

## Introduction

How do people achieve long-term goals in an uncertain environment, via repeated trials and noisy observations?

### Multi-armed bandit problem

- limited number of trials
- reward generated by a stationary, hidden Bernoulli
- the goal is to maximize total number of rewards

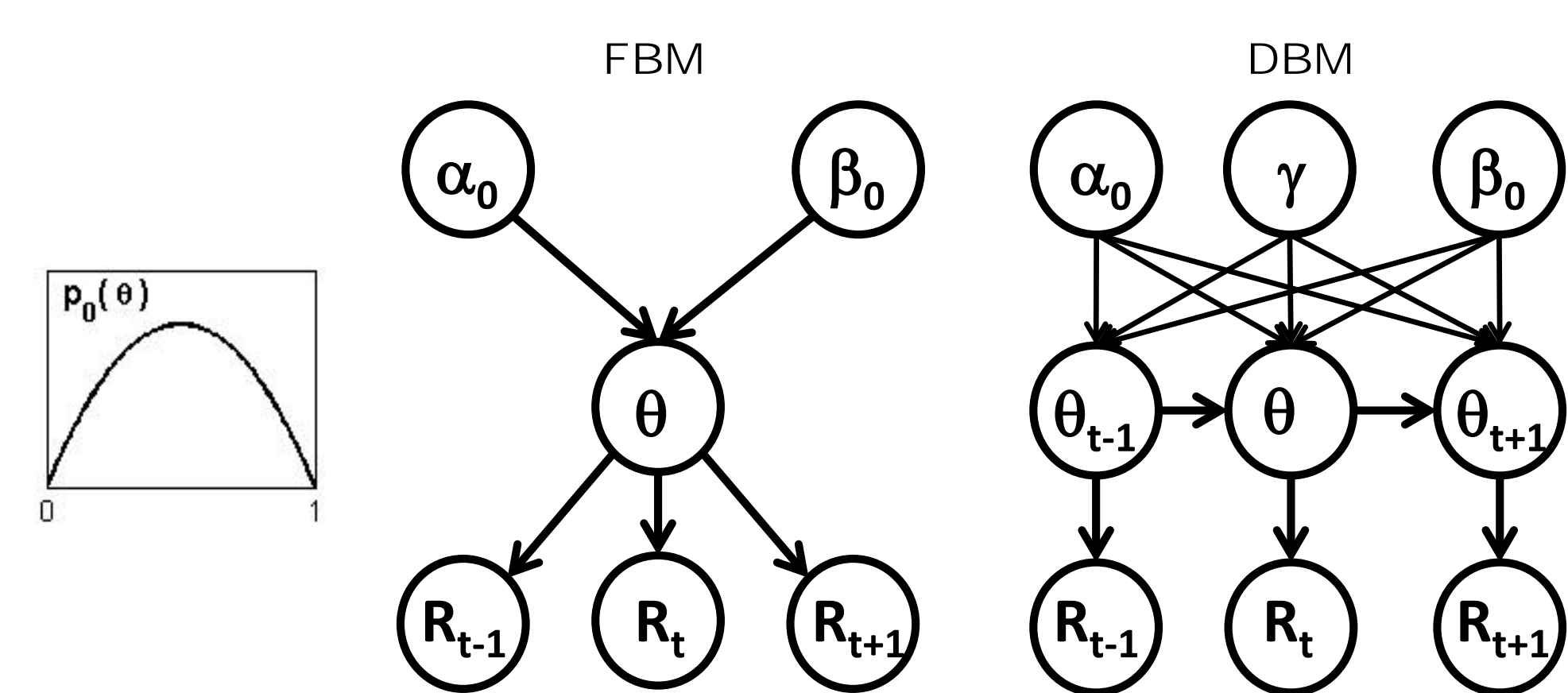
Balance between

- **exploration** selecting an arm to gain more information
- **exploitation** selecting an arm known to have relatively high expected reward

Our model

- learning component: **Dynamic Belief Model** [1]
  - Bayesian iterative inference assuming local patterns in a sequence of observations
- control component: **Knowledge Gradient** algorithm [2]
  - myopic approximation to the optimal policy
  - assuming next step being the last exploratory choice
  - chooses option that maximizes the future cumulative reward

## Temporal Dynamics



### Fixed Belief Model (FBM)

- assuming all reward rates fixed throughout the game
- $\Pr(\theta_k^t = \theta | S_k^{t-1}, F_k^{t-1}) = q_k^{t-1}$   
i.e. Beta( $\alpha_0 + S_k^{t-1}, \beta_0 + F_k^{t-1}$ )

### Dynamic Belief Model (DBM)

- assuming reward rates can undergo discrete changes according to a hidden Markov model
- updated belief is a weighted sum of last trial's belief and the generic prior

$$\Pr(\theta_k^t = \theta | S_k^{t-1}, F_k^{t-1}) = \gamma q_k^{t-1}(\theta) + (1 - \gamma)p_0(\theta)$$

where  $p_0 := \text{Beta}(\alpha_0, \beta_0)$

## Decision Policies

### The Optimal Algorithm

- computed via Bellman's dynamic programming principle

$$V^t(\mathbf{q}^t) = \max_k \hat{\theta}_k^t + \mathbb{E}[V^{t+1}(\mathbf{q}^{t+1})]$$

$$D^t(\mathbf{q}^t) = \text{argmax}_k \hat{\theta}_k^t + \mathbb{E}[V^{t+1}(\mathbf{q}^{t+1})]$$

- optimal policy computed off-line and then applied to data

### Knowledge Gradient (KG)

- expected increase in the best value among all arms if k were to be chosen, assuming next step is the last exploration
  - expectation taking over all possible outcomes incurred by pulling k

$$v_k^{\text{KG}, t} = \mathbb{E} \left[ \max_{k'} \hat{\theta}_{k'}^{t+1} | D^t = k, \mathbf{q}^t \right] - \max_{k'} \hat{\theta}_{k'}^t$$

- KG decision rule

$$D^{\text{KG}, t} = \text{argmax}_k \hat{\theta}_k^t + (T - t - 1) v_k^{\text{KG}, t}$$

### $\epsilon$ -Greedy

- with probability  $\epsilon$ , chooses randomly (exploration), otherwise chooses the greatest estimated value (exploitation)

$$\Pr(D^t = k | \epsilon, \hat{\theta}^t) = \begin{cases} (1 - \epsilon) / M^t & \text{if } k \in \text{argmax}_k \hat{\theta}_k^t \\ \epsilon / (K - M^t) & \text{otherwise} \end{cases}$$

$M^t$ : number of ties in greatest estimated value at  $t$

- maximizes the immediate gain with a constant rate, otherwise searching for information by random selection

### Win-stay, Lose-shift

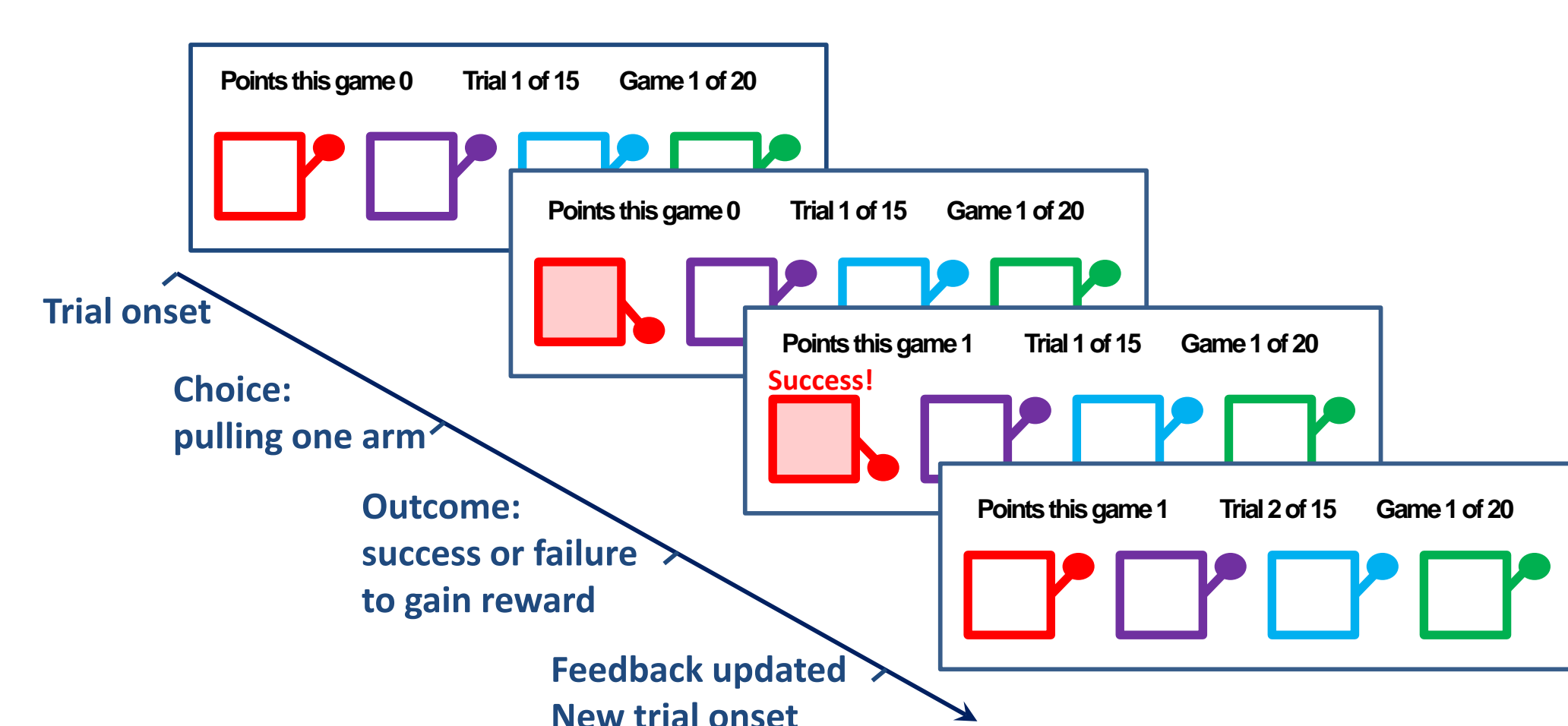
- chooses the same arm if it continues to produce a reward
- shifts to other arms (equal probabilities) following a failure

## Human Data [3]

### Behavioral experiment

- 451 subjects, 4 arms, 15 trials, 20 games
- reward rates fixed, iid from Beta(2, 2)
- participants played *same* set of games, order randomized
- participants were instructed the reward rates were drawn from the same environment, and fixed

### Schematic experimental design



## Bayesian Inference Model

For arm  $k$ , denote the outcome by  $R_k$  and reward rate by  $\theta_k$

$$R_k^t \stackrel{iid}{\sim} \text{Bernoulli}(\theta_k)$$

$$\theta_k \stackrel{iid}{\sim} \text{Beta}(\alpha_0, \beta_0)$$

The **belief state**, i.e. the posterior distribution of  $\theta_k$  given observations, is

$$q_k^t(\theta_k^t) := \Pr(\theta_k^t | S_k^t, F_k^t)$$

where  $S_k^t$  and  $F_k^t$  are the cumulative numbers of successes and failures from arm  $k$  at  $t$

The belief state can be calculated via Bayes' rule

$$q_k^t(\theta_k^t) \sim \Pr(R_k^t | \theta_k^t) \Pr(\theta_k^t | S_k^{t-1}, F_k^{t-1})$$

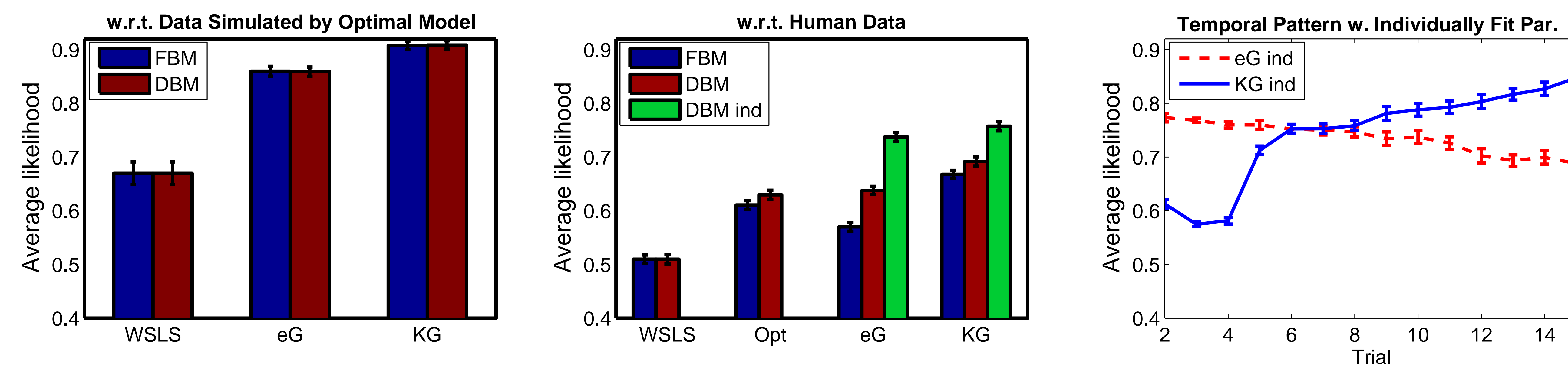
## Model Inference and Evaluation

MCMC to infer a posterior distribution of all the model parameters: Gibbs for  $\epsilon$ , Metropolis for the others

The **Average likelihood**: likelihood of consistent choices (between model and data, normalized by the set size of all consistent choices) of choices given each observed state of the game, based on MAP estimates of parameters

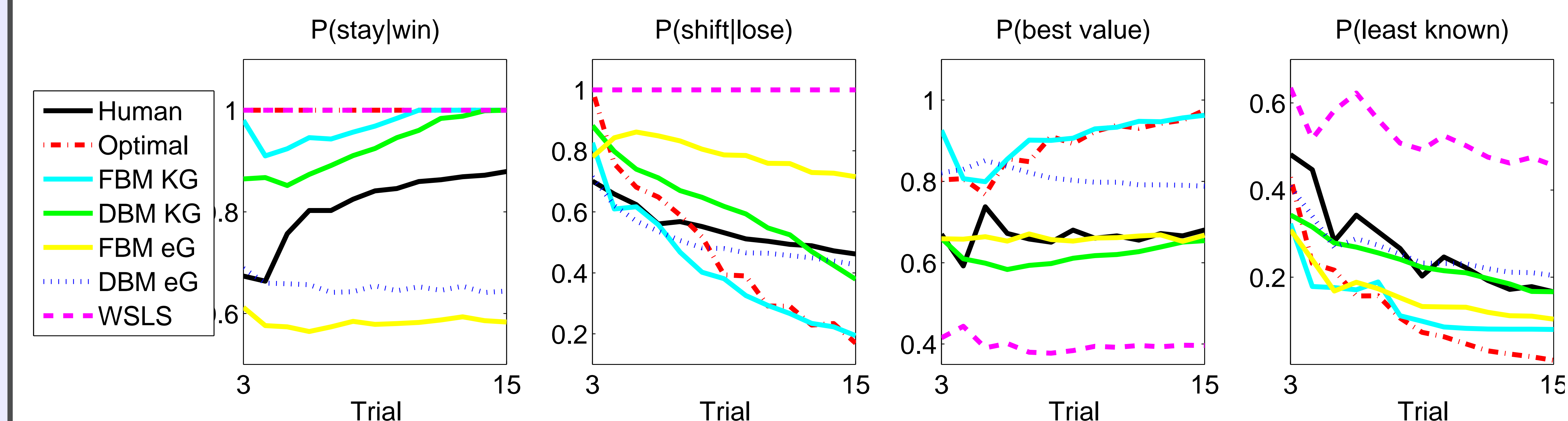
Leave-one(game)-out cross validation to calculate the average likelihood

### Model Comparison of Trialwise Choices: vs Optimal Model and Human Data



- Fitted average  $\gamma$  is .81 – people believe the world has changed in every 5 steps
- KG has the best fit to data simulated by the optimal algorithm and human data
- DBM outperforms FBM for human data for all policies
- Models may differ in their ability of capturing early/late trials

## Trialwise Behavioral Pattern from Model Simulations



- Compared to the optimal, human had more exploratory decisions, e.g. smaller probability of win-stay and choosing the best value
- Proportion of exploratory decisions decreases over trial
- DBM with KG captures the trial-wise patterns best among all models

## Summary

Human subjects tend to explore more often than policies that optimize the specific utility of the bandit problems

- DBM attributes it to the belief of a stochastically changing environment causing the sequential effects due to recent trial history
- KG decision rule favors uncertain options during early trials (to increase information gain)

## References

- [1] Yu, A. J. & Cohen, J. D. (2009). Sequential effects: Superstition or rational behavior? *NIPS*
- [2] Frazier, P.I. & Powell, W.B. (2008). The Knowledge-gradient stopping rule for ranking and selection, *Proc. of the 2008 Winter Simulation Conference*
- [3] Ryzhov, I., Powell, W & Frazier, P. (2012). The knowledge gradient algorithm for a general class of online learning problems, *Operations Research*
- [4] Steyvers, M., Lee, M. D. & Wagenmakers, E.-J. (2012). A Bayesian analysis of human decision-making on bandit problems, *Journal of Mathematical Psychology*