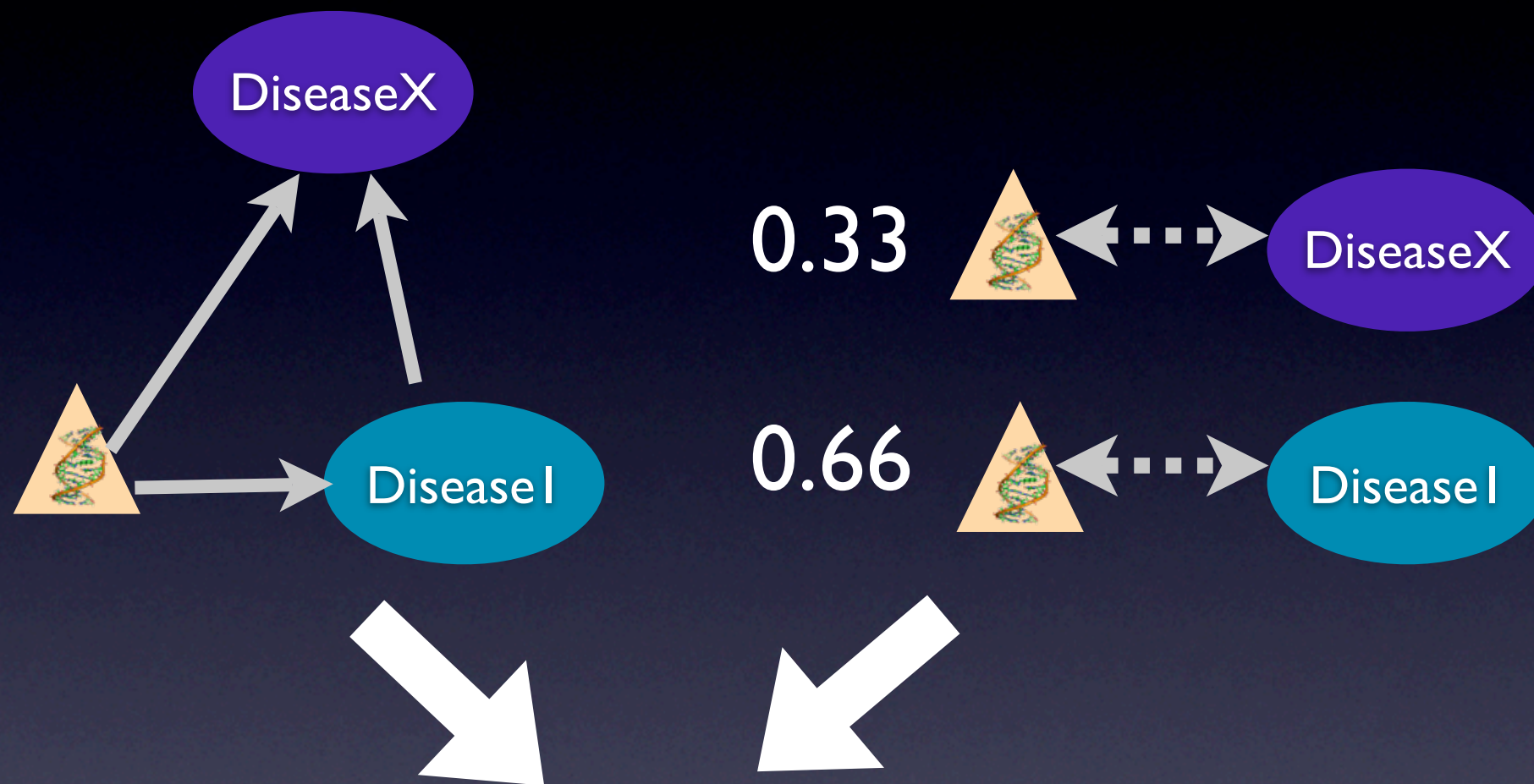
- Uncertainty how to explain the examples



A Solution

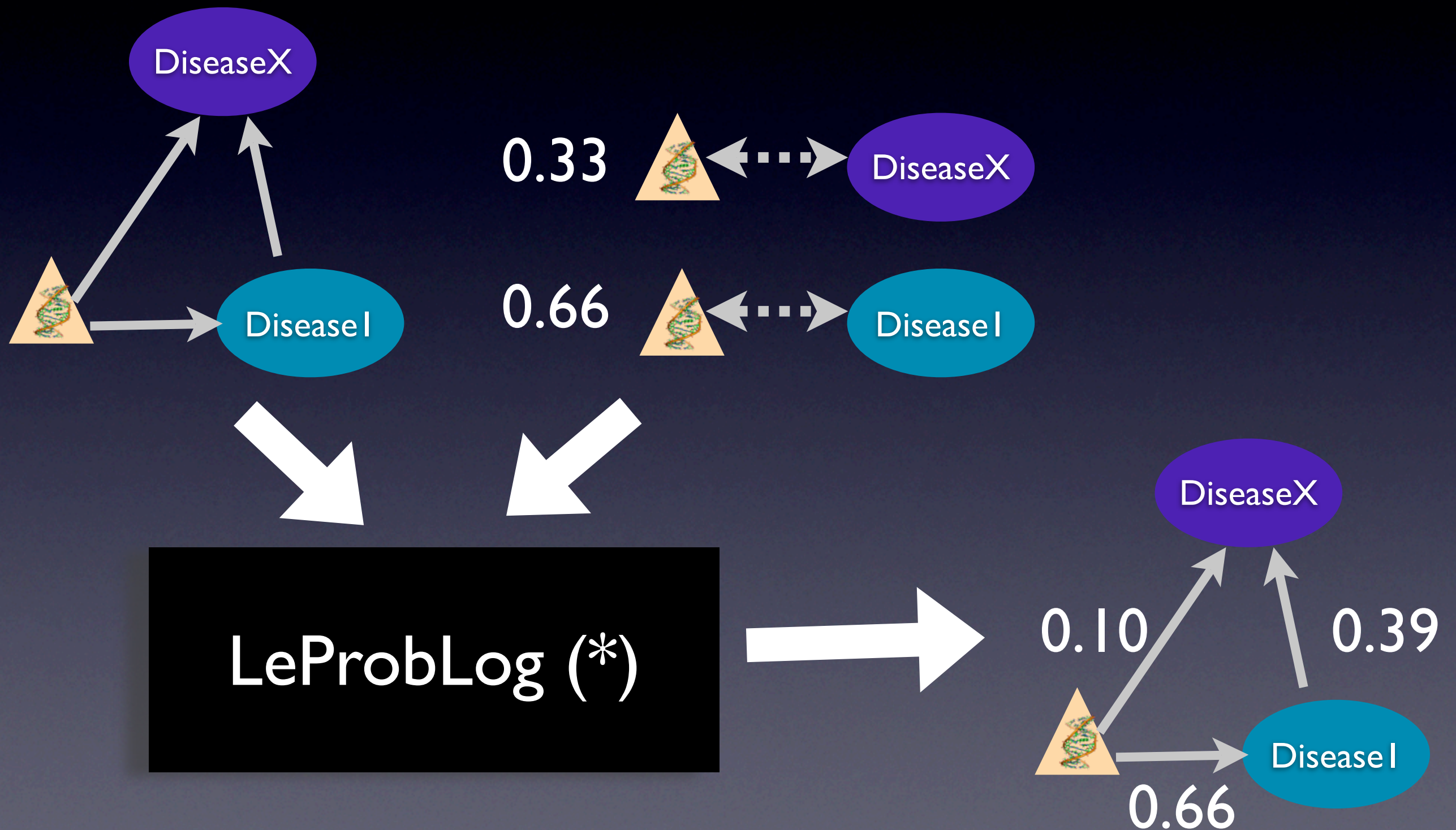


A Solution



LeProbLog (*)

A Solution

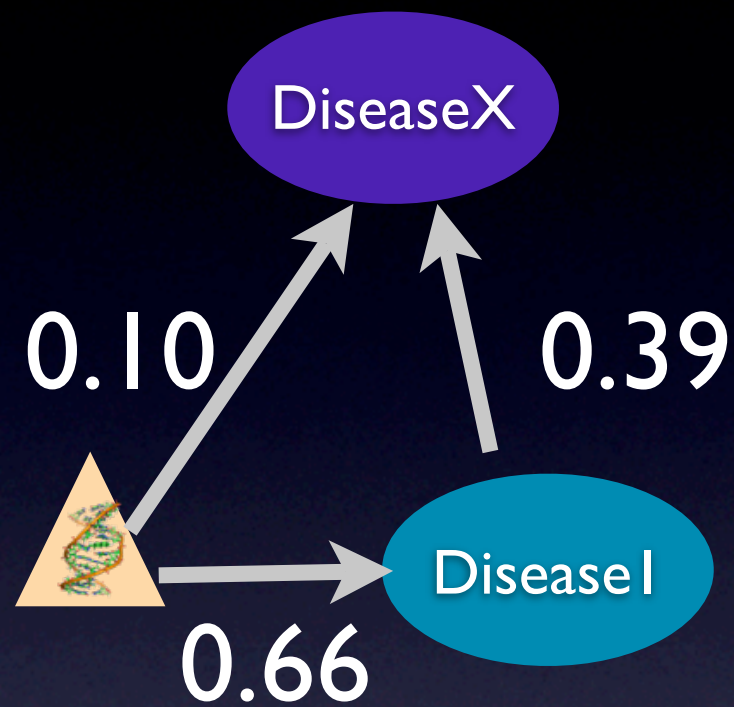


(*) Least Square Parameter Estimation for ProbLog ⁵

Outline

- Motivation
- ProbLog
- Parameter Learning for ProbLog
- Experiments

ProbLog



0.10 :: edge('x_gene', 'DiseaseX')

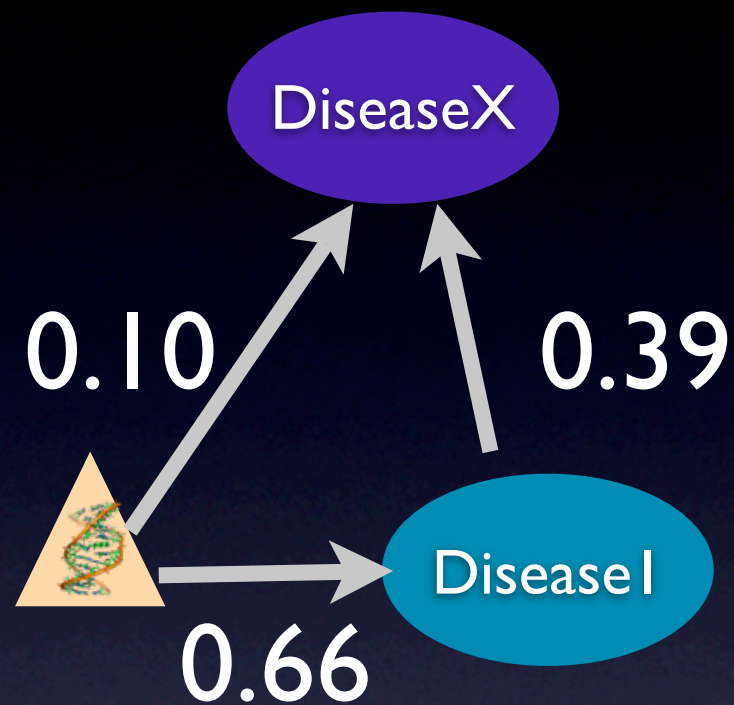
0.66 :: edge('x_gene', 'DiseaseI')

0.39 :: edge('DiseaseI', 'x_gene')

path(X,Y) :- edge(X,Y)

path(X,Y) :- edge(X,Z), path(Z,Y)

ProbLog

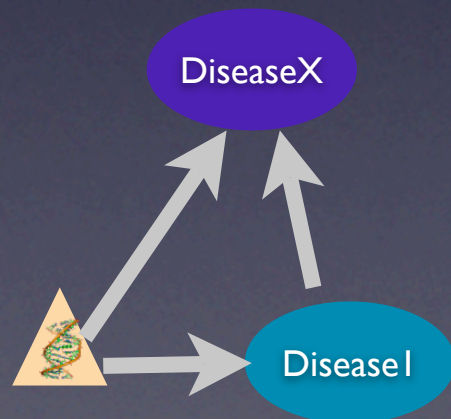


```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

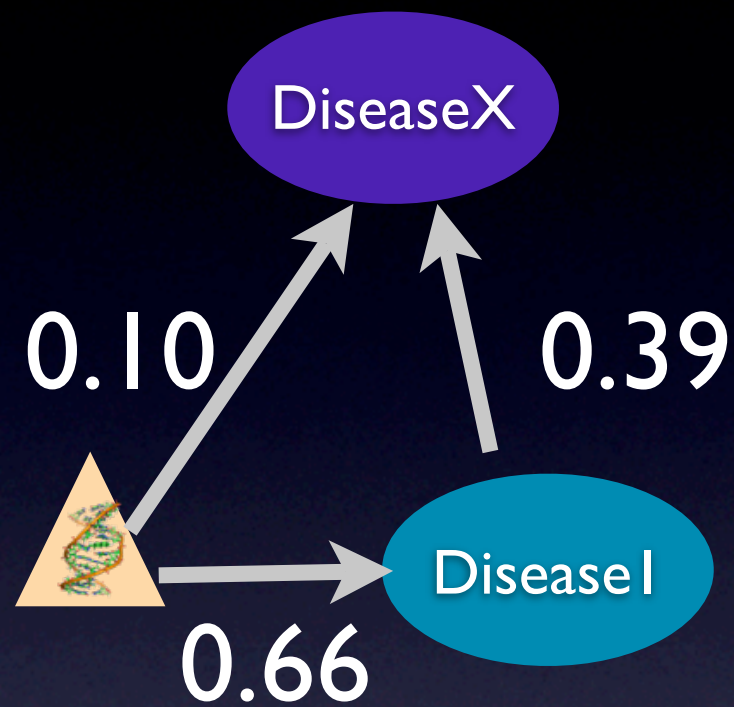
```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$$p = 0.10 * 0.66 * 0.39$$



ProbLog



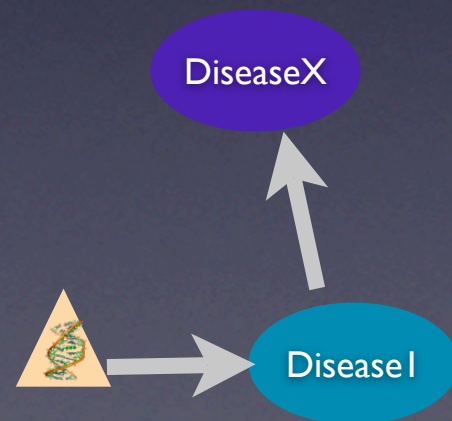
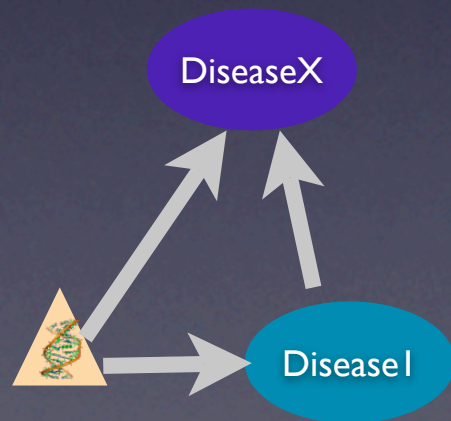
```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

```
path(X,Y) :- edge(X,Y)
```

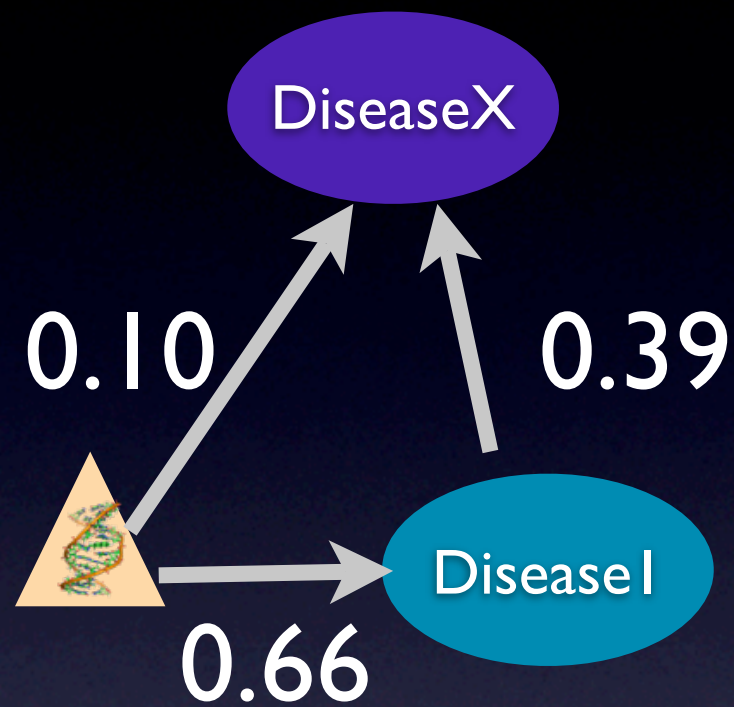
```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$$p=0.10 * 0.66 * 0.39$$

$$p=(1-0.10) * 0.66 * 0.39$$



ProbLog



```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

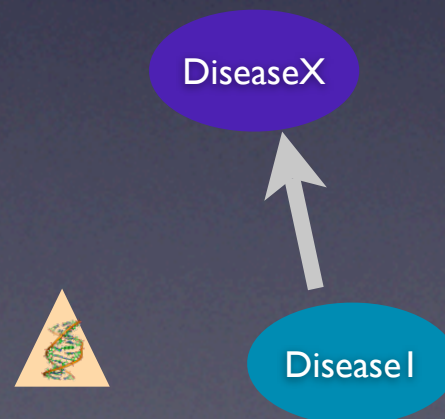
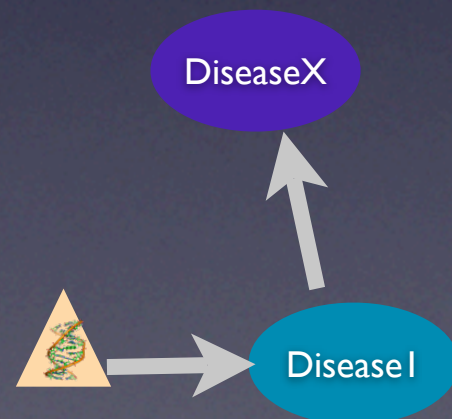
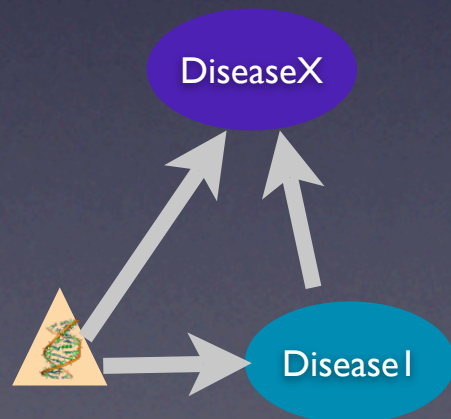
```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

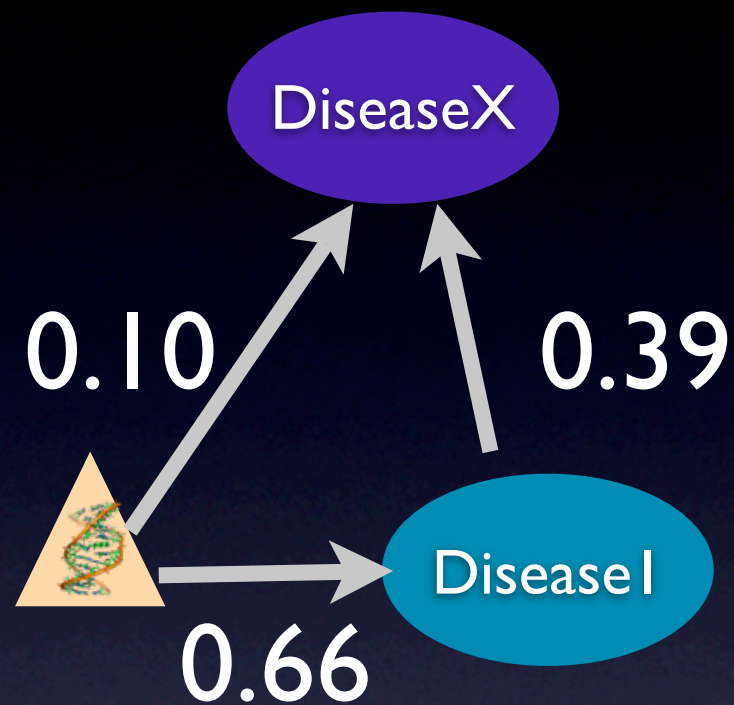
$$p=0.10 * 0.66 * 0.39$$

$$p=(1-0.10) * 0.66 * 0.39$$

$$p=(1-0.10) * (1-0.66) * 0.39$$



ProbLog



$0.10 :: \text{edge}(\text{'x_gene'}, \text{'DiseaseX'})$
 $0.66 :: \text{edge}(\text{'x_gene'}, \text{'DiseaseI'})$
 $0.39 :: \text{edge}(\text{'DiseaseI'}, \text{'x_gene'})$

$\text{path}(X,Y) :- \text{edge}(X,Y)$

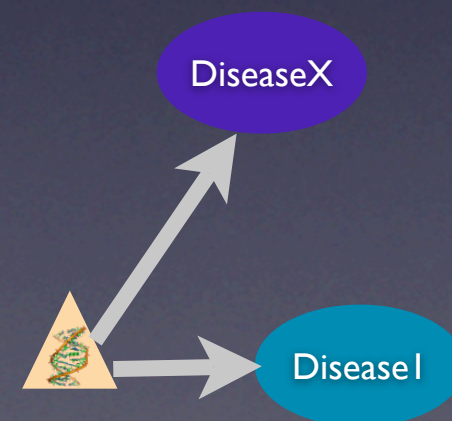
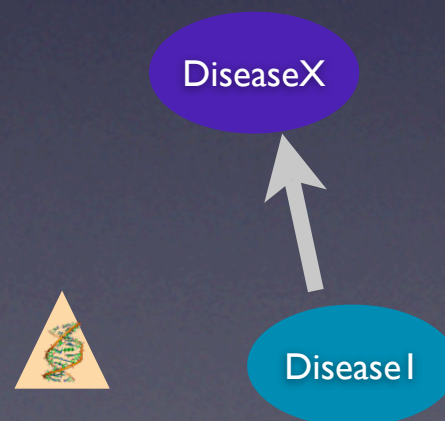
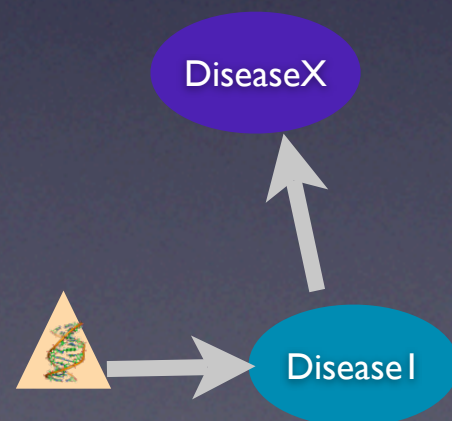
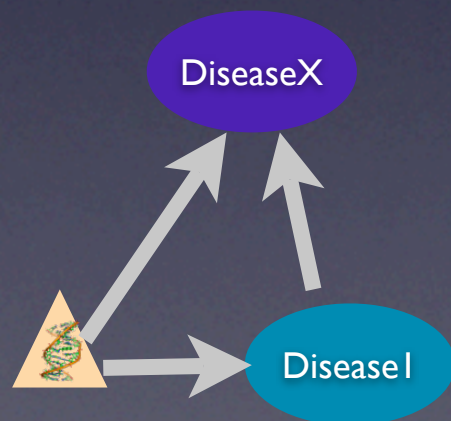
$\text{path}(X,Y) :- \text{edge}(X,Z), \text{path}(Z,Y)$

$$p = 0.10 * 0.66 * 0.39$$

$$p = (1 - 0.10) * 0.66 * 0.39$$

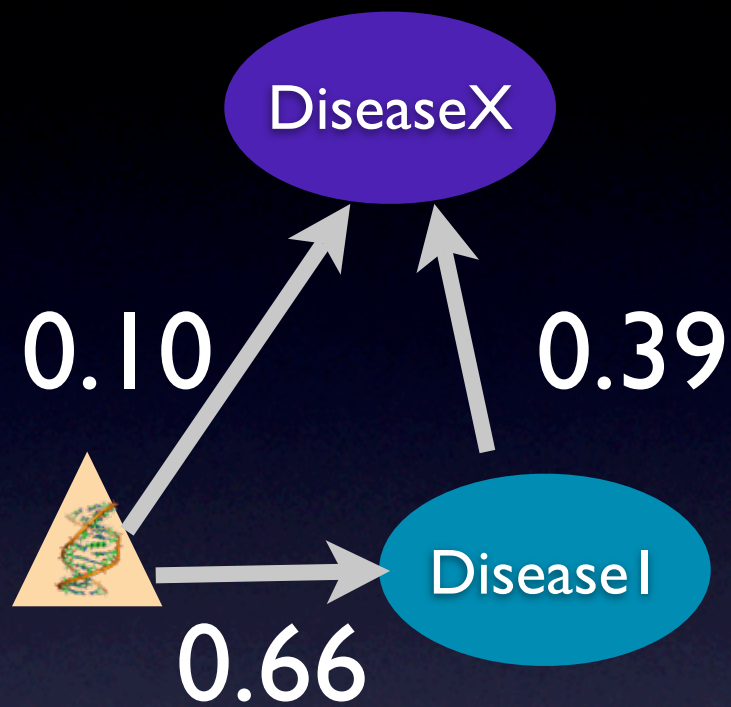
$$p = (1 - 0.10) * (1 - 0.66) * 0.39$$

$$p = 0.10 * 0.66 * (1 - 0.39)$$



...

ProbLog



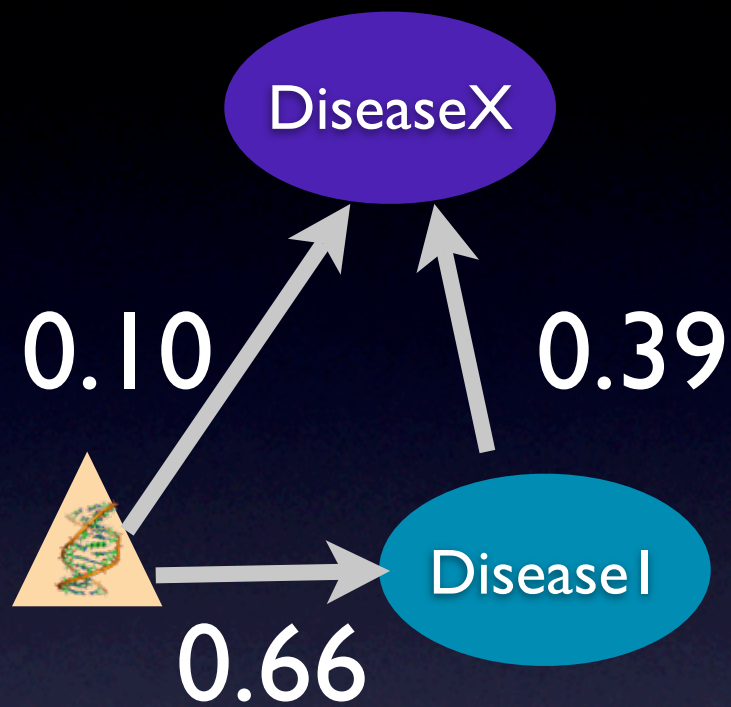
```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$P(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}))$

ProbLog

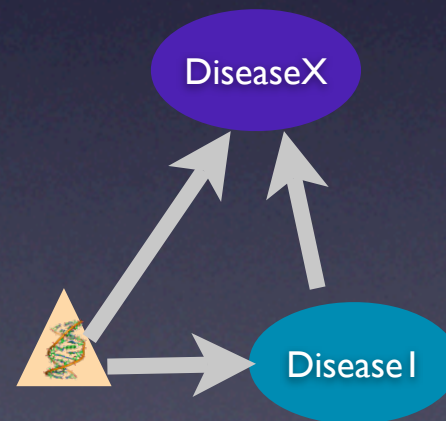


```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

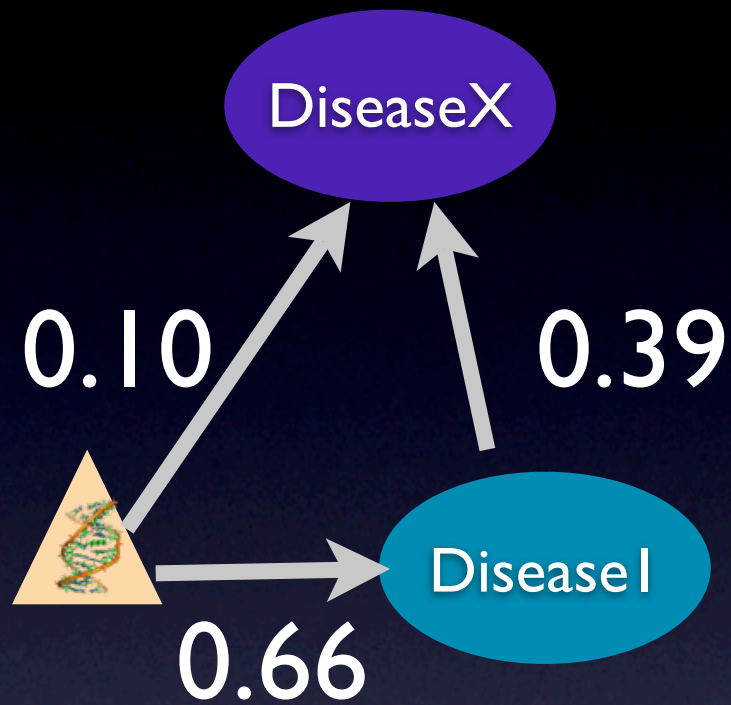
```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$P(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}))$



ProbLog



```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

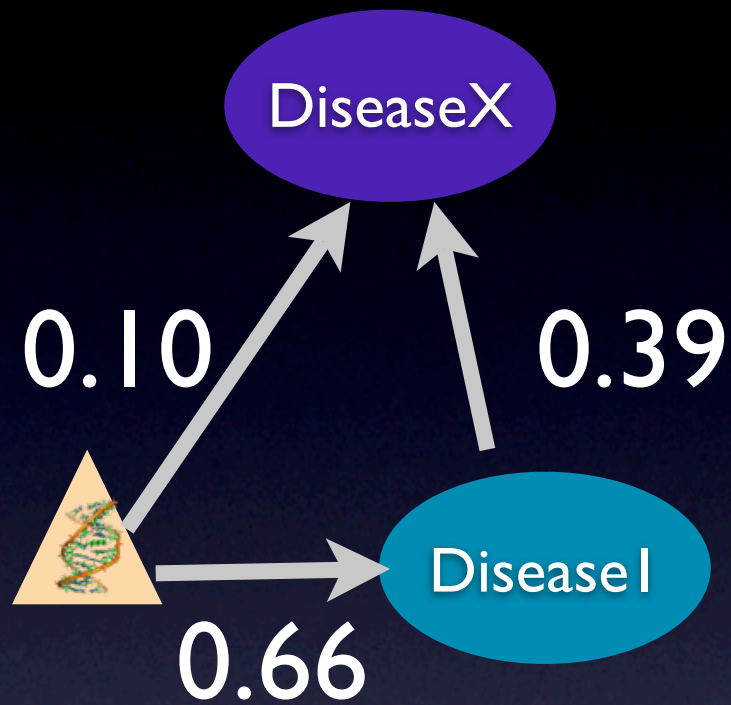
```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$P(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}))$



ProbLog

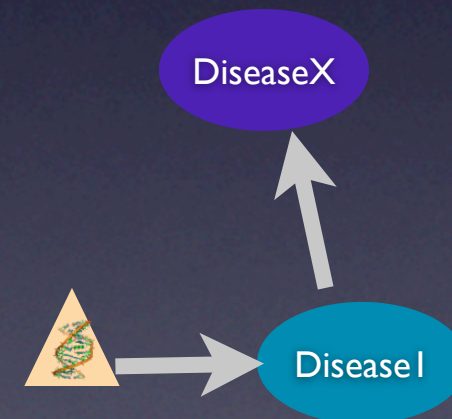
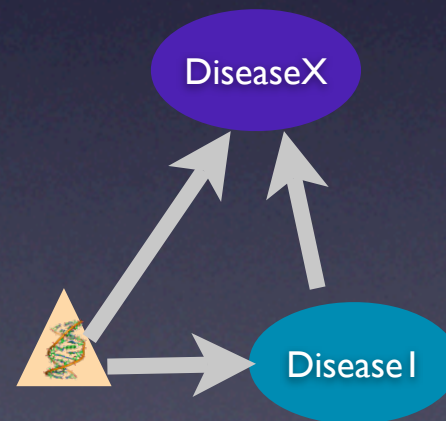


```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

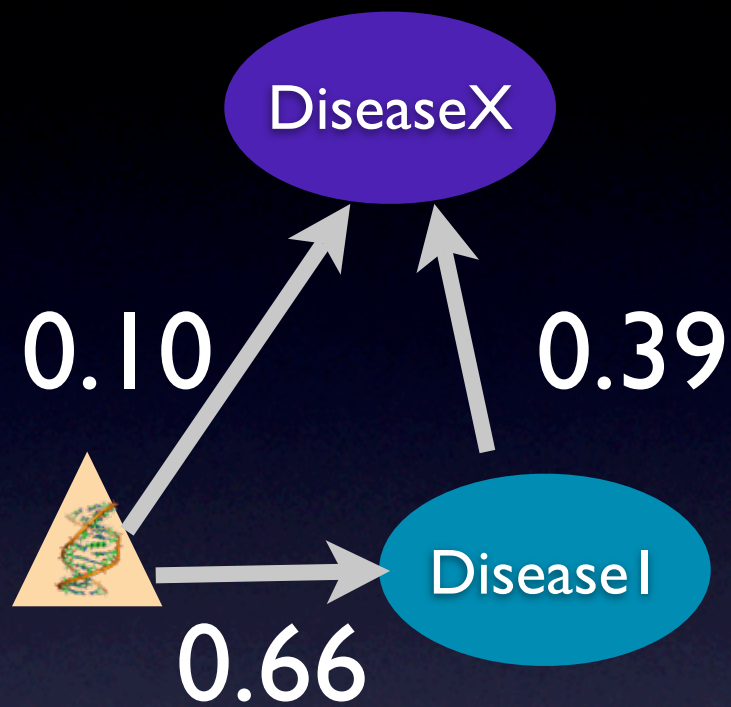
```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$P(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}))$



ProbLog

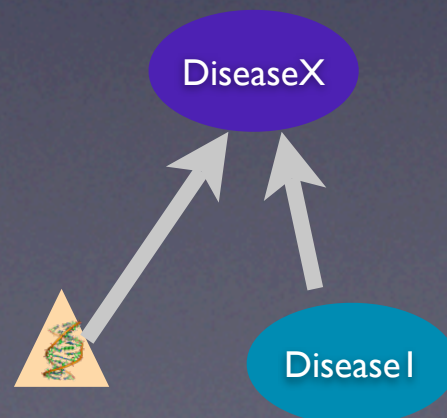
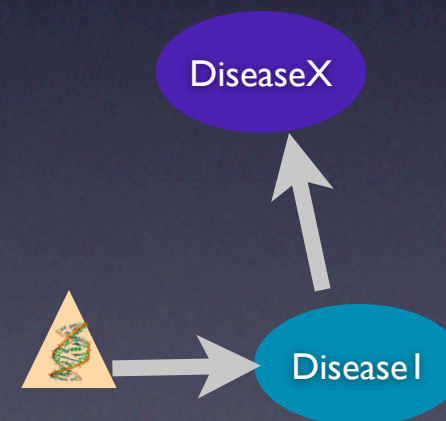
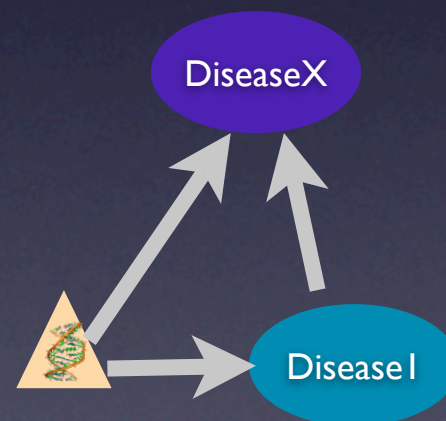


```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

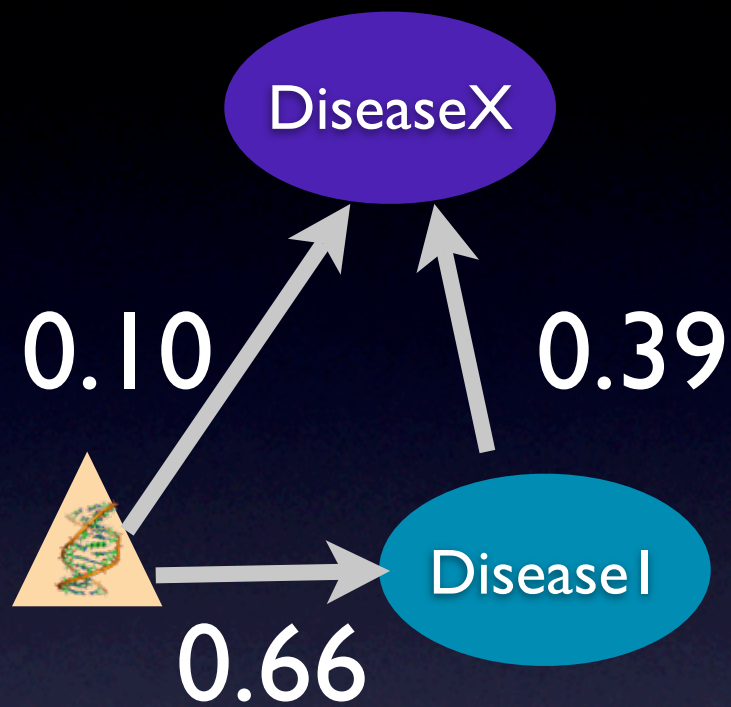
```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$P(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}))$



ProbLog

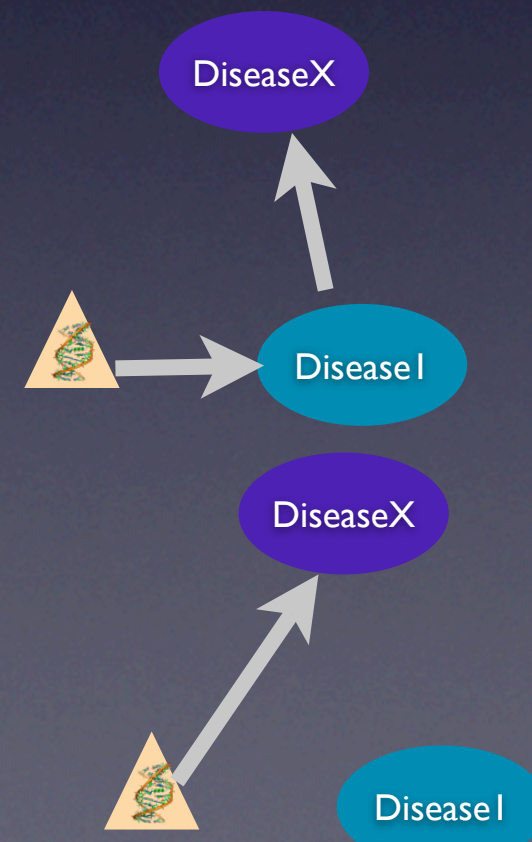
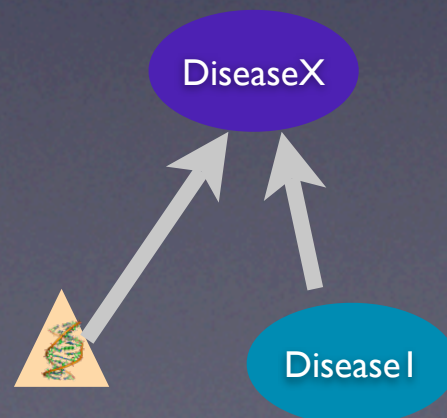
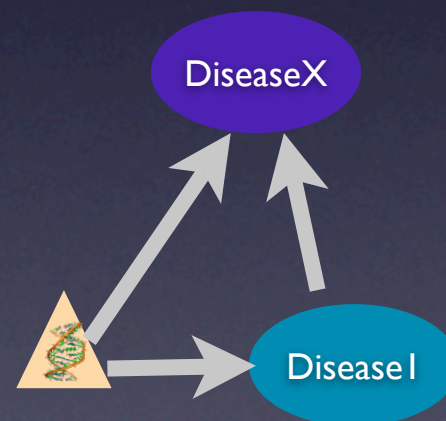


```
0.10 :: edge('x_gene', 'DiseaseX')
0.66 :: edge('x_gene', 'DiseaseI')
0.39 :: edge('DiseaseI', 'x_gene')
```

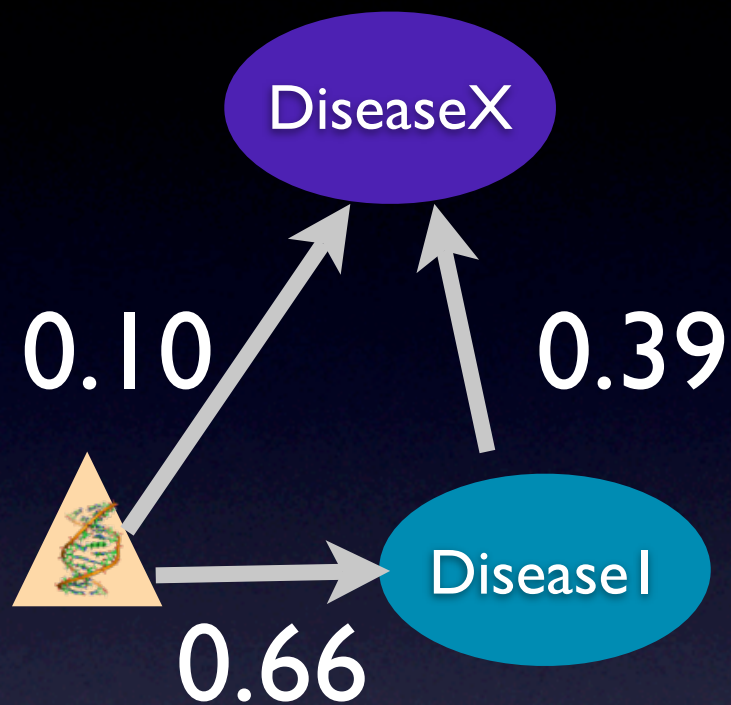
```
path(X,Y) :- edge(X,Y)
```

```
path(X,Y) :- edge(X,Z), path(Z,Y)
```

$P(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}))$



ProbLog

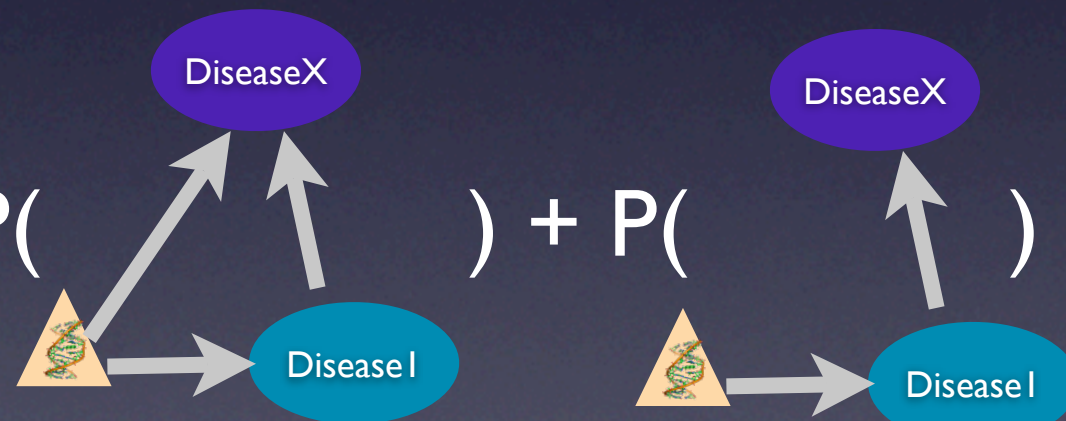


$0.10 :: \text{edge}(\text{'x_gene'}, \text{'DiseaseX'})$
 $0.66 :: \text{edge}(\text{'x_gene'}, \text{'DiseaseI'})$
 $0.39 :: \text{edge}(\text{'DiseaseI'}, \text{'x_gene'})$

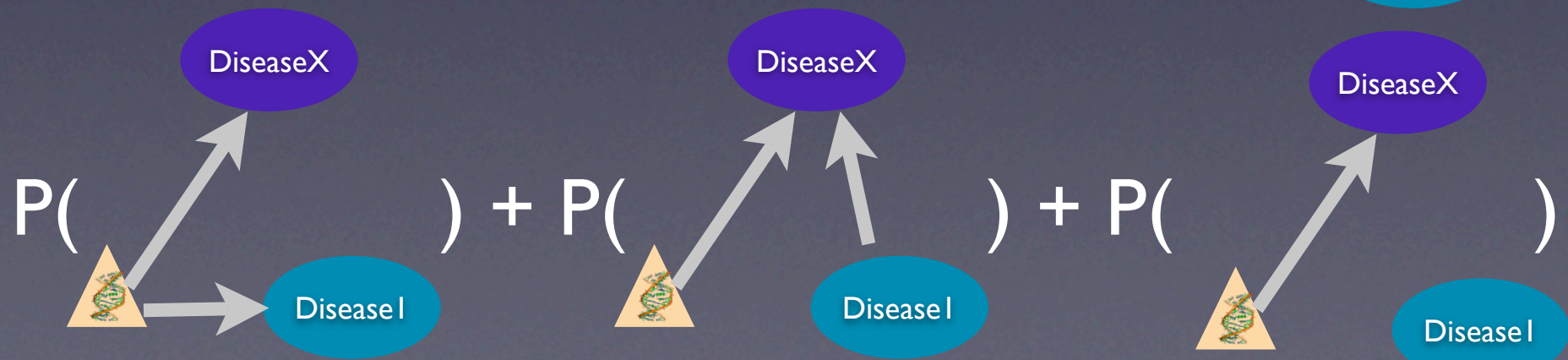
$\text{path}(X, Y) :- \text{edge}(X, Y)$

$\text{path}(X, Y) :- \text{edge}(X, Z), \text{path}(Z, Y)$

$P(\text{path}(\text{'x_gene'}, \text{'DiseaseX'})) = P($



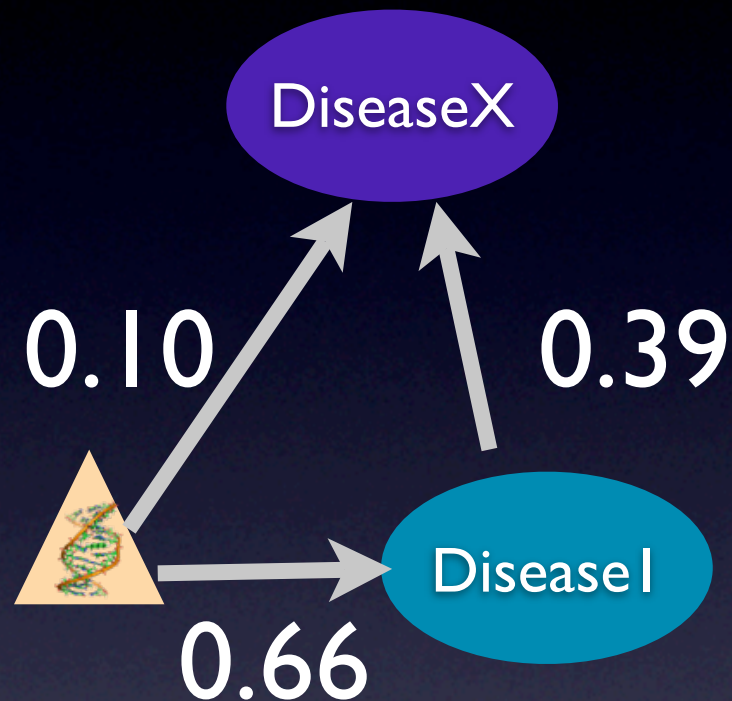
$) + P($



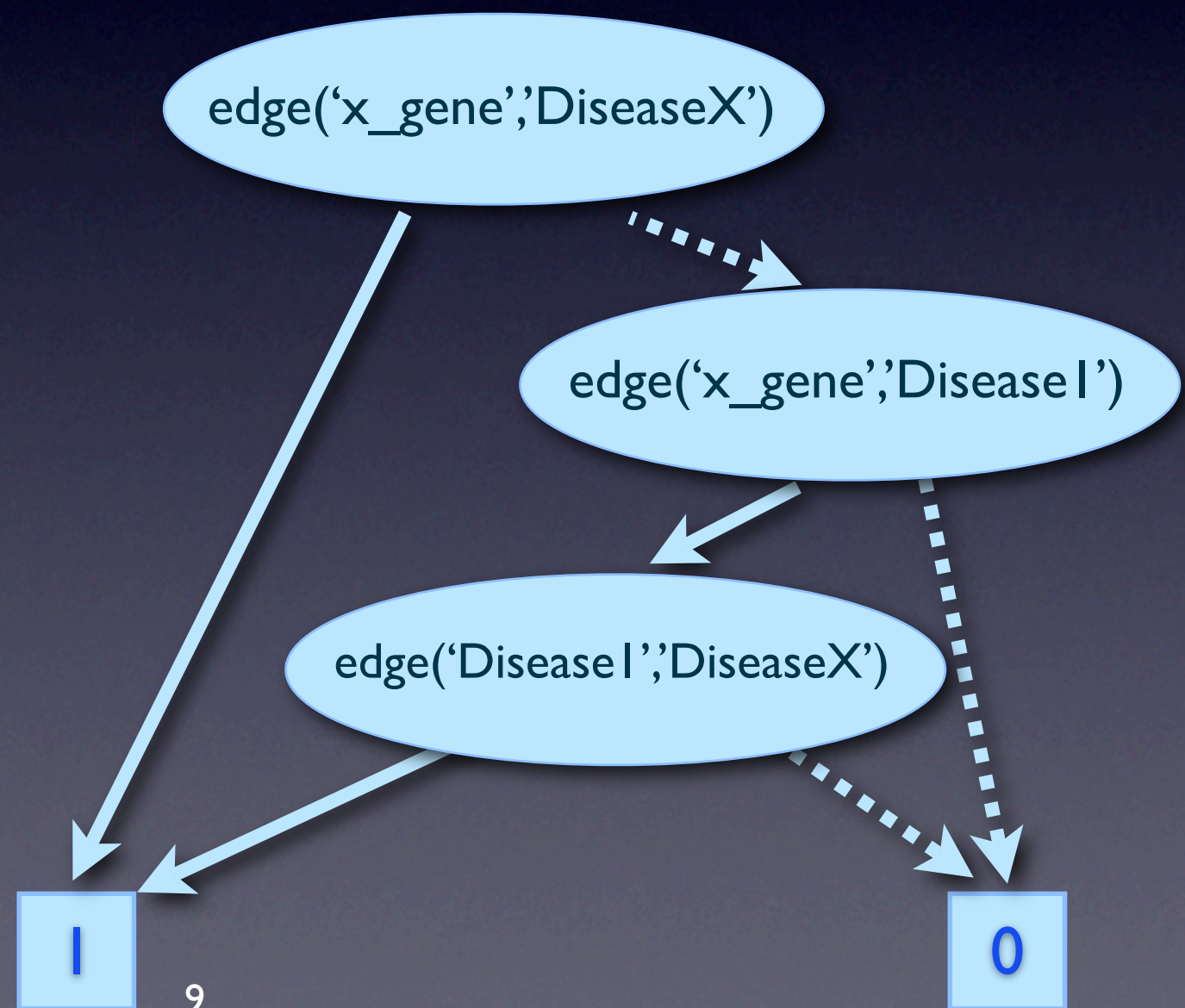
$) + P($

$) + P($

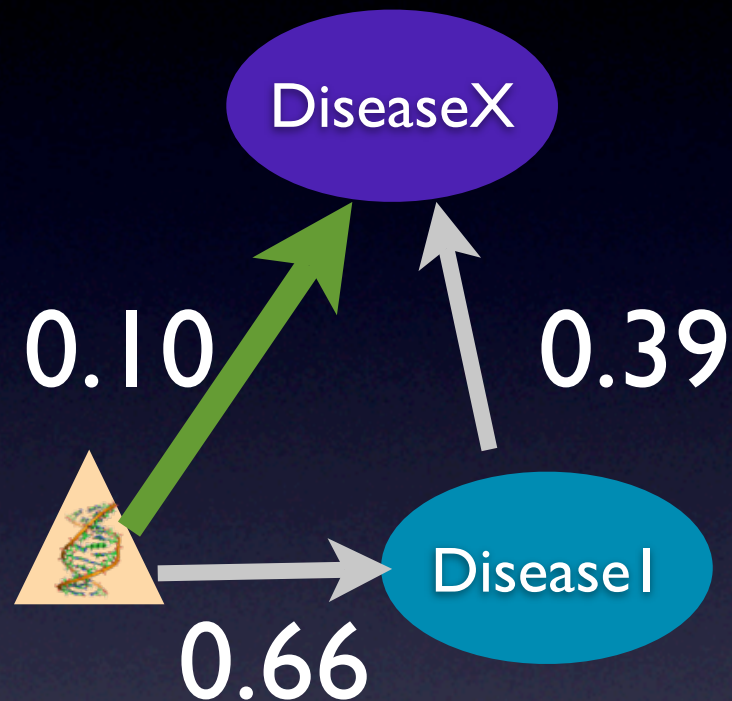
Calculating Probabilities



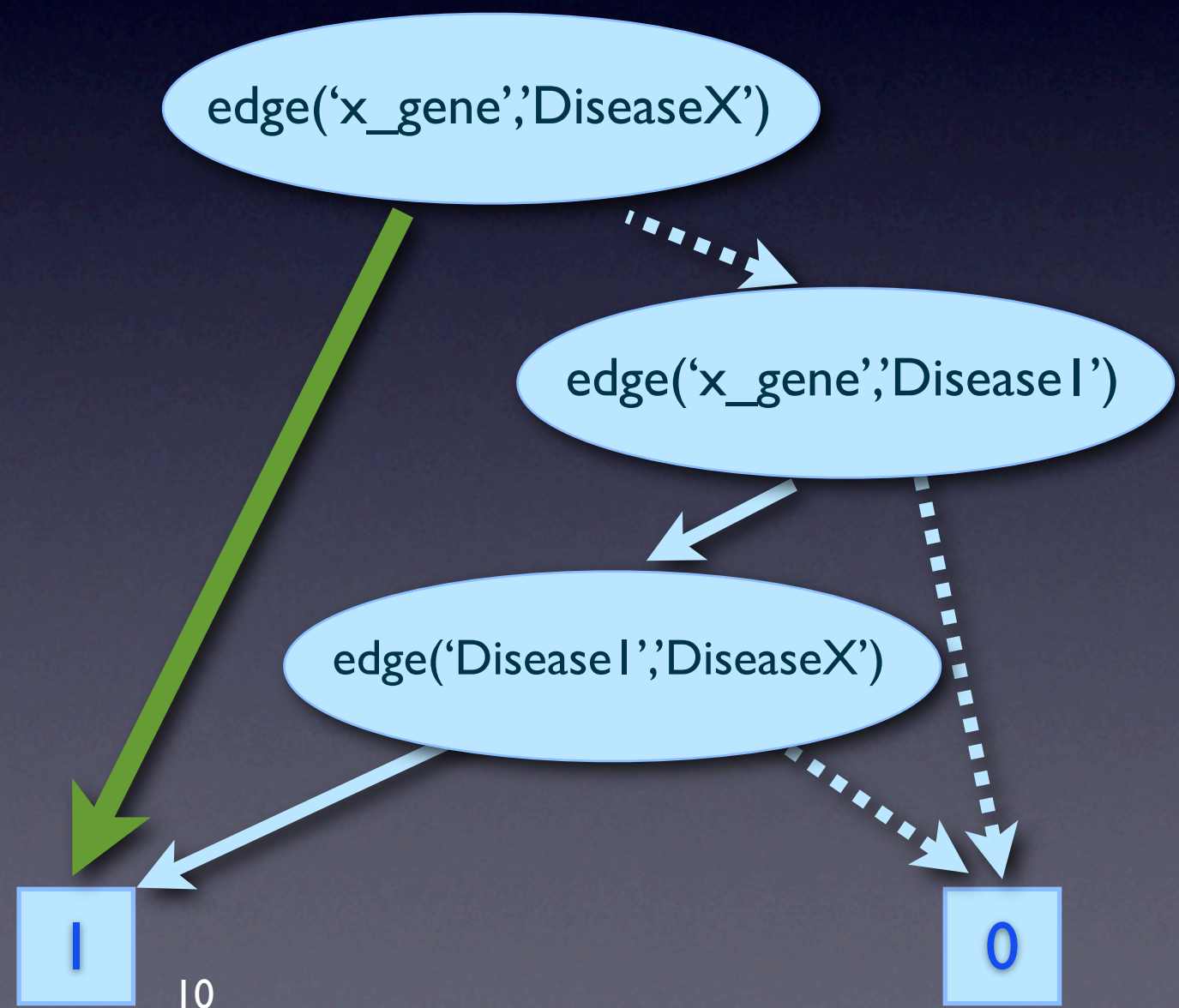
$$P_s('x_gene', 'DiseaseX'|L)=?$$



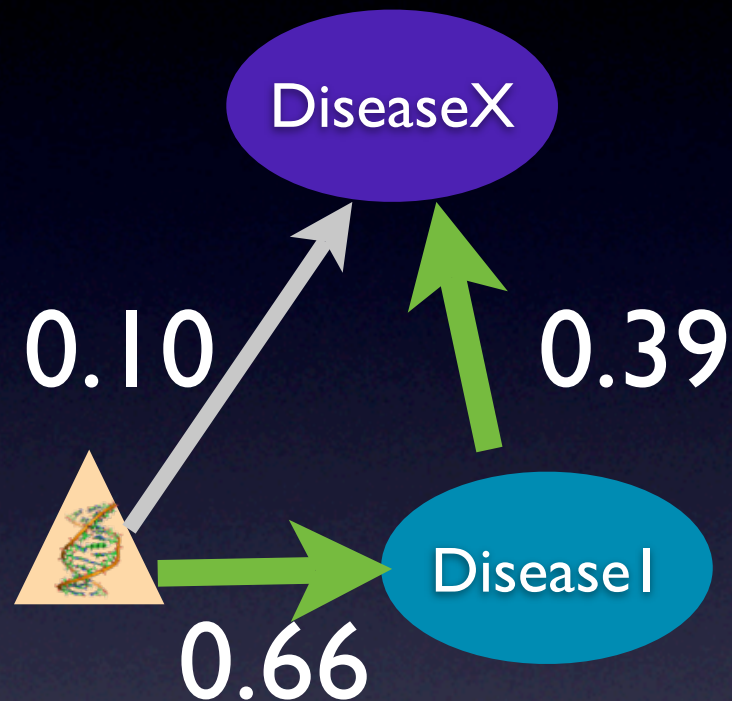
Calculating Probabilities



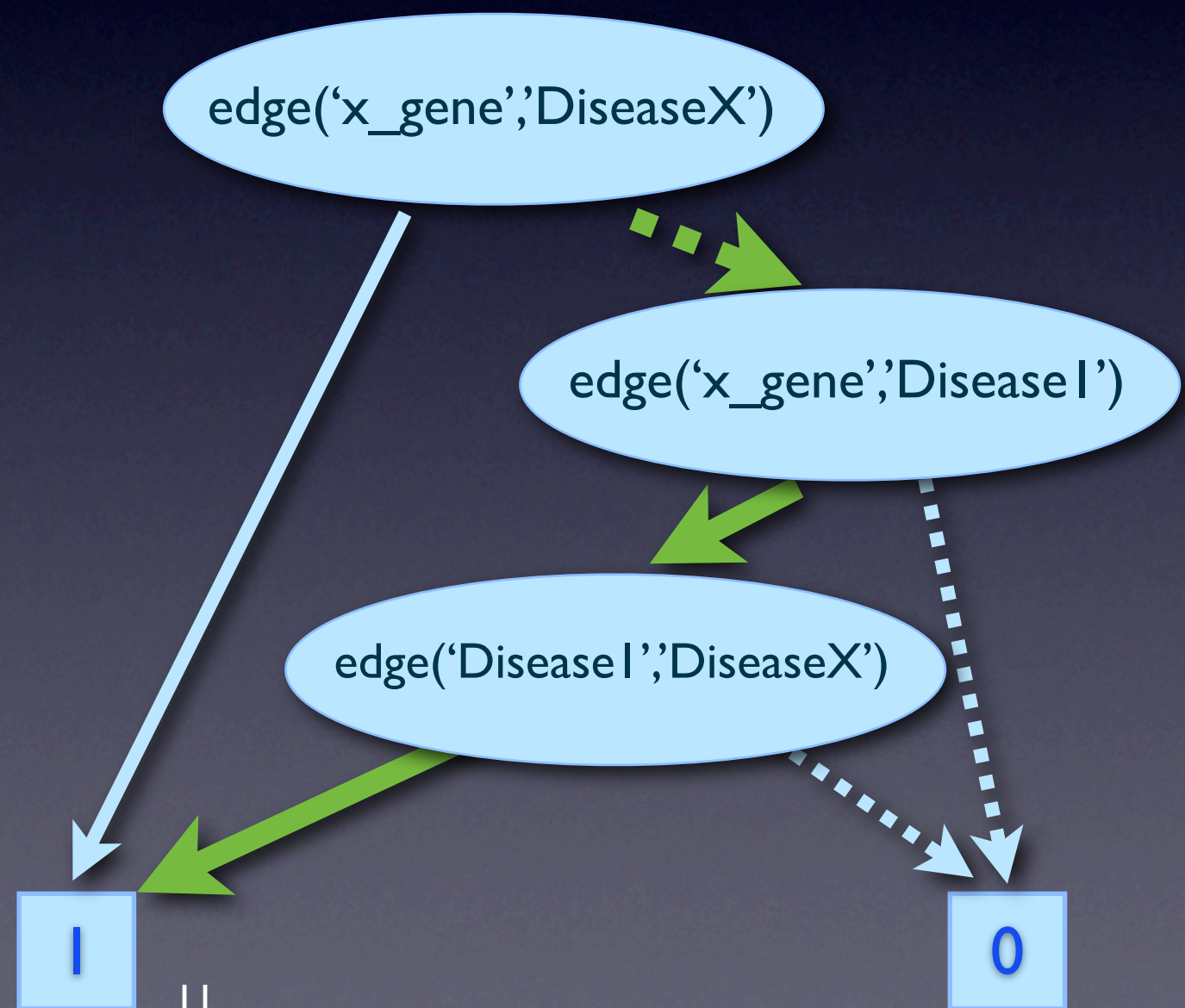
$$P_s('x_gene', 'DiseaseX'|L)=?$$



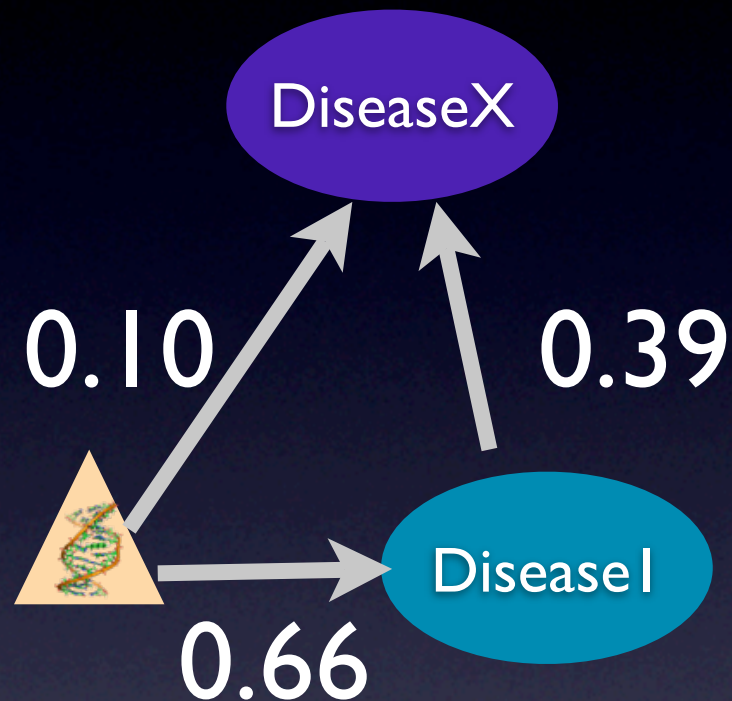
Calculating Probabilities



$$P_s('x_gene', 'DiseaseX'|L)=?$$



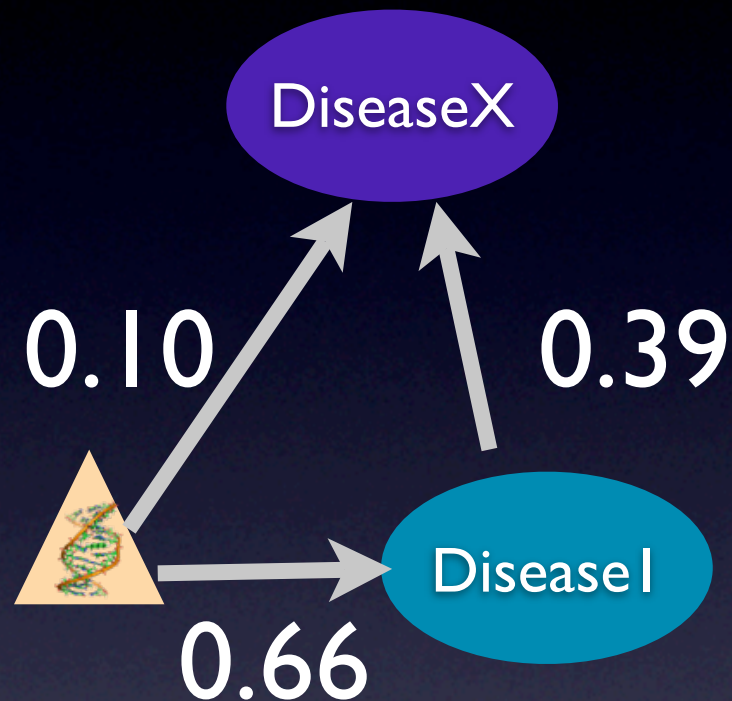
Calculating Probabilities



$$P_s('x_gene', 'DiseaseX'|L)=?$$



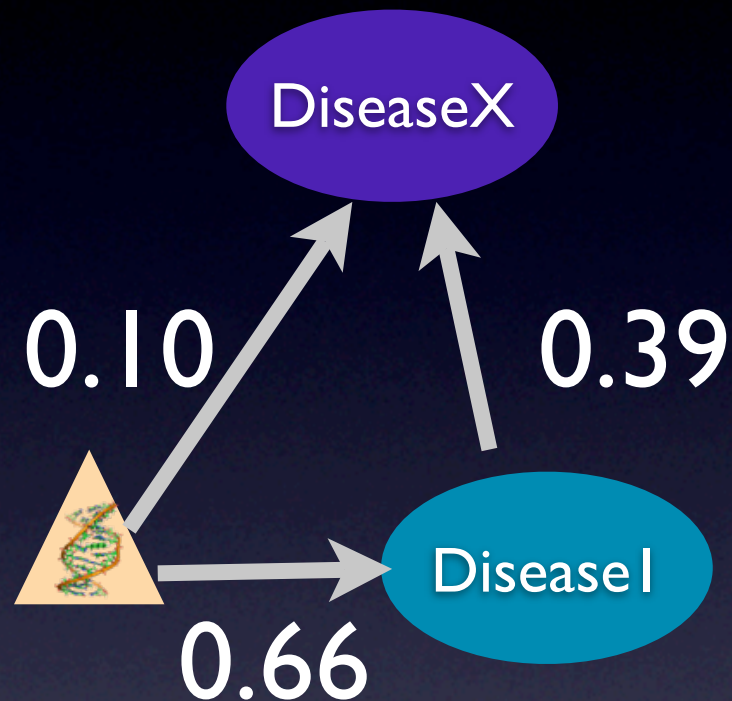
Calculating Probabilities



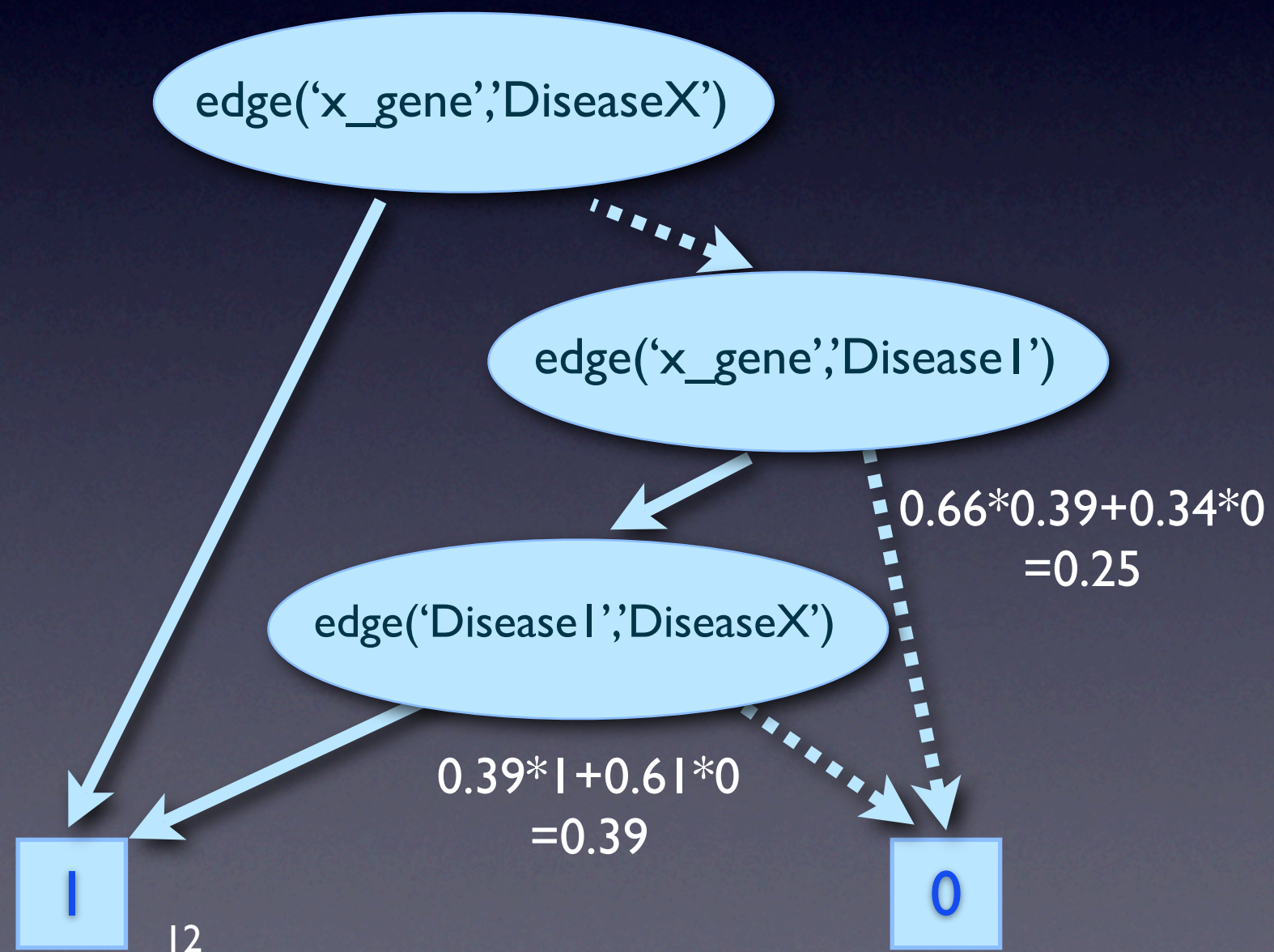
$$P_s('x_gene', 'DiseaseX'|L)=?$$



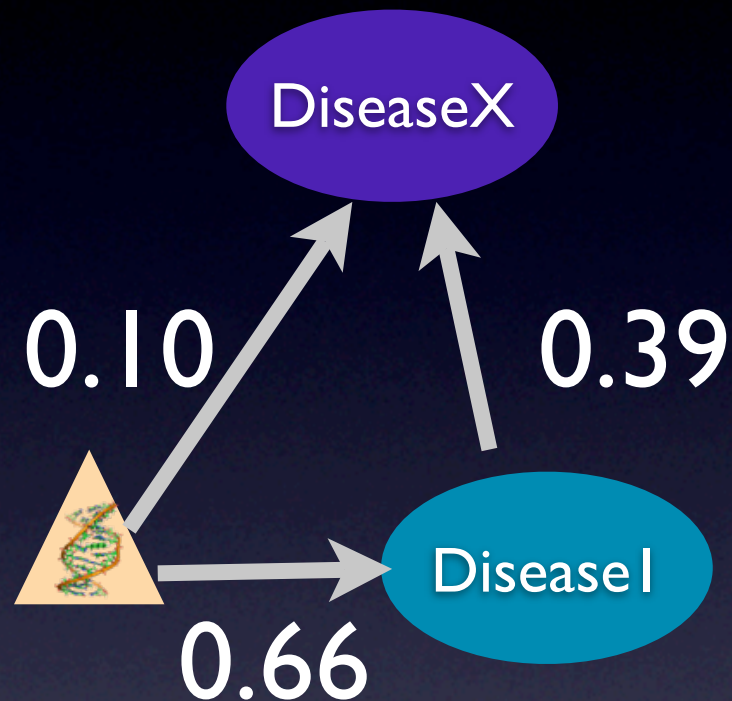
Calculating Probabilities



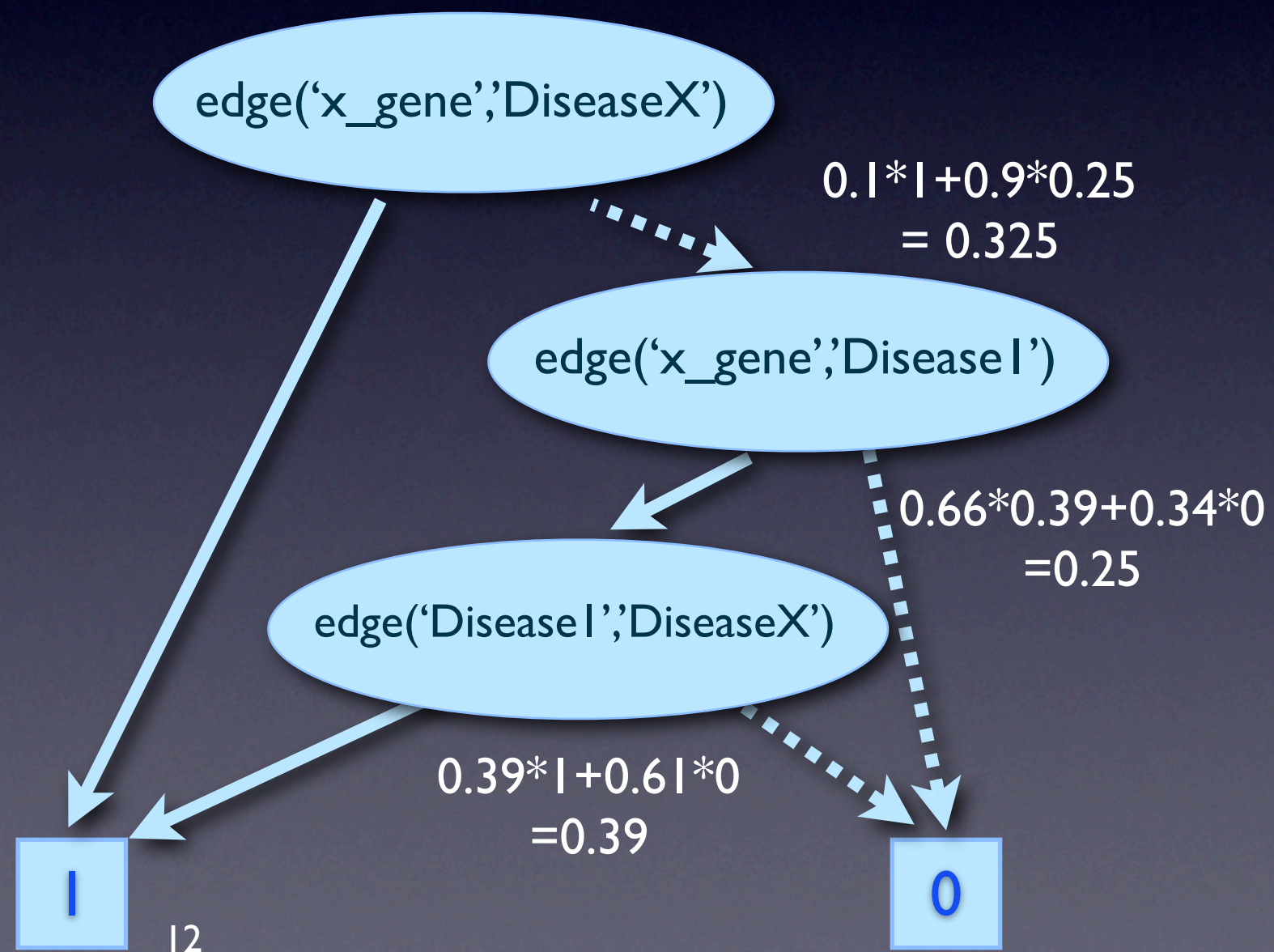
$$P_s('x_gene', 'DiseaseX'|L)=?$$



Calculating Probabilities



$$P_s('x_gene', 'DiseaseX'|L)=?$$



Parameter Learning

```
?? :: edge('x_gene', 'DiseaseX')  
?? :: edge('x_gene', 'DiseaseI')  
?? :: edge('DiseaseI', 'x_gene')
```

```
path(X,Y) :- edge(X,Y)  
path(X,Y) :- edge(X,Z), path(Z,Y)
```


Parameter Learning

```
?? :: edge('x_gene', 'DiseaseX')  
?? :: edge('x_gene', 'DiseaseI')  
?? :: edge('DiseaseI', 'x_gene')
```

```
path(X,Y) :- edge(X,Y)  
path(X,Y) :- edge(X,Z), path(Z,Y)
```

```
0.33 :: path('x_gene', 'DiseaseX')  
0.66 :: path('x_gene', 'DiseaseI')
```


Parameter Learning

?? :: edge('x_gene', 'DiseaseX')
?? :: edge('x_gene', 'DiseaseI')
?? :: edge('DiseaseI', 'x_gene')

path(X,Y) :- edge(X,Y)
path(X,Y) :- edge(X,Z), path(Z,Y)

$P_s(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}) \mid L) \iff 0.33 :: \text{path}(\text{'x_gene'}, \text{'DiseaseX'})$
 $P_s(\text{path}(\text{'x_gene'}, \text{'DiseaseI'}) \mid L) \iff 0.66 :: \text{path}(\text{'x_gene'}, \text{'DiseaseI'})$

Parameter Learning

?? :: edge('x_gene', 'DiseaseX')
?? :: edge('x_gene', 'DiseaseI')
?? :: edge('DiseaseI', 'x_gene')

path(X,Y) :- edge(X,Y)
path(X,Y) :- edge(X,Z), path(Z,Y)

$P_s(\text{path}(\text{'x_gene'}, \text{'DiseaseX'}) \mid L) \iff 0.33 :: \text{path}(\text{'x_gene'}, \text{'DiseaseX'})$
 $P_s(\text{path}(\text{'x_gene'}, \text{'DiseaseI'}) \mid L) \iff 0.66 :: \text{path}(\text{'x_gene'}, \text{'DiseaseI'})$

$$MSE(T) = \frac{1}{K} \sum_{1 \leq i \leq K} \left(P_s(q_i \mid T) - \tilde{p}_i \right)^2$$

Gradient Descent

Gradient Descent

- Init probabilities randomly

Gradient Descent

- Init probabilities randomly
- while MSE is too high

Gradient Descent

- Init probabilities randomly
- while MSE is too high
 - find proofs for queries and build BDDs

Gradient Descent

- Init probabilities randomly
- while MSE is too high
 - find proofs for queries and build BDDs
 - calculate gradient for MSE

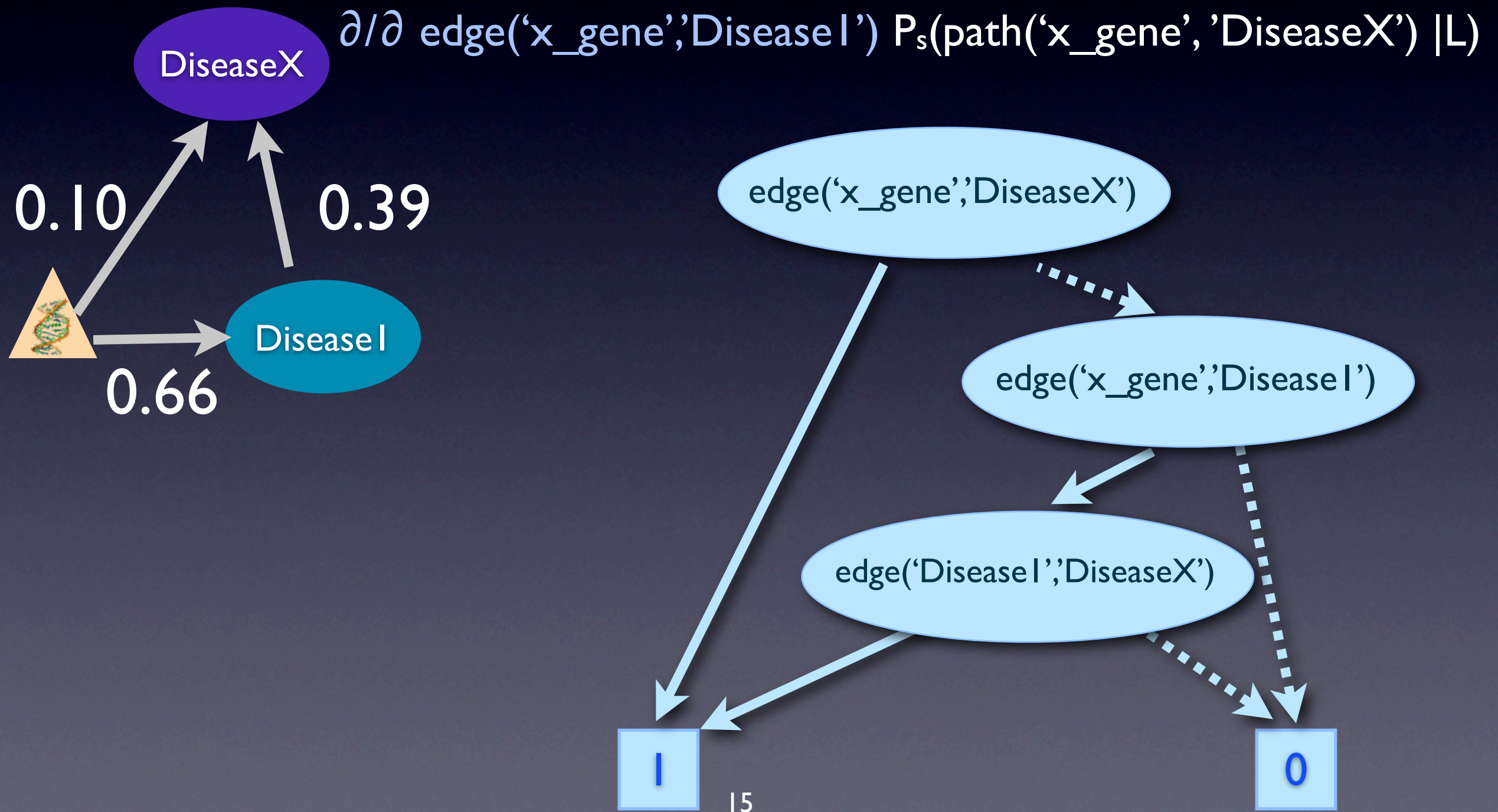
Gradient Descent

- Init probabilities randomly
- while MSE is too high
 - find proofs for queries and build BDDs
 - calculate gradient for MSE
 - add gradient

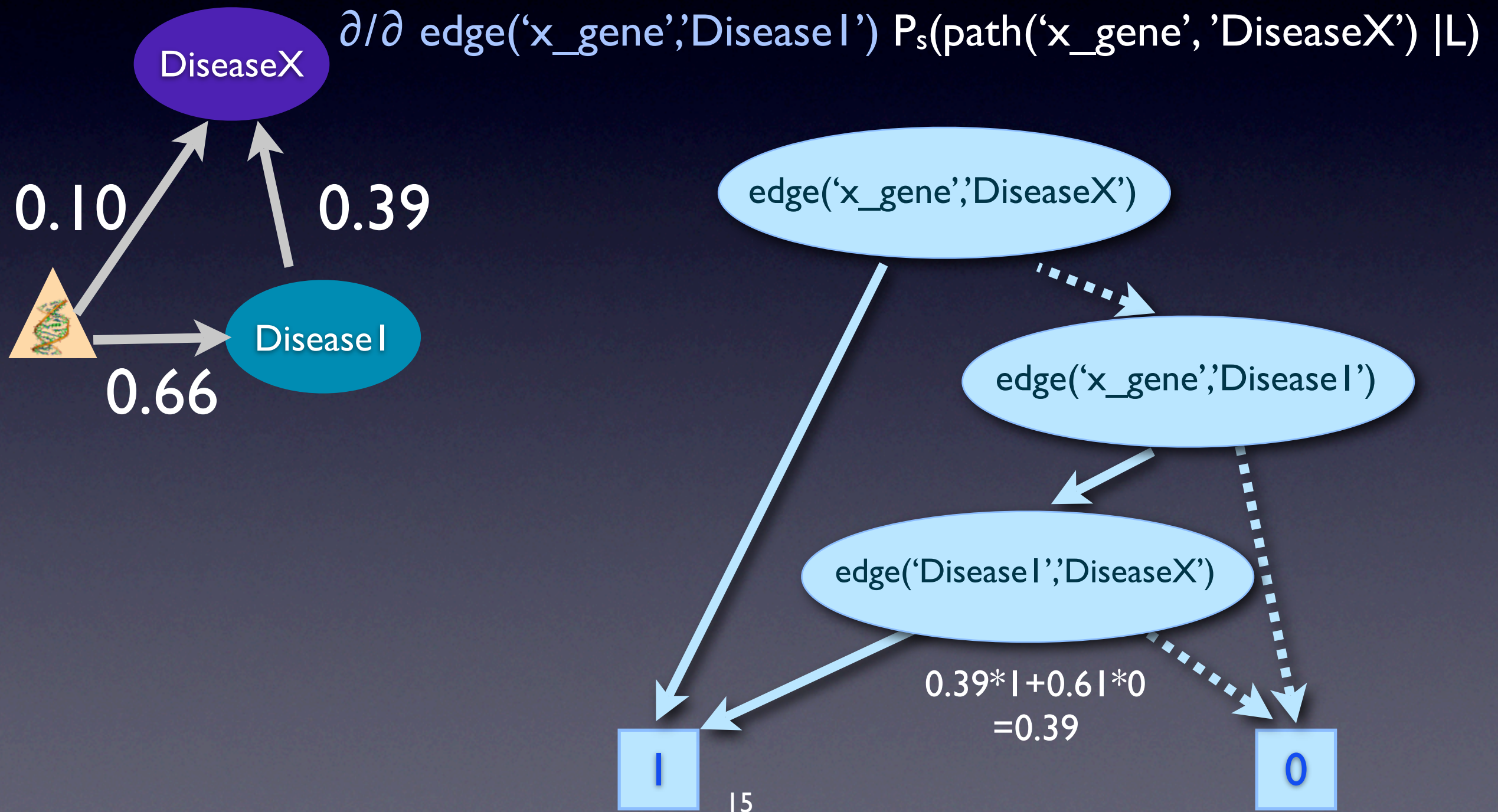
Gradient Descent

- Init probabilities randomly
- while MSE is too high
 - find proofs for queries and build BDDs
 - calculate gradient for MSE
 - add gradient
- return probabilities

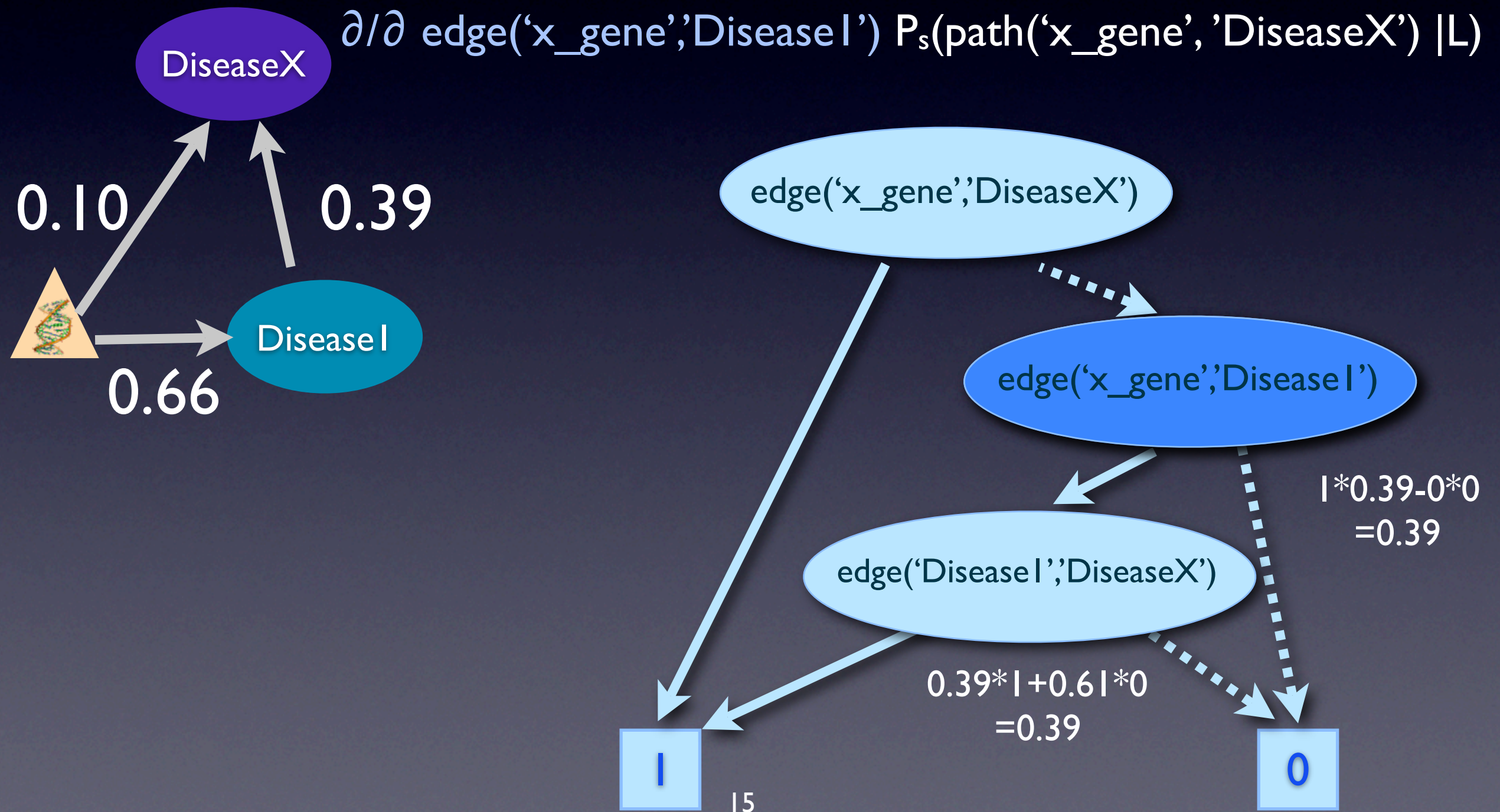
How to calculate the Gradient of the MSE?



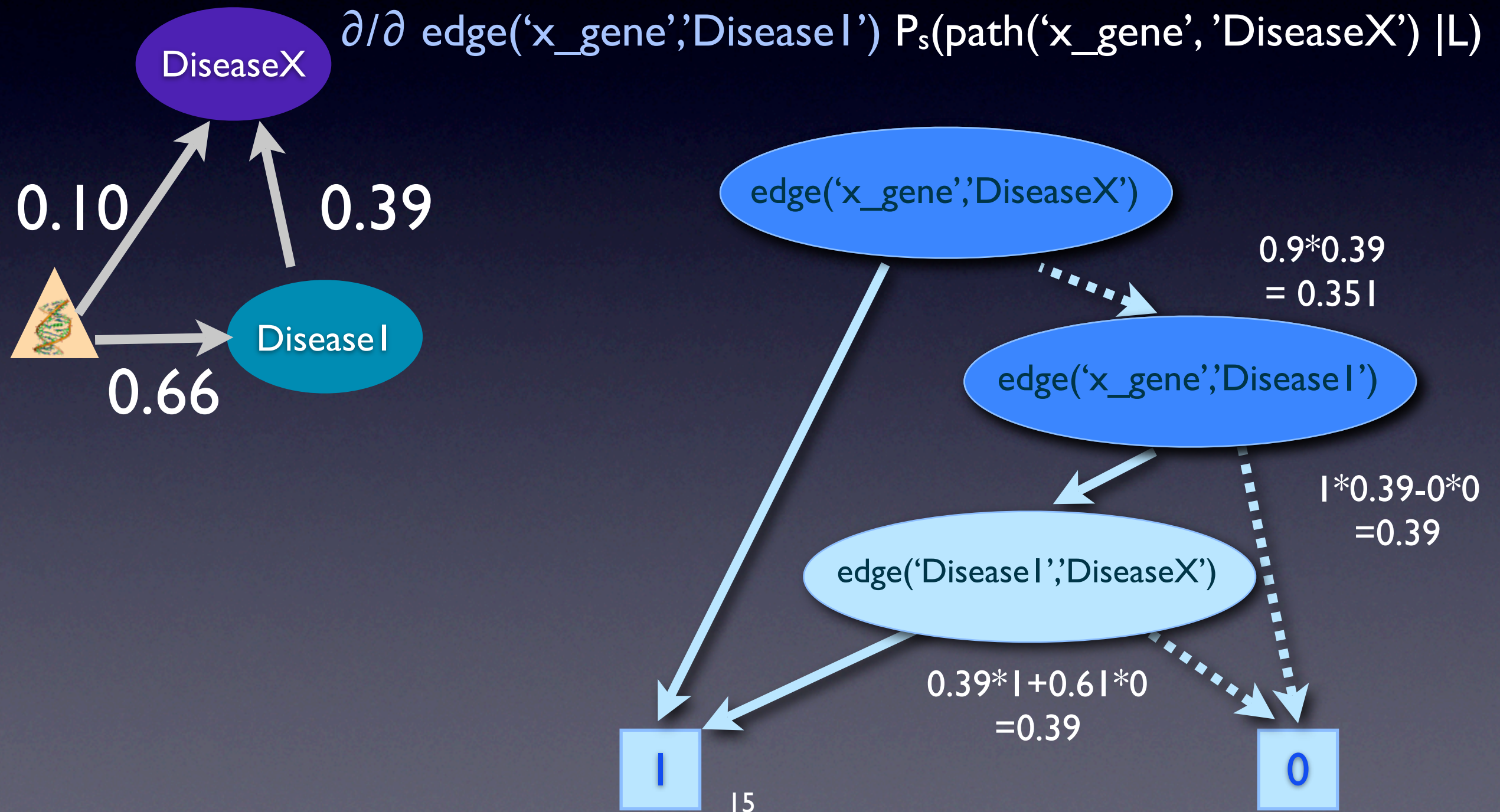
How to calculate the Gradient of the MSE?



How to calculate the Gradient of the MSE?



How to calculate the Gradient of the MSE?



Proofs are Queries

- queries explain that something is true, but they don't carry the concrete explanation

Proofs are Queries

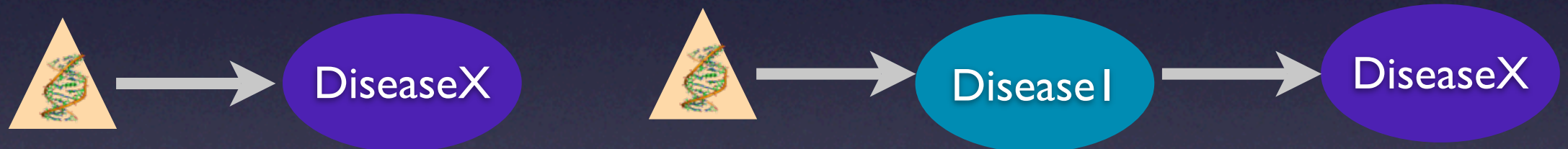
- queries explain that something is true, but they don't carry the concrete explanation

`path('x_gene', 'DiseaseX')`

Proofs are Queries

- queries explain that something is true, but they don't carry the concrete explanation

`path('x_gene', 'DiseaseX')`



Proofs are Queries

- queries explain that something is true, but they don't carry the concrete explanation

path('x_gene', 'DiseaseX')

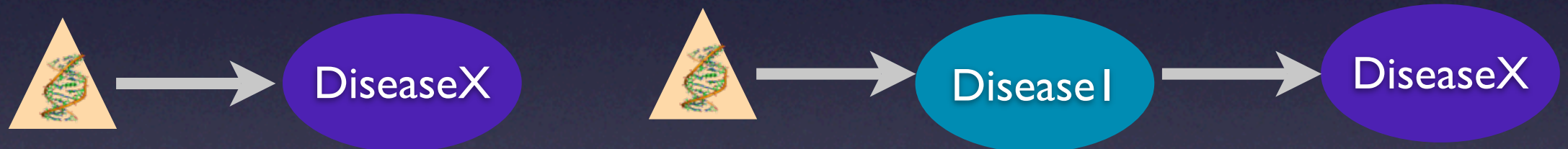


- a proof states, how something happens

Proofs are Queries

- queries explain that something is true, but they don't carry the concrete explanation

`path('x_gene', 'DiseaseX')`



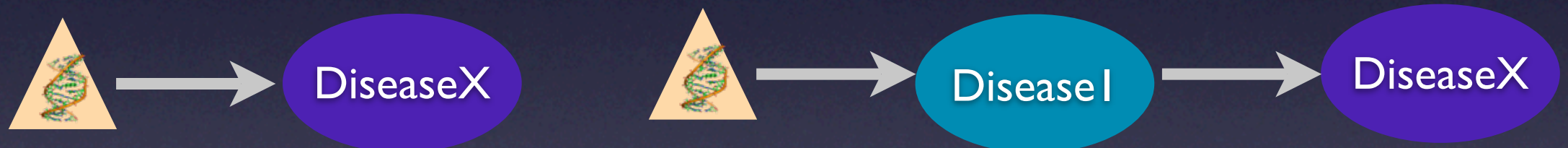
- a proof states, how something happens

`edge('x_gene', 'Disease I'), edge('Disease I', 'DiseaseX')`

Proofs are Queries

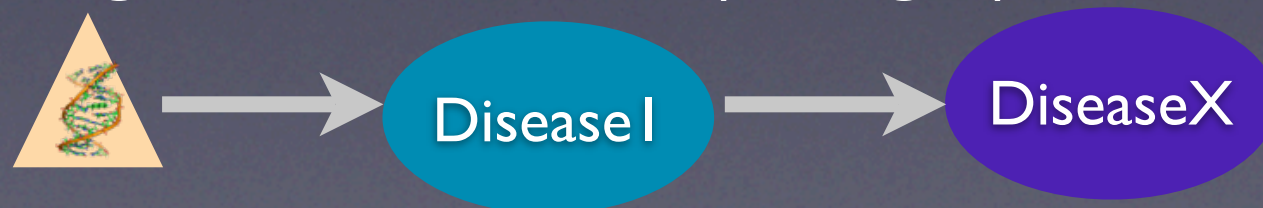
- queries explain that something is true, but they don't carry the concrete explanation

path('x_gene', 'DiseaseX')



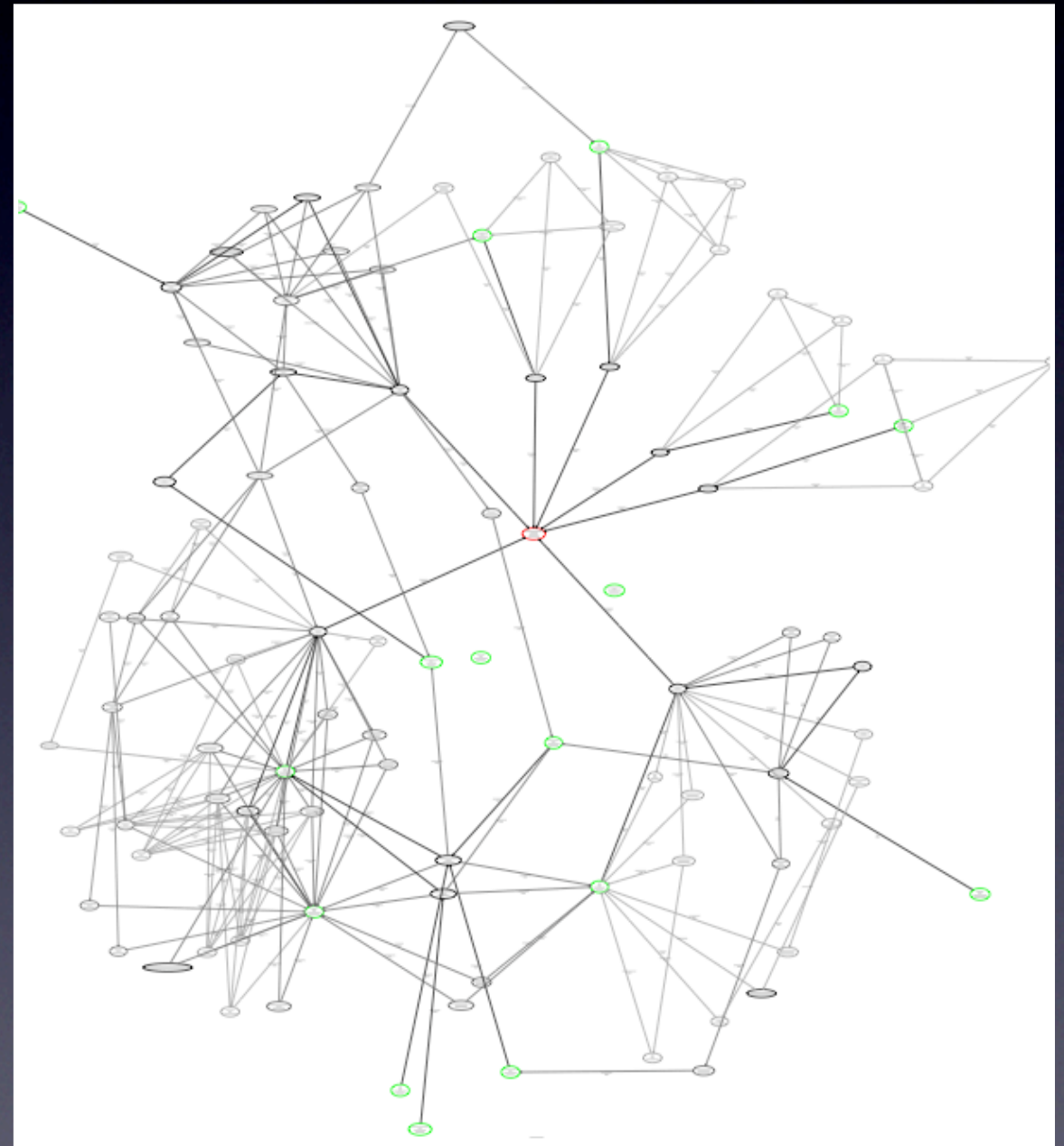
- a proof states, how something happens

edge('x_gene', 'DiseaseI'), edge('DiseaseI', 'DiseaseX')



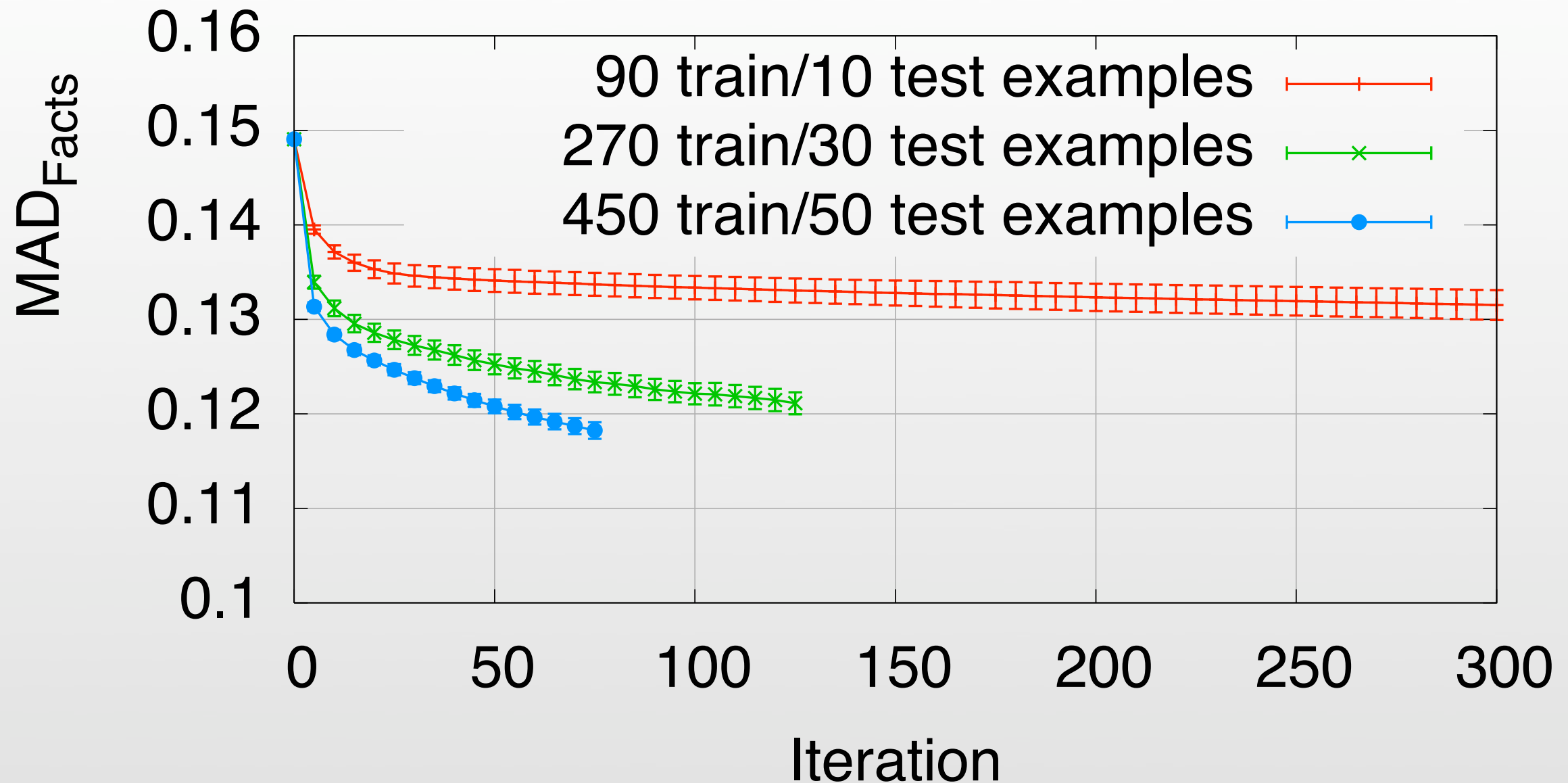
Experiments

- subgraph from the BIOMINE database around asthma genes
- 127 nodes, 241 edges
- sampled 500 random node pairs
- calculated $P(\text{path}(X,Y))$
- 10-fold cv



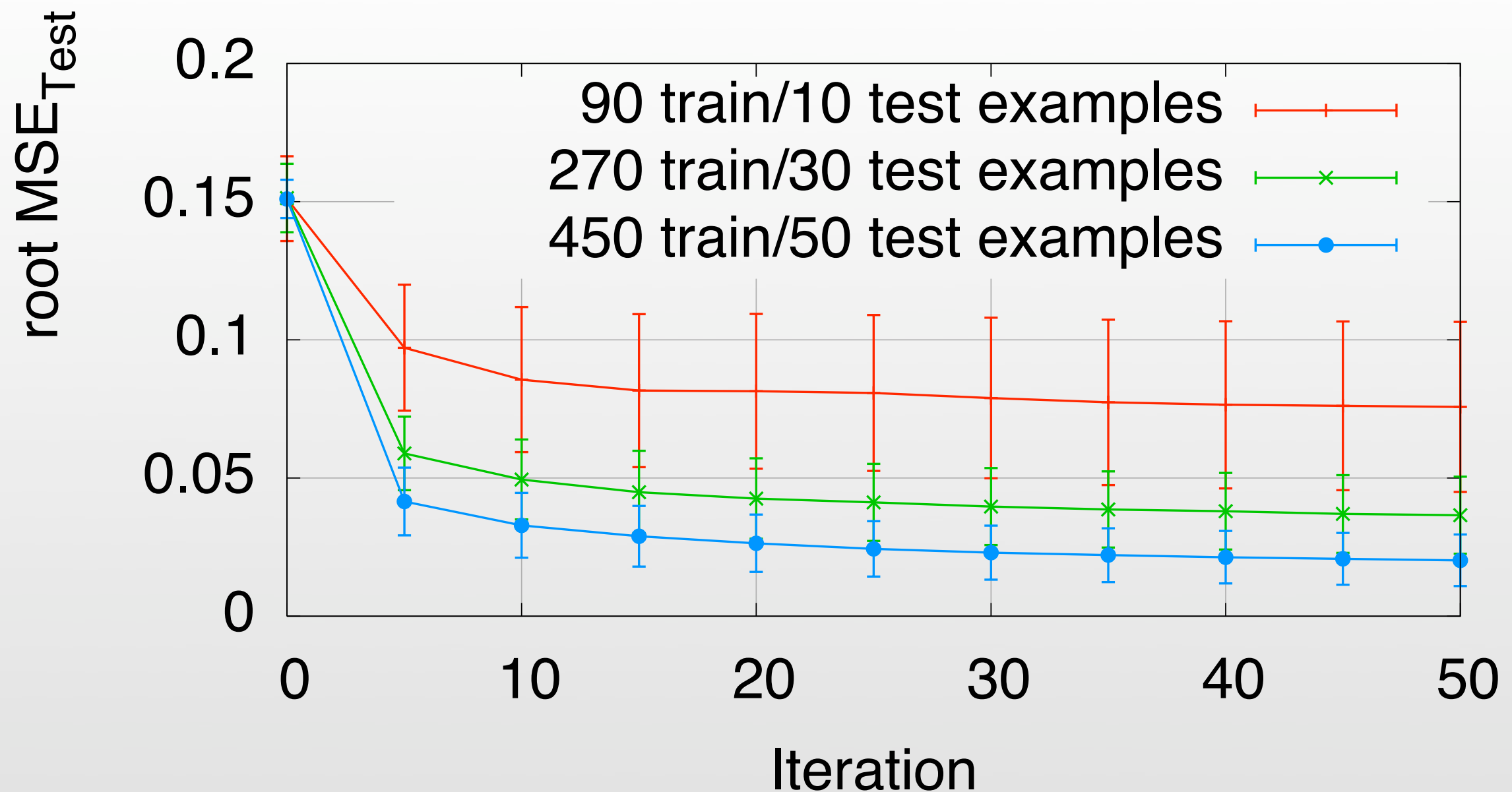
Recovering Probabilities

Asthma graph, update proofs every iteration, $k=5$



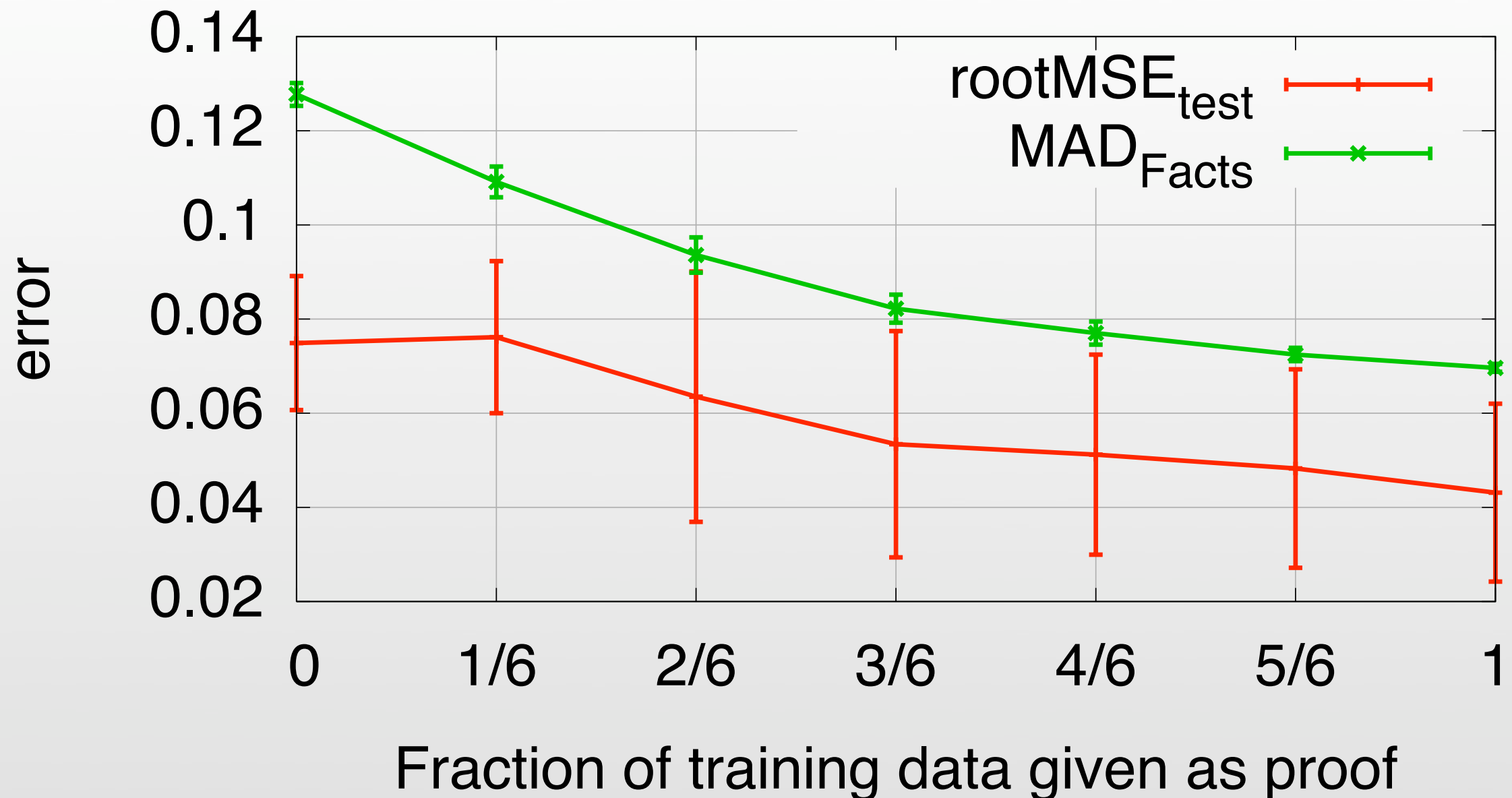
Reducing Test Error

Asthma graph, update proofs every iteration, $k=5$



Using Proofs

Asthma graph, update proofs every iteration, $k=1$



Summary

- it works
- it can naturally deal with proofs
- it is efficient due to the BDDs
- applicable for any prob. database where gradient of queries can be calculated

Questions?

More Details: See our ECML 2008 paper