# *Large Data Problems at the Long Tail:*
## The eBay Story
### Challenges and Opportunities

# Neel Sundaresan

**Research Labs**

**http://labs.ebay.com**

Neel Sundaresan

# What wags the Tail?

**eBay users trade about $1,400 worth of goods on the site every second.**

**On an average day on eBay…**

A vehicle sells every 2 minutes

A part or accessory sells every 3 seconds

Diamond jewelry sells every 83 seconds

A Timberland shoe sells every 10 minutes

A trading card sells every 6 seconds

Neel Sundaresan

# Then There was One…

$14.83

When asked if he understood that the laser pointer was broken, the buyer said "Of course, I'm a collector of broken laser pointers"

Neel Sundaresan

ebaY Research Labs

# Divine Reward!

**$28K**

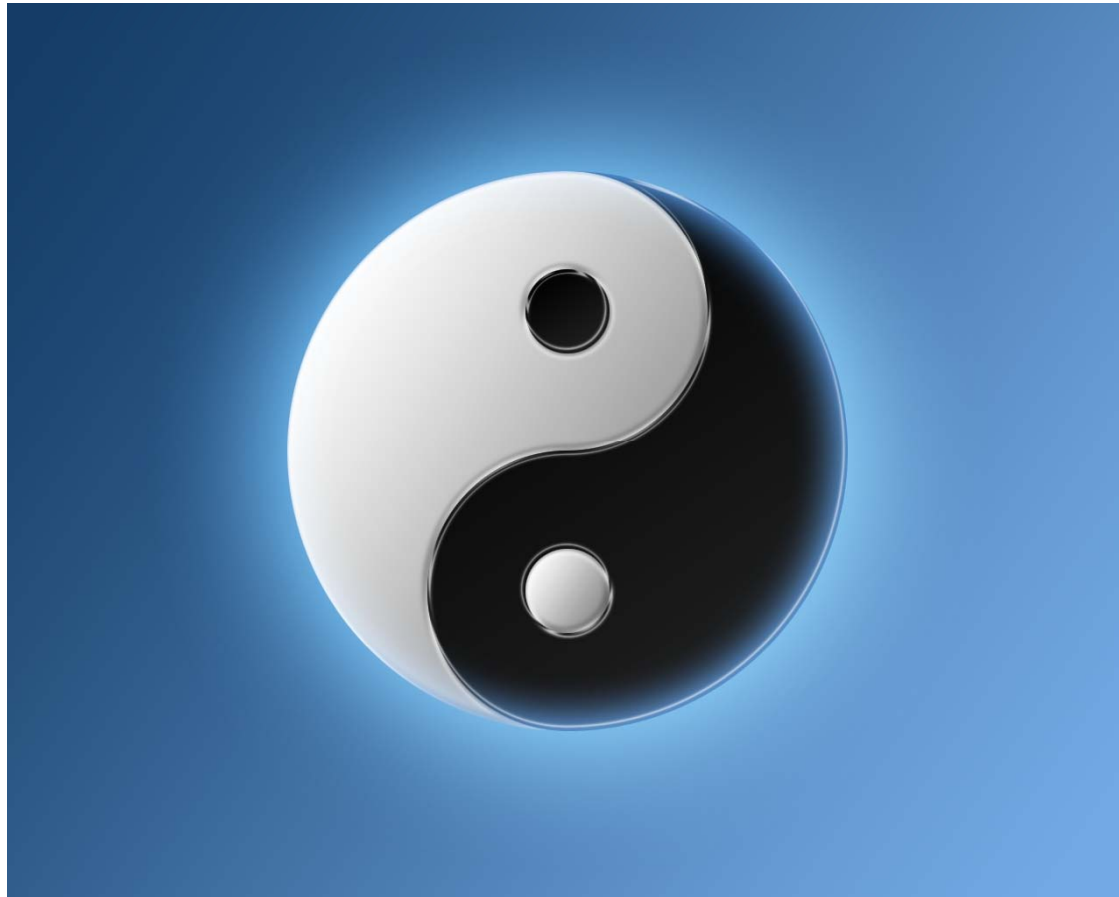*PetroliumJeliffe*

Neel Sundaresan

# The Long Tail Nature

- Buyers outnumber sellers 5:1
- Seller-items sold has a power law distribution
- Seller-revenue has a power law distribution
- Buyer-items bought has a power law distribution
  For a sample period over 1 month mean is 3 and median is 1
- Buyer-spent money has a power law distribution
  For a sample period of 1 month mean is 98$ and median was 45$
- Categories browsed
  mean 10.5, median 5
- ~10M new items a day, most items eventually sell, items last from a day to 30 days, most items not cataloged, some auction-some fixed price

Neel Sundaresan

# Opportunities at the Tail

- Vast majority of products appeal to small number of users
- Vast majority of products of this nature can only be carried by small number of sellers
- These account for sizable consumption
- "Selling less of more" becomes important

Neel Sundaresan

# Algorithm vs Data



Neel Sundaresan

# Search Challenge

- ## Near similar titles

  "Apple IPOD Nano 4GB Black NEW! Great Deal!"

  "Apple IPOD Nano 4GB Black NEW! Skin Great Deal!"

  What does someone querying for "ipod nano" look for?

Neel Sundaresan

# What's in a word

- **Spelling corrections for Swarowski?**
  64 of them!

swarvoski,swaroski,swavorski,swarovsky,swarski,**swarowski**,swarvorski,swarofski,
swarvski,sworovski,swarovksi,swarosky,svarovski,swarowsky,swarkovski,swarovki,
svarowski,swaraski,swaroviski,swarovoski,swaravski,swarorski,swartski,swarovsk,
swaovski,swarvosky,sworoski,swalovski,sworvski,swavroski,swrovski,
swaorvski,swavoski,sawrovski,swarovsi,swarovaski,svaroski,swarozski,swarouski,
swarokski,swarvaski,sworvoski,swarvowski,swarosvki,svarovsky,swaravoski,swarovky,
sawarski,sarovski,swarovzki,swarocski,swarovskl,swarovsku,swarovkski,swarrovski,
swarovske,swarowvski,swarvovski,saworvski,swarosvski,swarovrski,swarivski,swarsovski
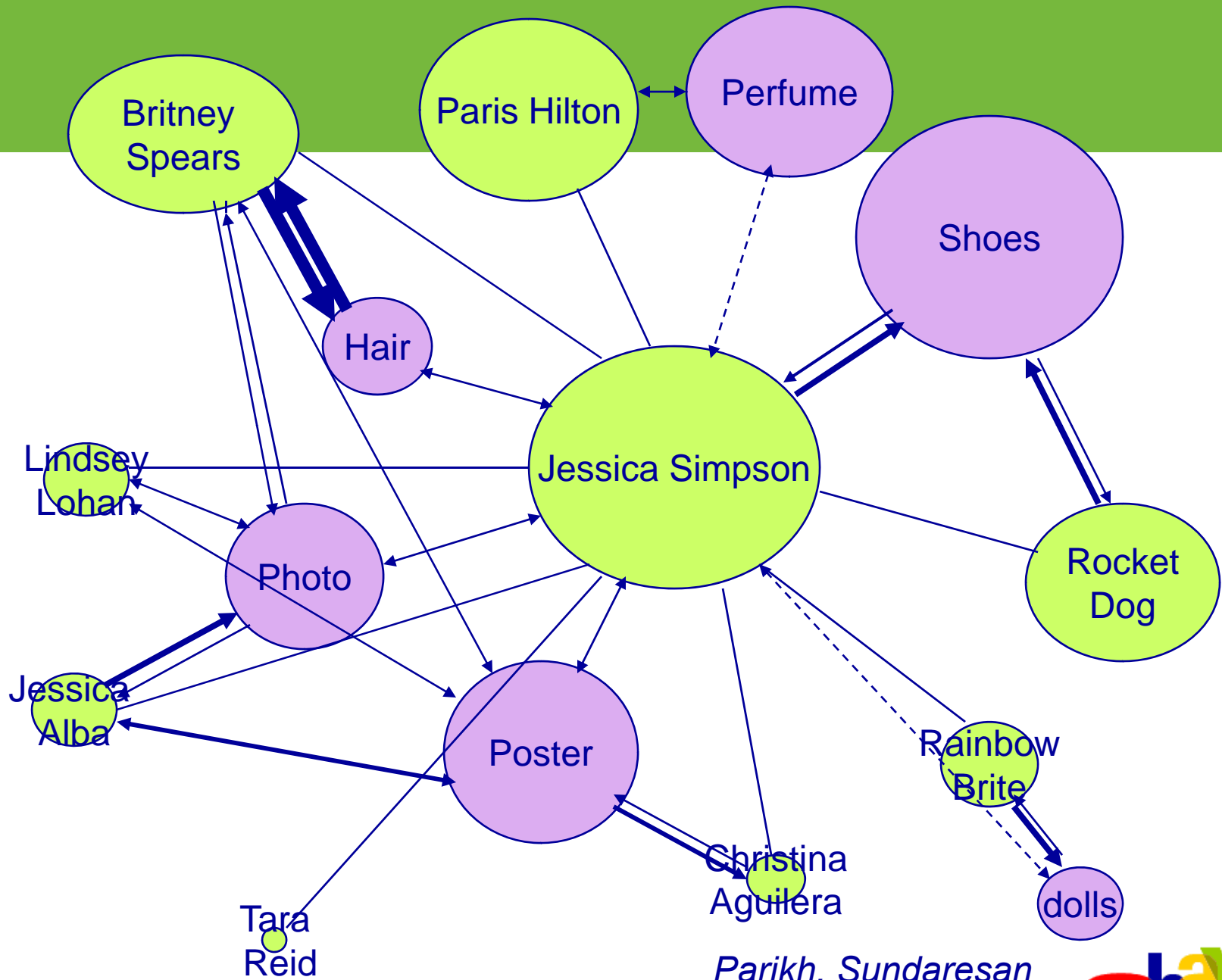
Courtesy K. Mauge'

Neel Sundaresan

# What's That in English?

- Avon SSS 5 OZ Shave Gel NIB FS GWP

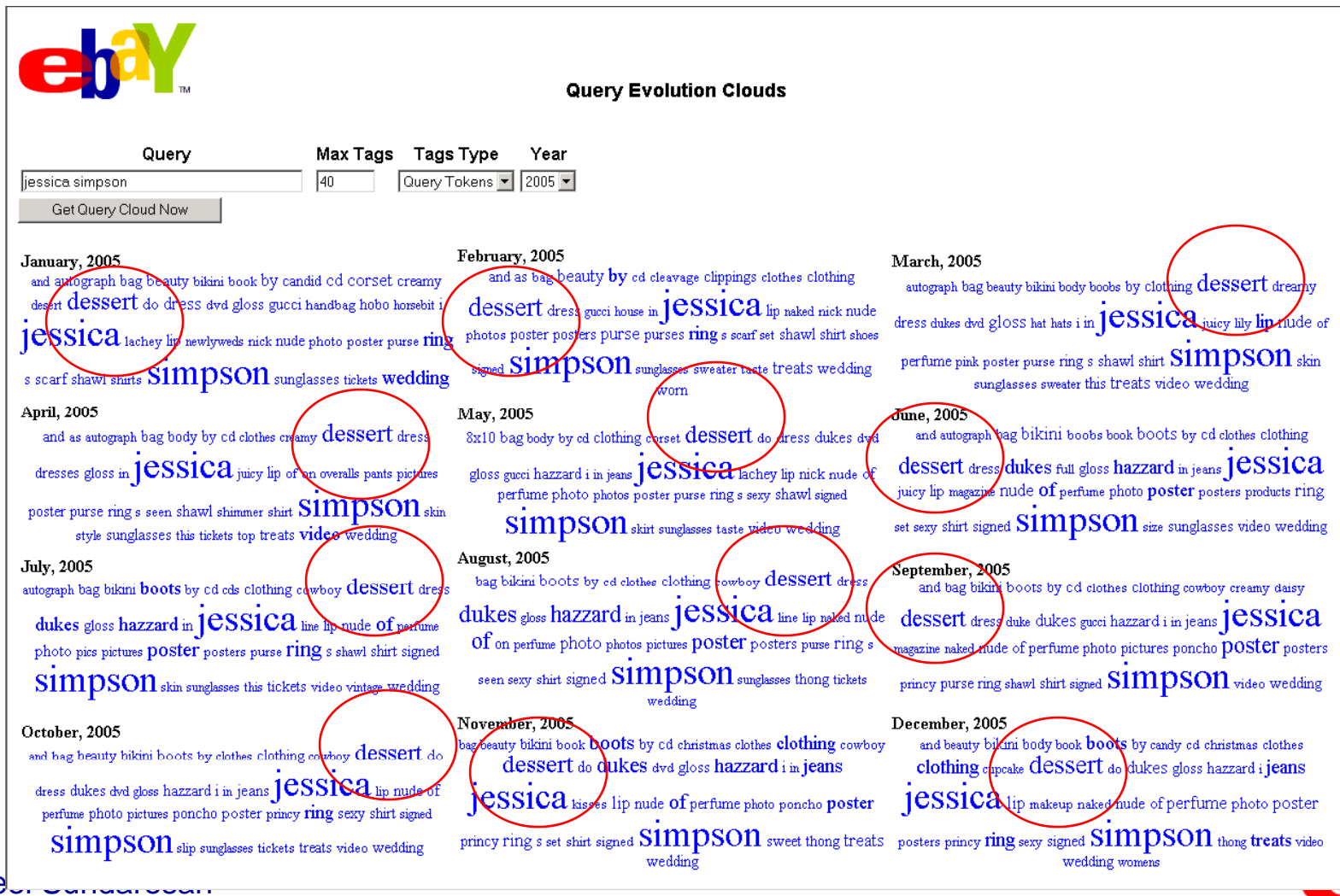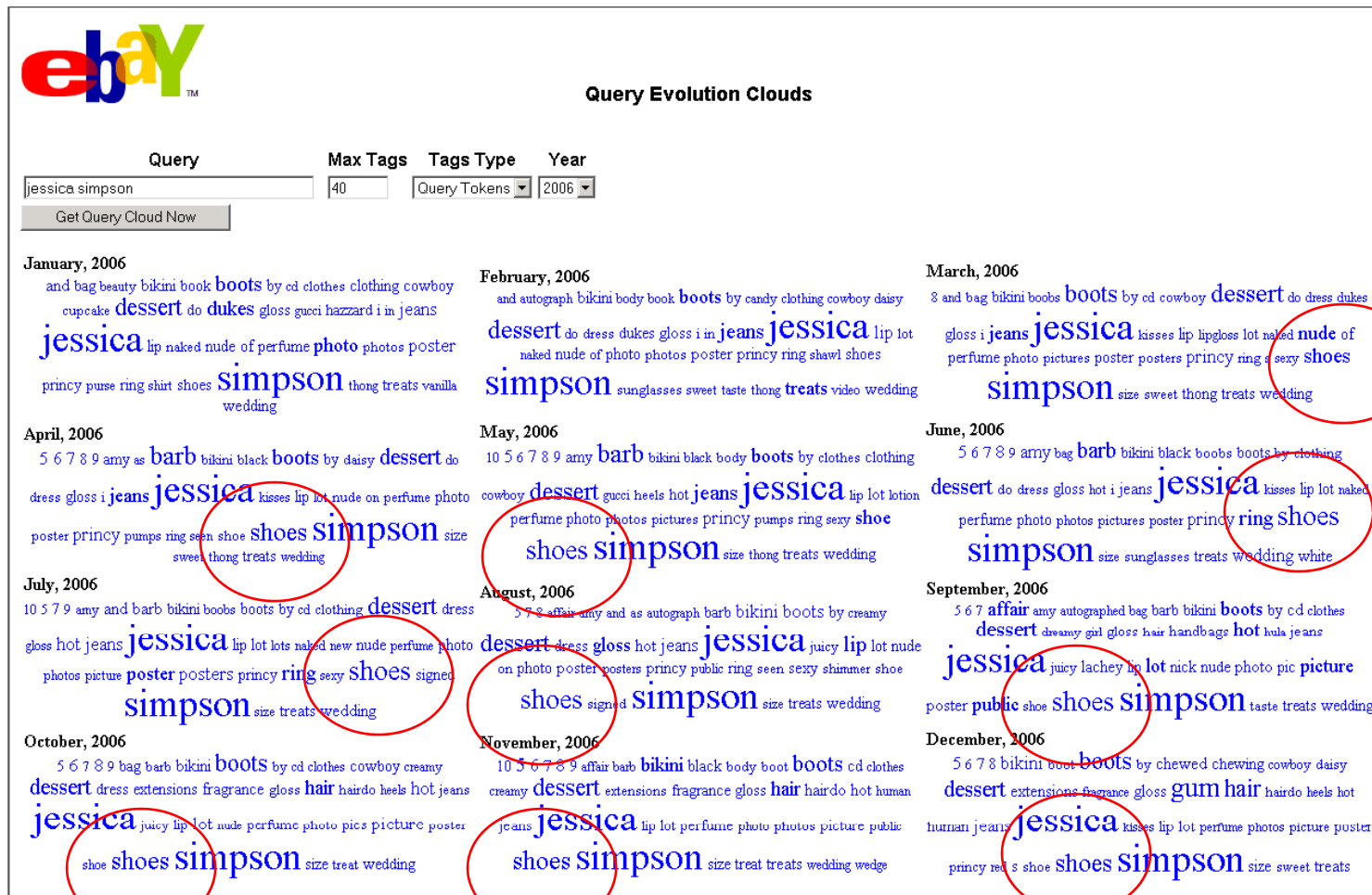| SKIN SO SOFT | NEW IN BOX | Free Shipping | Gift With Purchase |
|---|---|---|---|

Neel Sundaresan

eBaY®
**Research Labs**

Britney Spears

Paris Hilton

Perfume

Shoes

Hair

Lindsey Lohan

Jessica Simpson

Photo

Rocket Dog

Jessica Alba

Poster

Rainbow Brite

Tara Reid

Christina Aguilera

dolls

*Parikh, Sundaresan*

Neel Sundaresan

ebaY
Research Labs

# Evolution of Jessica Simpson

# JS (contd)



Neel Sundaresan

# JS (contd)



Neel Sundaresan

# Recommender Systems – At Scale

- MovieLens – 10M ratings, 10K movies, 70K users

- Netflix Prize Challenge – 100M ratings, 10K movies, 500K users

- eBay Challenge – 2+M txns/day, 10M+ new items, 200+M buyers and sellers

Neel Sundaresan

# Recommender Systems – Known Ways

- ### Collaborative  Filtering

    Neighborhood models

    - compute relationships between items or users

    Latent Factor Models

    -  Use Matrix Factorization. Work well with large users, item factors

Neel Sundaresan

# Matrix Factorization

- Users and items are mapped as a latent factor space (of dimensionality $\sigma$)

- User-item interactions are modeled as inner products

  Each item $v$ is represented by a vector $q_v$ in $R^\sigma$ representing the extent to which the <u>item has the different features</u> in $R^\sigma$

  Each user $u$ is represented by a vector $p_u$ in $R^\sigma$ representing the extent to which the <u>user is interested in different features</u> in $R^\sigma$

  Preference for item $v$ by user $u$ $r_{uv}$ is computed as the dot-product $q_v^T p_u$

  The major task is computing the mapping of each item and user to the corresponding factor vectors in $R^\sigma$

  *Ref. Koren et al. IEEE 09*

Neel Sundaresan

# Matrix Factorization (contd.)

- ## SVD raises difficulties due to high sparsity

  In Netflix sparsity is 1:100, eBay sparsity is even higher (1:10K)

  Instead of SVD, directly modeling on the observed data is preferred

  To learn the factor vectors $q_v$ and $p_u$

  We minimize the regularized square error on the known observations:

  $$\min_{q,p} \sum_{(u,v)\in\kappa} (r_{uv} - q_v^T p_u)^2 + \lambda(\|q_v\|^2 + \|p_u\|^2)$$

  where $r_{uv}$ is known for $\kappa$

  This minimization equation is solved using for e.g. Stochastic Gradient Descent or Alternating Least Squares

*Ref. Koren et al. IEEE 09*

Neel Sundaresan

# Matrix Factorization – the eBay Challenge

- The $q$ space of items is extremely large and volatile.

    There's very little overlap between buyer interests/purchases and actual items

- However, the query space is fairly static (even though in the millions)

    The user space can be mapped to queries and the item space can be mapped to queries

    One could use explicit feature vectors on the user and item space $q_v\, p_u$ then optimizing to discover a weight matrix W in $q_v\, W\, p_u$

    Alternatively, Latent Topic Models (e.g. Latent Dirichlet Allocation - LDA) can be used to map users to latent topics, and latent topics to search queries

Neel Sundaresan

# What Shape is the Universe of eBay?



Neel Sundaresan

# eBay - The Neck-Tie

Shen, Sundaresan

Cell Phones & PDAs 2005

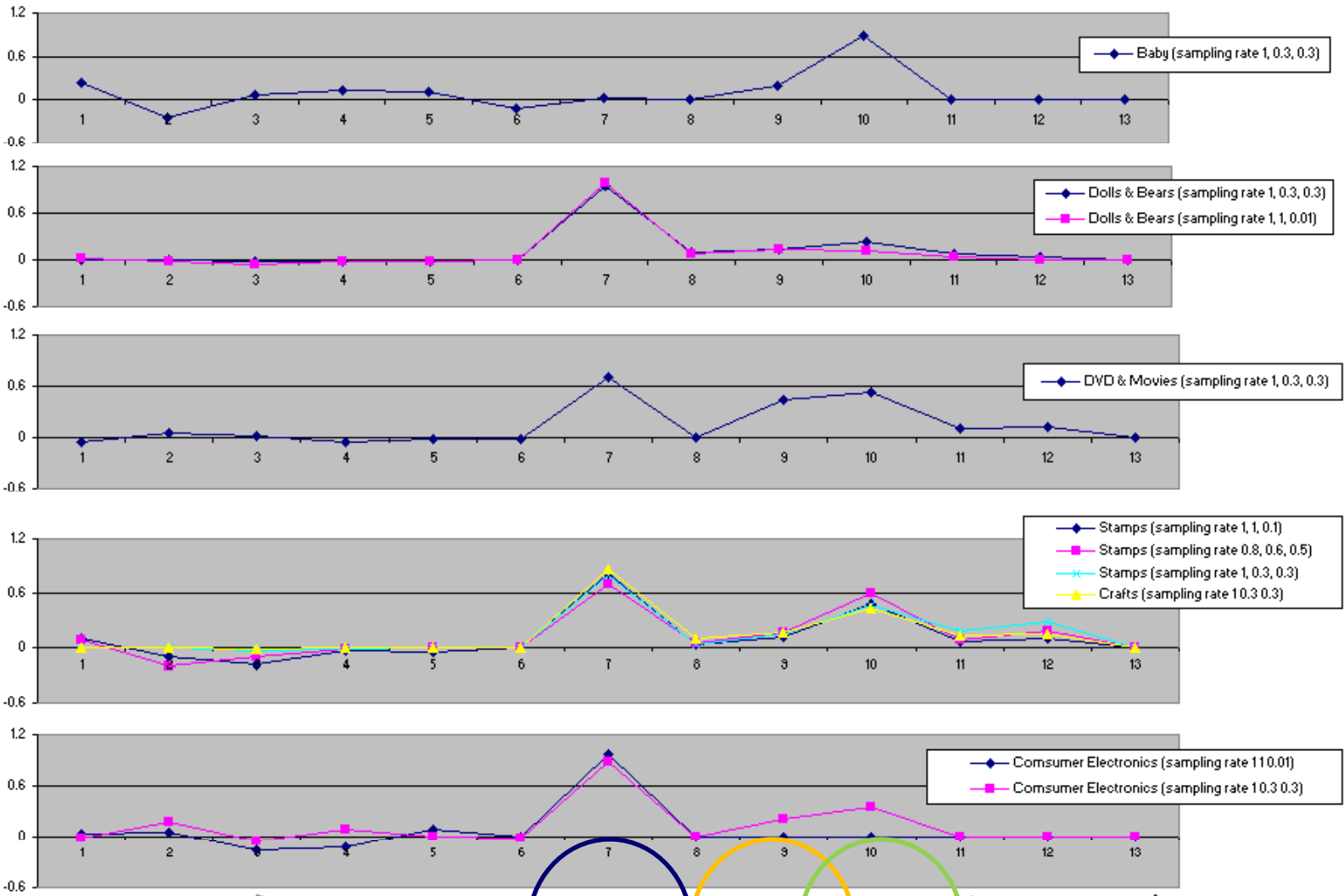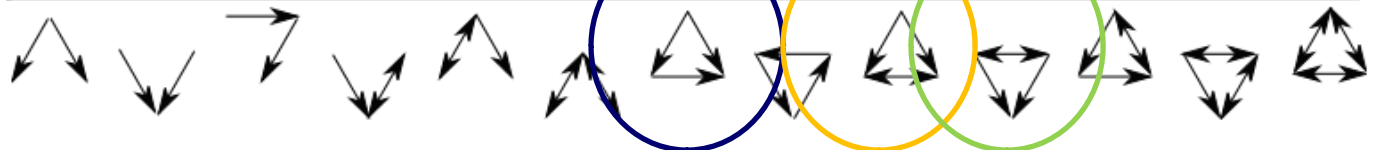Cell Phones & PDAs 2010

Neel Sundaresan

# Significance Profiles



Neel S
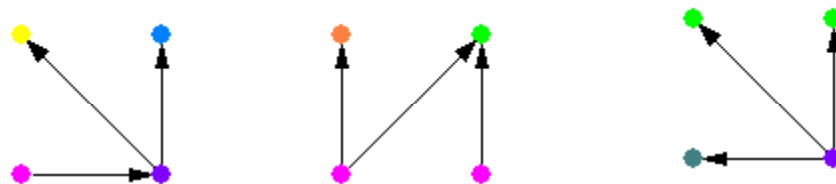
# Motifs: Feedback Enhanced Triad Distribution Across Categories

- Stamps



- Antiques



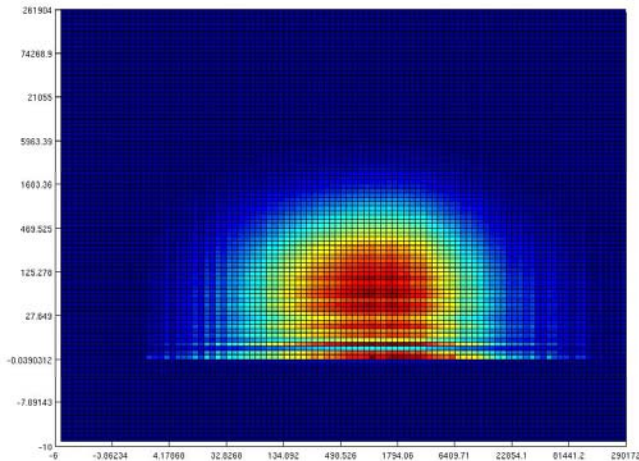| -1 | 0 | 5 | 20 | 100 | 500 | 1000 | 5000 | >5000 |
|----|---|---|----|-----|-----|------|------|-------|

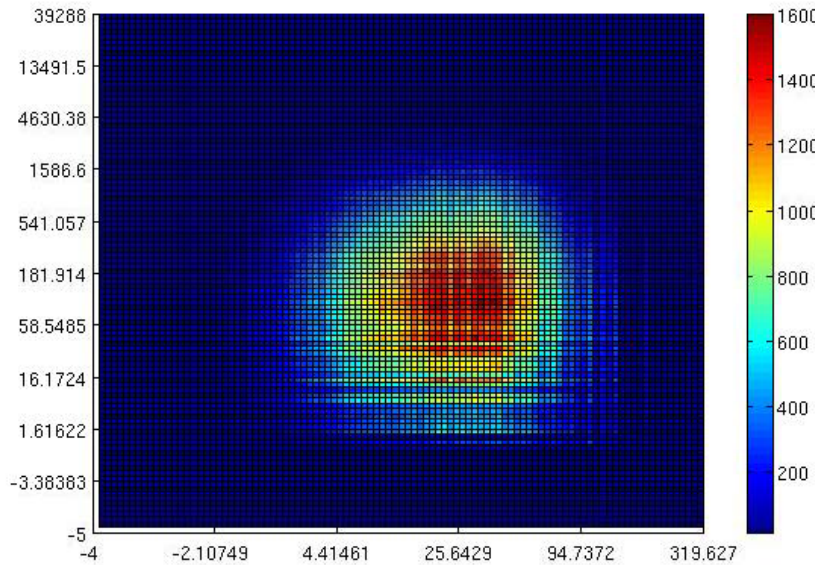Neel Sundaresan

Shen, Sundaresan

# Assortative Mining – Auroral Diagrams to Measure Preferential Attachment



Global
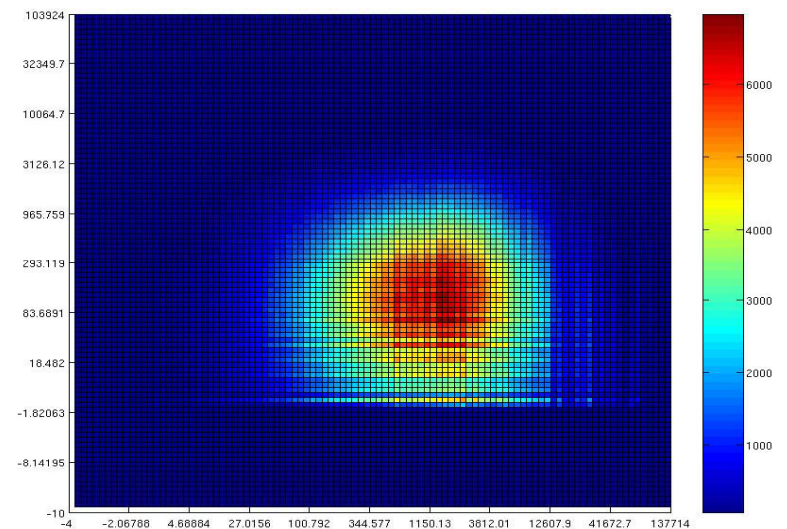
Arts and Crafts

Collectibles

Neel Sundaresan

Shen, Sundaresan

# Questions?



*ArtAmnesia*

Neel Sundaresan