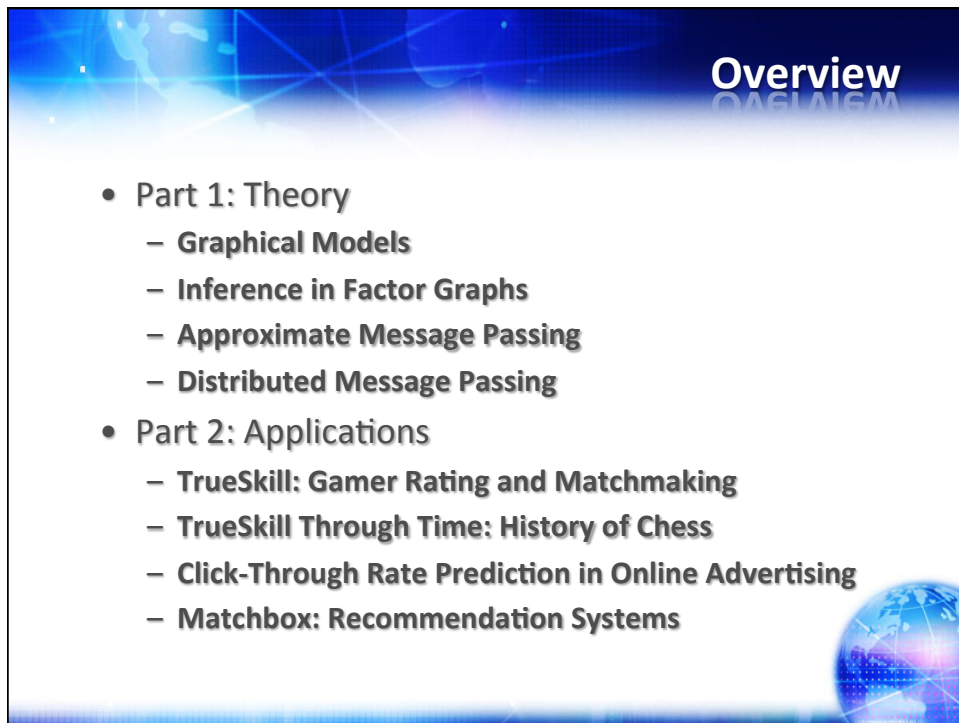# Distributed, Real-Time Bayesian Learning in Online Services

**Ralf Herbrich**
**Amazon**

---

# Overview

- Part 1: Theory
  - **Graphical Models**
  - **Inference in Factor Graphs**
  - **Approximate Message Passing**
  - **Distributed Message Passing**
- Part 2: Applications
  - **TrueSkill: Gamer Rating and Matchmaking**
  - **TrueSkill Through Time: History of Chess**
  - **Click-Through Rate Prediction in Online Advertising**
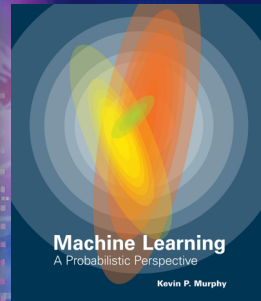  - **Matchbox: Recommendation Systems**

Part 1: Theory
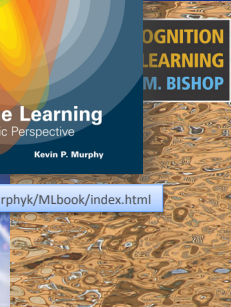
---

# Overview

- Graphical Models
- Inference in Factor Graphs
- Approximate Message Passing
- Distributed Message Passing

## Probabilities and Beliefs

- **Design**: System must assign degree of plausability **P**(A) to each logical statement A.
- **Axiom**:

**P** must be a probability measure!

3. P(A|C')>P(A|C) and P(B|AC')=P(B|AC) then
   P(AB|C') >= P(AB|C)



## Infer-Predict-Decide Cycle

**Inference**:
P(Parameters) + Data → **P**(Parameters|Data)

- Requires a (structural) model **P**(Data|Parameters)
- Allows to incorporate prior information **P**(Parameters|Data)

**Decision Making**:
Loss(Action,Data) + **P**(Data) → Action

- Business-loss not learning-loss!
- Often involves optimization!

**Prediction**:
**P**(Parameters) + Data → **P**(Data)

- Requires integration/ summation of parameter uncertainty
- Does not change state!

# Graphical Models

- **Definition:** Graphical representation of joint probability distribution
  - Nodes: ◯ = Variables
  - Edges: Relationship between variables
- **Variables**:
  - Observed Variables: Data
  - Unobserved Variables: 'Causes' + Temporary/Latent
- **Key Questions**:
  - (Conditional) *Dependency*: $p(a, b|c) \overset{?}{=} p(a|c) \cdot p(b|c)$
  - *Inference*/Marginalisation: $p(a, b) = \sum_c p(a, b, c)$
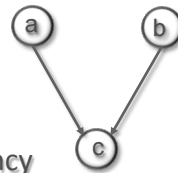
# Directed Models: Bayesian Networks

- **Definition:** Graphical representation of joint probability distribution (Pearl, 1988)
  - Nodes: ◯ = Variables
  - Directed Edges: Conditional probability distribution

- **Semantic:**

$$p(\mathbf{x}) = \prod_i p\left(x_i | \mathbf{x}_{\text{parents}(i)}\right)$$

  - Ancestral relationship of dependency

$$p(a, b, c) = p(a) \cdot p(b) \cdot p(c|a, b)$$

# Undirected Models: Markov Networks

- **Definition:** Graphical representation of joint probability distribution (Pearl, 1988)
  - Nodes: ◯ = Variables
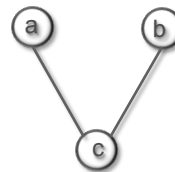  - Edges: Dependency between variables

- **Semantic:**

$$p(\mathbf{x}) = \frac{1}{Z} \cdot \prod_{\mathcal{C}} \phi(x_{\mathcal{C}}) \quad \phi \geq 0$$

  - Local potentials over cliques

$$p(a, b, c) = \frac{1}{Z} \cdot \phi_{ac}(a, c) \cdot \phi_{bc}(b, c)$$
$$Z = \sum_a \sum_b \sum_c \phi_{ac}(a, c) \cdot \phi_{bc}(b, c)$$

---

# Factor Graphs

- **Definition:** Graphical representation of product structure of a function (Wiberg, 1996)
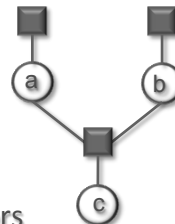  - Nodes: ■ = Factors   ◯ = Variables
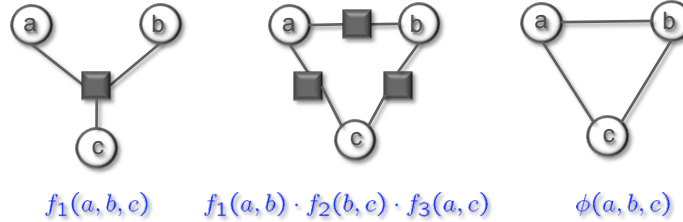  - Edges: Dependencies of factors on variables.

- **Semantic:**

$$p(\mathbf{x}) = \prod_f f\left(\mathbf{x}_{V(f)}\right)$$

  - Local variable dependency of factors

$$p(a, b, c) = f_1(a) \cdot f_2(b) \cdot f_3(a, b, c)$$

# Factor Graphs are Powerful!



$f_1(a, b, c)$     $f_1(a, b) \cdot f_2(b, c) \cdot f_3(a, c)$     $\phi(a, b, c)$

**Undirected graphical models can hide the factorisation within a clique!**

# Factor Graphs and Bayes' Law

- Bayes' law

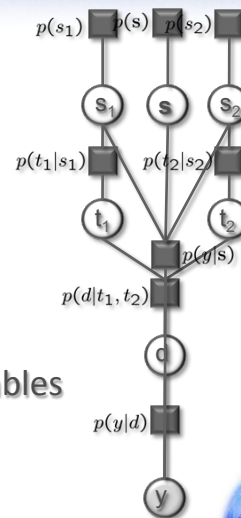$$p(\mathbf{s}|y) \propto p(y|\mathbf{s}) \cdot p(\mathbf{s})$$

- Factorising prior

$$p(\mathbf{s}) = p(s_1) \cdot p(s_2)$$

- Factorising likelihood

$$p(y, \mathbf{t}, d|\mathbf{s}) = \prod_i p(t_i|s_i) \cdot p(d|t_1, t_2) \cdot p(y|d)$$

- Inference: Sum out latent variables

$$p(y|\mathbf{s}) = \sum_{\mathbf{t}} \sum_d p(y, \mathbf{t}, d|\mathbf{s})$$

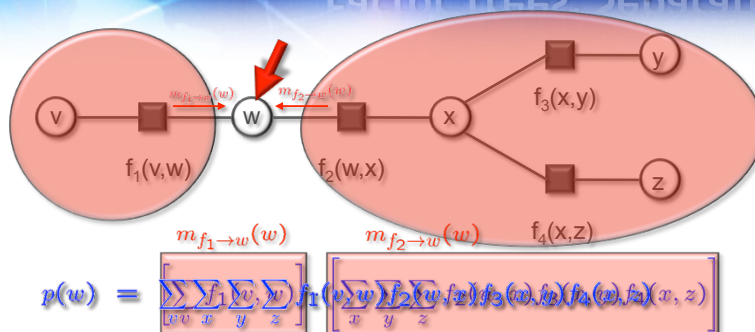|  | Dependency | Efficient Inference | Usage |
|---|---|---|---|
| Bayesian Networks | Yes | Somewhat | Ancestral Generative Process |
| Markov Networks | Yes | No | Local Couplings and Potentials |
| Factor Graphs | No | Yes | Efficient, distributed inference |

- Graphical Models
- Inference in Factor Graphs
- Approximate Message Passing
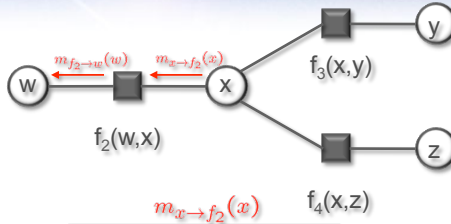- Distributed Message Passing

## Factor Graphs and Factor Trees

- **Factor Graphs:** Arbitrary functions
  - Bayesian Networks
  - Markov Networks
- **Factor Trees**: Functions where the variable indices never decrease from left to right
- **Factor Graph ➜ Factor Tree**:
  1. Pick an arbitrary node
  2. Build the spanning tree

## Factor Trees: Separation



**Observation:** Sum of products becomes product of sums of all messages from neighbouring factors to variable!
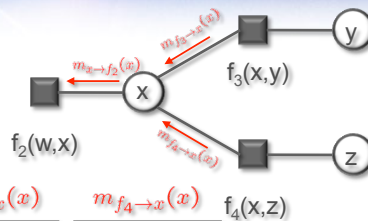
# Messages: From Factors To Variables



$$m_{f_2 \to w}(w) = \sum_x \sum_y \sum_z f_2(w,x) f_3(x,y) f_4(x,z)$$

**Observation:** Factors only need to sum out all their local variables!

# Messages: From Variables To Factors



$$m_{x \to f_2}(x) = \left[ \sum_y f_3(x,y) \right] \left[ \sum_z f_4(x,z) \right]$$

**Observation:** Variables pass on the product of all incoming messages!

# The Sum-Product Algorithm

- Three update equations (Aji & McEliece, 1997)

$$p(t) = \prod_{f \in F_t} m_{f \to t}(t)$$

$$m_{f \to t_1}(t_1) = \sum_{t_2} \sum_{t_3} \cdots \sum_{t_n} f(t_1, t_2, t_3, \ldots) \prod_{i>1} m_{t_i \to f}(t_i)$$

$$m_{t \to f}(t) = \prod_{f_j \in F_t \backslash \{f\}} m_{f_j \to t}(t)$$

- Update equations can be directly derived from the distributive law.

- Calculate all marginals at the same time!

- Only need to pass messages twice along each edge!

---

# Practical Considerations I

- **Log-Transform:** $\lambda_{f \to t}(t) := \log\left[m_{f \to t}(t)\right]$

$$\log[p(t)] = \sum_{f \in F_t} \lambda_{f \to t}(t)$$

$$\lambda_{f \to t_1}(t_1) = \sum_{t_2} \sum_{t_3} \cdots \sum_{t_n} f(t_1, t_2, t_3, \ldots) \exp\left[\sum_{i>1} \lambda_{t_i \to f}(t_i)\right]$$

$$\lambda_{t \to f}(t) = \sum_{f_j \in F_t \backslash \{f\}} \lambda_{f_j \to t}(t)$$

- **Exponential Family Messages:**

$$m(t) \propto \exp\left(\psi(t) \cdot \boldsymbol{\theta}\right)$$

- Message updates are just additions of the parameters $\boldsymbol{\theta}$ !

- (Univariate) Gaussian: $\theta := \left( \frac{\mu}{\sigma^2}, \frac{1}{\sigma^2} \right)$

- Bernoulli: $\theta := \log \left( \frac{p}{1-p} \right)$

- Binomial: $\theta := \log \left( \frac{p}{1-p} \right)$

- Beta: $\theta := (\alpha, \beta)$

- Gamma: $\theta := \left( \alpha, \frac{1}{\beta} \right)$

---

- **Redundant computations:**

$$p(t) = \prod_{f \in F_t} m_{f \to t}(t)$$

$$m_{t \to f}(t) = \prod_{f_j \in F_t \setminus \{f\}} m_{f_j \to t}(t)$$

$\Rightarrow$

$$p(t) = m_{t \to f}(t) \cdot m_{f \to t}(t)$$

- **Caching**: Only store $p(t)$ and $m_{f \to t}(t)$, then

$$m_{t \to f}(t) = \frac{p(t)}{m_{f \to t}(t)}$$

- Graphical Models
- Inference in Factor Graphs
- Approximate Message Passing
- Distributed Message Passing

- **Problem:** The exact messages from factors to variables may not be closed under products.

- **Solution:** Approximate *each* marginal as well as possible in using a divergence measure on beliefs.

- **General Idea:** Leave-one out approximation

$$\hat{p}(t) = \text{argmin}_{\hat{p}}, D\left[ m_{f \to t} \cdot \hat{m}_{t \to f}, \hat{p} \right]$$

$$\hat{m}_{f \to t}(t) = \frac{\hat{p}(t)}{\hat{m}_{t \to f}(t)}$$

$\widehat{m}_{t \to f}(t)$ $*$ $m_{f \to t}(t)$ $=$ $p(t)$

$\widehat{m}_{t \to f}(t)$ $*$ $\widehat{m}_{f \to t}(t)$ $=$ $\widehat{p}(t)$

---

# Divergence Measures

- **Kullback-Leibler Divergence:** Expected log-odd ratio between two distributions:

$$\mathsf{KL}(p, q) := \sum_t p(t) \log \left( \frac{p(t)}{q(t)} \right)$$

- **Minimizer for Exponential Families:** Matching the moments of the distribution $p(t)$!

- **General α-Divergence:**

$$D_\alpha(p, q) := \frac{1 - \sum_t \frac{p^{\alpha-1}(t)}{q^{\alpha-1}(t)}}{\alpha(1 - \alpha)}$$

- **Special Cases:**

$$D_0(p, q) = \mathsf{KL}(q, p)$$
$$D_1(p, q) = \mathsf{KL}(p, q)$$

α-Divergence in Pictures

$p(t)$
$\text{argmin}_q D_0(p, q)$
$\text{argmin}_q D_1(p, q)$

---

# Overview

- Graphical Models
- Inference in Factor Graphs
- Approximate Message Passing
- Distributed Message Passing

# Large-Data Challenge

- **Large Data (e.g. Facebook user actions)**
  - 500m daily users
  - 3 bln daily likes & comments
- **Two types of variables**
  - Observed → Data Factors
  - Latent → Model parameters
- **Discriminative Models**
  - Given the model parameters, data variables are independent

$$p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{Y}) \ \propto \ \prod_i p(y_i|\boldsymbol{\theta}, \mathbf{x}_i) \cdot \prod_j p(\theta_j)$$



# Distributed Message Passing

- **Idea**: Group variables and send messages across system boundaries

$$\prod_i p(y_i|\boldsymbol{\theta}, \mathbf{x}_i) \cdot p(\boldsymbol{\theta}) = \prod_k \underbrace{\prod_{j=1}^{n_k} p(y_{k,j}|\boldsymbol{\theta}, \mathbf{x}_{k,j})}_{f_k(\mathbf{X}_k, \mathbf{Y}_k, \boldsymbol{\theta})} \cdot \prod_l \underbrace{\prod_{r=1}^{m_l} p(\theta_{l,r})}_{g_l(\boldsymbol{\theta}_l)}$$

- **Data factors**: $f_k(\mathbf{X}_k, \mathbf{Y}_k, \boldsymbol{\theta})$
  - Know exactly which model parameter messages get updated
- **Parameter factors**: $g_l(\boldsymbol{\theta}_l)$
  - Need to keep track of which data factors need message update

# Distributed Conditional Models



**Belief Store ("Memory")**

**Message Passing ("Communicate")**

**Data Messages ("Compute")**

# A Systems Service View



Compute    Communicate    Store

Train Request

Train Request

Predict Request

$$p(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}) \;\propto\; \prod_k f_k(\mathbf{Y}_k|\boldsymbol{\theta}, \mathbf{X}_k) \cdot p(\boldsymbol{\theta})$$

- **Map-Reduce**
  - **Map**: Data nodes compute messages $m_{F_k \to \mu}$ from data $y_i$ and $m_{\mu \to F_k}$
  - **Reduce**: Combine messages $m_{F_k \to \mu}$ into $p_\mu$ by multiplication
  - Vanilla MR is a single pass only!
- **Caveats**:
  - Approximate data factors need all incoming message $m_{F_k \to \mu}$!
  - Each machine needs to be able to store the belief over $\mu$



---

$$p(y_i|\boldsymbol{\theta}, \mathbf{x}_i) \;=\; \Phi(y_i \boldsymbol{\theta}^{\mathrm{T}} \mathbf{x}_i)$$

$$p(\boldsymbol{\theta}) \;=\; \prod_j \mathcal{N}(\theta_j; \mu_j, \sigma_j^2)$$

**Sequential**

**Parallel**

# Approximation Quality

$$\mathbf{x} \quad = \quad [1; 1; \ldots; 1]^{\mathrm{T}}$$

**Single Bias Feature**                    **100 Bias Features**



# Solution : Dampening!

$$\lambda_{f \to \theta} \Rightarrow \alpha \cdot \lambda_{f \to \theta}$$

**First Step**                    **Second Step**

# Break!

Part 2: Applications

- TrueSkill: Gamer Rating and Matchmaking
- Click-Through Rate Prediction in Online Advertising
- Matchbox: Recommendation Systems

# TrueSkill™

Joint work with Thore Graepel, Tom Minka & Phillip Trelford

# Motivation

- Competition is central to our lives
  - Innate biological trait
  - Driving principle of many sports
- Chess Rating for fair competition
  - ELO: Developed in 1960 by Árpád Imre Élő
  - Matchmaking system for tournaments
- Challenges of online gaming
  - Learn from few match outcomes efficiently
  - Support multiple teams and multiple players per team

---

# The Skill Rating Problem

- **Given**:
  - Match outcomes: Orderings among k teams

- **Qu**

# Two Player Match Outcome Model

- Latent Gaussian performance model for fixed skills
- Possible outcomes: Player 1 wins over 2 (and vice versa)



$$\mathbf{P}(y_{12} = (1,2)|p_1, p_2) = \mathbb{I}(p_1 > p_2)$$

# Two Team Match Outcome Model

- Skill of a team is the sum of the skills of its members



$$\mathbf{P}(t_1|s_1, s_2) = \mathcal{N}\left(t_1; s_1 + s_2, 2 \cdot \beta^2\right)$$

## Multiple Team Match Outcome Model

- Possible outcomes: Permutations of the teams

$$\mathbf{P}(\boldsymbol{y}|t_1, t_2, t_3) = \mathbb{I}(\boldsymbol{y} = (i, j, k)) \text{ where } t_i > t_j > t_k$$



## Multiple Team Match Outcome Model

- But we are interested in the (Gaussian) posterior!

$$\mathbf{P}(s_i|\boldsymbol{y} = (1, 2, 3)) = \mathcal{N}(s_i; \mu_i, \sigma_i^2)$$

$$y_{12} = (1, 2) \qquad y_{23} = (2, 3)$$

Efficient Approximate Inference

Gaussian Prior Factors

$s_1$ $s_2$ $s_3$ $s_4$

Fast and efficient approximate message passing using Expectation Propagation

Ranking Likelihood Factors



Applications to Online Gaming

- **Leaderboard**
  - Global ranking of all players

$$\mu_i - 3 \cdot \sigma_i$$

- **Matchmaking**
  - For gamers: Most uncertain outcome

$$\mathbf{P}(p_i \approx p_j | \mu_i = \mu_j, \sigma_i^2 + \sigma_j^2)$$

$$\mathbf{P}(p_i \approx p_j | \mu_i - \mu_j = 0, \sigma_i^2 + \sigma_j^2 = 0)$$

## Experimental Setup

- **Data Set: Halo 2 Beta**
  - 3 game modes
    - Free-for-All
    - Two Teams
    - 1 vs. 1
  - > 60,000 match outcomes
  - ≈ 6,000 players
  - 6 weeks of game play
  - Publically available

## Convergence Speed

## Convergence Speed (ctd.)

Winning probability (y-axis): 100%, 80%, 60%, 40%, 20%, 0%

- char wins
- SQLWildman wins
- Both players draw

5/8 games won by char

Number of games played (x-axis): 0, 100, 200, 300, 400, 500

---

## Xbox 360 & Halo 3

- **Xbox 360 Live**
  - Launched in September 2005
  - Every game uses TrueSkill™ to match players
  - > 10 million players
  - > 2 million matches per day
  - > 2 billion hours of gameplay
- **Halo 3**
  - Launched on 25[th] September 2007
  - Largest entertainment launch in history
  - > 200,000 player concurrently (peak: 1,000,000)

# Halo 3 in Action



# Halo 3 Public Beta Analysis

## Skill Distributions of Online Games

**Golf (18 holes)**: 60 levels

**Car racing (3-4 laps)**: 40 levels

**UNO (chance game)**: 10 levels

---

## TrueSkill™ Through Time: Chess

- Model time-series of skills by smoothing across time
- History of Chess
  - 3.5M game outcomes (ChessBase)
  - 20 million variables (each of 200,000 players in each year of lifetime + latent variables)
  - 40 million factors

ChessBase Analysis: 1850 - 2006



Online Advertising

Joint work with Thore Graepel, Joaquin Quiñonero Candela, Onno Zoeter, Tom Borchert , Phillip Trelford

## Why Predict Probability-of-Click?



$$b_1 \cdot p_1 \geq b_2 \cdot p_2 \geq \cdots$$

$$c_i = b_{i+1} \cdot \frac{p_{i+1}}{p_i}$$

---

## The Scale of Things

- **Several weeks of data in training**:

    7,000,000,000 impressions

- **2 weeks of CPU time during training**:

    2 wks × 7 days × 86,400 sec/day =

    1,209,600 seconds

- **Learning algorithm speed requirement**:

    – 5,787 impression updates / sec

    – 172.8 μs per impression update

# The Flow of Information

User interaction → Raw Logs → Structured Data

- Why structured data?
  - Data validation and cleaning
  - Principled feature transformations

# Uncertainty: Bayesian Probabilities

**Client IP**
- 102.34.12.201
- 15.70.165.9
- 221.98.2.187
- 92.154.3.86

**Match Type**
- Exact Match
- Broad Match

**Position**
- ML-1
- SB-1
- SB-2

+ → $\mathbf{p}$(pClick)

# Principled Exploration



average: 25% (3 clicks out of 12 impressions)
average: 30% (30 clicks out of 100 impressions)

# Training Algorithm in Action



Prediction
Training Update

# Inference: An Optimization View

$$\mu_i \leftarrow \mu_i + \frac{\sigma_i^2}{s} \cdot h \left[ \frac{\sum_{j=1}^d \mu_j}{s} \right] \quad \sigma_i^2 \leftarrow \sigma_i^2 \left( 1 - \frac{\sigma_i^2}{s^2} \cdot g \left[ \frac{\sum_{j=1}^d \mu_j}{s} \right] \right)$$

$$s^2 = \beta^2 + \sum_{j=1}^d \sigma_j^2$$

$$h(t) = \frac{\mathcal{N}(t;\, 0, 1)}{\Phi(t)} \qquad g(t) = h(t) \cdot [h(t) + t]$$



# Client IP: Mean & Variance

UserAgent: Mean Posterior Effects


Accuracy

# MatchBox

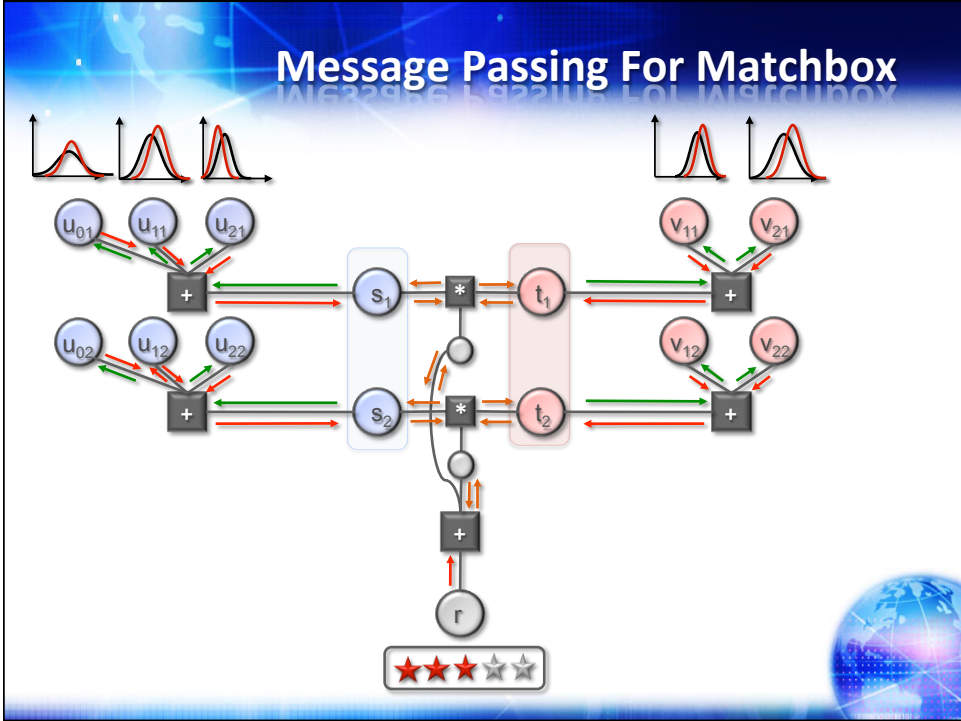Joint work with Thore Graepel, Joaquin Quiñonero Candela, David Stern
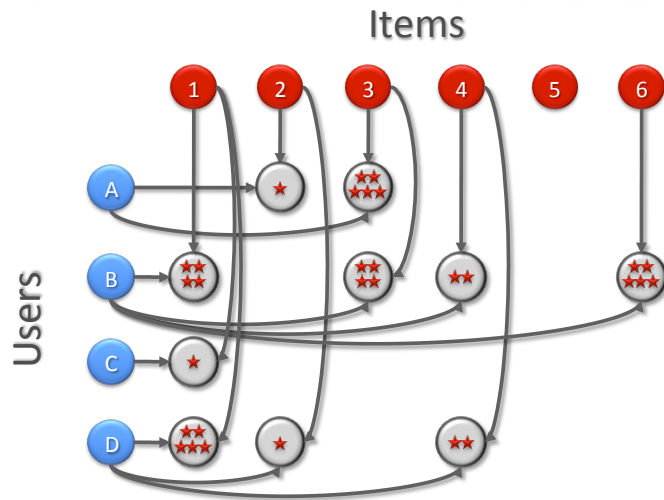
# Matchbox With Metadata

User $\quad \mathbf{s} = \mathbf{U}\mathbf{x}$

Item $\quad \mathbf{t} = \mathbf{V}\mathbf{y}$

Rating potential $\sim \mathcal{N}(\mathbf{s}^{\top}\mathbf{t}, \beta^2)$



# Recommender System: MatchBox

Message Passing For Matchbox



User/Item Trait Space

# Message Passing Iteration 2

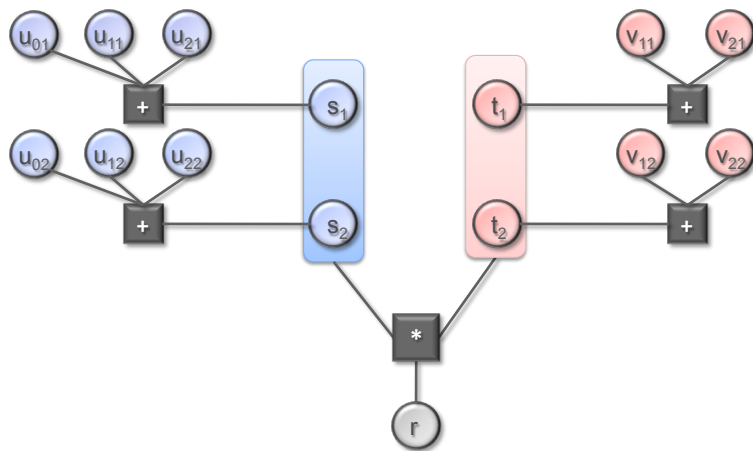# Message Passing Iteration 3

feedback models

Feedback Models



Message Passing: Compositionality

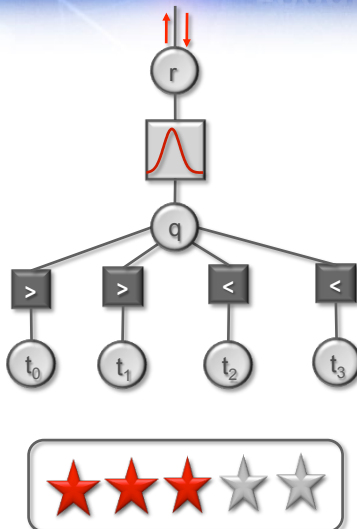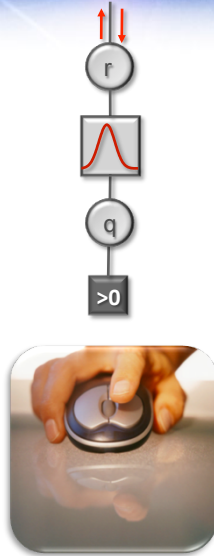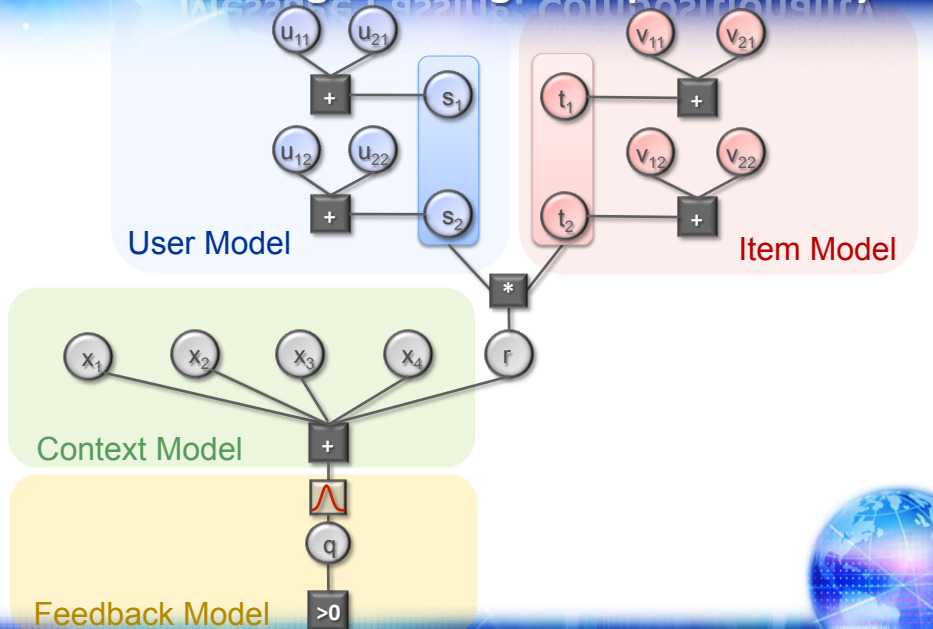# accuracy

---

## Performance and Accuracy



MovieLens Data

- 1 million ratings
- 3,900 movies / 6,040 users
- User / movie metadata

# MovieLens – 1,000,000 ratings

## 6,040 users

| User ID | |
|---|---|

| User Job | | User Age |
|---|---|---|

| User Job | | User Age |
|---|---|---|
| Other | Lawyer | <18 |
| Academic | Programmer | 18-25 |
| Artist | Retired | 25-34 |
| Admin | Sales | 35-44 |
| Student | Scientist | 45-49 |
| Customer Service | Self-Employed | 50-55 |
| | | >55 |
| Health Care | Technician | |
| Managerial | Craftsman | **User Gender** |
| Farmer | Unemployed | Male |
| Homemaker | Writer | Female |

## 3,900 movies

| Movie ID | |
|---|---|

| Movie Genre | |
|---|---|
| Action | Horror |
| Adventure | Musical |
| Animation | Mystery |
| Children's | Romance |
| Comedy | Thriller |
| Crime | Sci-Fi |
| Documentary | War |
| Drama | Western |
| Fantasy | Film Noir |

---

# MovieLens with Thresholds Model

## (ADF), Training Time= 1 Minute

# MovieLens Error with Thresholds



# Recommendation Speed
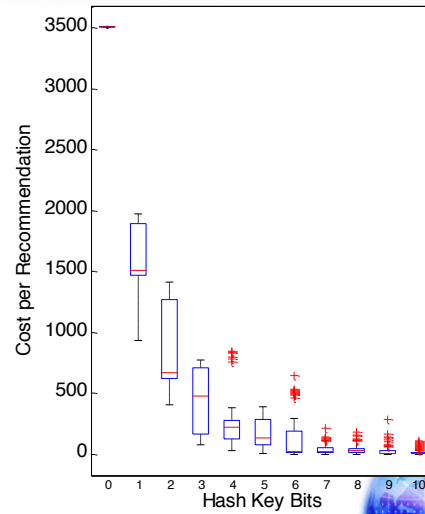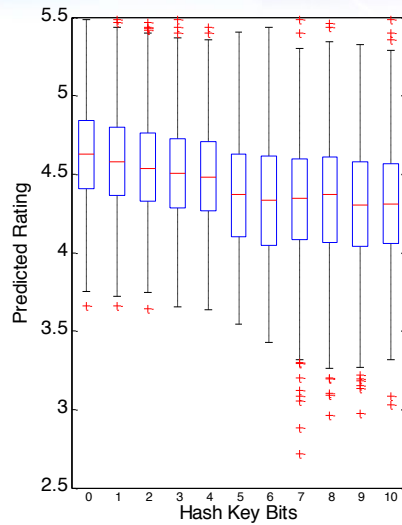
# Recommendation Speed

- **Goal:**
  find N items with highest predicted rating.
- **Challenge:**
  potentially have to consider all items.
- Two approaches to make this faster:
  - Locality Sensitive Hashing
  - KD Trees
- Locality Sensitive Hash:

$$P(h(x) = h(y)) = \mathsf{sim}(x, y)$$

# Random Projection Hashing

- Random Projections:
  - Generate random hyper planes
    (m random vectors, $\mathbf{a}_i$).
  - Gives m bit hash, $\{x_0, x_1, \cdots, x_m\}$ , by:
    $$x_i = \mathbf{1}[\mathbf{a}_i \cdot \mathbf{t} > 0]$$
- p(all bits match) $\propto$ cosine similarity.
- Store items in buckets indexed by keys.
- Given a user trait vector:

  1. Generate key, q.

  2. Search buckets by hamming distance
     from q until find N items.

## Accuracy and Speedup