

Machine Learning for Finance

Neal Parikh

Cornell University

Spring 2018

k-means

k -means

given $\mathcal{D} = \{x_1, \dots, x_N\}$, $x_i \in \mathbf{R}^n$, group data into a few 'clusters'

- 1 randomly initialize cluster centroids $\mu_1, \dots, \mu_k \in \mathbf{R}^n$
- 2 repeat until convergence
 - 1 find cluster assignment for x_i

$$c_i := \operatorname{argmin}_j \|x_i - \mu_j\|_2^2$$

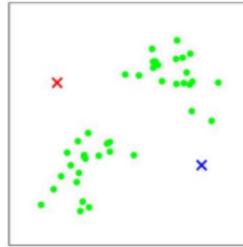
- 2 recompute cluster centroids using these assignments

$$\mu_j := \frac{\sum_{i=1}^N [c_i = j] x_i}{\sum_{i=1}^N [c_i = j]}$$

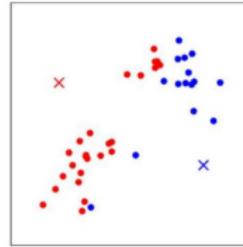
k -means



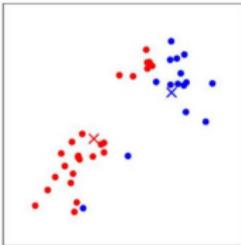
(a)



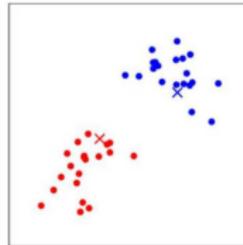
(b)



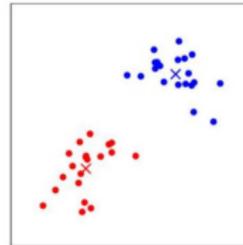
(c)



(d)



(e)



(f)

Alternating minimization

- k -means can also be viewed as alternating minimization on the (biconvex) ‘distortion function’

$$J(c, \mu) = \sum_{i=1}^N \|x_i - \mu_{c_i}\|_2^2$$

- results dependent on initialization, so do random restarts and pick one with lowest distortion
- can also derive k -means as a limit of a probabilistic model

Mixture models and the EM algorithm

Mixture of Gaussians

- probabilistic model for clustering / density estimation
- consider data $\mathcal{D} = \{x_1, \dots, x_N\}$
- generative model

$$z \sim \text{Multinomial}(\phi)$$
$$x | z = k \sim \text{N}(\mu_k, \Sigma_k)$$

- *i.e.*, each x_i generated by sampling a **unobserved** (hidden, latent) $z_i \in [K]$ and then drawing x_i from the corresponding Gaussian
- presence of these latent variables is the key new wrinkle
- model parameters are ϕ, μ_k, Σ_k

Maximum likelihood estimation

- model parameters are ϕ, μ_k, Σ_k
- as usual, write down likelihood for $w = (\phi, \mu_k, \Sigma_k)$

$$\begin{aligned}\ell(w) &= \sum_{i=1}^N \log p(x_i; w) \\ &= \sum_{i=1}^N \log \sum_{z_i=1}^K p(x_i | z_i) p(z_i)\end{aligned}$$

- this function is *nonconvex* due to sum over values of z_i

Maximum likelihood estimation

- if z_i were known, problem is easy and becomes

$$\ell(w) = \sum_{i=1}^N \log p(x_i | z_i) + \sum_{i=1}^N \log p(z_i)$$

- maximizing with respect to ϕ, μ, Σ gives

$$\phi_j = \frac{1}{N} \sum_{i=1}^N [z_i = j], \quad \mu_j = \frac{\sum_{i=1}^N [z_i = j] x_i}{\sum_{i=1}^N [z_i = j]}$$

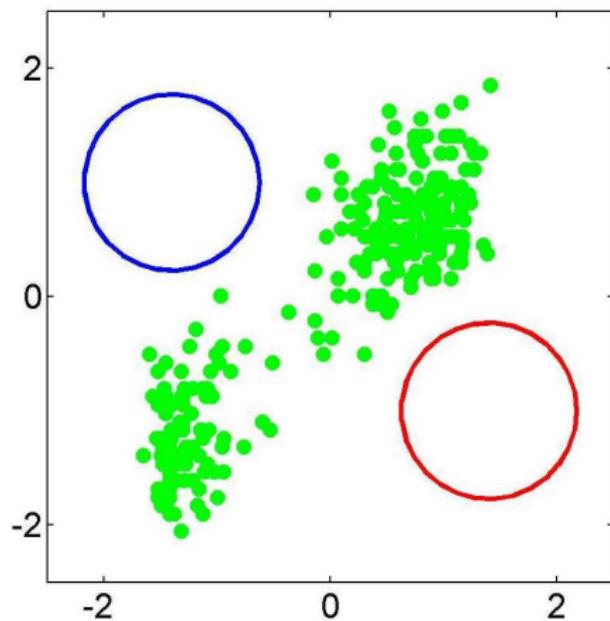
similar expression for Σ

- *i.e.*, if z_i were known, nearly identical to maximum likelihood estimates in GDA (with z_i 's as class labels)

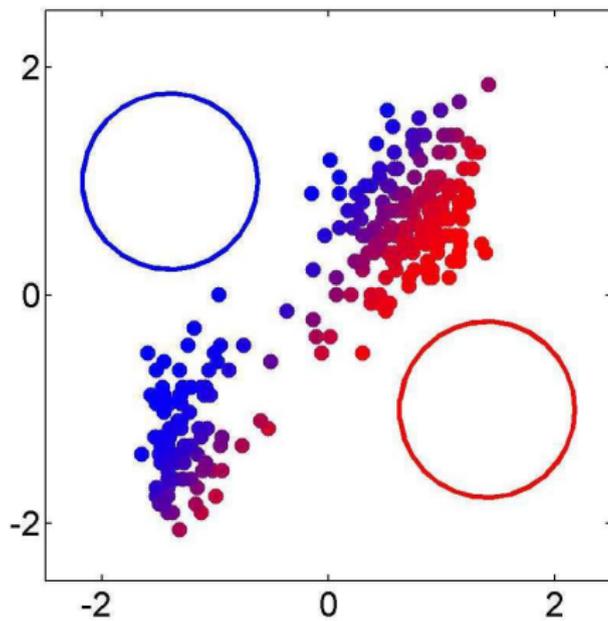
EM algorithm

- **idea:** iteratively guess the z_i and then use formulas above:
 - ① E-step: compute $\rho_{ij} = p(z_i = j | x_i; \theta, \mu, \Sigma)$
 - ② M-step: use formulas above with ρ_{ij} in place of $[z_i = j]$
- E-step computes posterior probability of z_i 's, given data and current setting of parameters; 'soft guesses' for values of z_i
- M-step is maximum likelihood estimation, but there is uncertainty around the value of the z_i and that's incorporated in estimates
- a 'soft' version of k -means in this context

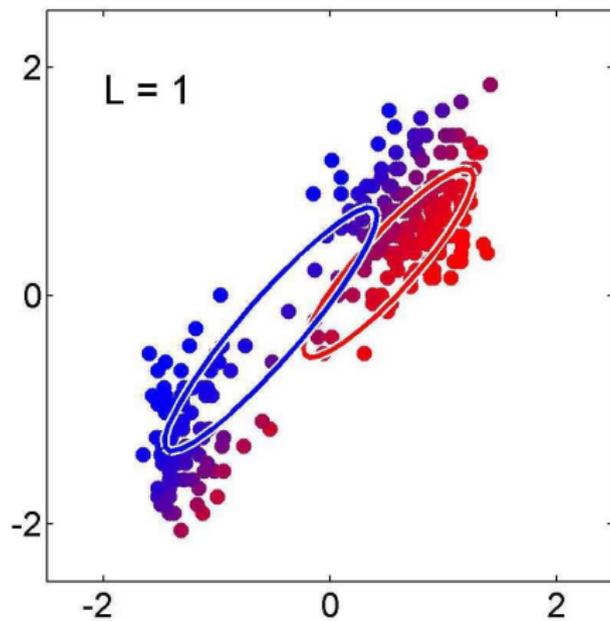
Gaussian mixture model



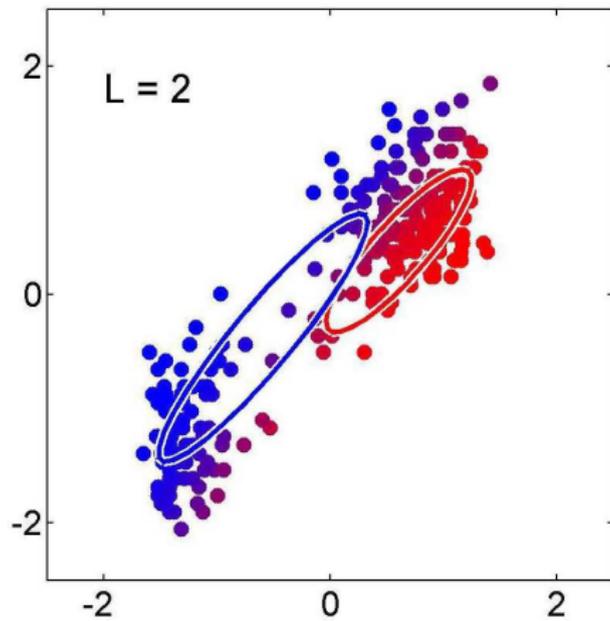
Gaussian mixture model



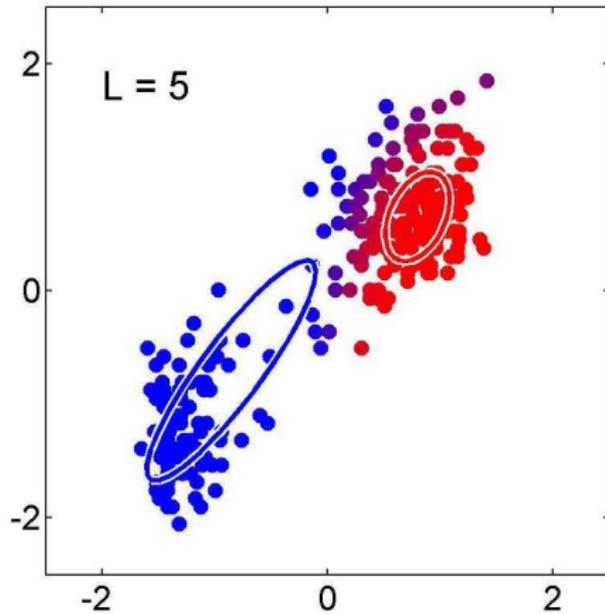
Gaussian mixture model



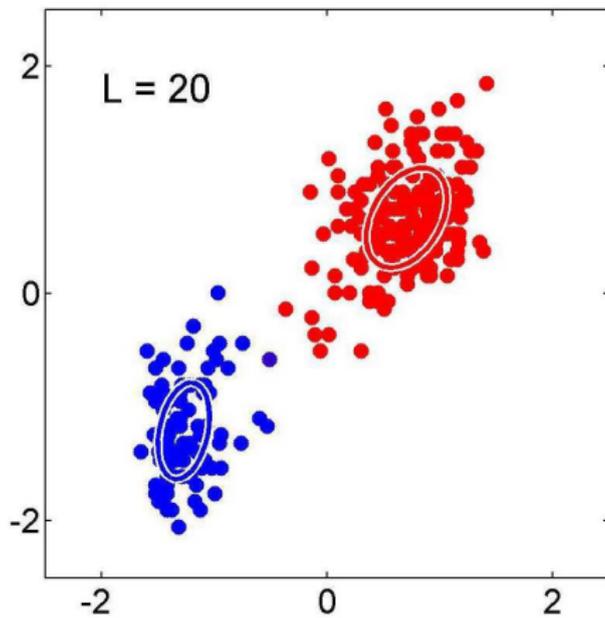
Gaussian mixture model



Gaussian mixture model



Gaussian mixture model



EM algorithm

- in general, EM algorithm is standard approach to maximum likelihood estimation with latent variable models
- data $\mathcal{D} = \{x_1, \dots, x_N\}$
- want to fit model $p(x, z)$ with z hidden
- likelihood is given by

$$\ell(w) = \sum_{i=1}^N \log p(x; w) = \sum_{i=1}^N \log \sum_z p(x, z; w)$$

- often the case that maximum likelihood estimation of x would be easy if z were known, so alternate the two steps

EM algorithm

- EM algorithm can be motivated and analyzed in various ways
- iteratively lower bound ℓ , then maximize that lower bound
- for each i , let q_i be a distribution over z 's

$$\begin{aligned}\sum_{i=1}^N \log p(x_i) &= \sum_{i=1}^N \log \sum_{z_i} p(x_i, z_i) \\ &= \sum_{i=1}^N \log \sum_{z_i} q_i(z_i) \frac{p(x_i, z_i)}{q_i(z_i)} \\ &\geq \sum_{i=1}^N \sum_{z_i} q_i(z_i) \log \frac{p(x_i, z_i)}{q_i(z_i)}\end{aligned}$$

by Jensen's inequality

EM algorithm

- previous formula gives lower bound for *any* q_i ; ideally, have the lower bound be tight (inequality holds with equality) for current value of w
- can show that this is the case when $q_i(z_i) = p(z_i | x_i; w)$; it suffices that $q_i(z_i) \propto p(x_i, z_i; w)$, so

$$\begin{aligned}q_i(z_i) &= \frac{p(x_i, z_i; w)}{\sum_z p(x_i, z; w)} \\ &= \frac{p(x_i, z_i; w)}{p(x_i; w)} \\ &= p(z_i | x_i; w)\end{aligned}$$

- E-step (above): obtain lower bound (has form of an expectation) on ℓ
- M-step: maximize this lower bound with respect to w

EM algorithm

- can show that this algorithm converges because it monotonically improves the log likelihood
- *i.e.*, can show $\ell(w^k) < \ell(w^{k+1})$

EM algorithm

- EM algorithm can also be viewed as coordinate ascent on

$$J(q, w) = \sum_{i=1}^N \sum_{z_i} q_i(z_i) \log \frac{p(x_i, z_i; w)}{q_i(z_i)}$$

- E-step: maximization with respect to q
- M-step: maximization with respect to w

Factor analysis

- fitting Gaussian mixture model to data $x_1, \dots, x_N \in \mathbf{R}^n$ assumes enough data ($N \gg n$) to discern this structure
- if $n \gg N$, cannot even fit a single Gaussian
- here, data points span low-dimensional subspace of \mathbf{R}^n , so MLEs of the parameters result in degenerate Gaussian (singular covariance matrix) that puts all mass in affine space spanned by the data
- consider models that explicitly handle low rank structure

Factor analysis

- consider generative model $p(x, z)$ given by

$$\begin{aligned}z &\sim \text{N}(0, I) \\x | z &\sim \text{N}(\mu + \Lambda z, \Psi)\end{aligned}$$

where $\mu \in \mathbf{R}^n$, $\Lambda \in \mathbf{R}^{n \times k}$, $\Psi \in \mathbf{R}^{n \times n}$ diagonal

- x observed, z latent
- low dimensional structure: $k < n$, i.e., data is generated by affine transformation of k -dimensional Gaussian (plus noise)

Factor analysis

- $p(x, z)$ is Gaussian, and need to find its mean and covariance from the generative model
- ideally, would want to maximize \log (marginal) likelihood of data, using marginal distribution of x , but this function is hard to optimize
- so, use EM
 - E-step: compute $q_i(z_i) = p(z_i | x_i)$ (also Gaussian)
 - M-step: maximize lower bound

$$\sum_{i=1}^N \int_{z_i} q_i(z_i) \log \frac{p(x_i, z_i)}{q_i(z_i)} dz_i$$

- involves some messy algebra, but can obtain closed form solutions for all these subproblems (matrix computations)

Principal components analysis

Dimensionality reduction

- model data $x \in \mathbf{R}^n$ as approximately lying in some k -dimensional subspace, with $k \ll n$
- has many different use cases
 - data compression
 - data visualization
 - noise reduction
 - preprocessing for supervised learning
 - feature discovery
 - structure discovery

Principal components analysis

- let $\mathcal{D} = \{x_1, \dots, x_N\}$, with $x_i \in \mathbf{R}^n$, $n < N$
- rescale data to have mean zero and unit variance
 - ① replace x_i with $x_i - (1/N) \sum_i x_i$
 - ② replace x_i^j with x_i^j / σ_j , where $\sigma_j = (1/N) \sum_i (x_i^j)^2$

Principal components analysis

- several ways to motivate PCA
- select directions on which to project points to maximize variance
- compute top k eigenvectors of empirical covariance matrix
- pick k -dimensional basis so approximation error of projecting data onto it is minimized

Principal components analysis

- given \mathcal{D} , find unit vector u such that projection of \mathcal{D} onto direction u has maximum variance
- length of projection of x onto u is $x^T u$, so solve

$$\begin{aligned} & \text{maximize} && (1/N) \sum_{i=1}^N (x_i^T u)^2 \\ & \text{subject to} && \|u\|_2 = 1 \end{aligned}$$

- objective can be rewritten as quadratic form $u^T \Sigma u$, where

$$\Sigma = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$$

is **empirical covariance matrix** of (preprocessed) \mathcal{D}

- so solution of problem above is computing principal eigenvector of Σ

Principal components analysis

- in general, find top k eigenvectors u_1, \dots, u_k of Σ
- these give an orthonormal basis for \mathbf{R}^k
- compute rank k approximation to x_i as

$$y_i = (u_1^T x_i, \dots, u_k^T x_i)$$

- choice of u_i maximizes $\sum_i \|y_i\|_2^2$

Topic models

Topic models

- topic models: methods for automatically organizing, understanding, searching, and summarizing large electronic archives
 - discover hidden themes that pervade the collection
 - annotate documents with those themes
 - use annotations to organize, summarize, and search texts
- unsupervised generative latent variable models of document structure
- originally introduced by Blei, Ng, and Jordan (2003); much subsequent work by Blei and collaborators, among many others

Topic models

- idea: documents composed of multiple topics
- each topic is a distribution over words
- each document is a mixture of corpus-wide topics
- each word is drawn from one of these topics

Latent Dirichlet allocation

- generative model $p(\theta, z, w \mid \alpha, \beta)$

$$\theta \sim \text{Dirichlet}(\alpha)$$

$$z_n \sim \text{Multinomial}(\theta), \quad n = 1, \dots, N$$

$$w_n \sim \text{Multinomial}(\beta_{z_n}), \quad n = 1, \dots, N$$

- estimate parameters by, e.g., maximizing log-likelihood

$$\ell(\alpha, \beta) = \sum_{d=1}^D \log p(\mathbf{w} \mid \alpha, \beta)$$

where $\mathbf{w}_1, \dots, \mathbf{w}_D$ are documents (training set)

- want to compute posterior of latent variables
- conceptually, use EM (but need approximations here)

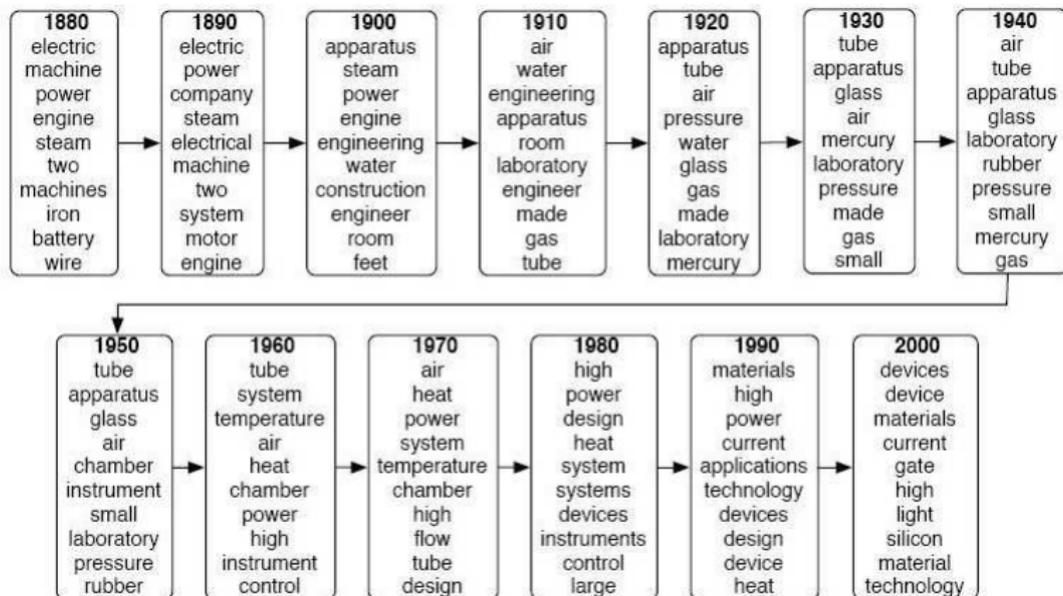
A 100 topic model of Science 1980-2000

sound speech acoustic language sounds	quantum laser light optical electron	brain memory human visual cognitive	computer data information problem computers	ice climate ocean sea temperature
stars universe galaxies astronomers star	research national science new funding	materials organic molecules molecular polymer	fossil species evolution birds evolutionary	volcanic years fig deposits rocks

Variants and applications

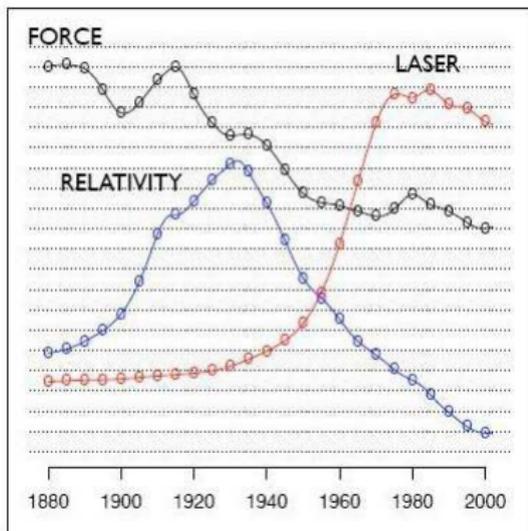
- finding similar documents
- measuring scholarly impact (detect influential articles)
- discover evolution of topics over time
- discover correlations between topics
- annotate images with captions
- characterizing political decisions
- organize and browse large document collections

Model Evolution of Topics over Time

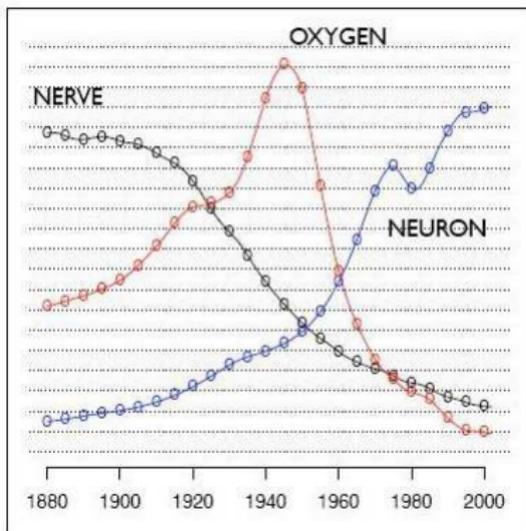


Visualizing Trends Within Topics

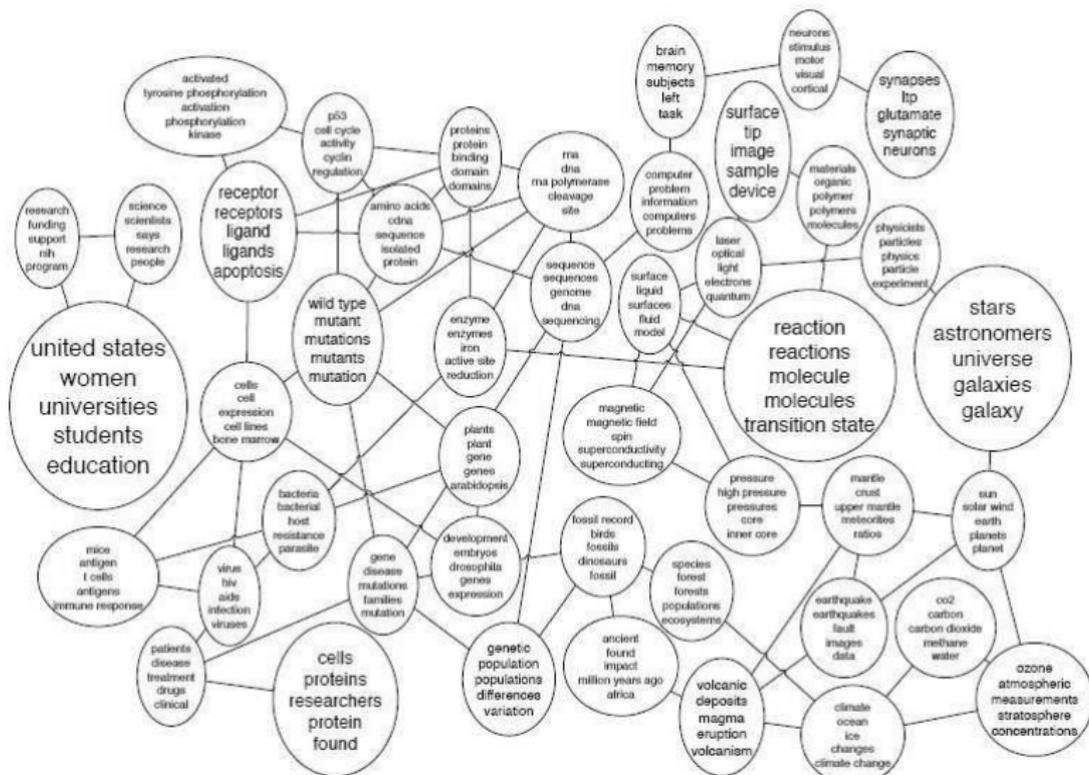
"Theoretical Physics"



"Neuroscience"

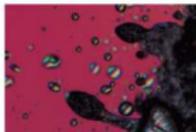
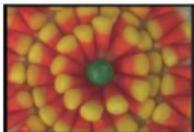


Model Connections Between Topics

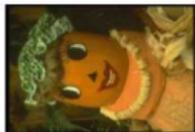


Matching Words and Pictures

Candy



Sunset



People
& Fish



Matching Words and Pictures



True caption
market people

Corr-LDA
people market pattern textile display



True caption
scotland water

Corr-LDA
scotland water flowers hills tree



True caption
bridge sky water

Corr-LDA
sky water buildings people mountain



True caption
sky tree water

Corr-LDA
tree water sky people buildings



True caption
birds tree

Corr-LDA
birds nest leaves branch tree



True caption
fish reefs water

Corr-LDA
fish water ocean tree coral



True caption
mountain sky tree water

Corr-LDA
sky water tree mountain people



True caption
clouds jet plane

Corr-LDA
sky plane jet mountain clouds