

**DFCI Boston: Using the Weighted
Histogram Analysis Method (WHAM) in
cancer biology and the Yeast Protein
Databank (YPD); Latent Dirichlet
Analysis (LDA) for biological sequences
and structures**

**Russell Hanson
DFCI – April 24, 2009**

Outline

■ Part I

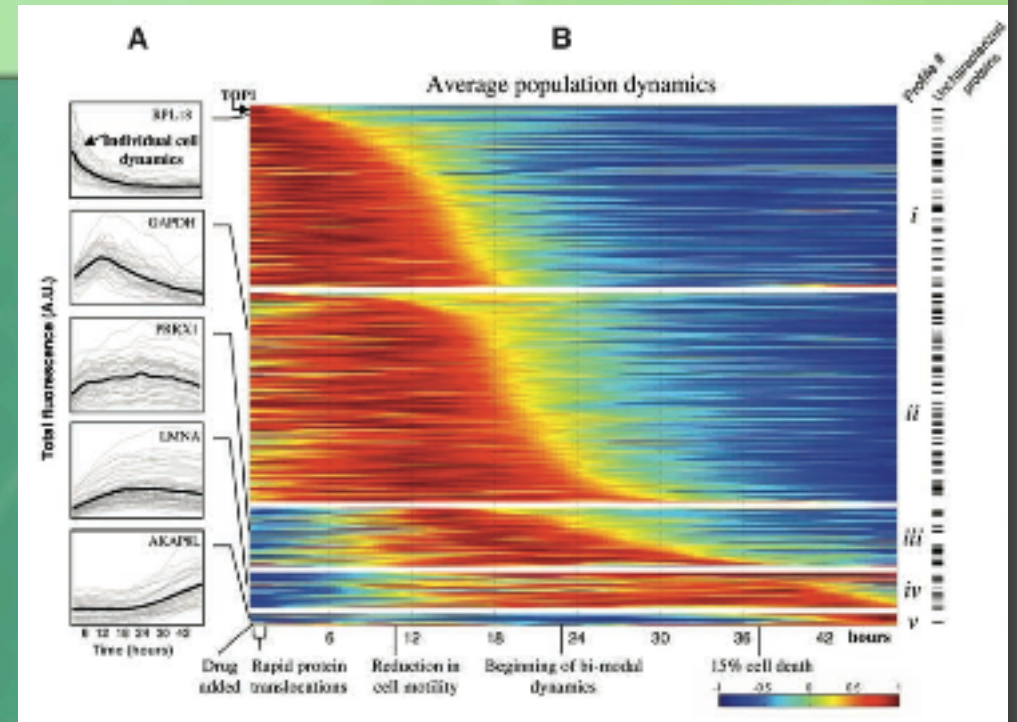
- Uri Alon's cell fate modeling
- Weighted histogram analysis method

■ Part II

- Gene function in cell
- Weighted histogram analysis method

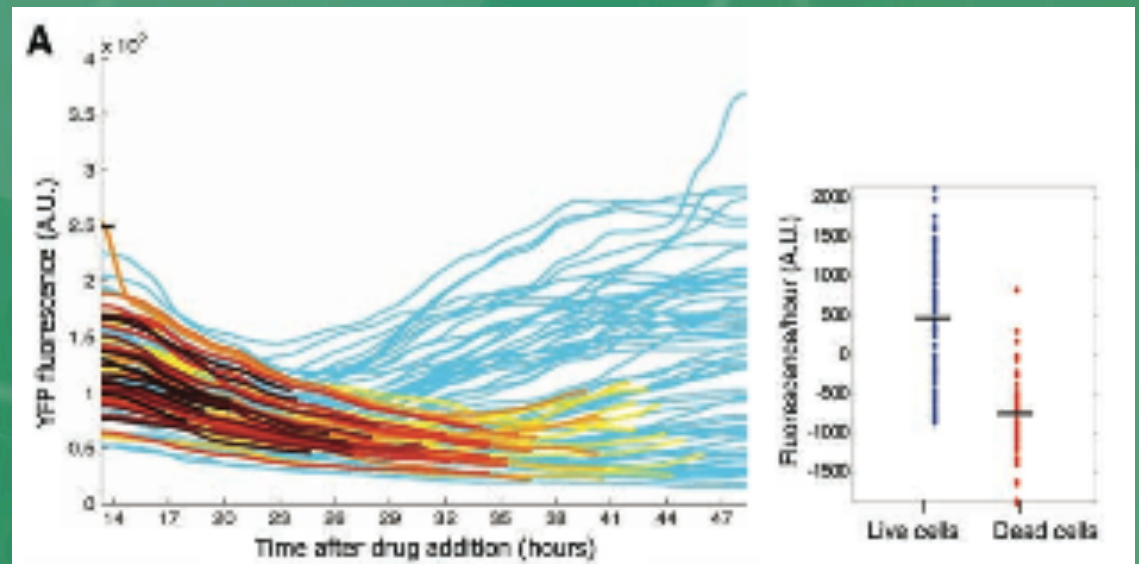
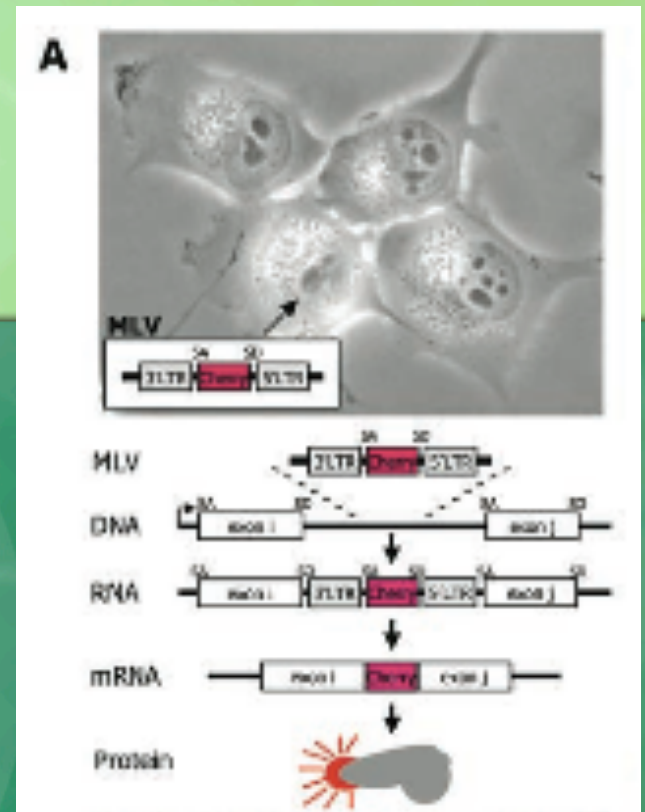
■ Part III

- Corpora clustering
- Latent Dirichlet Analysis for topical clustering

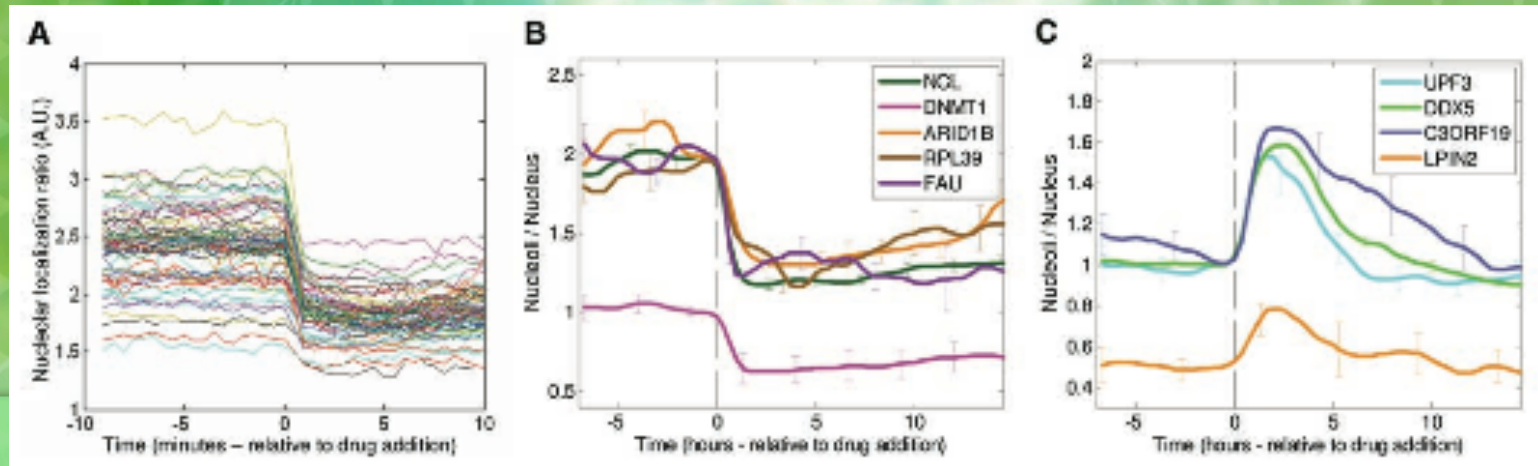


Uri Alon's talk (p1)

- Question: Why do some cells live and some die?
- Procedure: Add chemotherapy drugs; monitor cell fate



Uri Alon's talk (p2)



- “Superposition” idea:

$$P_{1,2}(t) = a \cdot P_1(t) + b \cdot P_2(t)$$

s.t. $a+b=1$.

- N.B. non-linearities; non-linear responses; non-linear systems; 5 coupled differential equations in biochemical processes; simulation methods in biochemical pathways (Plectix); system doesn't saturate to 1 upon additional drug addition

Some ideas from statistical mechanics - WHAM

- “Superposition” idea:

$$P_{1,2}(t) = a*P_1(t)+b*P_2(t)$$

$$\text{s.t. } a+b=1.$$

The density of states is related to the histogram by

$$W(S) = \overline{N_n(S)n_n^{-1}} \exp[f_n - K_n S] \quad (1)$$

For a set of values $\{K_n|n = 1, R\}$, these can be recombined as

$$W(S) = \sum_{n=1}^R p_n(S) N_n(S) n_n^{-1} \exp[f_n - K_n S] \quad (2)$$

with

$$\sum_{n=1}^R p_n(S) = 1 \quad (3)$$

Inserting the actual histograms from $W(S)$ above, and minimize the error, obtain

$$p_n(S) = \frac{n_n g_n^{-1} \exp[K_n S - f_n]}{\sum_{m=1}^R n_m g_m^{-1} \exp[K_m S - f_m]} \quad (4)$$

Define:

$$P(S, K) = W(S) \exp[KS] \quad (5)$$

And finally the essential multiple-histogram equations are obtained as

$$P(S, K) = \frac{\sum_{n=1}^R g_n^{-1} N_n(S) \exp[KS]}{\sum_{m=1}^R n_m g_m^{-1} \exp[K_m S - f_m]} \quad (6)$$

Part II – Yeast functional classes

One seeks to determine the network of interacting proteins, the “interactome”, by choosing one protein and hypothesizing that if it were in a multi-protein complex, it would have 4–5 interacting in-membrane neighbors.

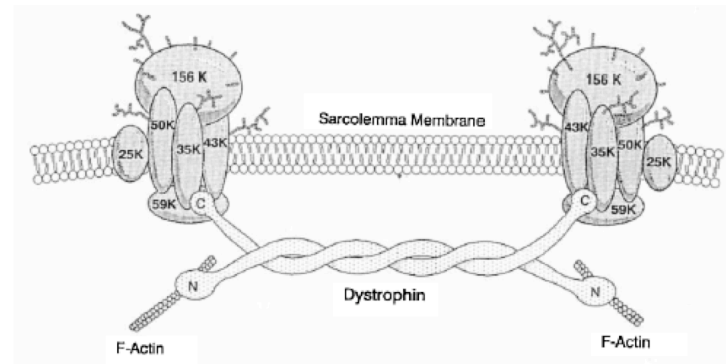


Figure 1. Proteins interacting in a trans-membrane protein complex, the dystrophin-glycoprotein complex from (Ervasti, 1991).

YPD data

- Yeast Protein Data Bank has the following data:

| Classes | |
|---------|--|
| 1 | Metabolism |
| 2 | Energy |
| 3 | Cell cycle and DNA processing |
| 4 | Transcription |
| 5 | Protein synthesis |
| 6 | Protein fate |
| 7 | Cellular transportation and transportation mechanism |
| 8 | Cell rescue, defense and virulence |
| 9 | Interaction with cell environment |
| 10 | Cell fate |
| 11 | Control of cell organization |
| 12 | Transport facilitation |
| 13 | Others |

Table 1. The first 12 YPD functional classes plus a container class.

A generic weighting function is:

$$\begin{aligned} \text{Total weight} = & A * \text{metabolism} + B * \text{energy} + C * \text{DNA processing} \\ & + D * \text{transcription} + \dots \end{aligned} \quad (1)$$

Heuristic methods for the solution of equations of the form:

$$f(\cdot) = A * g(\cdot) + B * h(\cdot) + C * i(\cdot) + D * j(\cdot) + [E * k(\cdot) + F * l(\cdot)] + \dots \quad (2)$$

DFCI for constants A, B, C, D, \dots . A variety of heuristic methods for finding solutions of marginally ill-posed problems may be referenced in Michalewicz et al. (2000).

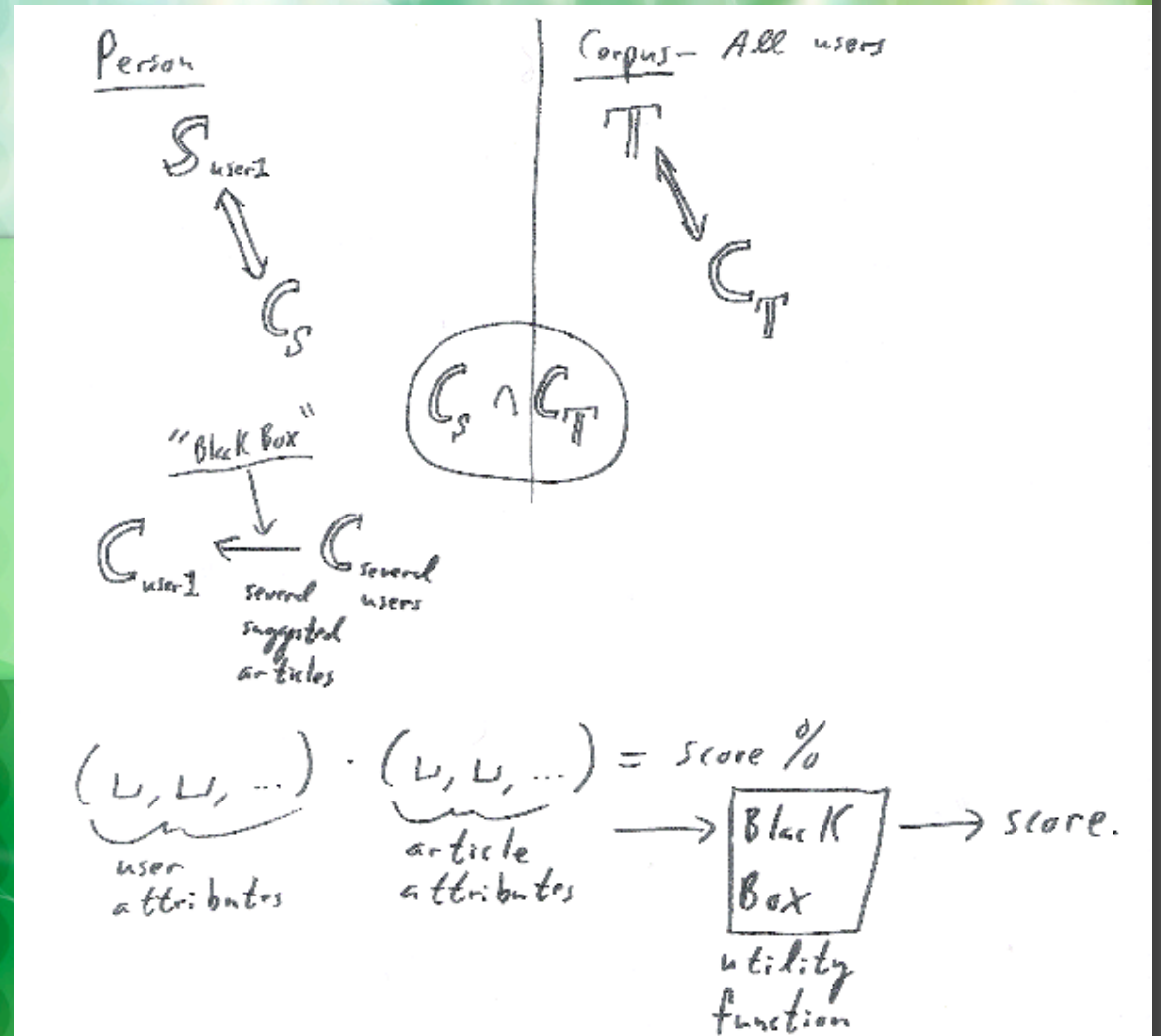
Matrix representation for YPD classifications

| Protein index no. | f | g | h | i | ... |
|-------------------|-------|------|-----|-----|-----|
| 1 | True | 0.8 | 0.7 | 0.4 | |
| 2 | False | 0.5 | 0.7 | 0.5 | |
| 3 | True | 0.28 | 0.6 | 0.6 | |
| ... | | | | | ... |
| 10 | False | 0.3 | 0.3 | 0.4 | ... |

Table 2. Matrix representation used for the YPD. Column *f* represents the true or false classification annotation and columns *g*, *h*, *i*, ... represent decimal scores for the functional classes.

Part III – Corpora and thematic clustering

- Shift gears slightly: *making recommendations for news websites, etc.* will return shortly
- Want to have more fine-grained recommendations than connectivity in user network -- weight in a given thematic cluster.



Attributized Bayesian Choice Modeling

Attributized content items, i , are stored as vectors in the choice-set database

such that:

$$\mathbf{A}_i = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{iN}),$$

where $\alpha_{i1}, \dots, \alpha_{iN}$ are the scores on the N attributes for item i .

Summary of tastes, T :

$$\mathbf{T}_u = (\beta_{u1}, \beta_{u2}, \dots, \beta_{uN}),$$

where $\beta_{u1}, \dots, \beta_{uN}$ are user u 's weights on N attributes.

- Collaborative Filtering for text and “news”:
 - Cold Start Problem (it isn't collaborative until it's collaborative)
 - Past Experience: Some people want the most popular (“Dodgers make offer to Manny Ramirez - Boston.com”); some don't (“Non-Abelian Anyons and Topological Quantum Computation”)
 - By weight in whole network; by weight in user's network; by weight in thematic cluster

Method (p1)

Document Universe. Let Ω be the *document universe*, i. e. the set of all documents known to the IR system:

$$\Omega = \{D_i \mid i \in \mathbf{N}\}. \quad (1.1)$$

\mathbf{N} is the set of natural numbers.

Document Set. Let \mathcal{S} be the *document set*, i. e. those n documents that form the input to the clustering process:

$$\mathcal{S} = \{D_1, D_2, \dots, D_n\} \subseteq \Omega. \quad (1.2)$$

Note: For many—but not all—applications \mathcal{S} equals Ω .

Feature Set. Let

$$\mathcal{F} = \{f_1, f_2, \dots, f_m\}, \quad (1.3)$$

with \mathcal{F} a *set of m features* and f_i an individual feature i . Each feature stands for a concrete or abstract document property.

Document Vector. Let

$$\mathbf{d}_i = (d_{i1}, d_{i2}, \dots, d_{im}), \quad (1.4)$$

with \mathbf{d}_i the *document vector* of document D_i in an m -dimensional feature space \mathcal{F} . The j th component of \mathbf{d}_i (written as d_{ij}) corresponds to the *value* or *strength* of feature f_j in document D_i . d_{ij} is usually a non-negative real number:

$$d_{ij} \in \mathbf{R}_0^+. \quad (1.5)$$

Document Feature Matrix. Let

$$H = \begin{pmatrix} \mathbf{d}_1 \\ \mathbf{d}_2 \\ \vdots \\ \mathbf{d}_n \end{pmatrix}, \quad (1.6)$$

with H the *document feature matrix*, defined by the individual vector representations \mathbf{d} of all $D \in \mathcal{S}$. The document feature matrix is usually the input to the clustering algorithm.

Document Vectorisation Function. Let τ be a function which transforms a text document into an m -dimensional vector representation in feature space \mathcal{F} :

$$\mathbf{d}_i = \tau(D_i), \text{ with } \tau : \Omega \rightarrow \mathbf{R}^m, \quad (1.7)$$

with \mathbf{R} the set of real numbers.

Feature Transformation Function. Let ϕ be a function which transforms a document vector from one feature space (\mathcal{F}_1) into another (\mathcal{F}_2), sometimes making use of additional information from the document feature matrix H :

$$\mathbf{d}'_i = \phi(\mathbf{d}_i, H), \text{ with } \phi : \mathbf{R}^{m_{\mathcal{F}_1}}, \mathbf{R}^{n \times m} \rightarrow \mathbf{R}^{m_{\mathcal{F}_2}}. \quad (1.8)$$

Document Frequency. Let

$$df(j, H) = \sum_i |\text{sgn}(h_{ij})|, \quad (1.9)$$

with the *document frequency* $df(j, H)$ the number of documents with a non-zero value for feature f_j .

ering

Metho (p2)

ring

Cluster. Let a *cluster* C_i be an subset of \mathcal{S} :

$$C_i \subseteq \mathcal{S}, \quad (1.10)$$

and let n_i be the number of objects in cluster C_i :

$$n_i = |C_i|. \quad (1.11)$$

Cluster Solution. Let

$$\mathcal{C} = \{C_1, C_2, \dots, C_k \mid C_i \subseteq \mathcal{S} \ \forall i \in 1 \dots k\}. \quad (1.12)$$

A *cluster solution* \mathcal{C} is thus defined as a set of k clusters.

Cluster Algorithm. Let

$$\mathcal{C} = \kappa(H), \quad \text{with } \kappa: \mathbf{R}^{n \times m} \rightarrow \mathcal{P}(\mathcal{S}) \quad (1.13)$$

and with κ denoting the cluster algorithm, $\mathcal{P}(\mathcal{S})$ the power set of \mathcal{S} and \mathbf{R} the set of real numbers.

Cluster Representative. Let

$$\mathbf{r}_i = (r_{i1}, r_{i2}, \dots, r_{im}), \quad (1.14)$$

with \mathbf{r}_i a *representative vector* for cluster C_i in an m -dimensional feature space \mathcal{F} .

Individual Cluster Criterion Function. An individual *criterion function* $E(C)$ measures the quality of a single cluster:

$$E: \mathcal{P}(\mathcal{S}) \rightarrow \mathbf{R}, \quad (1.15)$$

with $\mathcal{P}(X)$ the power set of X and \mathbf{R} the set of real numbers.

Overall Cluster Criterion Function. An overall *criterion function* $\Psi(\mathcal{C})$ measures the quality of an entire cluster solution:

$$\Psi: \mathcal{P}(\mathcal{P}(\mathcal{S})) \rightarrow \mathbf{R}, \quad (1.16)$$

with $\mathcal{P}(\mathcal{P}(\mathcal{S}))$ the set of all possible cluster solutions.

Type and Token. Within documents it is common to refer to word *types* and word *tokens*. The former refer abstractly to features in a document or a corpus, while the latter refer to individual occurrences. Formally speaking, the tokens of a document are a *bag* (which allows multiple occurrences of the same element). The types are the *set* created by eliminating all duplicates from the token bag.

Recall and Precision. In IR two widespread performance measures are defined by the set of documents in a collection that are *relevant* to a particular query (\mathcal{A}) and those documents that are actually *retrieved* by the system (\mathcal{B}):

$$\text{Recall } (R) = \frac{|\mathcal{A} \cap \mathcal{B}|}{|\mathcal{B}|}, \quad (1.17)$$

$$\text{Precision } (P) = \frac{|\mathcal{A} \cap \mathcal{B}|}{|\mathcal{A}|}. \quad (1.18)$$

Their weighted arithmetic mean, the so-called *F-Measure* is also used frequently (see Equation 2.77 for an example).

Latent Dirichlet Allocation/Analysis (p1)

We describe *latent Dirichlet allocation* (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying set of topics. Each topic is, in turn, modeled as an infinite mixture over an underlying set of topic probabilities. In the context of text modeling, the topic probabilities provide an explicit representation of a document. We present efficient approximate inference techniques based on variational methods and an EM algorithm for empirical Bayes parameter estimation. We report results in document modeling, text classification, and collaborative filtering, comparing to a mixture of unigrams model and the probabilistic LSI model.

Latent Dirichlet allocation (LDA) is a generative probabilistic model of a corpus. The basic idea is that documents are represented as random mixtures over latent topics, where each topic is characterized by a distribution over words.¹

LDA assumes the following generative process for each document w in a corpus D :

1. Choose $N \sim \text{Poisson}(\xi)$.
2. Choose $\theta \sim \text{Dir}(\alpha)$.
3. For each of the N words w_n :
 - (a) Choose a topic $z_n \sim \text{Multinomial}(\theta)$.
 - (b) Choose a word w_n from $p(w_n | z_n, \beta)$, a multinomial probability conditioned on the topic z_n .

Several simplifying assumptions are made in this basic model, some of which we remove in subsequent sections. First, the dimensionality k of the Dirichlet distribution (and thus the dimensionality of the topic variable z) is assumed known and fixed. Second, the word probabilities are parameterized by a $k \times V$ matrix β where $\beta_{lj} = p(w^j = 1 | z^l = 1)$, which for now we treat as a fixed quantity that is to be estimated. Finally, the Poisson assumption is not critical to anything that follows and more realistic document length distributions can be used as needed. Furthermore, note that N is independent of all the other data generating variables (θ and z). It is thus an ancillary variable and we will generally ignore its randomness in the subsequent development.

Latent Dirichlet Allocation/Analysis (p2)

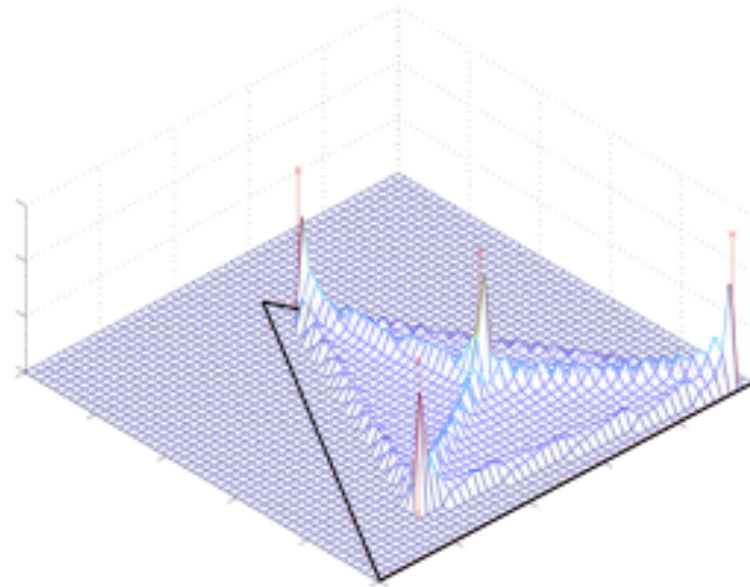


Figure 2: An example density on unigram distributions $p(w|\theta, \beta)$ under LDA for three words and four topics. The triangle embedded in the x-y plane is the 2-D simplex representing all possible multinomial distributions over three words. Each of the vertices of the triangle corresponds to a deterministic distribution that assigns probability one to one of the words; the midpoint of an edge gives probability 0.5 to two of the words; and the centroid of the triangle is the uniform distribution over all three words. The four points marked with an x are the locations of the multinomial distributions $p(w|z)$ for each of the four topics, and the surface shown on top of the simplex is an example of a density over the $(V - 1)$ -simplex (multinomial distributions of words) given by LDA.

Latent Dirichlet Allocation/Analysis (p3)

| “Arts” | “Budgets” | “Children” | “Education” |
|---------|------------|------------|-------------|
| NEW | MILLION | CHILDREN | SCHOOL |
| FILM | TAX | WOMEN | STUDENTS |
| SHOW | PROGRAM | PEOPLE | SCHOOLS |
| MUSIC | BUDGET | CHILD | EDUCATION |
| MOVIE | BILLION | YEARS | TEACHERS |
| PLAY | FEDERAL | FAMILIES | HIGH |
| MUSICAL | YEAR | WORK | PUBLIC |
| BEST | SPENDING | PARENTS | TEACHER |
| ACTOR | NEW | SAYS | BENNETT |
| FIRST | STATE | FAMILY | MANIGAT |
| YORK | PLAN | WELFARE | NAMPHY |
| OPERA | MONEY | MEN | STATE |
| THEATER | PROGRAMS | PERCENT | PRESIDENT |
| ACTRESS | GOVERNMENT | CARE | ELEMENTARY |
| LOVE | CONGRESS | LIFE | HAITI |

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Figure 8: An example article from the AP corpus. Each color codes a different factor from which the word is putatively generated.

Latent Dirichlet Allocation/Analysis for proteins/genes/cancer targets/clinical data

- Just as in previous slide, topics were clustered around “Arts,” “Budgets,” “Children,” and “Education.”
 - Cluster around patient type
 - Cluster around drug type
 - Cluster around drug-interaction type
 - Cluster around phenotype/genotype
 - Cluster around structure/sequence types

Selected references:

- A. A. Cohen, et al. Response to a Drug Dynamic Proteomics of Individual Cancer Cells in DOI: 10.1126/Science.1160165, 1511 (2008); 322
- Ferrenberg AM, Swendsen RH. Optimized Monte-Carlo data-analysis Physical Review Letters volume: 63 issue: 12 pages: 1195-1198 Sep 18 1989.
- Document Clustering in Large German Corpora Using Natural Language Processing, Richard Forster (2006), University of Zurich
- Latent Dirichlet Allocation, Blei, Ng, and Jordan, Journal of Machine Learning Research 3 (2003) 993-1022
- LobeLink.com
- A. Ben-Hur and W.S. Noble "Kernel methods for predicting protein-protein interactions." Bioinformatics (Proceedings of the Intelligent Systems for Molecular Biology Conference). 21:38-46, 2005.
- D.-S. Huang, K. Li, and G.W. Irwin (Eds.): ICIC 2006, LNBI 4115, pp. 514–524, 2006. "Prediction of Protein Complexes Based on Protein Interaction Data and Functional Annotation Data Using Kernel Methods." Springer-Verlag Berlin Heidelberg 2006.
- B. Schwikowski, P. Uetz, S. Fields. A network of protein-protein interactions in yeast. Nature Biotechnology 18 (12): 1257-1261 Dec. 2000.
- K. Tsuda, H.J. Shin, B. Schölkopf. Fast protein classification with multiple networks. Bioinformatics 2005 21(Suppl 2):ii59-ii65.
- J.-P. Vert, K. Tsuda and B. Schölkopf, A primer on kernel methods, in Kernel Methods in Computational Biology, B. Schölkopf, K. Tsuda and J.-P. Vert (Eds.), MIT Press, p.35-70, 2004.