# Overview of Machine Learning Research

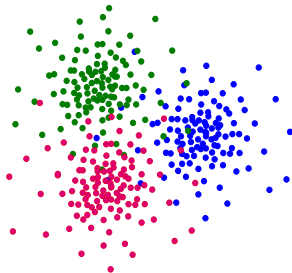**Anima Anandkumar**

U.C. Irvine

# Learning with Big Data



High Dimensional Regime

- Missing observations, gross corruptions, outliers, ill-posed problems.
- Needle in a haystack: finding low dimensional structures in high dimensional data.

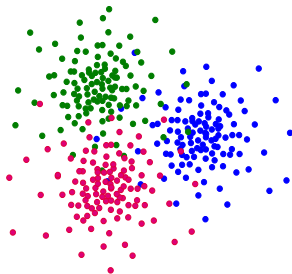  Principled approaches for finding low dimensional structures?

# Application 1: Clustering

- Basic operation of grouping data points.
- Hypothesis: each data point belongs to an unknown group.

# Application 1: Clustering

- Basic operation of grouping data points.
- Hypothesis: each data point belongs to an unknown group.



Probabilistic/latent variable viewpoint

- The groups represent different distributions. (e.g. Gaussian).
- Each data point is drawn from one of the given distributions. (e.g. Gaussian mixtures).

# Application 2: Topic Modeling

Document modeling

- Observed: words in document corpus.
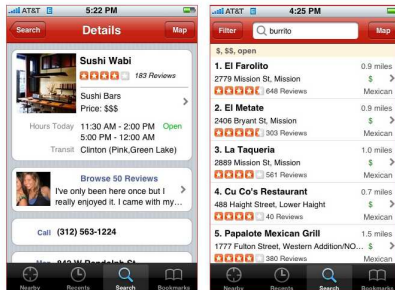- Hidden: topics.
- Goal: carry out document summarization.

# Application 3: Understanding Human Communities



**Social Networks**

- Observed: network of social ties, e.g. friendships, co-authorships
- Hidden: groups/communities of actors.

# Application 4: Recommender Systems



Recommender System

- Observed: Ratings of users for various products, e.g. yelp reviews.
- Goal: Predict new recommendations.
- Modeling: Find groups/communities of users and products.
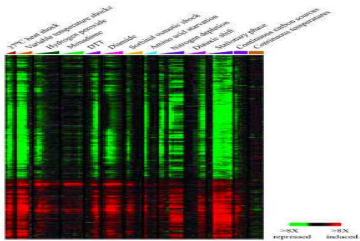
# Application 5: Feature Learning



| Label | Features |
|---|---|
| 0 | 2.1  5.2  0   0 |
| 1 | 0    0    2   1 |
| 1 | 1.1  0    0   0 |
| 0 | 0    0    7   0 |

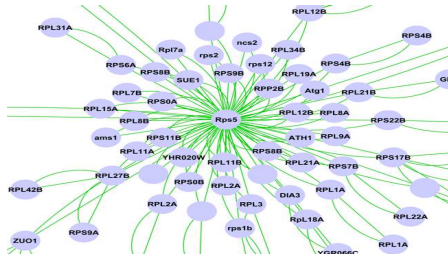## Feature Engineering

- Learn good features/representations for classification tasks, e.g. image and speech recognition.
- Sparse representations, low dimensional hidden structures.

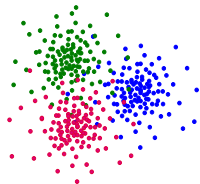# Application 6: Computational Biology



Gasch *et al.* Mol Biol Cell 2000.

- Observed: gene expression levels
- Goal: discover gene groups
- Hidden variables: regulators controlling gene groups

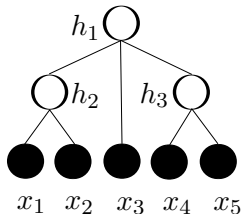# How to model hidden effects?

Basic Approach: mixtures/clusters
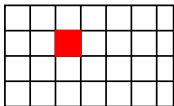- Hidden variable $h$ is categorical.
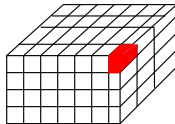


Advanced: Probabilistic models
- Hidden variable $h$ has more general distributions.
- Can model mixed memberships.

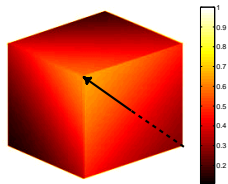# Learning Algorithms through Tensor Factorization



vs.

- Co-occurrence of three-words in a document, e.g. [apple, orange, banana].

## Tensor Eigenvectors

- Can learn the hidden topics by finding tensor eigenvectors.
- Common friends (neighbors) of triplets of nodes in a social networks.
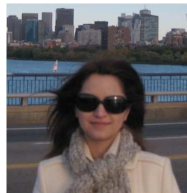
# My Research Group

Furong Huang

Majid Janzamin

Hanie Sedghi

Niranjan UN

Forough Arabshahi

# Resources and Course Information

## Resources and Course Information

- ML summer school lectures available at
  `http://newport.eecs.uci.edu/anandkumar/MLSS.html`
- Publications at `http://newport.eecs.uci.edu/anandkumar/`

# Resources and Course Information

- ML summer school lectures available at
  `http://newport.eecs.uci.edu/anandkumar/MLSS.html`
- Publications at `http://newport.eecs.uci.edu/anandkumar/`

Courses

- EECS 298: Large scale ML: theory and practice.
  - ▶ Cloud-based programming, spectral and tensor methods. Hadoop framework.
- EECS 298 (formerly 251B): Statistical learning theory
  - ▶ Theoretical course. Non-parametrics, optimization, regularization, concentration bounds.