# Bayesian Knowledge Tracing and Other Predictive Models

## Zachary A. Pardos

## PSLC Summer School 2012

**Ph.D. Committee**
Dr. Neil Heffernan, Advisor, WPI - Computer Science
Dr. Ryan S.J.d. Baker, WPI - Social Science and Policy Studies
Dr. Gabor Sarkozy, WPI - Computer Science
Dr. Kenneth Koedinger, CMU - Human Computer Interaction Institute

zpardos@gmail.com
http://wpi.edu/~zpardos

# Outline of Talk

- Introduction to Knowledge Tracing
  - History
  - Intuition
  - Generative example
  - Influence of parameters
    - Demo?
  - Prior Per Student model
  - Variations (and other models)
  - MATLAB Code demo

History in the literature

- Introduced in 1995 (Corbett & Anderson)

- Four parameter simplification of ACT-R theory of skill acquisition (Anderson 1993)

- Computations based on a variation of Bayesian calculations proposed in 1972 (Atkinson)

- Formalized as equivalent to a Dynamic Bayesian Network (Rye, 2004) "Student modeling based on belief networks"
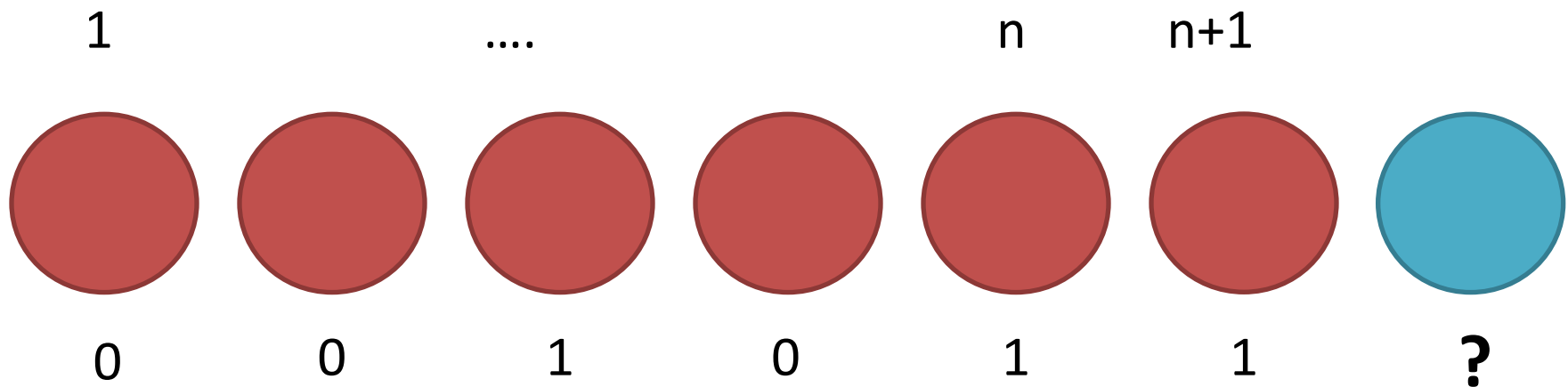
Real world deployment

- Used in the Cognitive Tutors (Carnegie Learning) to determine when a student has mastered a skill and can move on in the curriculum

- Replies on a skill model (tagging of skills to items)

- Parameters of the model can be learned with Expectation Maximization (EM) or grid search

# For some Skill K:

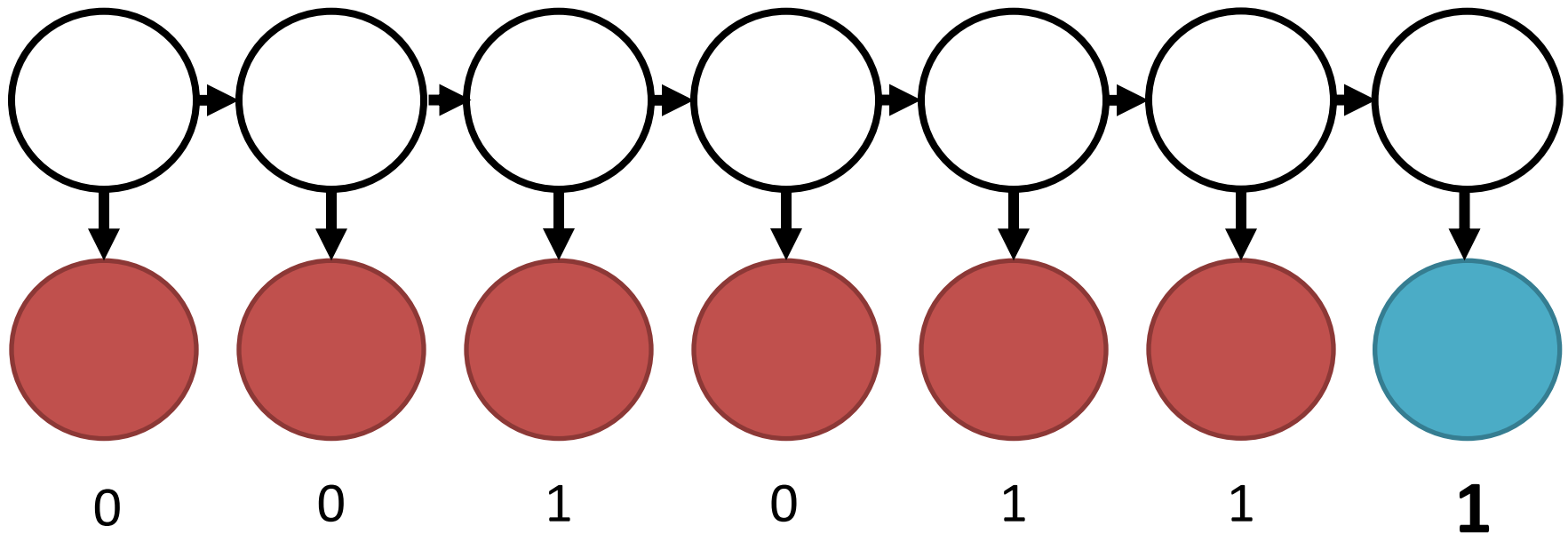Given a student's response sequence 1 to n, predict n+1

| 1 | | .... | | | n | n+1 |

| 0 | 0 | 1 | 0 | 1 | 1 | ? |

Chronological response sequence for student *Y*
[ 0 = Incorrect response     1 = Correct response]

**Track knowledge over time**
*(model of <u>learning</u>)*



0　　　0　　　1　　　0　　　1　　　1　　　**1**

# Intro to Knowledge Tracing

Knowledge Tracing (KT) can be represented as a simple HMM



Latent

Observed

Node representations
K = Knowledge node
Q = Question node

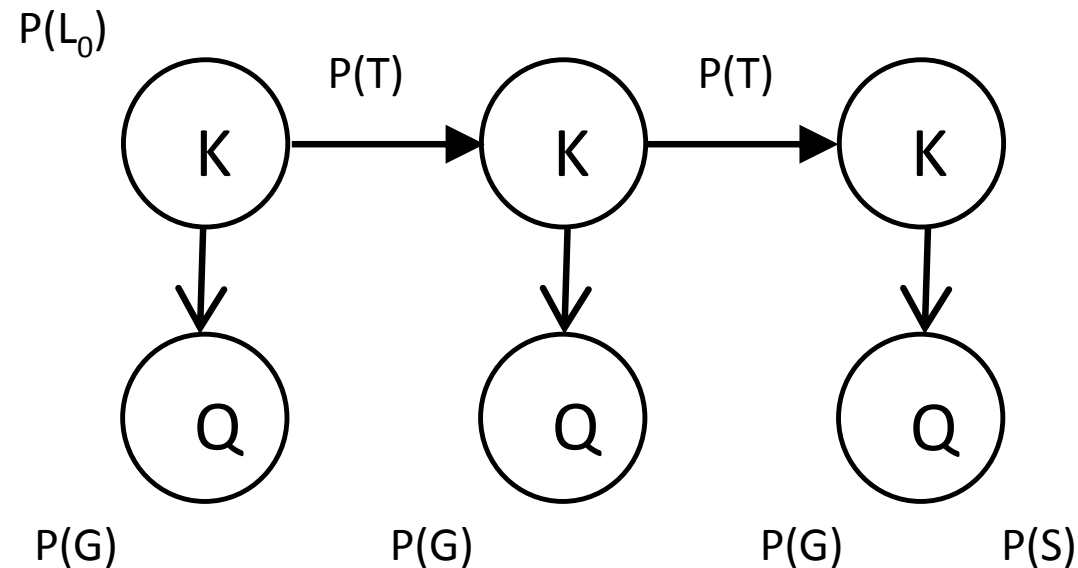Node states
K = Two state (0 or 1)
Q = Two state (0 or 1)

Four parameters of the KT model:

$P(L_0)$ = Probability of initial knowledge
$P(T)$ =  Probability of learning
$P(G)$ = Probability of guess
$P(S)$ = Probability of slip

$P(L_0)$

$P(T)$          $P(T)$

K → K → K

↓        ↓        ↓

Q        Q        Q

$P(G)$          $P(G)$          $P(G)$          $P(S)$

Probability of forgetting assumed to be zero (fixed)

# Formulas for inference and prediction

$If\ Correct_n$

$$P(L_{n-1}) = \frac{P(L_{n-1})*(1-P(S))}{P(L_{n-1})*(1-P(S))+ (1-P(L_{n-1}))*(P(G))} \qquad (1)$$

$Incorrect_n$

$$P(L_{n-1}) = \frac{P(L_{n-1})*P(S)}{P(L_{n-1})*P(S)+ (1-P(L_{n-1}))*(1-P(G))} \qquad (2)$$

$$P(L_n) =\ P((L_{n-1}) * (1 - P(F)) + ((1 - P(L_{n-1})) * P(T)) \qquad (3)$$

- Derivation (Reye, JAIED 2004):

$$p(L_{n-1} \mid C_n) = \frac{p(C_n \mid L_{n-1})p(L_{n-1})}{p(C_n \mid L_{n-1})p(L_{n-1}) + p(\neg L_{n-1})p(C_n \mid \neg L_{n-1})}$$

- Formulas use Bayes Theorem to make inferences about latent variable
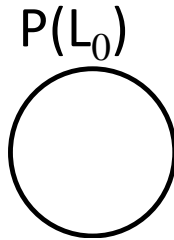
# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Generative - Example
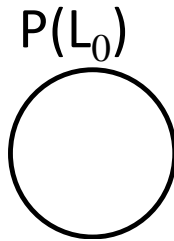
Zach Pardos | Predictive Models of Student Learning | July 20, 2012

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40

$P(L_0)$

How a Bayesian Knowledge Tracing World Works

Prior = 0.40

$P(L_0)$

How a Bayesian Knowledge Tracing World Works

Prior = 0.40

knowledge

$P(L_0)$

0

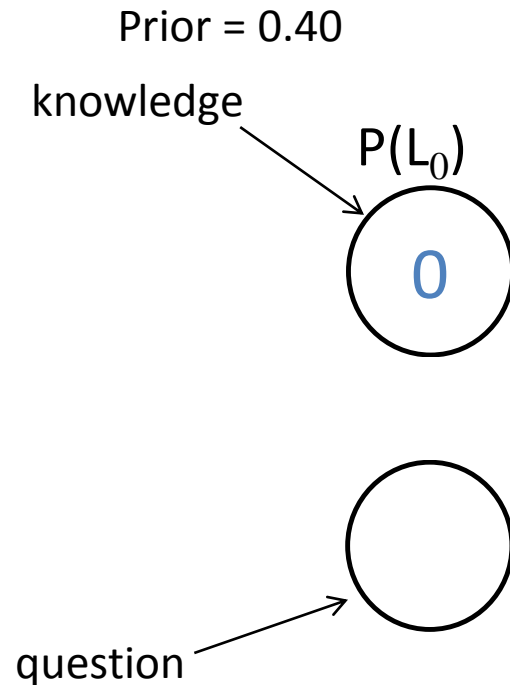How a Bayesian Knowledge Tracing World Works
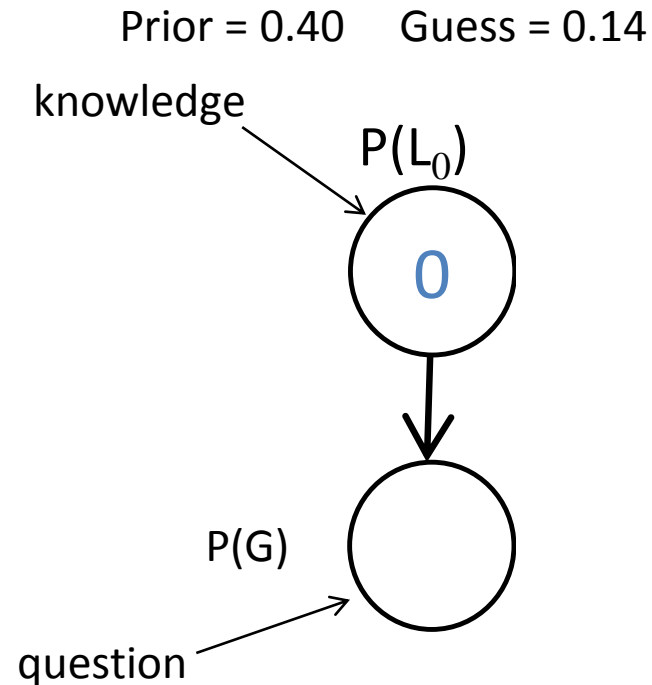
Prior = 0.40

knowledge

$P(L_0)$

0

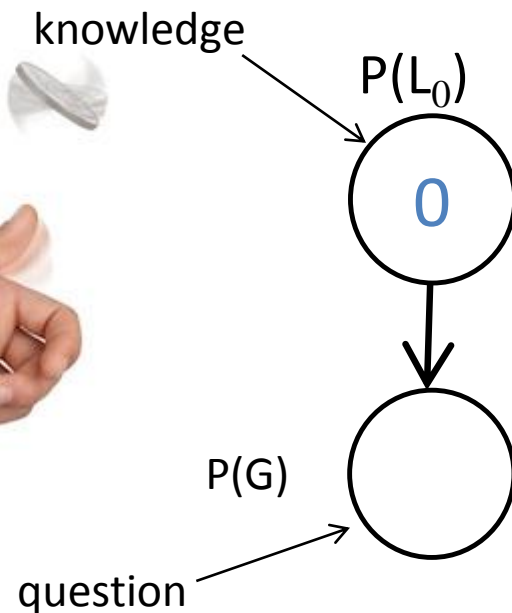question

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14

knowledge

$P(L_0)$

0

$P(G)$

question

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14

knowledge

$P(L_0)$

0

$P(G)$

question

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14

knowledge

$P(L_0)$

0

$P(G)$

0

question

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20

Zach Pardos                    Predictive Models of Student Learning                    July 20, 2012

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20

knowledge

$P(L_0)$          $P(T)$

0

P(G)

0

question

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20

knowledge

$P(L_0)$     $P(T)$



0 → 1

0

P(G)

question

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$    $P(T)$

$P(G)$    $P(S)$

question

## How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05



knowledge

$P(L_0)$    $P(T)$

0    1

P(G)    0    P(S)

question

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20     Slip = 0.05

knowledge

$P(L_0)$     $P(T)$



question

$P(G)$     $P(S)$
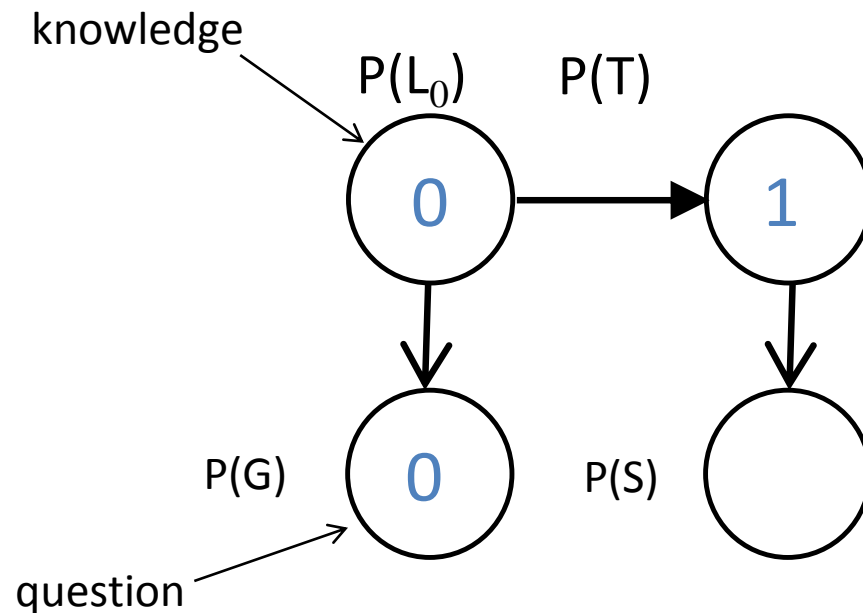
# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20     Slip = 0.05



Generalization of the response prediction calculation:

$$P(Correct_n) = P(L_n)\big(1 - P(S)\big) + (1 - P(L_n))P(G)$$

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20     Slip = 0.05



Generalization of the probability of learning calculation:

$$P(L_{n+1}) = P(L_n) + (1 - P(L_n))P(T)$$

Zach Pardos          Predictive Models of Student Learning          July 20, 2012

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20     Slip = 0.05



knowledge

$P(L_0)$     $P(T)$

You want to infer $P(L_n)$ from the student's responses

question

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$        $P(T)$

You want to infer $P(L_n)$ from the student's responses

question

First, infer the knowledge at the first opportunity:

$$P(Knowledge|Response = 0) = \frac{P(L_0)P(S)}{P(L_0)P(S) + (1 - P(L_0))(1 - P(G))}$$

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

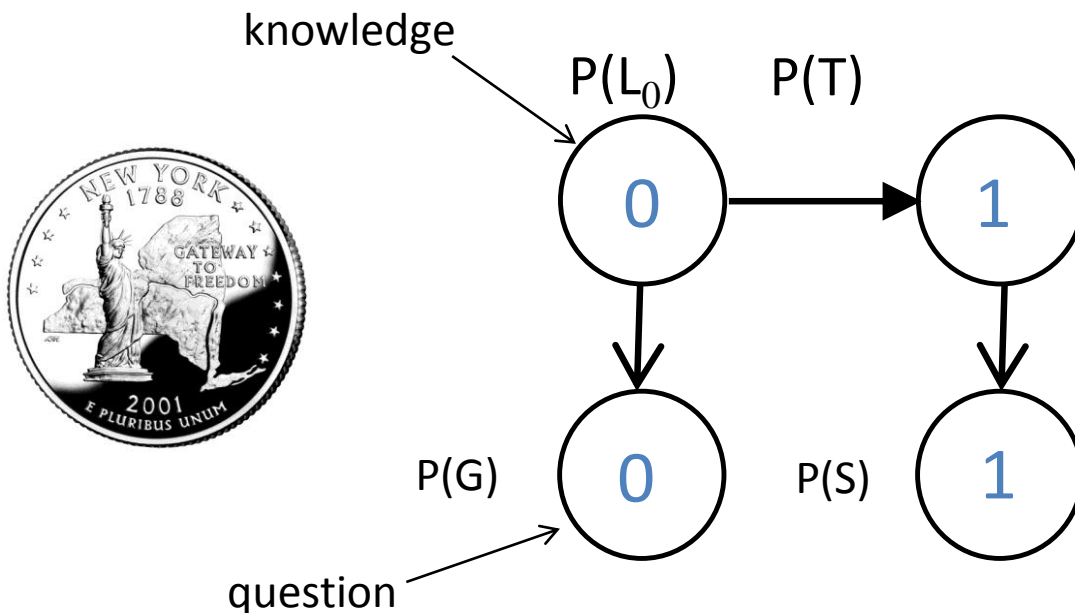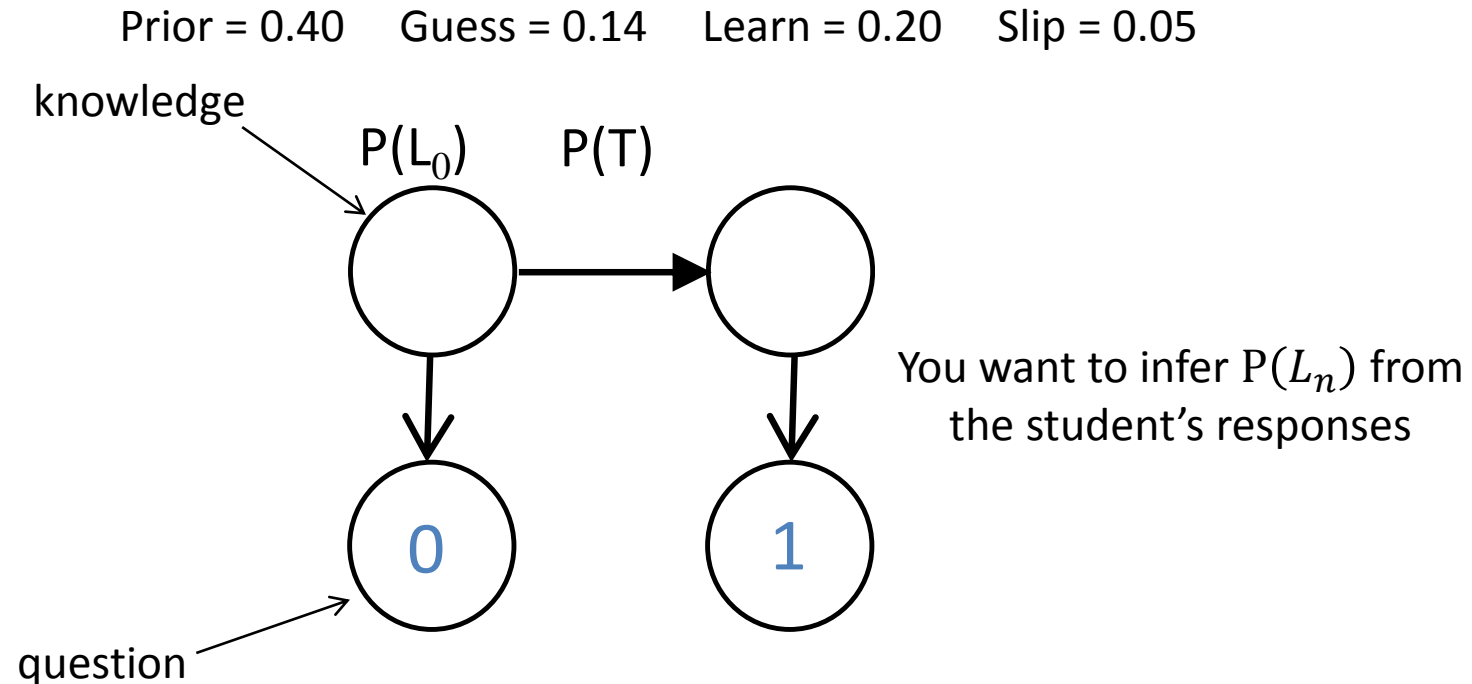Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$      $P(T)$

You want to infer $P(L_n)$ from the student's responses

question

First, infer the knowledge at the first opportunity:

$$P(Knowledge|Response = 0) = \frac{0.40 \cdot P(S)}{0.40 \cdot P(S) + (1 - 0.40)(1 - P(G))}$$

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$    $P(T)$

You want to infer $P(L_n)$ from the student's responses

0          1

question

First, infer the knowledge at the first opportunity:

$$P(Knowledge|Response = 0) = \frac{0.40 \cdot 0.05}{0.40 \cdot 0.05 + (1 - 0.40)(1 - P(G))}$$

29

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$    $P(T)$

You want to infer $P(L_n)$ from the student's responses

question
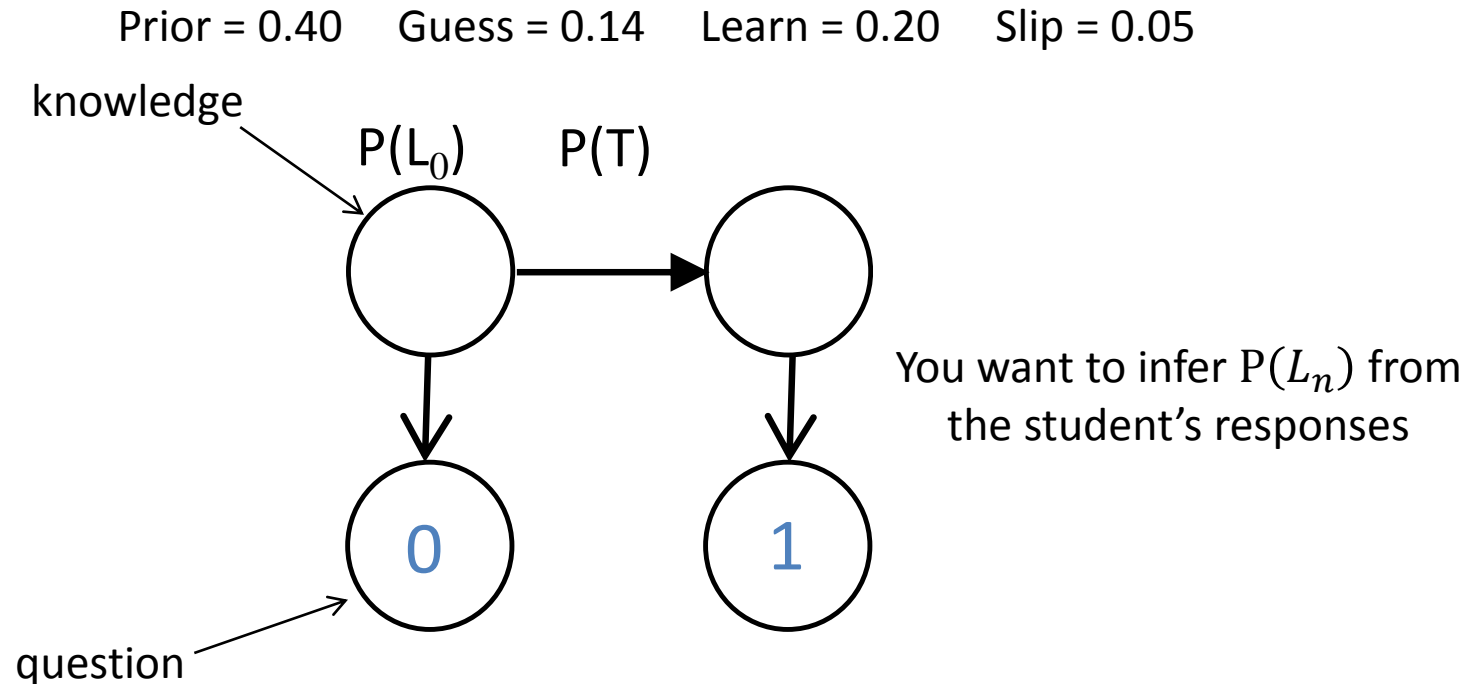
First, infer the knowledge at the first opportunity:

$$P(Knowledge|Response = 0) = \frac{0.40 \cdot 0.05}{0.40 \cdot 0.05 + (1 - 0.40)(1 - 0.14)}$$

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20     Slip = 0.05

knowledge

$P(L_0)$     $P(T)$

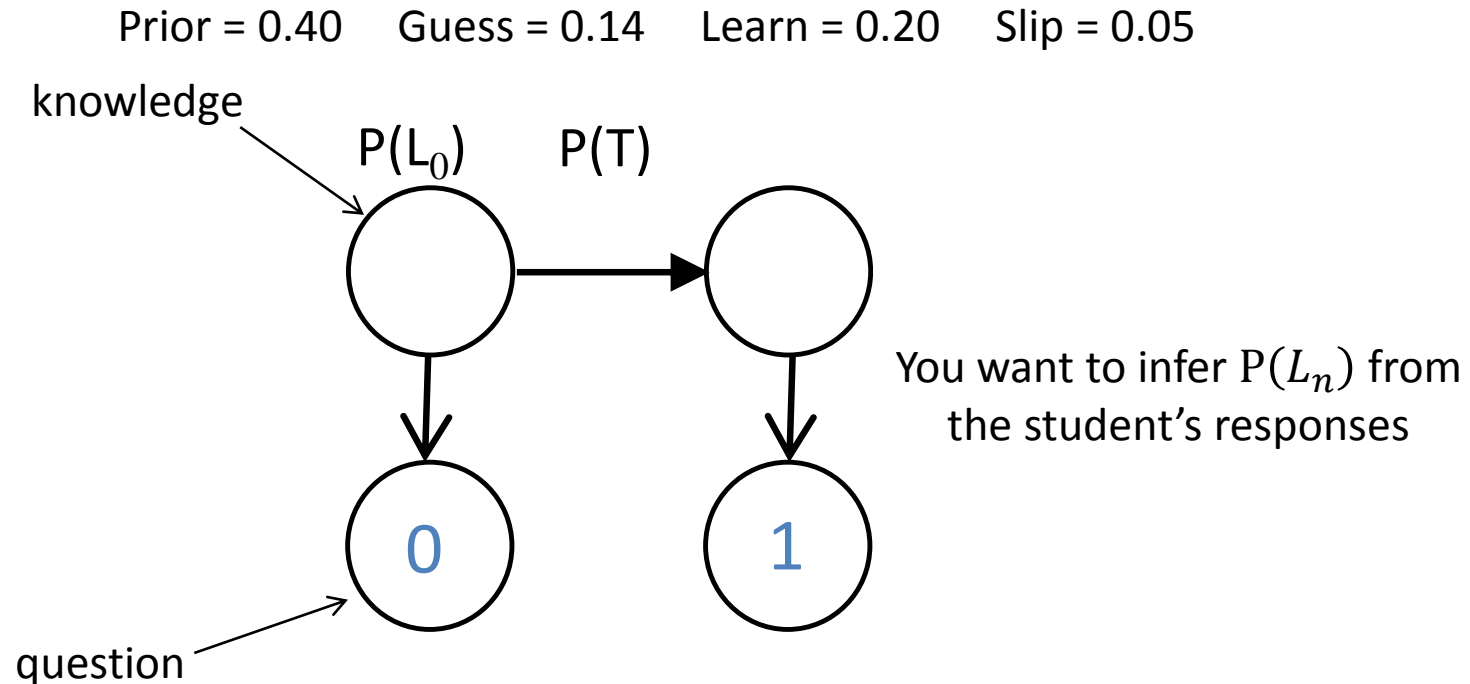You want to infer $P(L_n)$ from the student's responses

0     1

question
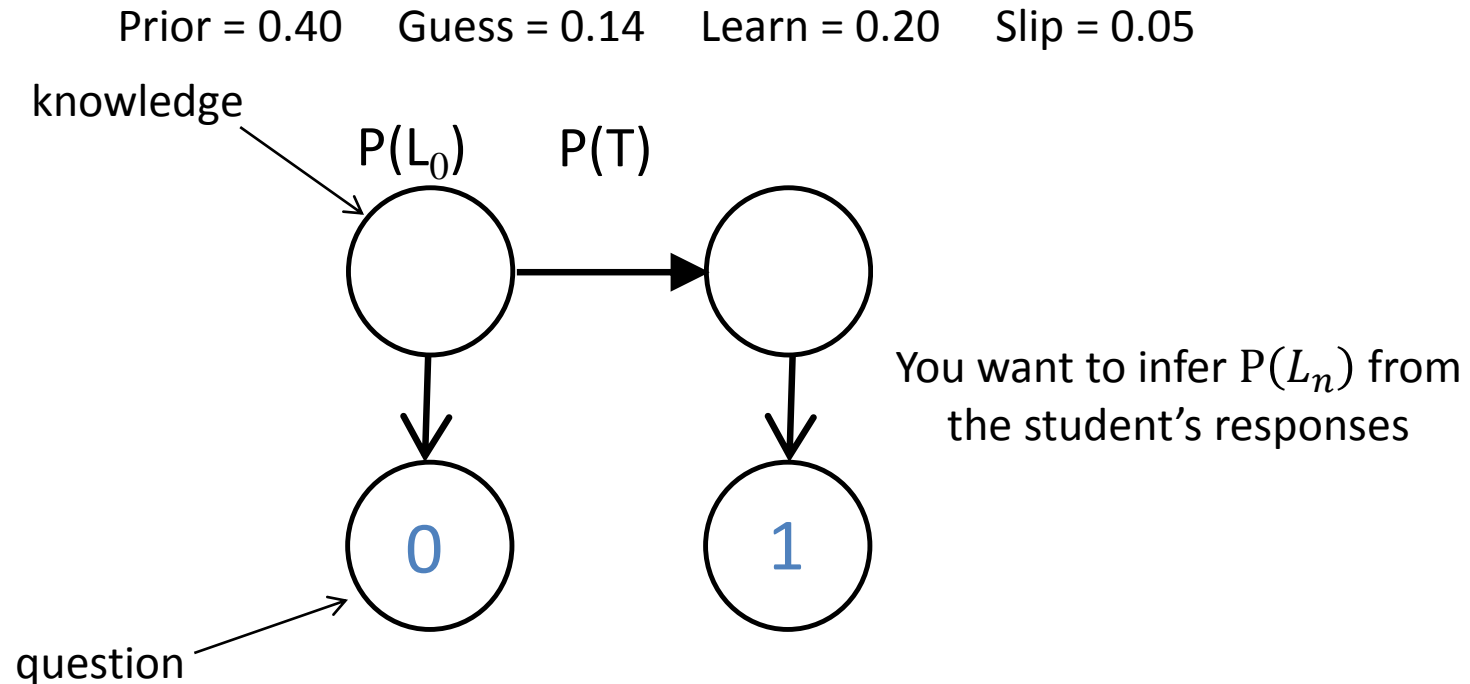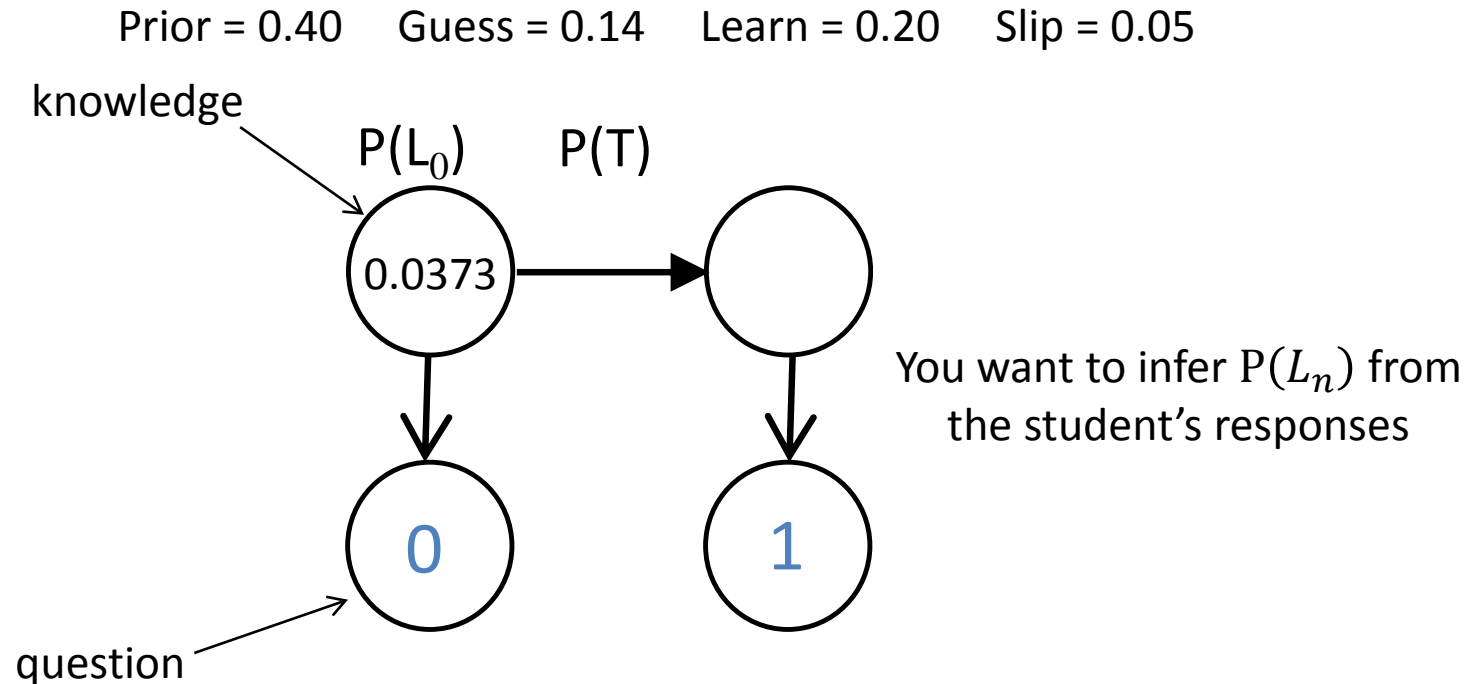
First, infer the knowledge at the first opportunity:

$$P(Knowledge|Response = 0) = \frac{0.40 \cdot 0.05}{0.40 \cdot 0.05 + (1 - 0.40)(1 - 0.14)} =$$

31

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$    $P(T)$

0.0373

You want to infer $P(L_n)$ from the student's responses

0

1

question
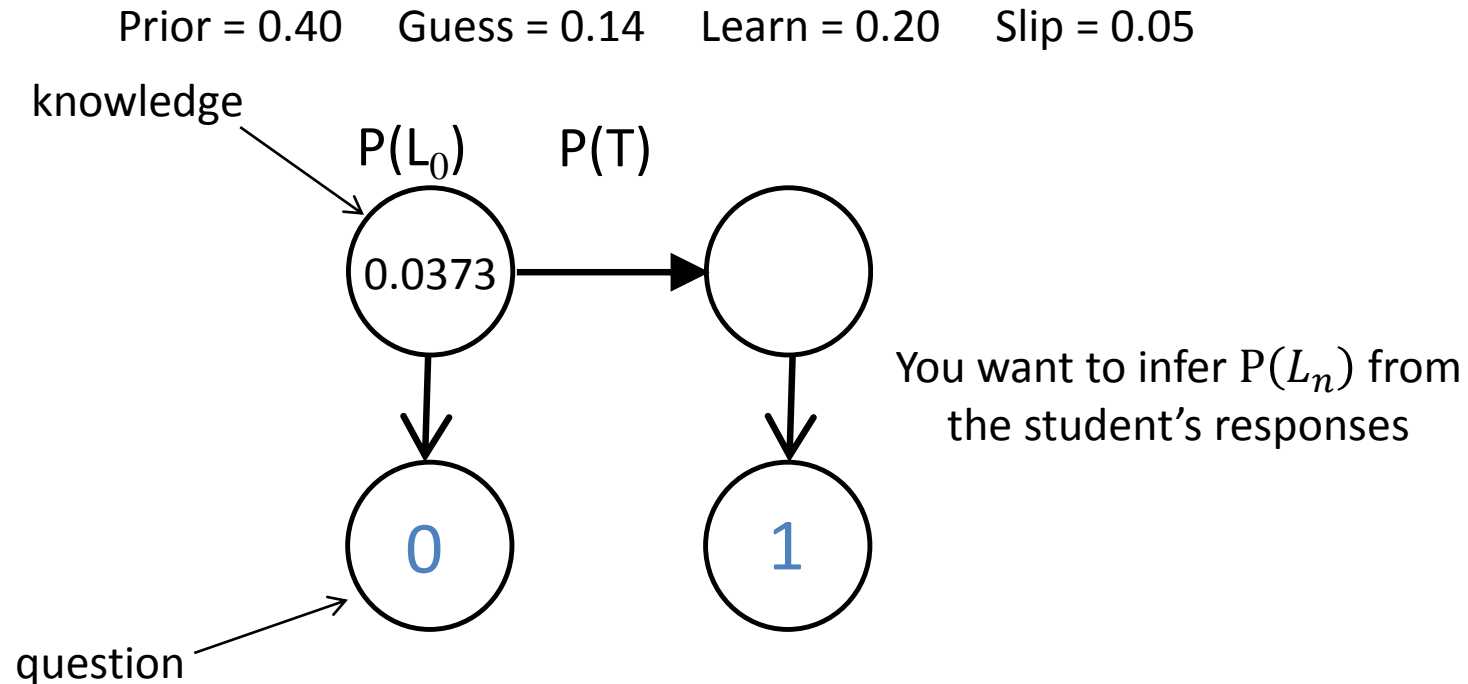
First, infer the knowledge at the first opportunity:    Posterior probability of knowledge

$$P(Knowledge|Response = 0) = \frac{0.40 \cdot 0.05}{0.40 \cdot 0.05 + (1 - 0.40)(1 - 0.14)} = 0.0373$$

# Knowledge Tracing

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20     Slip = 0.05

knowledge

$P(L_0)$         $P(T)$



0.0373

You want to infer $P(L_n)$ from the student's responses

0

1

question

Next, apply the learning transition formula:

$$P(L_{n+1}) = 0.0373 + (1 - 0.0373)(0.20) =$$

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$    $P(T)$    0.2298

0.0373

question

You want to infer $P(L_n)$ from the student's responses
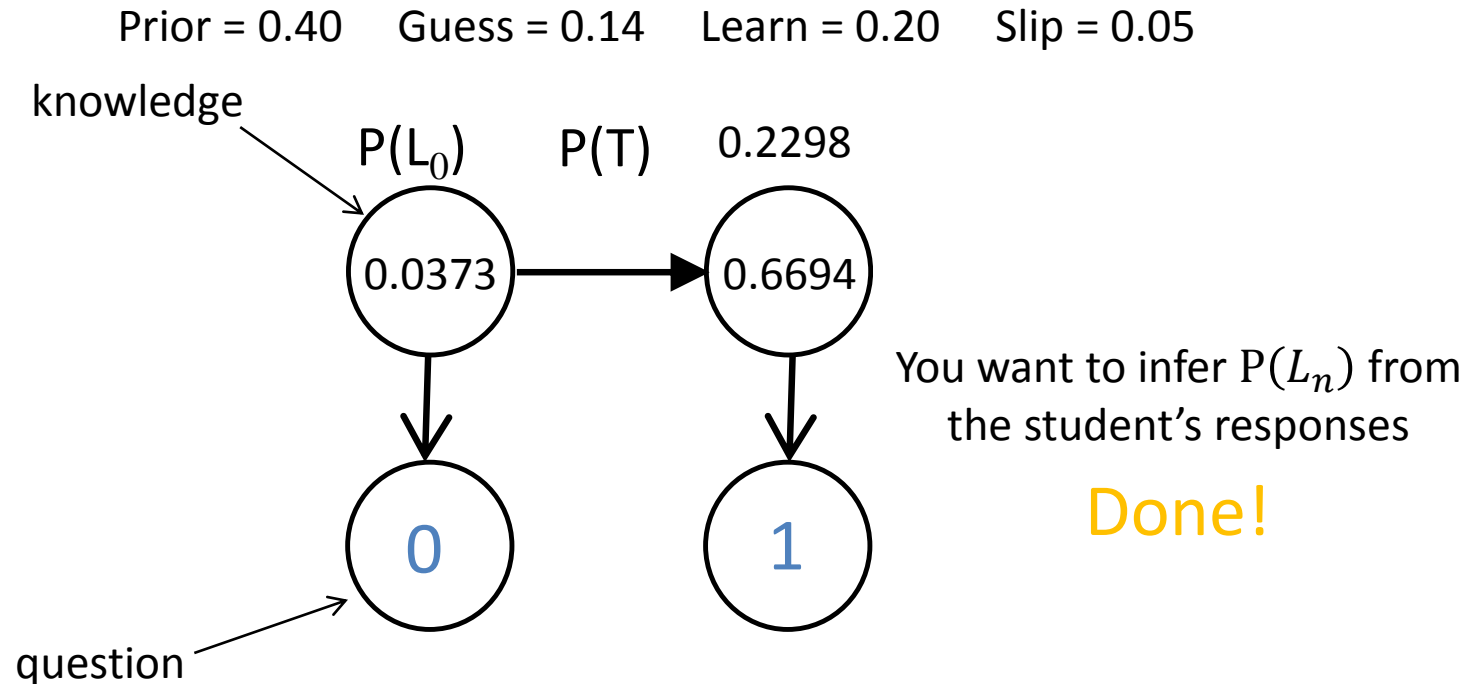
Next, apply the learning transition formula:

$$P(L_{n+1}) = 0.0373 + (1 - 0.0373)(0.20) = 0.2298$$  New prior for $L_{n+1}$

How a Bayesian Knowledge Tracing World Works

Prior = 0.40     Guess = 0.14     Learn = 0.20     Slip = 0.05

knowledge

$P(L_0)$     $P(T)$     0.2298

( 0.0373 ) → ( 0.6694 )

You want to infer $P(L_n)$ from the student's responses

( 0 )     ( 1 )

Done!

question
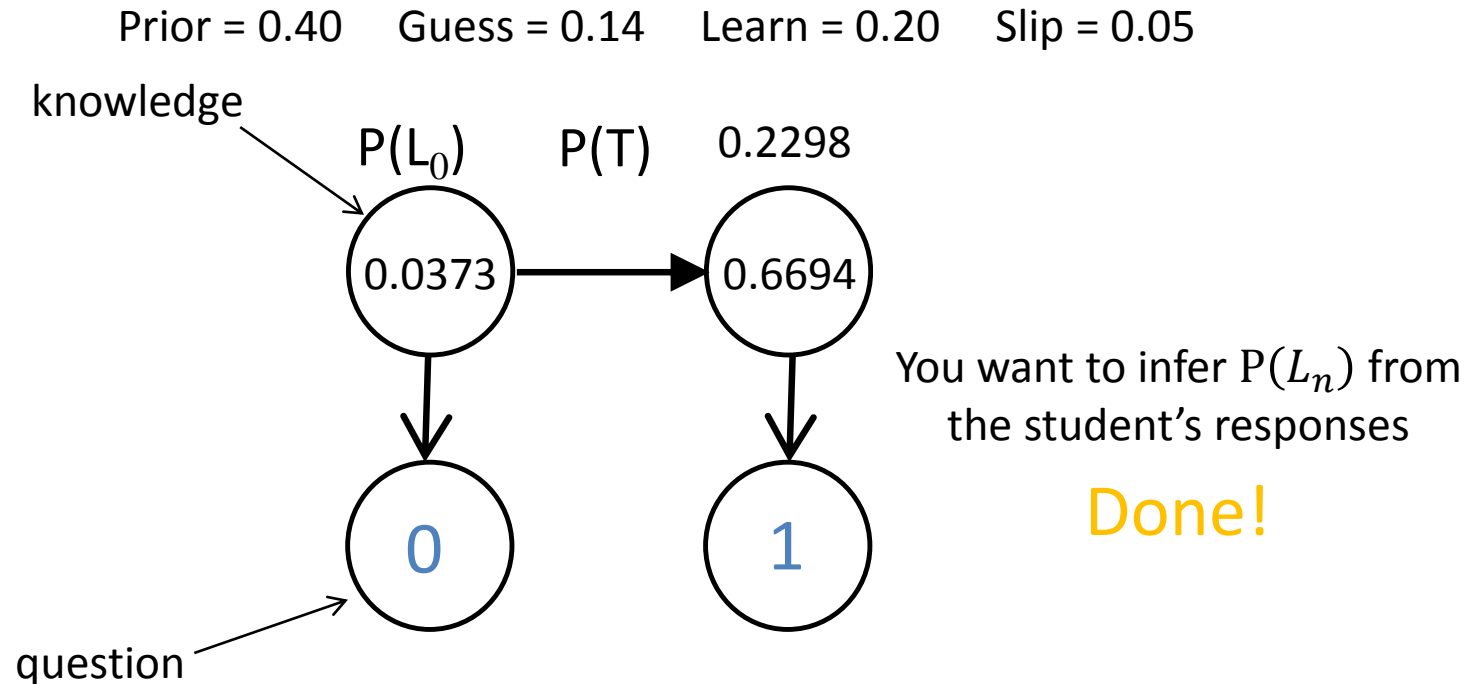
Lastly, infer the knowledge at the second opportunity:

$$P(Knowledge|Response = 1) = \frac{0.2298 \cdot (1 - 0.05)}{0.2298 \cdot (1 - 0.05) + (1 - 0.2298) \cdot 0.14} = 0.6694$$

35

How a Bayesian Knowledge Tracing World Works

Prior = 0.40    Guess = 0.14    Learn = 0.20    Slip = 0.05

knowledge

$P(L_0)$    $P(T)$    0.2298

0.0373 ———▶ 0.6694

You want to infer $P(L_n)$ from the student's responses

0    1

Done!

question
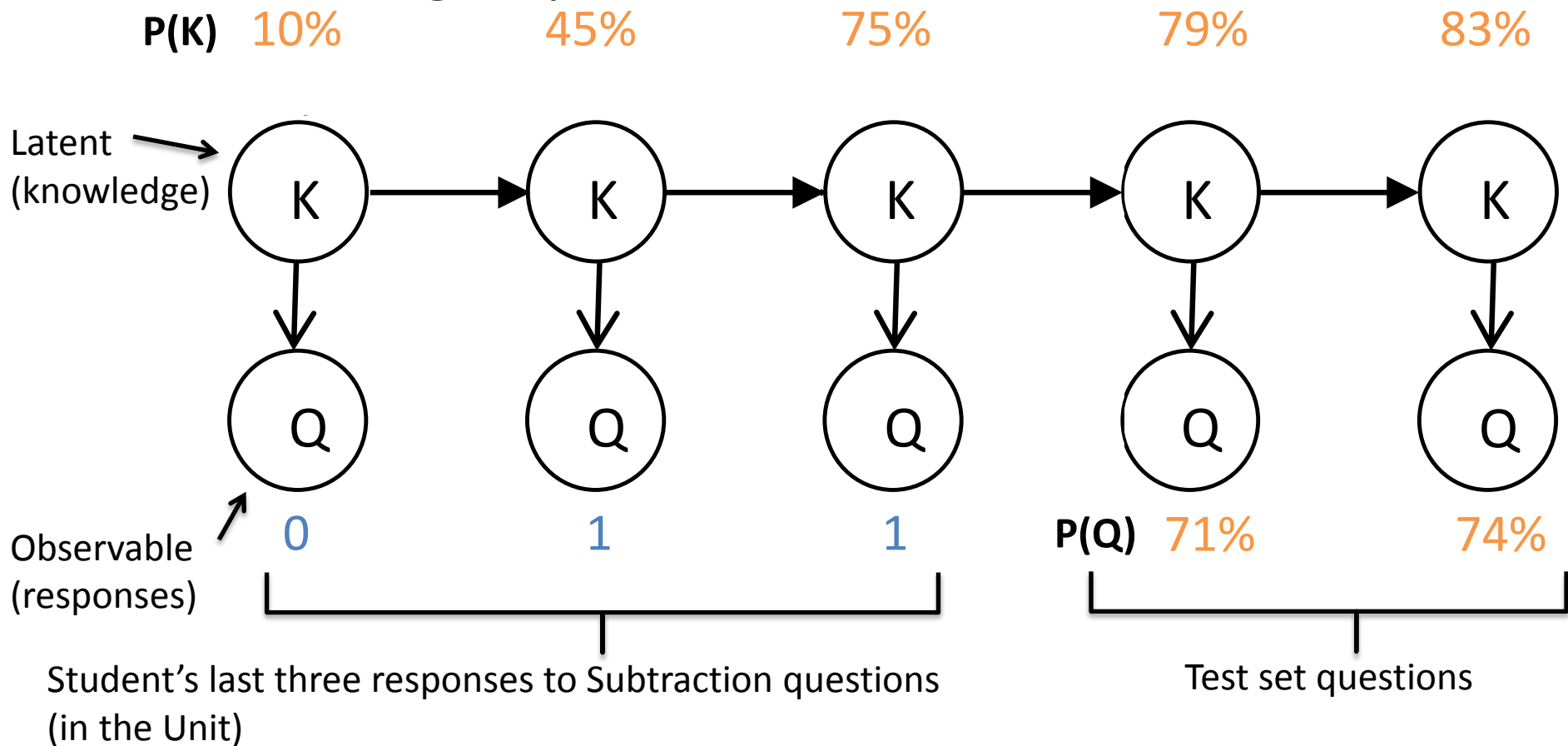
Inference calculations are applications of Bayes theorem: $P(K|Q) = \dfrac{P(Q|K)P(K)}{P(Q)}$

**Model Prediction**

Model Tracing Step – Skill: Subtraction

**P(K)** 10%        45%        75%        79%        83%

Latent
(knowledge)

K → K → K → K → K

Q   Q   Q   Q   Q

0        1        1     **P(Q)** 71%        74%

Observable
(responses)

Student's last three responses to Subtraction questions
(in the Unit)

Test set questions

# Influence of parameter values

Estimate of knowledge for student with response sequence: 0 1 1 1 1 1 1 1 1 1



Temporal Bayes Net Experimenter
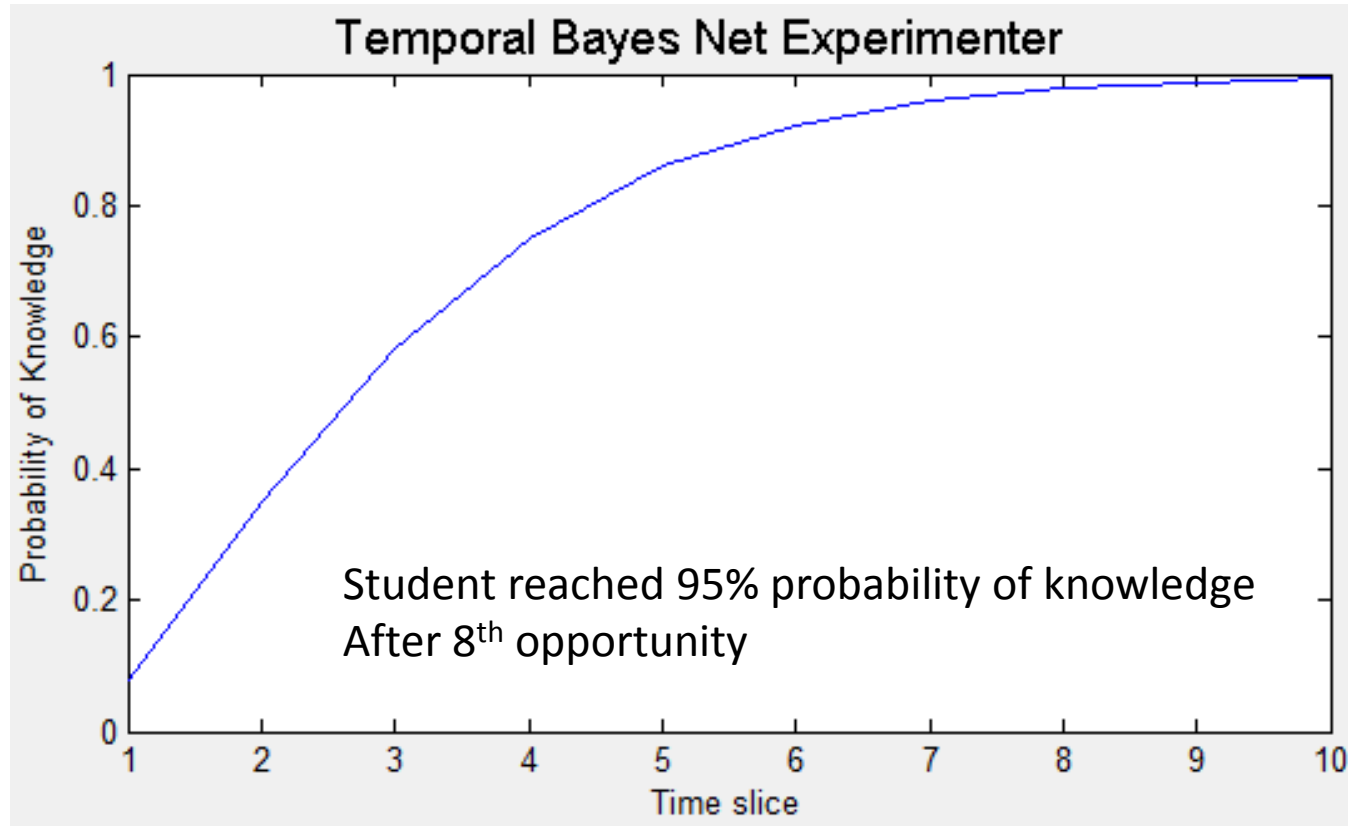
Student reached 95% probability of knowledge
After 4th opportunity

$P(L_0)$: 0.50  $P(T)$: 0.20  $P(G)$: 0.14  $P(S)$: 0.09

# Influence of parameter values

Estimate of knowledge for student with response sequence: 0 1 1 1 1 1 1 1 1 1



**Temporal Bayes Net Experimenter**

Student reached 95% probability of knowledge
After 8th opportunity

$P(L_0)$: 0.50  $P(T)$: 0.20  $P(G)$: 0.14  $P(S)$: 0.09
$P(L_0)$: 0.50  $P(T)$: 0.20  **$P(G)$: 0.64  $P(S)$: 0.03**

( Demo )

# Parameter fitting

## -EM, Grid-search, Spectral DS (Gordon)
## -1st workshop on Parameter fitting (ITS 2012)



Standard Knowledge Tracing

Prior Per Student (cold start heuristic)

Pardos, Z. A., Heffernan, N. T. In Press (2010**) Navigating the parameter space of Bayesian Knowledge Tracing models: Visualizations of the convergence of the Expectation Maximization algorithm**. In *Proceedings of the 3rd International Conference on Educational Data Mining. Pittsburg*
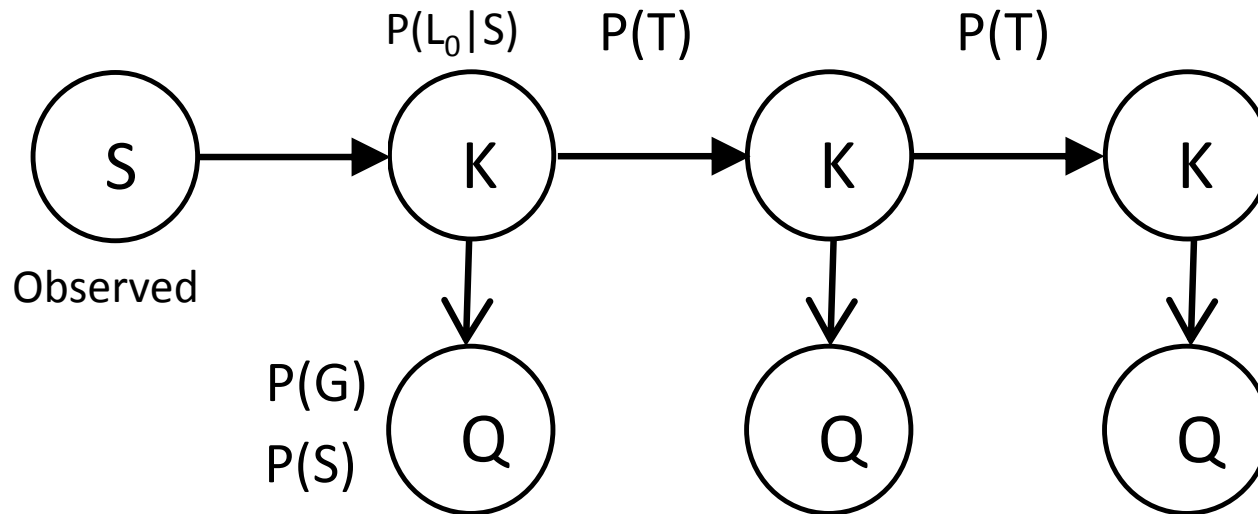
# Prior Per Student Model

- Knowledge Tracing, the current state of the art in knowledge assessment
  - Has no student specific parameters
    - Individual prior knowledge
    - Individual learn rates
  - Research objective is to add individualization to improve knowledge assessment and prediction accuracy.

# Prior Individualization Approach

Do all students enter a lesson with the same background knowledge?



**Node representations**
K = Knowledge node
Q = Question node
**S = Student node**

**Node states**
K = Two state (0 or 1)
Q = Two state (0 or 1)
**S = Multi state (1 to N)**

# Prior Individualization Approach

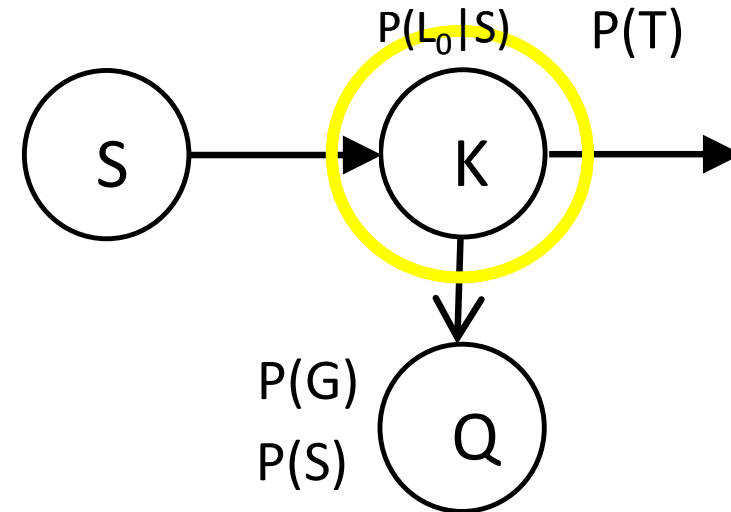Conditional Probability Table of Student node and Individualized Prior node

• Now that the model enables a prior parameter per student, how are these parameters going to be learned?

CPT of Individualized Prior node

| S value | $P(L_0|S)$ |
|---------|-----------|
| 1 | 0.05 |
| 2 | 0.30 |
| | .95 |
| | : |

### Several strategies tried

| | Most accurate predictor (of 42) | | Avg. Correlation | |
|---|---|---|---|---|
| $P(L_0)$ Strategy | PPS | KT | PPS | KT |
| Percent correct heuristic | 33 | 8 | 0.3515 | 0.1933 |
| Cold start heuristic | 30 | 12 | 0.3014 | 0.1726 |
| Random parameter values | 26 | 16 | 0.2518 | 0.1726 |

$P(L_0|S)$    $P(T)$

S → K →

$P(G)$
$P(S)$
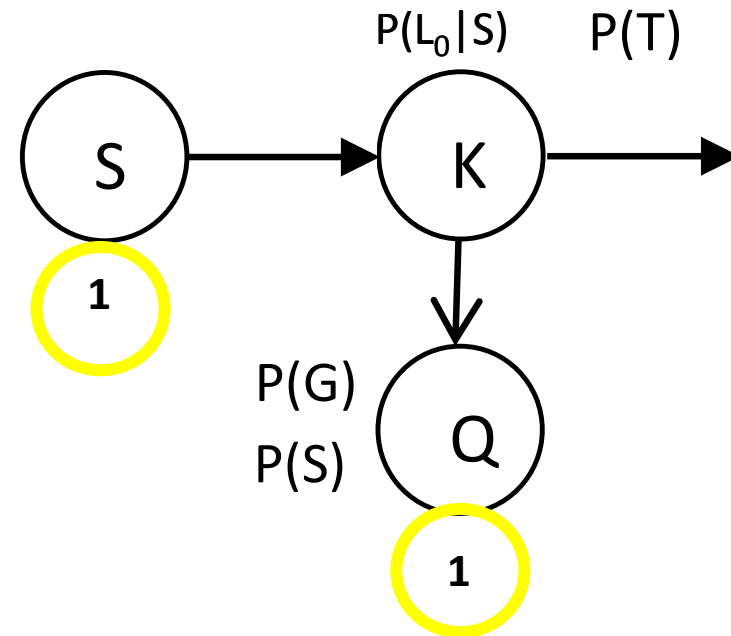
K → Q

(Pardos & Heffernan, 2010a)

# Prior Individualization Approach

### Conditional Probability Table of Student node and Individualized Prior node

• Cold Start Heuristic

CPT of Individualized Prior node

| S value | $P(L_0|S)$ |
|---------|-----------|
| 0 | 0.05 |
| 1 | 0.30 |

$P(L_0|S)$    $P(T)$

**S** → **K** →

**1**

$P(G)$
$P(S)$    **Q**

**1**

# Prior Individualization Approach

What values to use for the two priors?

CPT of Individualized Prior node

What values to use for the two priors?

| S value | $P(L_0|S)$ |
|---------|------------|
| 0 | 0.05 |
| 1 | 0.30 |

$P(L_0|S)$     $P(T)$

S → K →

**1**

P(G)
P(S)     Q

**1**

# Prior Individualization Approach

What values to use for the two priors?

CPT of Individualized Prior node

1. **Use ad-hoc values**

| S value | $P(L_0|S)$ |
|---------|------------|
| 0 | 0.10 |
| 1 | 0.85 |

$P(L_0|S)$   $P(T)$

S → K →

**1**

$P(G)$
$P(S)$   Q

K → Q

**1**

# Prior Individualization Approach

What values to use for the two priors?

1. Use ad-hoc values
2. **Learn the values**

CPT of Individualized Prior node

| S value | $P(L_0|S)$ |
|---------|------------|
| 0 | EM |
| 1 | EM |

$P(L_0|S)$    $P(T)$

S → K →

**1**

P(G)
P(S)    Q

**1**
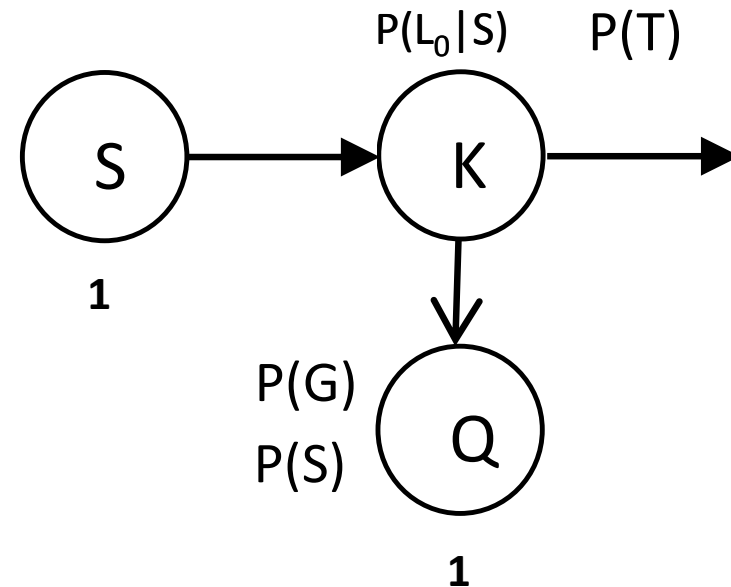
# Prior Individualization Approach

What values to use for the two priors?

CPT of Individualized Prior node

1. Use ad-hoc values
2. Learn the values
3. **Link with the guess/slip CPT**

| S value | $P(L_0\|S)$ |
|---------|-------------|
| 0 | *Slip* |
| 1 | 1-*Guess* |



$P(L_0|S)$  P(T)

S  K

**1**

P(G)
P(S)

Q

**1**

# Prior Individualization Approach

## What values to use for the two priors?

CPT of Individualized Prior node

1. Use ad-hoc values
2. Learn the values
3. Link with the guess/slip CPT

| S value | P($L_0$\|S) |
|---------|-----------|
| 0 | *Slip* |
| 1 | 1-*Guess* |

$P(L_0|S)$    $P(T)$

S → K →

**1**

P(G)
P(S)   Q

**1**

Algebra (development)

| | Strategy | RMSE |
|---|----------|------|
| 1 | adjustable | 0.3659 |
| 2 | guess/slip | 0.3660 |
| 3 | *Ad-hoc* | 0.3662 |

Bridge to Algebra (development)

| | Strategy | RMSE |
|---|----------|------|
| 1 | guess/slip | 0.3227 |
| 2 | adjustable | 0.3228 |
| 3 | *Ad-hoc* | 0.3236 |

(Pardos & Heffernan, JMLR In Press)

With an ASSISTments Platform dataset, PPS (ad-hoc) achieved an $R^2$ of 0.301 (0.176 with KT)

(Pardos & Heffernan, UMAP 2010)

# Prior Per Student Model

<u>Cold start heuristic was a success</u>
- Performed well, improvement in prediction over KT in 30/42 problem sets

- Requires no extra information outside of the responses in the problem set being predicted

- Reduces the free parameters to three instead of four
    - Faster parameter training time with more accurate prediction

- The most simple individualization technique to add to existing KT models
    - One binary node addition and one arc

- Parameters can be learned from one population of students to predict another
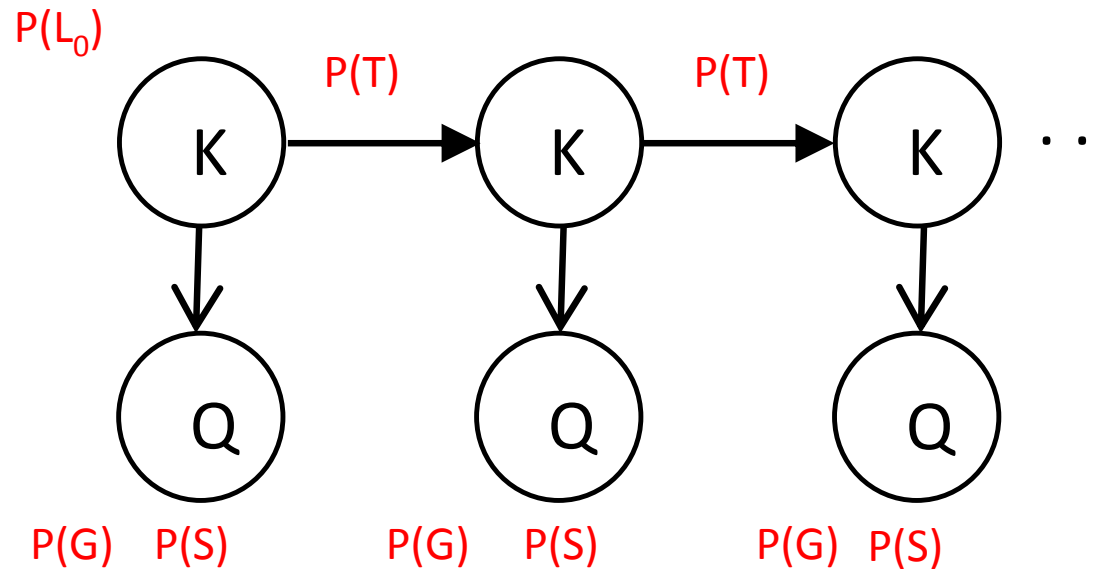
# Variations on Knowledge Tracing (and other models)

## 1. BKT-BF

Learns values for these parameters by performing a grid search (0.01 granularity) and chooses the set of parameters with the best squared error

$P(L_0)$ = Probability of initial knowledge
$P(T)$ = Probability of learning
$P(G)$ = Probability of guess
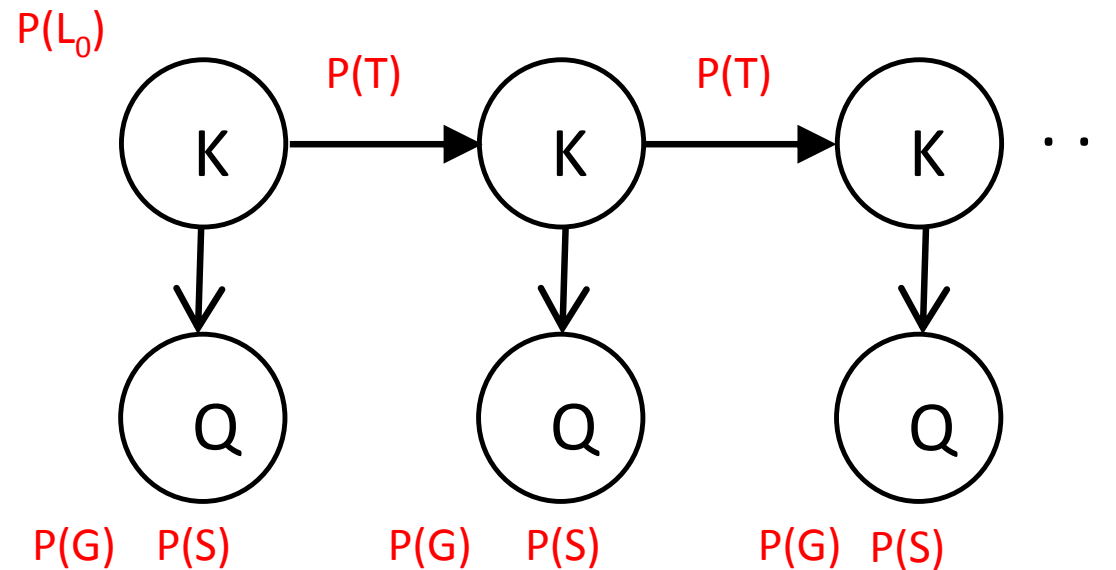$P(S)$ = Probability of slip



(Baker et al., 2010)

## 2. BKT-EM

Learns values for these parameters with Expectation Maximization (EM). Maximizes the log likelihood fit to the data

$P(L_0)$ = Probability of initial knowledge
$P(T)$ = Probability of learning
$P(G)$ = Probability of guess
$P(S)$ = Probability of slip



(Chang et al., 2006)
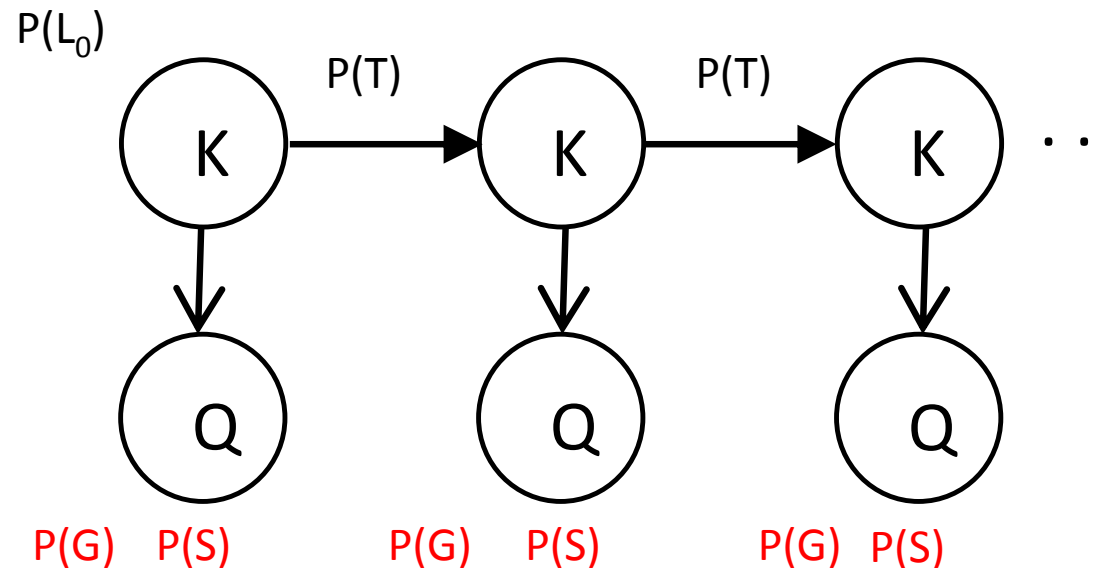
## 3. BKT-CGS

Guess and slip parameters are assessed <u>contextually</u> using a regression on features generated from student performance in the tutor

$P(L_0)$ = Probability of initial knowledge
$P(T)$ = Probability of learning
$P(G)$ = Probability of guess
$P(S)$ = Probability of slip
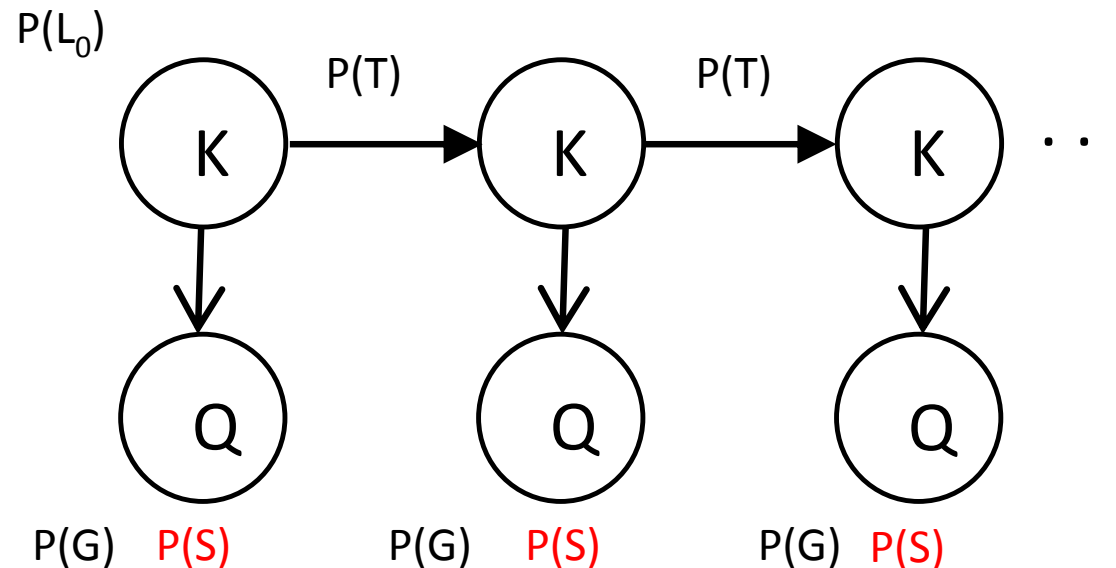


(Baker, Corbett, & Aleven, 2008)

## 4. BKT-CSlip

Uses the student's averaged <u>contextual</u> Slip parameter learned across all incorrect actions.

$P(L_0)$ = Probability of initial knowledge
$P(T)$ = Probability of learning
$P(G)$ = Probability of guess
$P(S)$ = Probability of slip


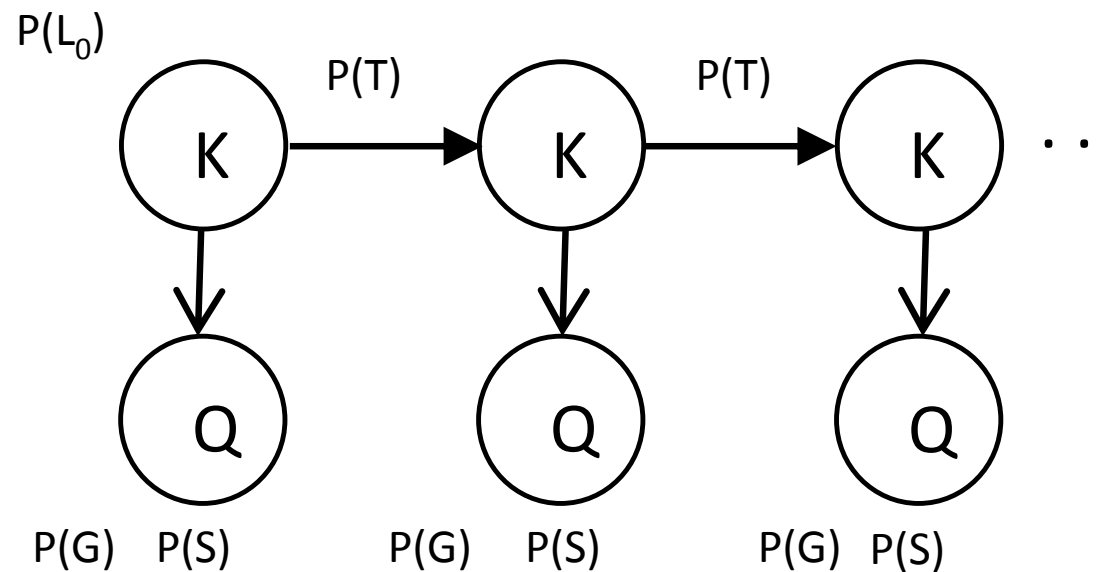
(Baker, Corbett, & Aleven, 2008)

## 5. BKT-LessData

<u>Limits</u> students response sequence length to the most recent 15 during EM training.

$P(L_0)$ = Probability of initial knowledge
$P(T)$ = Probability of learning
$P(G)$ = Probability of guess
$P(S)$ = Probability of slip



Most recent 15 responses used (max)

(Nooraiei et al, 2011)

## 6. BKT-PPS

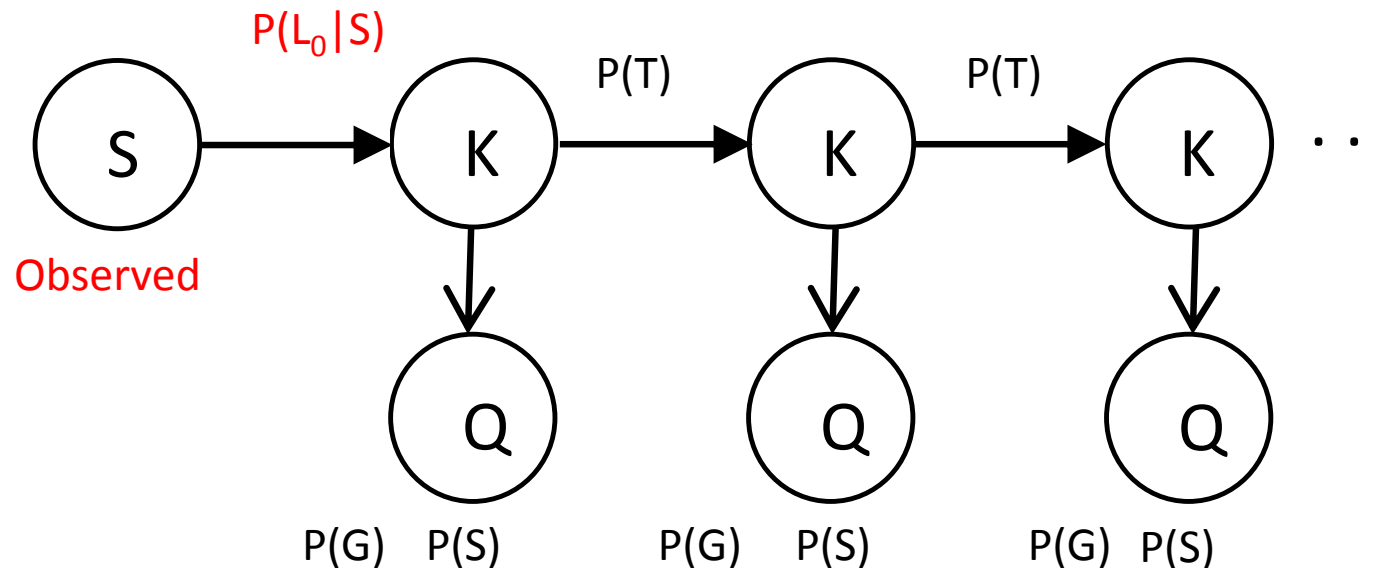Prior per student (PPS) model which <u>individualizes</u> the prior parameter. Students are assigned a prior based on their response to the first question.

$P(L_0)$ = Probability of initial knowledge
$P(T)$ = Probability of learning
$P(G)$ = Probability of guess
$P(S)$ = Probability of slip

$P(L_0|S)$

Observed

P(T)      P(T)

P(G)   P(S)      P(G)   P(S)      P(G)   P(S)

(Pardos & Heffernan, 2010)

## 7. CFAR

Correct on First Attempt Rate (CFAR) calculates the student's <span style="color:red">percent correct</span> on the <u>current skill</u> up until the question being predicted.

Student responses for Skill X: <span style="color:red">0 1 0 1 0 1</span>

Predicted next response would be 0.50

(Yu et al., 2010)

## 8. Tabling

Uses the student's response sequence (max length 3) to predict the next response by looking up the <span style="color:red">average next response</span> among student with the <u>same sequence</u> in the training set

Training set

Student A: 0 1 1 <span style="color:red">0</span>
Student B: 0 1 1 <span style="color:red">1</span>
Student C: 0 1 1 <span style="color:red">1</span>

Max table length set to 3:
Table size was $2^0+2^1+2^2+2^3=15$

Test set student: 0 0 1 _

Predicted next response would be 0.66

(Wang et al., 2011)

## 9. IRT

Item response theory (IRT) the standard assessment tool used for GRE testing.

$$p(+|v,i) = \frac{\exp(\theta_v - \sigma_i)}{1 + \exp(\theta_v - \sigma_i)} = \frac{1}{1 + e^{-(\theta_v - \sigma_i)}}$$

Where $p(+|v,i)$ is the probability of positive performance of student v on test item i and $\sigma\_i$ is the difficulty of item i.

Extension which breaks down an item into cognitive operations (Scheiblechner, 1972)

$$\sigma_i = \sum_{j=1}^{m} q_{ij} n_j + c$$

Where $j$ is a cognitive operation, $q_{ij}$ is the number of times the operation occurs in item $i$, $n_j$ is the difficulty of cognitive operation $j$ and $c$ is a scaling constant

Learning from the test addition:
$$\sigma = \sum_{j=1}^{m} q_{ij} n_j - q_{uj} h^*_{ij} \beta_j) + c$$

62

## 10. PFA

Performance Factors Analysis (PFA). <u>Logistic regression</u> model which elaborates on the Rasch IRT model. Predicts performance based on <span style="color:red">the count of student's prior failures and successes</span> on the current skill.

An overall difficulty parameter $^\beta$ is also fit for each skill or each item in this formula the variant of PFA that fits $^\beta$ for each skill is shown. The PFA equation is:

$$m(i, j \in KCs, s, f) = \beta_j + \sum(\gamma_j S_{ij} + \rho_j F_{ij})$$

(Pavlik et al., 2009)

# Conclusion

Time Left? If yes  then KT-IDEM / IEM
Else
  next slide

Bayesian Knowledge Tracing
MATLAB Demo / code

Questions?

Thesis

Pardos, Z.A., Heffernan, N.T. (2012) Tutor Modeling vs. Student Modeling. In *Proceedings of the 25th annual Florida Artificial Intelligence Research Society Conference *Invited paper*

zpardos@gmail.com
http://wpi.edu/~zpardos