# *Is the Whole Greater Than the Sum of Its Parts?*
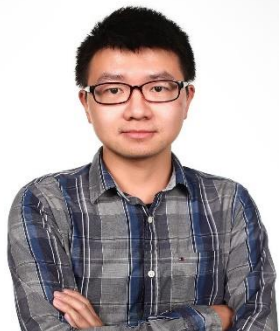
Presenter: **Liangyue Li**

Joint work with

**Liangyue Li**
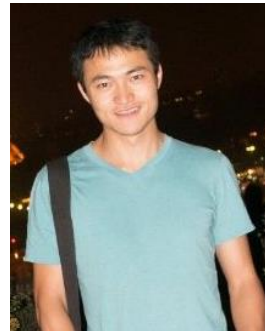(ASU)

**Hanghang Tong**
(ASU)

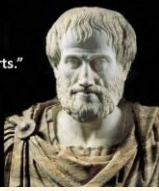**Yong Wang**
(HKUST)

**Conglei Shi**
(IBM->Airbnb)

**Nan Cao**
(Tongji)

**Norbou Buchler**
(US ARL)

**DATA Lab**

**Arizona State University**

# From the Ancient Philosophy

*The whole is greater than the sum of its parts.* -- **Aristotle**

- **Whole**: a collection of parts

- **Parts**: individual elements

- **Aristotle's hypothesis**:

  – whole > sum of parts

**DATA Lab**

# Part-Whole in Team Science



Research Team



Sports Team



Film Crew



Sales Team

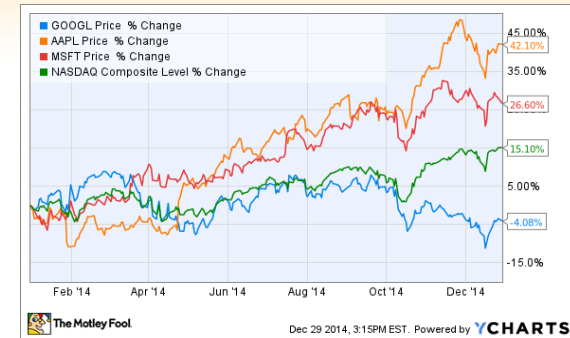Whole – Team
Parts – Team members

DATA Lab

# Part-Whole Beyond Teams
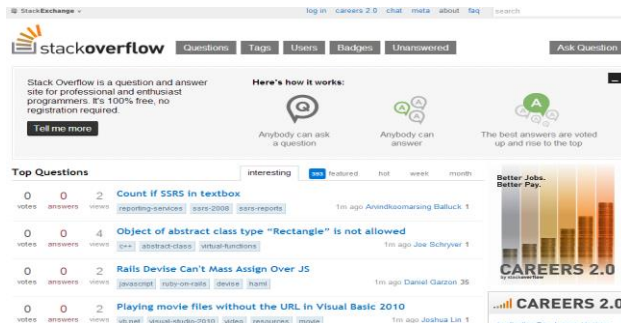


**Autonomous System**
Whole: system
Parts: individual drones
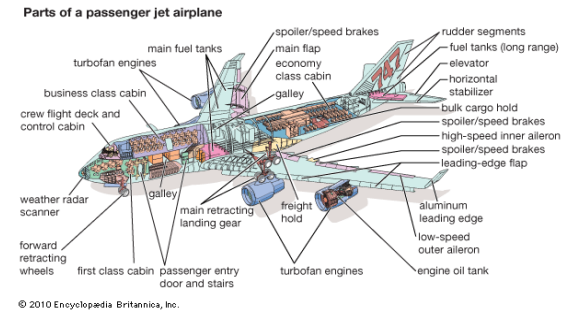


**Stock Market**
Whole: DJIA
Parts: individual stock



**Community Question Answering**
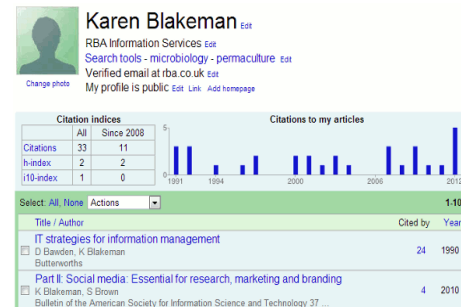Whole: question
Parts: individual answers



Parts of a passenger jet airplane

**System Reliability**
Whole: system
Parts: individual component

DATA Lab

# Outcome of Part-Whole



**Whole**: Team
**Part**: Members

**Whole outcome**: Team's performance
**Part outcome**: each member's performance

**Whole**: Researcher
**Part**: Publications

**Whole outcome**: h-index
**Part outcome**: #citations of publications

Question: how can we predict the outcome of whole/parts?

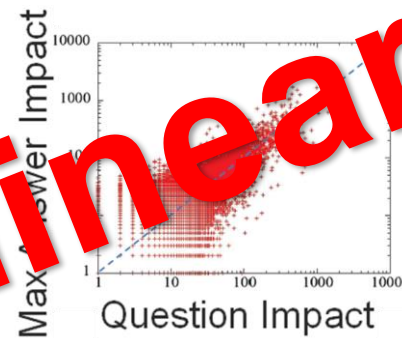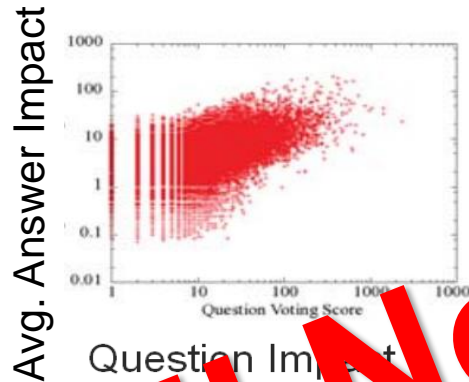DATA Lab

# Predict the Part-Whole Outcomes

- **Existing Algorithmic Work**

  - Separate models for parts and whole

  - Joint <span style="color:red">linear</span> models

- **Aristotle's hypothesis: whole>sum(parts)**

- **Question: how to jointly predict part/whole**

  - by leveraging the part-whole relationship *beyond* the linear models?
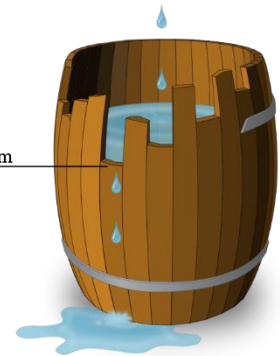
DATA
Lab

# Challenges -- Modeling

- **Nonlinear** Part-whole Relationship

  - **Max**: impact of a question is strongly correlated with that of the *best* answer



  - **Min**: classic Wooden Bucket Theory

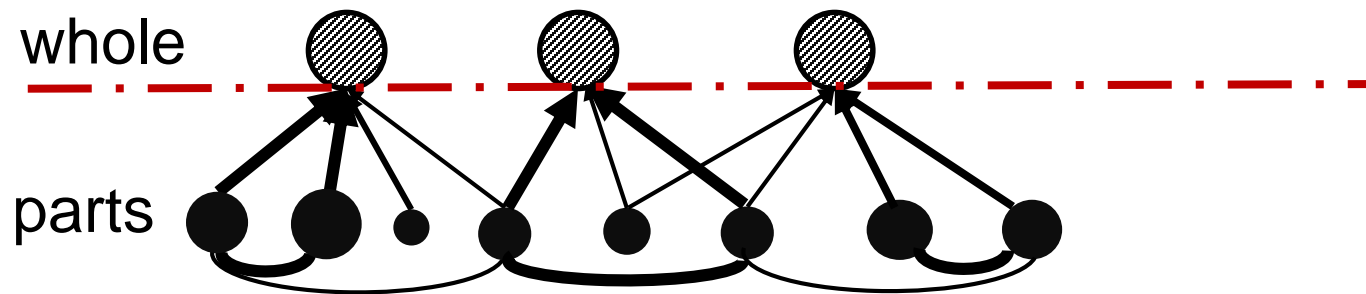  - **Sparsity**: team performance often dominated by a few top-performing team members

**DATA Lab**

**Arizona State University**

# Challenges – Modeling (con't)

- **Part-part Interdependency**

  - Parts are connected via underlying network

  - Impact of such interdependency on outcomes

Hypothesis-1: similar parts -> similar contribution to whole
Hypothesis-2: similar parts -> similar parts outcome



whole

parts

Question: how can we utilize
      1. nonlinear part-whole relationship
      2. part-part interdependency

# Challenges -- Algorithm

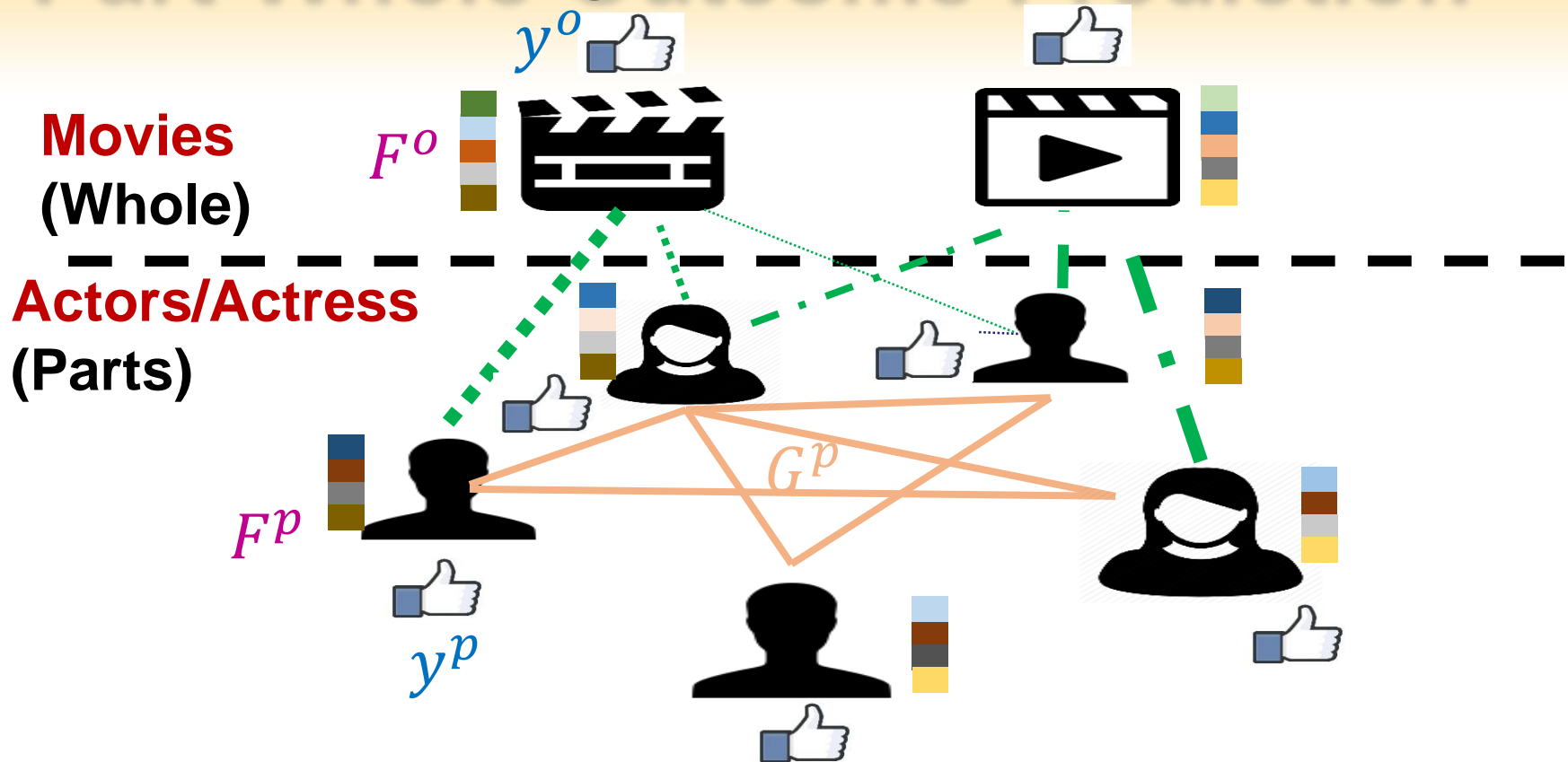**Non-linearity**

**+**

**Interdependency**

} **high complexity**

Question: how to scale up the computation?

**DATA Lab**

# Part-Whole Outcome Prediction



**Movies (Whole)**

$y^o$

$F^o$

**Actors/Actress (Parts)**

$F^p$

$G^p$

$y^p$

**Given**: 1. feature matrix for whole/part $F^o/F^p$
2. outcome vector for whole/part $y^o/y^p$
3. whole to part mapping $\phi$
4. parts' network $G^p$ (optional)

**Predict**: outcome of new whole/parts

DATA Lab

# Roadmap

- **Motivations**

- **PAROLE -- Modeling**

  - **Generic Framework**

  - Modeling Part-Whole Relationship

  - Modeling Part-Part Interdependency

- **PAROLE -- Optimization**

- **Empirical Evaluations**

- **Conclusions**

**DATA Lab**

**Arizona State University**

- Formulation

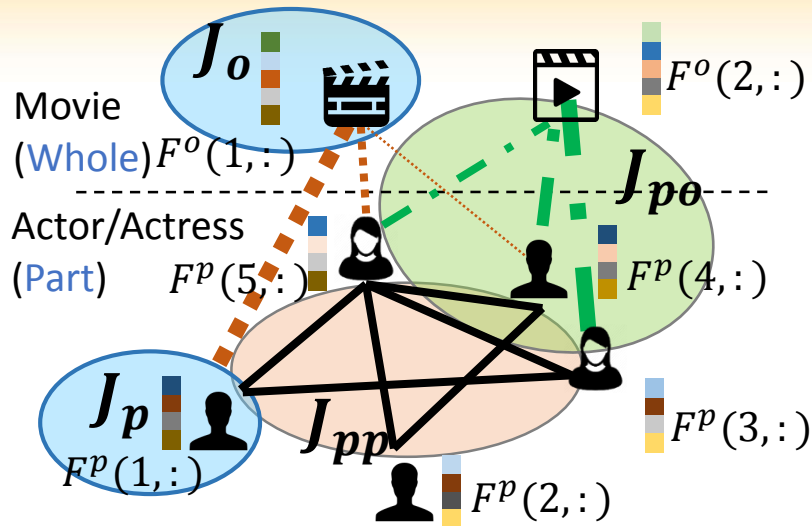$$\min J = J_o + J_p + J_{po} + J_{pp} + J_r$$

$$= \frac{1}{n_o} \sum_{i=1}^{n_o} L[f(F^o(i,:), w^o), y^o(i)]$$

$$+ \frac{1}{n_p} \sum_{i=1}^{n_p} L[f(F^p(i,:), w^p), y^p(i)]$$

$$+ \frac{\alpha}{n_o} \sum_{i=1}^{n_o} h(f(F^o(i,:), w^o), Agg(\phi(o_i)))$$

$$+ \frac{\beta}{n_p} \sum_{i=1}^{n_p} \sum_{j=1}^{n_p} G_{ij}^p g(f(F^p(i,:), w^p), f(F^p(j,:), w^p))$$

$$+ \gamma(\Omega(w^o) + \Omega(w^p))$$



Movie (Whole) $F^o(1,:)$
$F^o(2,:)$
$J_o$
$J_{po}$
Actor/Actress (Part)
$F^p(5,:)$
$F^p(4,:)$
$J_p$
$F^p(1,:)$
$J_{pp}$
$F^p(3,:)$
$F^p(2,:)$

$J_o$: **Predictive Model for Whole**

$J_p$: **Predictive Model for Part**

$J_{po}$: **Part-whole Relationship**

$J_{pp}$: **Part-part Interdependency**

$J_r$: **parameter regularizer**

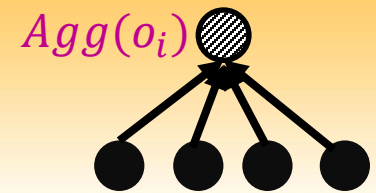**DATA Lab**

# Roadmap



- Motivations

- **PAROLE -- Modeling**

  - Generic Framework

  - **Modeling Part-Whole Relationship**

  - Modeling Part-Part Interdependency

- PAROLE -- Optimization

- Empirical Evaluations

- Conclusions

**DATA Lab**

**Arizona State University**
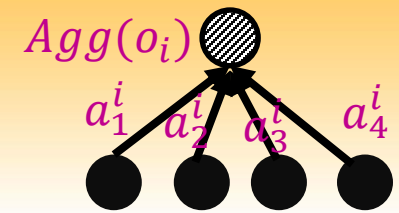
# Modeling Part-Whole Relationship

■ **Overview**: for each whole entity $o_i$, define

$$\underline{e_i} = \mathbf{F^o}(i,:)\mathbf{w^o} - \mathrm{Agg}(o_i)$$

– $e_i$: Measure the difference between

- predicted whole outcome using whole feature
- predicted whole outcome using aggregated parts outcome

■ **Key idea**: model part-whole relations by

- Different loss functions on $e_i$
- Different aggregation functions $Agg(\cdot)$

**DATA Lab**

# Overview

$Agg(o_i)$

$a_1^i \quad a_2^i \quad a_3^i \quad a_4^i$

- **Intuition:** whole ← (weighted) sum of parts

- **Details:**

$$e_i = F^o(i,:)w^o - Agg(o_i)$$

$$Agg(o_i) = \sum_{j \in \phi(o_i)} a_j^i F^p(j,:)w^p$$

  – $a_j^i$: weight of part $j$'s contribution to the whole
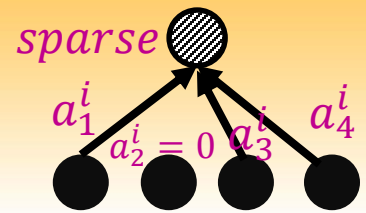
  $o_i$'s outcome

- **Remark:**

  – Characterize part-whole relationships

    • Use different loss functions on $e_i$

    • Use different norms on $a_i$

DATA Lab

# Linear Part-Whole Relation



- **Intuition**: Whole ← linear combination of parts

  - some parts play more important roles than the others in contributing to the whole outcome

- **Details**: $J_{po} = \frac{\alpha}{2n_o} \sum_{i=1}^{n_o} e_i^2$

- **Remark**:

  - $a_j^i = 1$: the whole is the sum of its parts
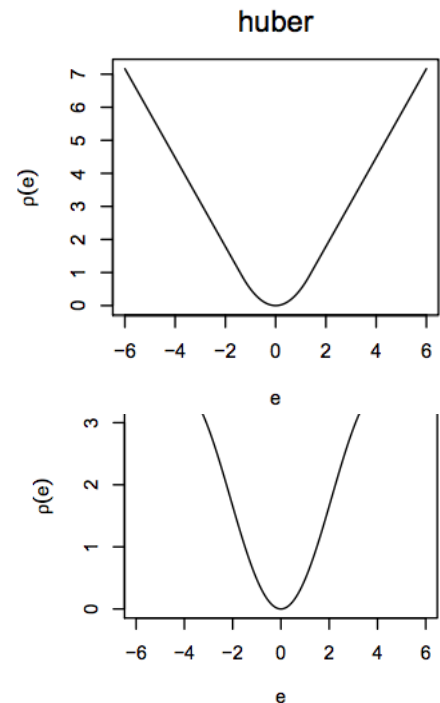
  - $a_j^i = \frac{1}{|o_i|}$: average coupling

**DATA Lab**

**Arizona State University**

# **Sparse Part-Whole Relation**



- **Intuition**: Whole ← a few parts

  - some parts have little or no effect on the whole outcome

- **Details**: $J_{po} = \frac{\alpha}{n_o} \sum_{i=1}^{n_o} (\frac{1}{2} e_i^2 + \lambda |\mathbf{a}_i|_1)$

- **Remark**:

  - The $l_1$ norm can shrink some part contributions $a_j^i$ to exactly zero

  - **Nonlinear** part-whole relation

DATA Lab

# Ordered Sparse Part-Whole Relation

- **Intuition**: Whole ← a few top parts

  - team performance is determined by not only a few key members, but also the structural hierarchy between them

- **Details:** $J_{po} = \frac{\alpha}{n_o} \sum_{i=1}^{n_o} (\frac{1}{2} e_i^2 + \lambda \Omega_w(\mathbf{a}_i))$

  - $\Omega_w(x) = \sum_{i=1}^{n} |x|_{[i]} w_i = \mathbf{w}^T |\mathbf{x}|_{\downarrow}$: ordered weighted $l_1$ norm

  - $w \in \mathcal{K}_{m+}$: vector of non-increasing non-negative weights

DATA Lab

# Robust Part-Whole Relation

- **Intuition**: Whole ← parts that are not outliers

  – squared loss is sensitive to outliers.

- **Solution**: robust regression model

- **Details**: $J_{po} = \frac{\alpha}{n_o} \sum_{i=1}^{n_o} \rho(e_i)$

  – $\rho(\cdot)$ is robust estimator

huber

| Case \ Method | $|e| \leq t$ | $|e| > t$ |
|---|---|---|
| Huber $\rho_H(e)$ | $\frac{1}{2}e^2$ | $t|e| - \frac{1}{2}t^2$ |
| Bisquare $\rho_B(e)$ | $\frac{t^2}{6}\left\{1 - [1 - (\frac{e}{t})^2]^3\right\}$ | $\frac{t^2}{6}$ |

DATA Lab

# Maximum Part-Whole Relation

- **Intuition**: Whole ← max(parts)

  - team performance dominated by the best team member/leader

- **Details**:

  - $Agg(o_i) = \max(parts'\ outcome)$ [not differentiable]

  - Soft max function: $\max(x_1, x_2, \dots, x_n) \approx$ $\ln(\exp(x_1) + \exp(x_2) + \cdots + \exp(x_n))$

  - Aggregation: $Agg(o_i) = \ln(\sum_{j \in \phi(o_i)} \exp(F^p(j,:)w^p))$

$$J_{po} = \frac{\alpha}{2n_o} \sum_{i=1}^{n_o} e_i^2$$

**DATA Lab**

# Summarize Part-Whole Relations

| Name | $Agg(o_i)$<br>Aggregation of parts | $J_{po}$<br>Sub-objective | Remark |
|------|-----------------------------------|---------------------------|--------|
| **Maximum** | $\ln(\sum \exp(F^p(j,:)w^p))$ | $\frac{\alpha}{2n_o}\sum e_i^2$ | Nonlinear<br>Whole ← max(parts) |
| **Linear** | $\sum a_j^i F^p(j,:)w^p$ | $\frac{\alpha}{2n_o}\sum e_i^2$ | Linear<br>Whole ← linear<br>combination of parts |
| **Sparse** | $\sum a_j^i F^p(j,:)w^p$ | $\frac{\alpha}{n_o}\sum(\frac{1}{2}e_i^2 + \lambda|a_i|_1)$ | Nonlinear<br>Whole ← a few parts |
| **Ordered Sparse** | $\sum a_j^i F^p(j,:)w^p$ | $\frac{\alpha}{n_o}\sum(\frac{1}{2}e_i^2 + \lambda\Omega_w(a_i))$ | Nonlinear<br>Whole ← a few top parts |
| **Robust** | $\sum a_j^i F^p(j,:)w^p$ | $\frac{\alpha}{n_o}\sum \rho(e_i)$ | Nonlinear<br>Whole ← parts that are not outliers |

**DATA Lab**

# Roadmap



- **Motivations**

- **PAROLE -- Modeling**

  – Generic Framework

  – Modeling Part-Whole Relationship

  – **Modeling Part-Part Interdependency**

- PAROLE -- Optimization

- Empirical Evaluations

- Conclusions

**DATA Lab**

**Arizona State University**

# Modeling Part-Part Interdependency

- **Effect on the whole outcome**

  - **Intuition**: closely connected parts might have similar contribution to the whole outcome

  - **Details**:

  $$\mathcal{J}_{po} = \frac{\alpha}{n_o} \sum_{i=1}^{n_o} \left[ \frac{1}{2} e_i^2 + \lambda |\mathbf{a}_i|_1 + \frac{1}{2} \sum_{k,l \in \phi(o_i)} G_{kl}^p (a_k^i - a_l^i)^2 \right]$$

  

  $a_1^i \quad a_2^i \quad a_3^i \quad a_4^i$

  $G_{12}^p \qquad G_{14}^p$

  - Similar parts (large $G_{kl}^p$)

  $\rightarrow$ similar contributions ($a_k^i \approx a_l^i$)

DATA Lab

# Modeling Part-Part Interdependency

- **Effect on the parts outcome**

  - **Intuition**: closely connected parts might share similar outcomes themselves

  - **Details**:

$$\mathcal{J}_{pp} = \frac{\beta}{2n_p} \sum_{i=1}^{n_p} \sum_{j=1}^{n_p} G_{ij}^p (\mathbf{F}^p(i,:)\mathbf{w}^p - \mathbf{F}^p(j,:)\mathbf{w}^p)^2$$



$F^p(4,:)w^p$

$F^p(1,:)w^p$

$G_{12}^p$　$G_{14}^p$

- Similar parts (large $G_{ij}^p$)

→ similar predicted outcomes ($F^p(i,:)w^p \approx F^p(j,:)w^p$)

DATA Lab

# Roadmap

- Motivations

- PAROLE -- Modeling

- **PAROLE -- Optimization**

- Empirical Evaluations

- Conclusions

**DATA Lab**

**Arizona State University**

# Optimization Solution



- **Formulation**:

  - $J = J_o(w^o) + J_p(w^p) + J_{po}(w^o, w^p, a_j^i) + J_{pp}(w^p) + J_r(w^o, w^p)$

- **Observation**:

  - not jointly convex w.r.t. $w^o, w^p, a_i^j$

  - Convex w.r.t. to one block while fixing others

- **Solution**: block coordinate descent

**DATA Lab**

**Arizona State University**

# Block Coordinate Descent

- Three coordinate blocks: $w^o, w^p, a_j^i$

- Update one block while fixing others

- Update each block

  - (proximal) gradient descent

| | $\dfrac{\partial J_{po}}{\partial w^o}$ | $\dfrac{\partial J_{po}}{\partial w^p}$ | $\dfrac{\partial J_{po}}{\partial a_i}$ or proximal gradient update |
|---|---|---|---|
| Maximum Agg | $\dfrac{\alpha}{n_o}\sum_{i=1}^{n_o} e_i (F^o(i,:))'$ | $\dfrac{\alpha}{n_o}\sum_{i=1}^{n_o} e_i \dfrac{\sum_{j\in\phi(o_i)}(F^p(j,:))'\tilde{y}_i^p}{\sum_{j\in\phi(o_i)}\tilde{y}_i^p}$ | N/A |
| Linear Agg | $\dfrac{\alpha}{n_o}(F^o)'(F^o w^o - MF^p w^p)$ | $-\dfrac{\alpha}{n_o}(F^p)'M'(F^o w^o - MF^p w^p)$ | $e_i(-F^p(\phi(o_i),:)w^p) + L_i^p a_i$ |
| Sparse Agg | $\dfrac{\alpha}{n_o}(F^o)'(F^o w^o - MF^p w^p)$ | $-\dfrac{\alpha}{n_o}(F^p)'M'(F^o w^o - MF^p w^p)$ | $z = a_i - \tau\big[e_i(-F^p(\phi(o_i),:)w^p) + L_i^p a_i\big]$ $a_i \leftarrow prox_{\lambda\tau l_1}(z)$ |
| Order Sparse Agg | $\dfrac{\alpha}{n_o}(F^o)'(F^o w^o - MF^p w^p)$ | $-\dfrac{\alpha}{n_o}(F^p)'M'(F^o w^o - MF^p w^p)$ | $z = a_i - \tau\big[e_i(-F^p(\phi(o_i),:)w^p) + L_i^p a_i\big]$ $a_i \leftarrow prox_{\lambda\tau\Omega_w}(z)$ |
| Robust Agg | $\dfrac{\alpha}{n_o}\sum_{i=1}^{n_o} \dfrac{\partial\rho(e_i)}{\partial e_i} F^o(i,:)'$ | $\dfrac{\alpha}{n_o}\sum_{i=1}^{n_o} \dfrac{\partial\rho(e_i)}{\partial e_i}\Big(-\sum_{j\in\phi(o_i)} a_j F^p(j,:)'\Big)$ | $\dfrac{\alpha}{n_o}\Big[\dfrac{\partial\rho(e_i)}{\partial e_i}(-F^p(\phi(o_i),:)w^p) + L_i^p a_i\Big]$ |

# Optimization Properties

details

- **Convergence and Optimality**
  - Under mild conditions, the optimization alg converges to a coordinate-wise minimum point

- **Complexity**
  - The alg scales linearly w.r.t. the size of part-whole graph in both time and space

Whole

Parts

Complexity: $O(n_o d_o + n_p d_p + m_{po} + m_{pp})$
$n_o$: #whole entities
$n_p$: #part entities
$m_{po}$: #links from whole to parts
$m_{pp}$: #links in part-part network
$d_o, d_p$: feature dimension of whole, parts

DATA Lab

**Arizona State University**

# Roadmap

- Motivations

- PAROLE -- Modeling

- PAROLE -- Optimization

- **Empirical Evaluations**

- Conclusions

**DATA Lab**

**Arizona State University**

# Datasets

| Data | Whole | Part | #Whole | #Part |
|------|-------|------|--------|-------|
| Math | Question (#votes) | Answer (#votes) | 16,638 | 32,876 |
| SO | Question (#votes) | Answer (#votes) | 1,966,272 | 4,282,570 |
| DBLP | Author (h-index) | Paper (#citation) | 234,681 | 129,756 |
| Movie | Movie (# 👍 ) | Actors/Actress (# 👍 ) | 5,043 | 37,365 |

- **Setup**: sort whole in chronological order, gather first $x$ percent and corresponding parts as training, test on last 10%
- **Metric**: root mean squared error (RMSE)

**DATA Lab**

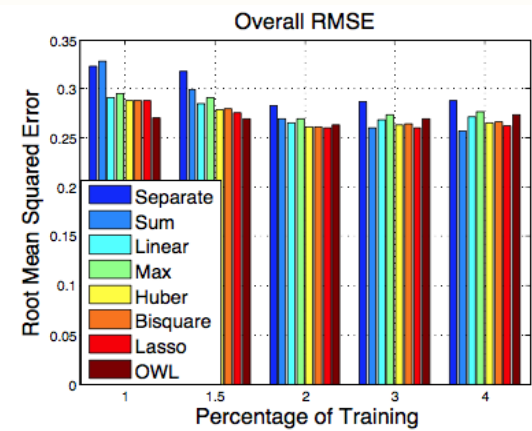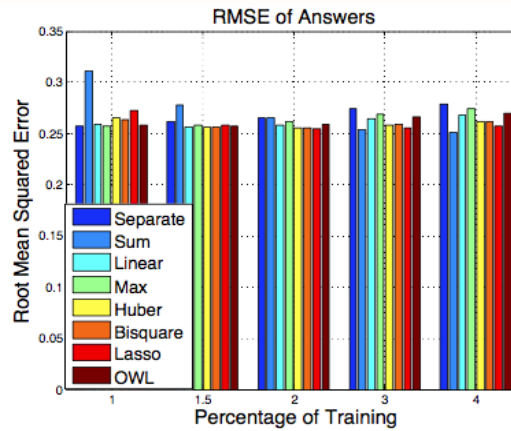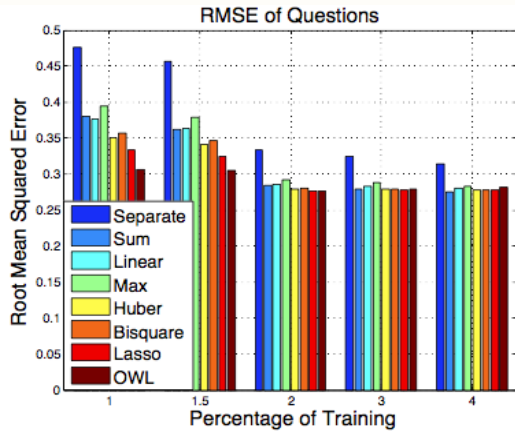# Outcome Prediction Performance



Math

## Observations

1. Joint prediction models > separate models
2. Non-linear part-whole relationships (max, Huber, Bisquare, Lasso, OWL) > linear relationships (Sum, Linear)
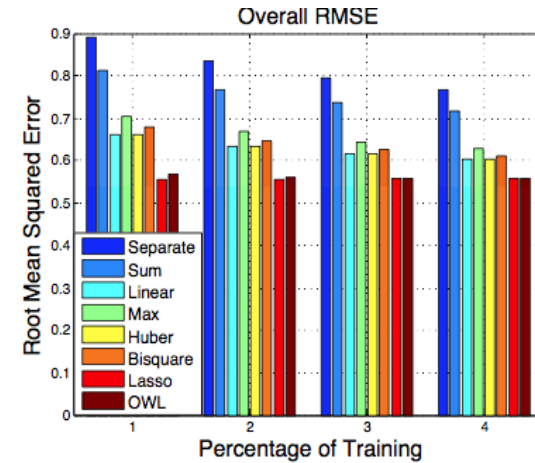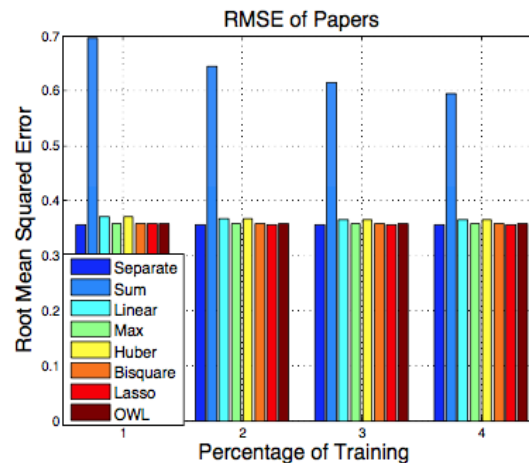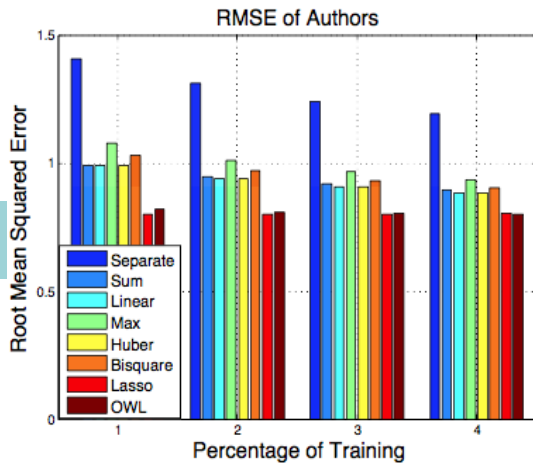3. Lasso and OWL are the best methods in most cases

Whole

Parts

Overall

SO



(a) RMSE of question outcome prediction.  (b) RMSE of answer outcome prediction.  (c) Overall RMSE.

DBLP



(a) RMSE of author outcome prediction.  (b) RMSE of paper outcome prediction.  (c) Overall RMSE.
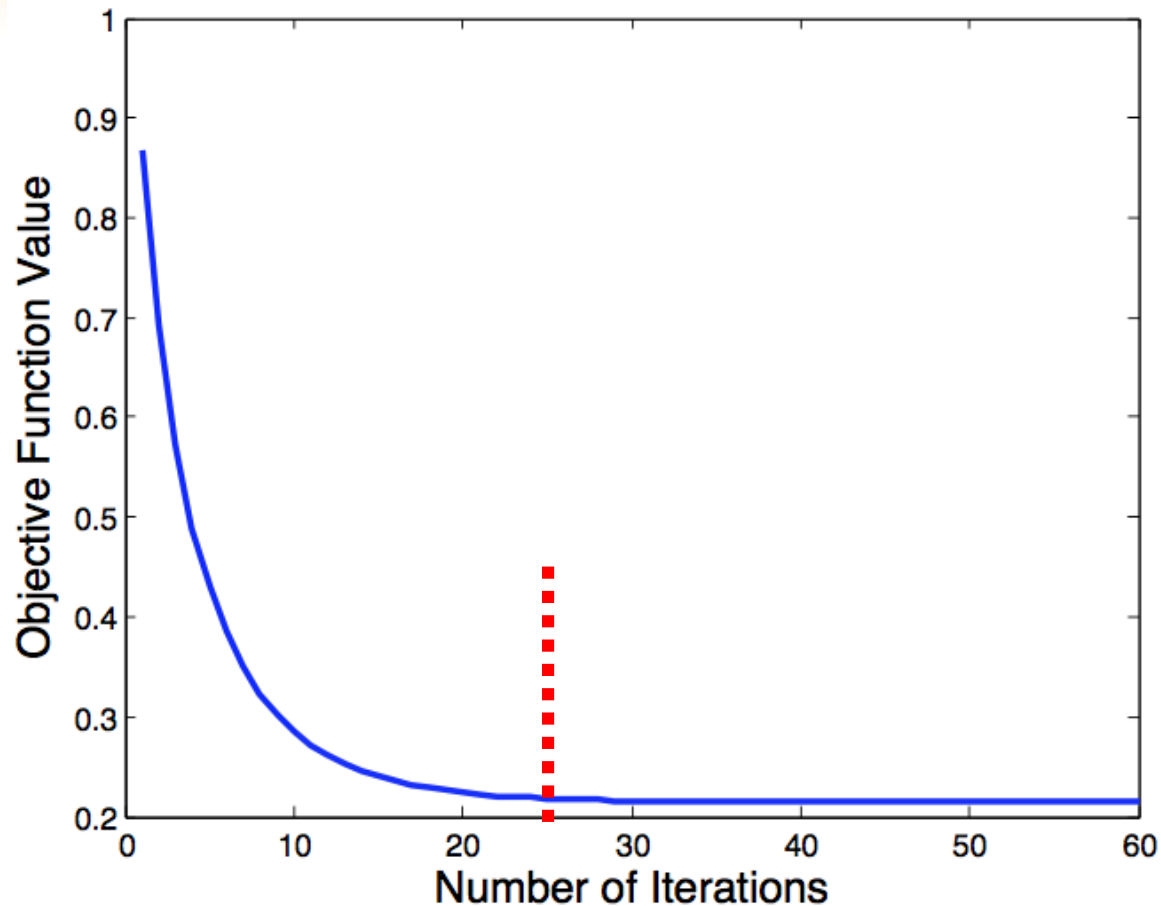
**DATA Lab**

# Effect of part-part interdependency

Movie



- **PAROLE-Basic** – no network information
- Part-part interdependency on whole outcome and parts outcome both boost the performance
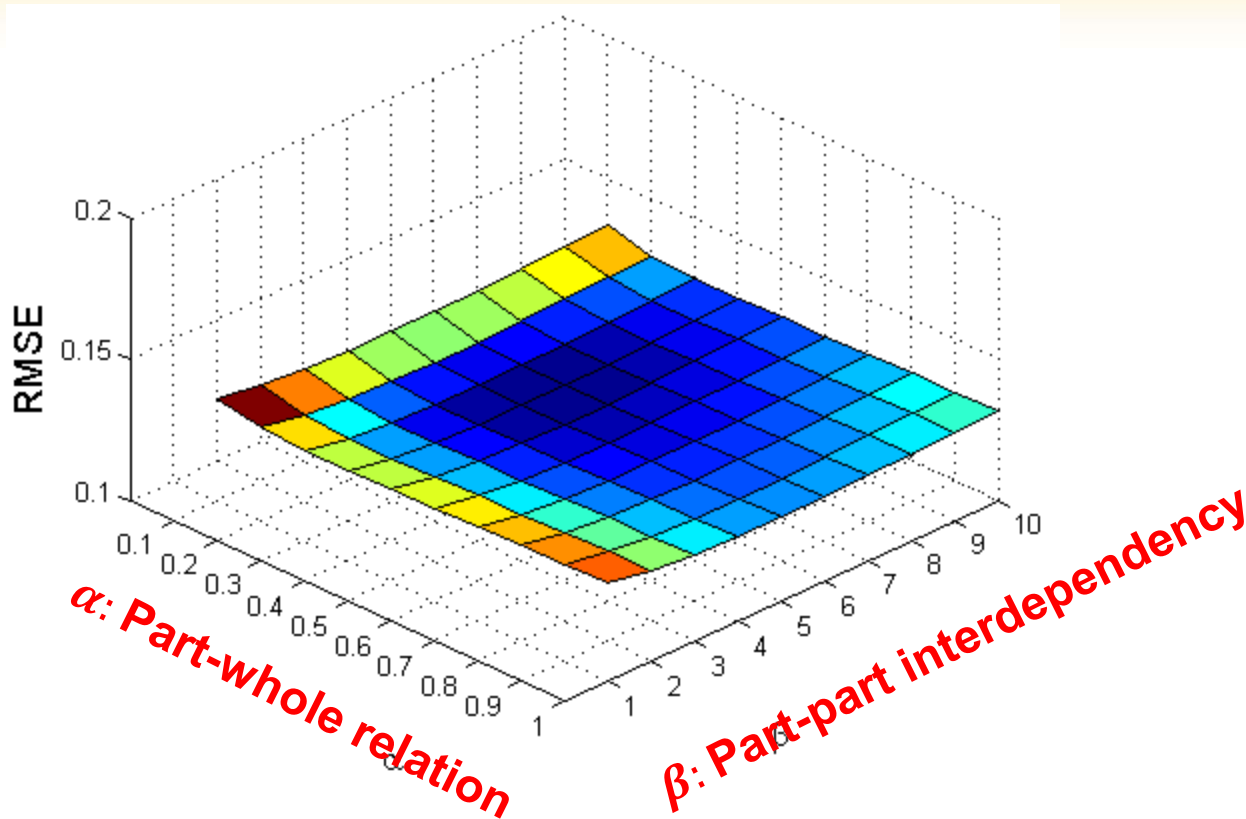
DATA Lab

# Convergence Analysis

SO



- PAROLE converges fast (25-30 iterations)

**DATA Lab**

# Parameter Sensitivity



Movie

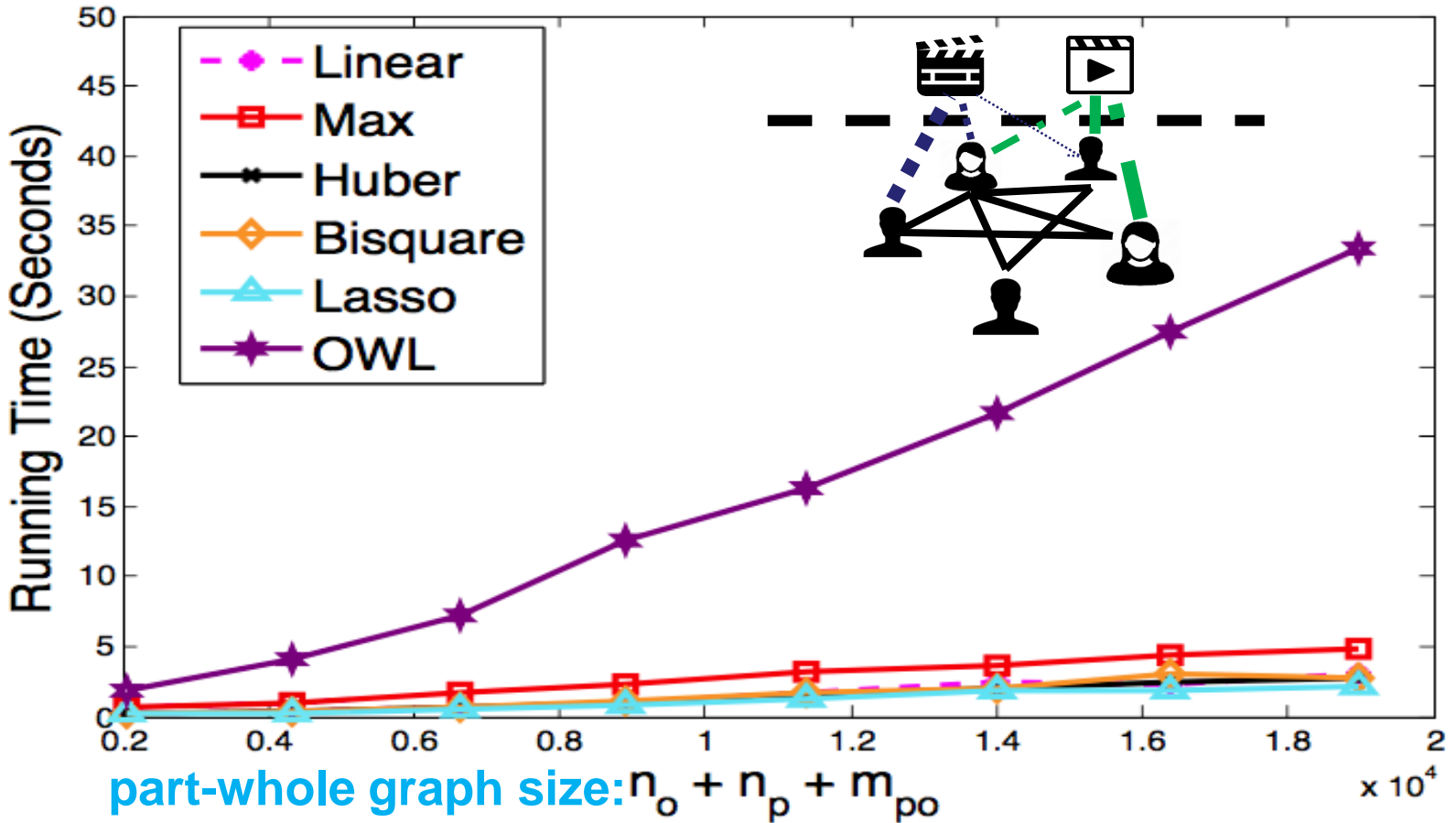- $\alpha$ controls importance of part-whole relation
- $\beta$ controls importance of part-part interdependency
- Stable in a relatively large parameter space

DATA Lab

**Arizona State University**

# Scalability of PAROLE



SO

- PAROLE scales linearly w.r.t. part-whole graph size

DATA Lab

# Roadmap

- Motivations

- PAROLE -- Modeling

- PAROLE -- Optimization

- Empirical Evaluations

- **Conclusions**

**DATA Lab**

**Arizona State University**

# Conclusions -- PAROLE

- **Goals**: leverage potentially non-linear part-whole relationships for outcome prediction

- **Solutions**: PAROLE

  - **Modeling**
    - Part-whole relationship
    - Part-part interdependency

  - **Optimization**
    - Block coordinate descent
    - Converges to a coordinate-wise minimum point
    - Scales linearly w.r.t. the part-whole graph size

**DATA Lab**