

# Experimental Projects on Web Algorithms

Yury Lifshits  
<http://yury.name>

CalTech, Fall'07  
Invited lecture at CS141a

1 / 19

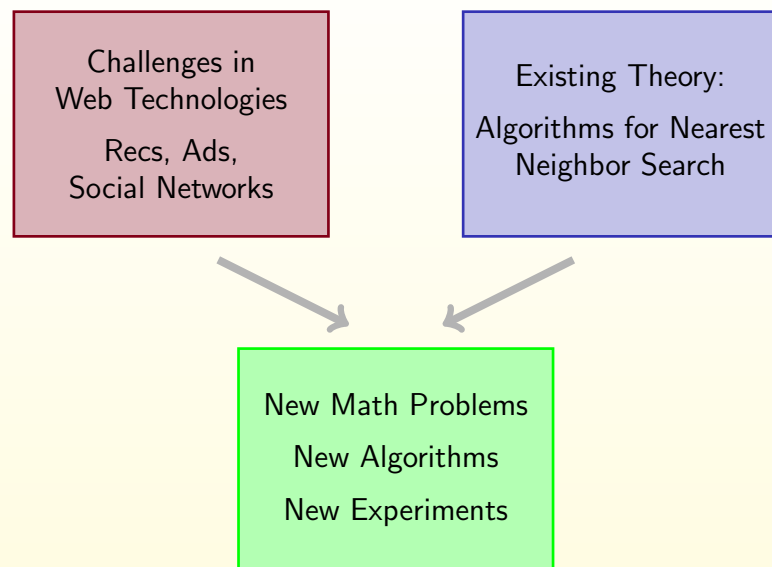
## Invitation to CS101.2

New Caltech course  
Algorithmic Problems Around the Web:

- <http://yury.name/algoweb.html>
- MW 11:00-11:55, Jorgensen 287
- Lectures: algorithms for nearest neighbor search
- Projects: adjusting above algorithms to web technologies
- Datasets: friendship graph, users-ads graph

2 / 19

## Course Philosophy



3 / 19

## Outline

- 1 Challenges in Web Technologies
- 2 Existing Theory: Nearest Neighbors
- 3 Topics for Experimental Projects

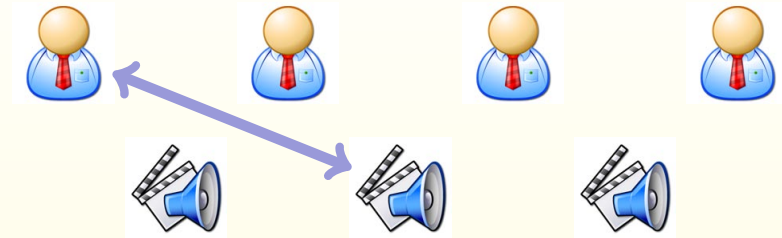
4 / 19

# Part I

## Challenges in Web Technologies

5 / 19

## Recommendation Systems



### Approaches:

- Content-based
- Collaborative filtering

6 / 19

## Behavioral Targeting

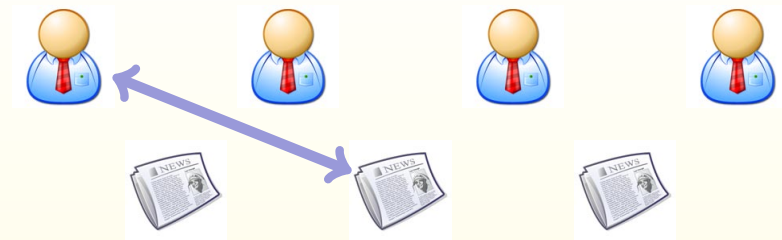


### Ad targeting:

- Ancient: broadcasting
- Current: contextual
- Future: behavioral

7 / 19

## Personalized News Aggregation



### Factors to take into account:

- Friendship
- Content
- Feedback (previous ratings)
- Popularity (votes, comments, hyperlinks)

8 / 19

# Social Networks Analysis

Social network:

Nodes

Edges

Examples of relations: financial exchange, friends, dislike, conflict, trade, web links, sexual relations, disease transmission, airline routes, etc.

## Our focus

Community discovery

Burst detection

9 / 19

# Part II Theory of Nearest Neighbors

10 / 19

# Nearest Neighbors Informally

To preprocess a database of  $n$  objects so that given a query object, one can effectively determine its nearest neighbors in database

11 / 19

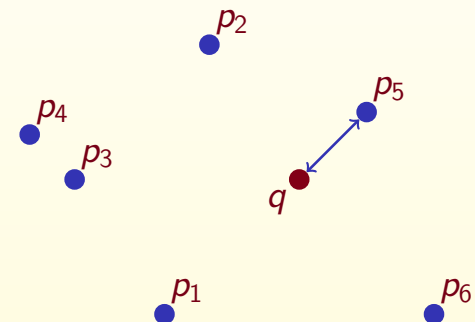
# More Formally

**Search space:** object domain  $\mathbb{U}$ , similarity function  $\sigma$

**Input:** database  $S = \{p_1, \dots, p_n\} \subseteq \mathbb{U}$

**Query:**  $q \in \mathbb{U}$

**Task:** find  $\operatorname{argmax}_{p_i} \sigma(p_i, q)$



12 / 19

## Some Solutions for NN Problem

Sphere Rectangle Tree Orchard's Algorithm LAESA  
k-d-B tree Geometric near-neighbor access tree  
Excluded middle vantage point forest.mvp-tree Fixed-height  
fixed-queries tree AESA Vantage-point  
tree R\*-tree Burkhard-Keller tree BBD tree  
Navigating Nets Voronoi tree Balanced aspect ratio tree Metric tree  
vp<sup>s</sup>-tree M-tree Locality-Sensitive Hashing  
SS-tree R-tree Spatial approximation tree Multi-vantage  
point tree Bisector tree mb-tree  
Generalized hyperplane tree  
Hybrid tree Slim tree Spill Tree Fixed queries tree X-tree k-d  
tree Balltree Quadtree Octree Post-office tree

13 / 19

## Part III Topics for Experimental Projects

14 / 19

## E1 Recommendations for Blog Posts

### Available information:

Friendship graph  
Comments, hyperlinks  
Keywords of interests, post content

**Task:** For every user recommend 10 posts from last day that seems to be the most interesting for him/her

15 / 19

## E2 CTR Prediction

### Available information:

Click-or-not bipartite graph

**Task:** Predict click-through rate for given pair "user-ad"

16 / 19

## E3 Social Networks Visualization

### Input:

Friendship graph

### Similarity:

Number of joint friends  
Length of shortest path

### Task:

Construct embedding into 2D  
that put similar people close to each other

17 / 19

## E4 Disorder Analysis

**Disorder inequality** for some constant  $D$ :

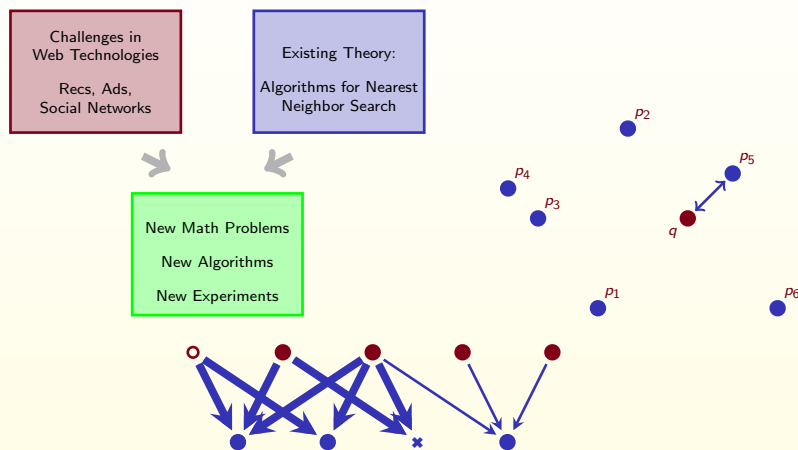
$$\forall p, r, s \in \{q\} \cup S : \text{rank}_r(s) \leq D \cdot (\text{rank}_p(r) + \text{rank}_p(s))$$

### Tasks:

- Compute disorder values for various datasets
- Implement disorder-based algorithms for NNS
- Study their performance

18 / 19

## Last Slide



Thanks for your attention! Questions?

19 / 19