

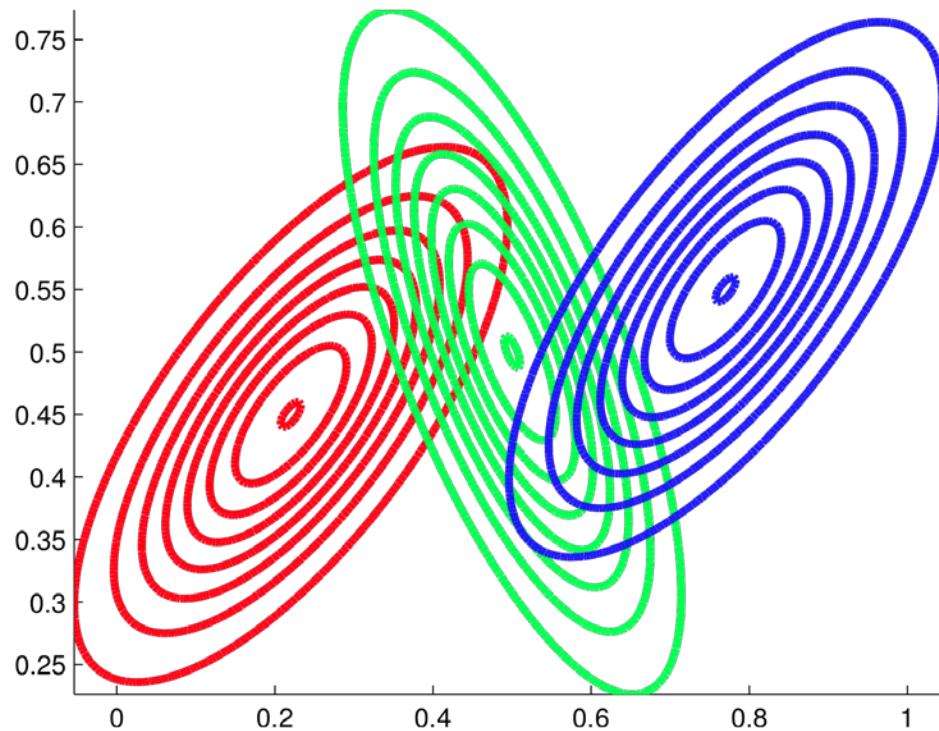
Introduction to Machine Learning

Brown University CSCI 1950-F, Spring 2011
Prof. Erik Sudderth

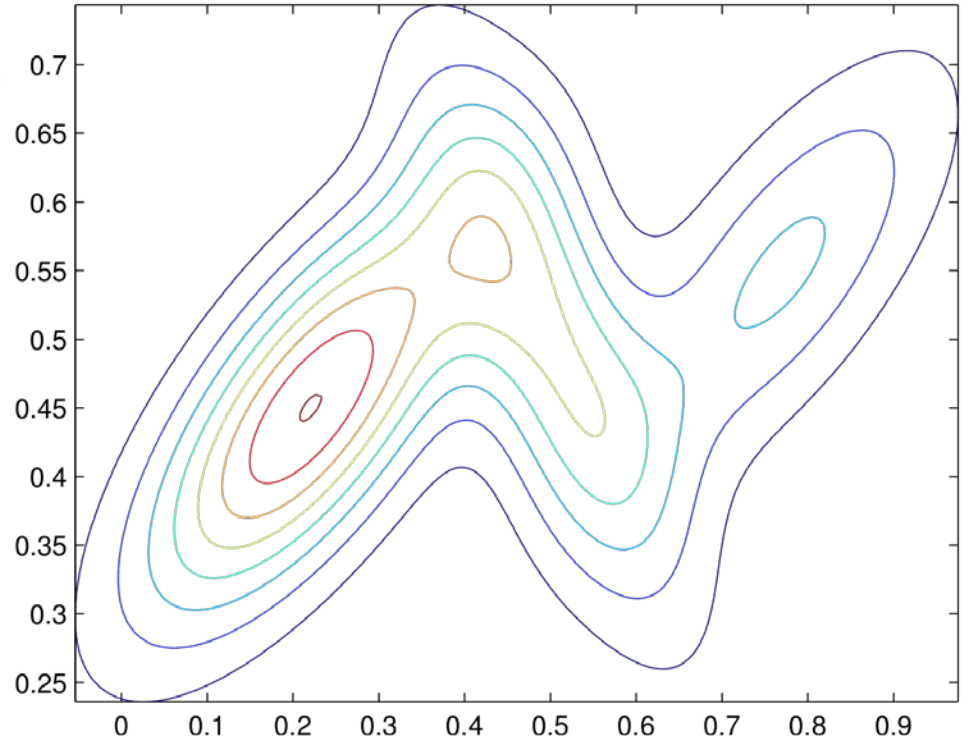
Lecture 22: Topic Models,
Hidden Markov Models (HMMs)

Many figures courtesy Kevin Murphy's textbook,
Machine Learning: A Probabilistic Perspective

Gaussian Mixture Models

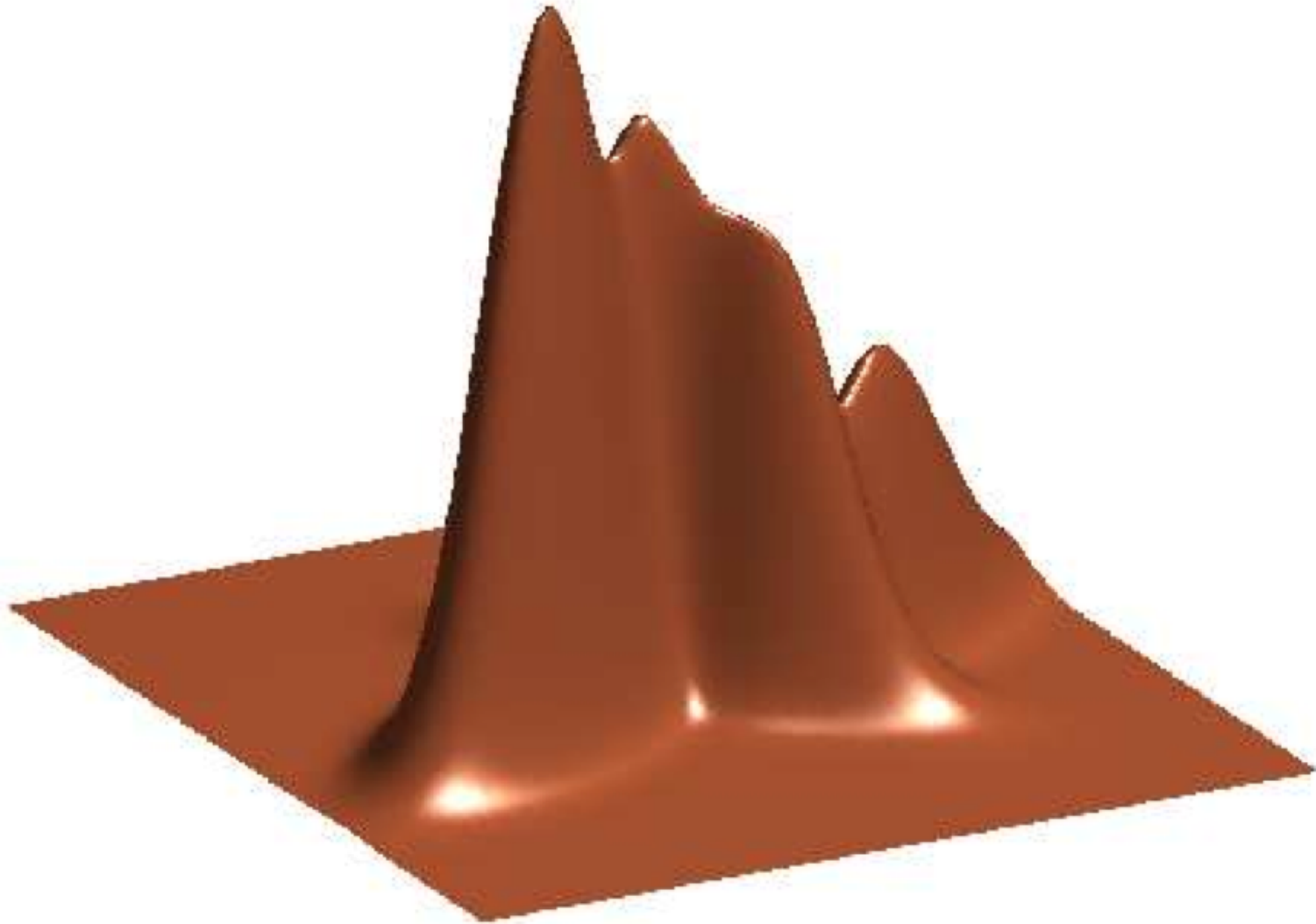


*Mixture of 3 Gaussian
Distributions in 2D*



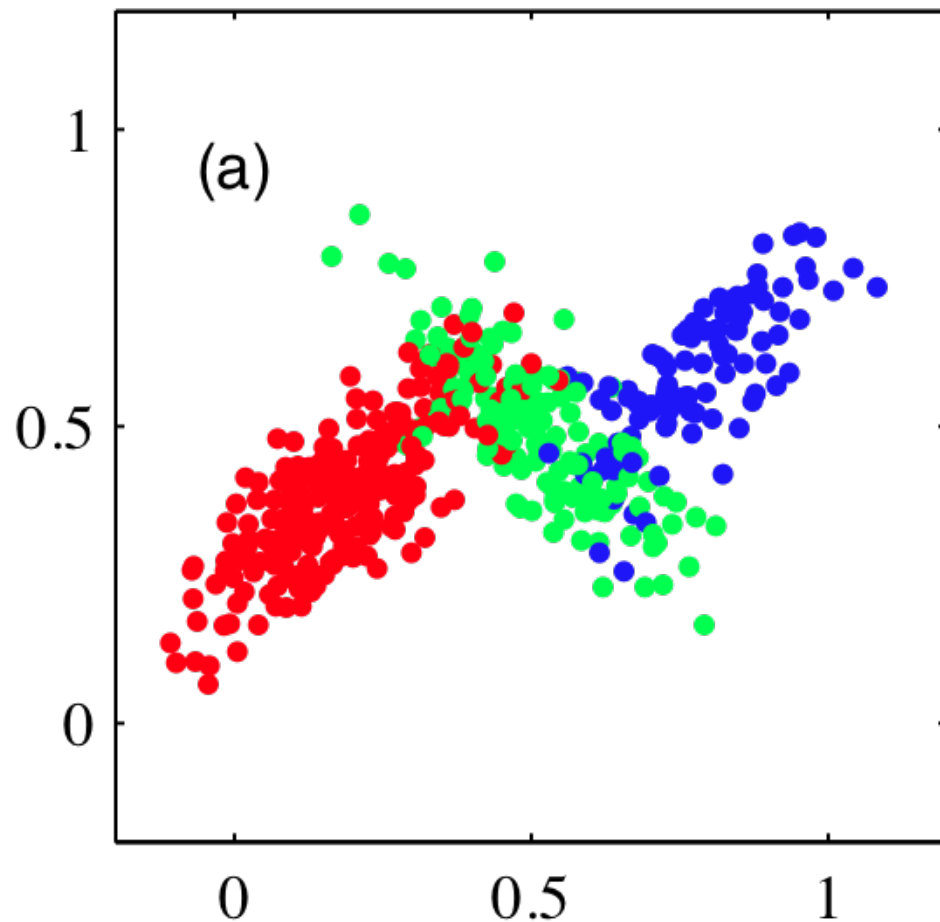
*Contour Plot of Joint Density,
Marginalizing Cluster Assignments*

Gaussian Mixture Models

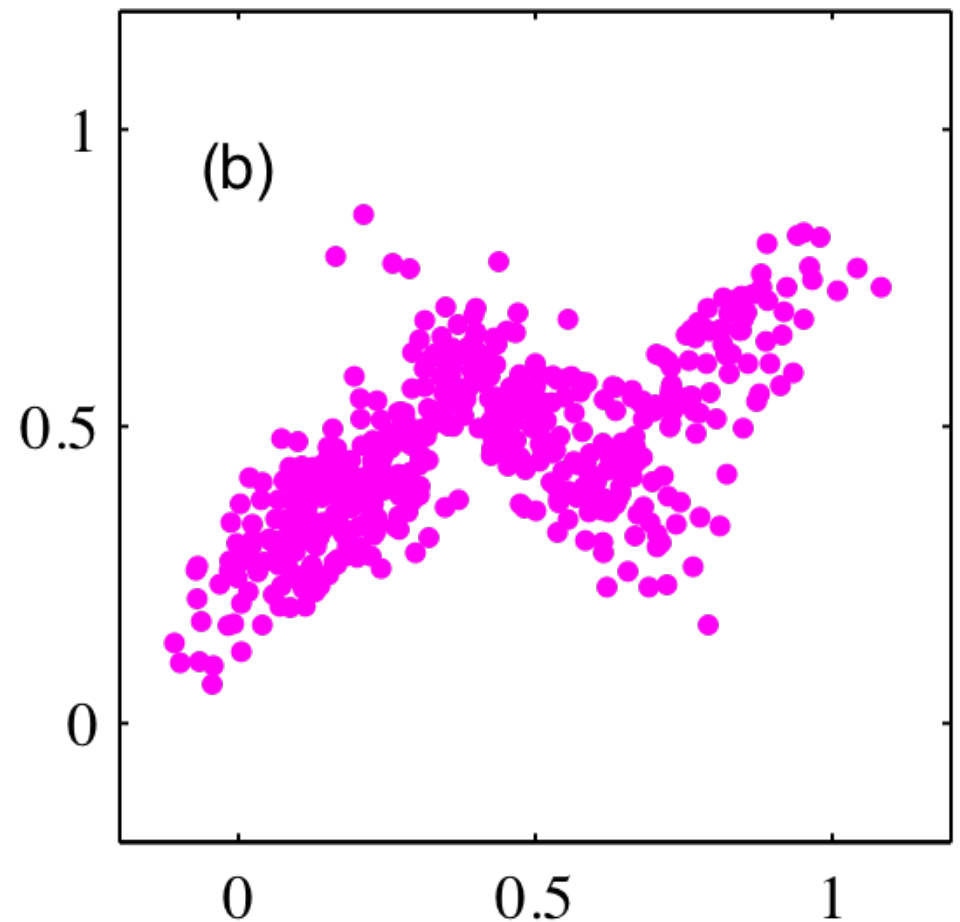


*Surface Plot of Joint Density,
Marginalizing Cluster Assignments*

Fitting Gaussian Mixtures



*Complete Data Labeled
by True Cluster Assignments*



*Incomplete Data:
Points to be Clustered*

Collections of Mixture Models

- Standard mixture models assume a single, “flat” dataset
- Many applications involve multiple related, but distinct, “groups” of data
 - Multiple documents in a text corpus
 - Multiple images in a photo repository
 - Multiple users with their own spam filtering decisions
 - Multiple hospitals in a clinical trial
 - Multiple companies in a financial market
- How can we jointly model this data?
- Lumping into single large dataset ignores group differences
- Modeling groups independently can be ineffective, especially when limited data about any one group
- Hierarchical Bayesian models share between groups

Multiple Gaussian Mixtures

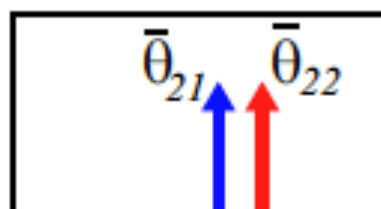
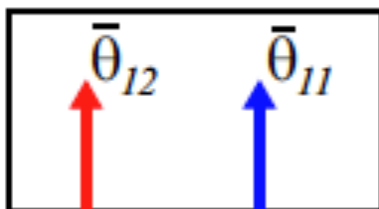


Use data to learn set of shared mixture identities, and their frequencies across groups

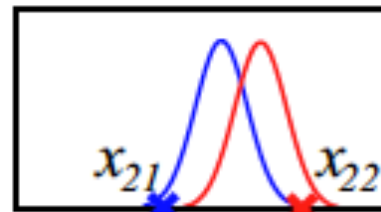
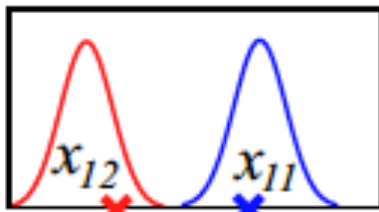


Each group has its own weight on shared mixture parameters

$$\pi_j \sim \text{Dir}(\alpha)$$



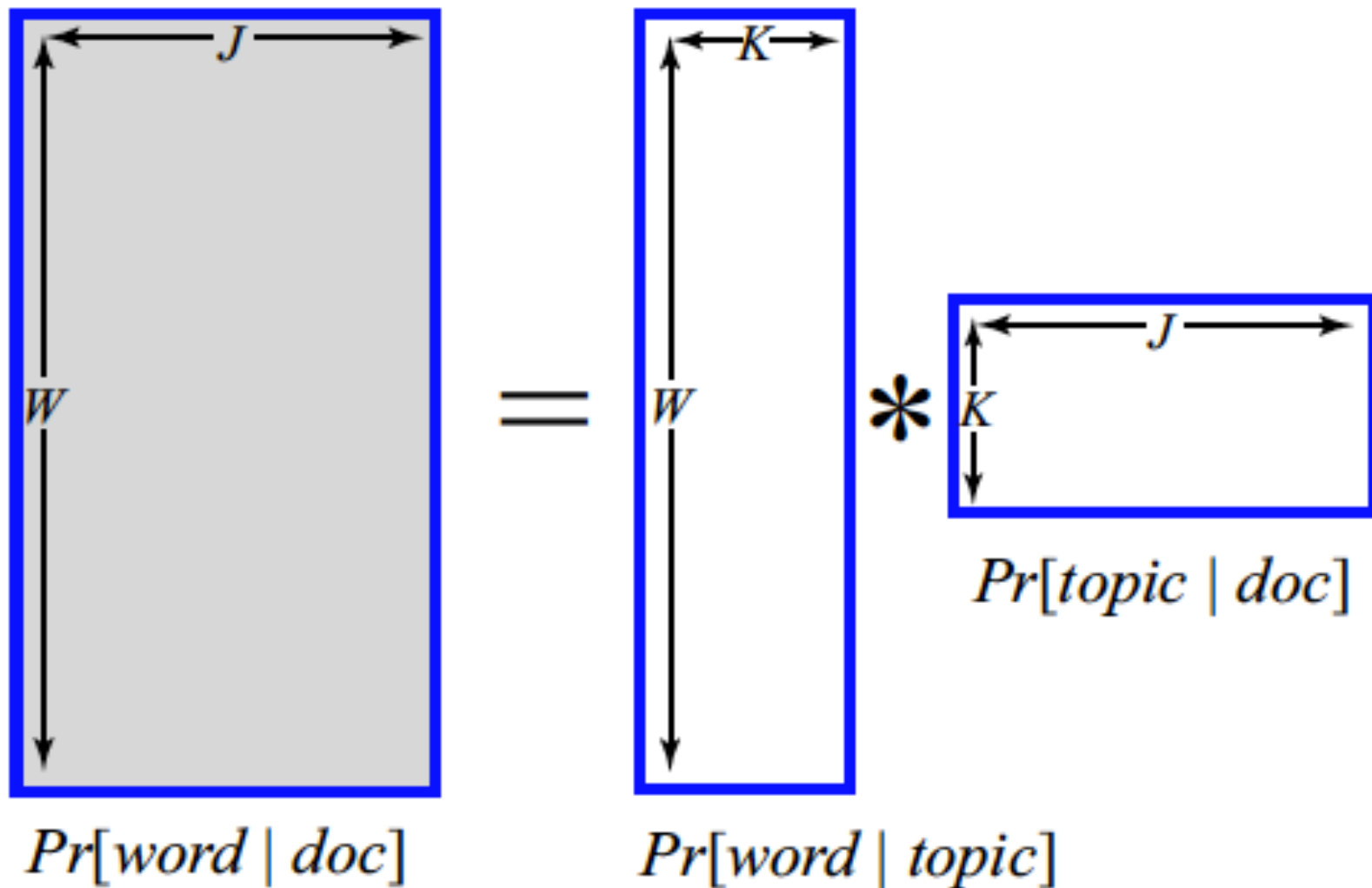
Each observation comes from some mixture component



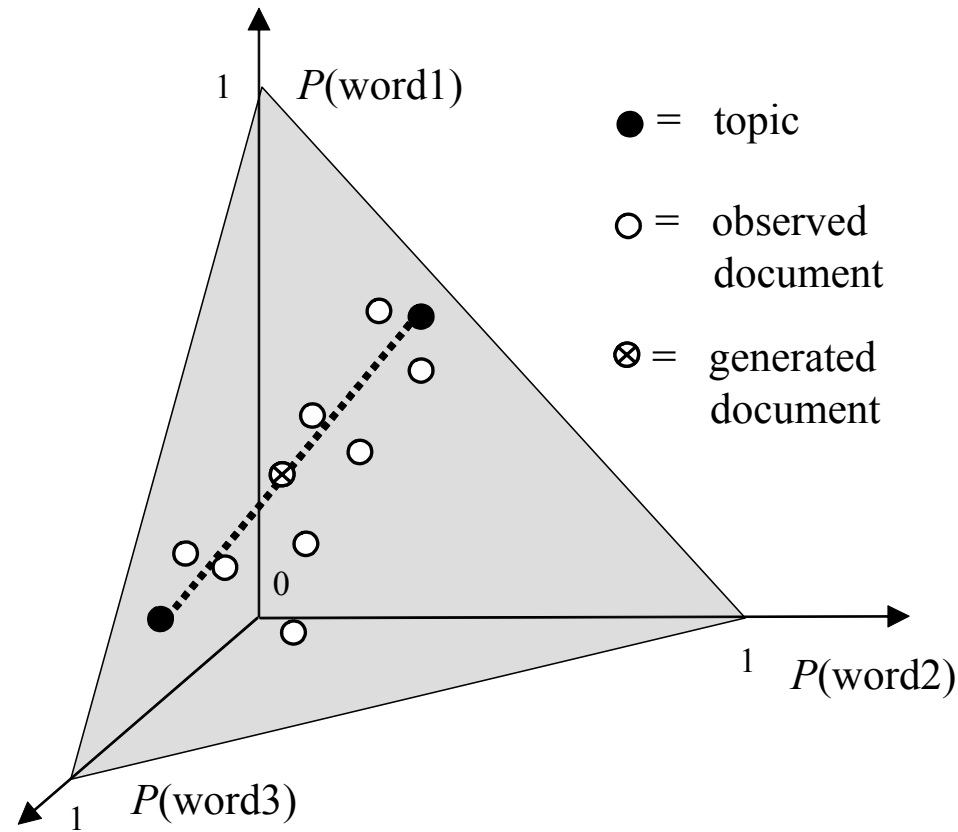
J=2 Groups of Data

$$p(x_{ji} \mid \pi_j, \theta_1, \dots, \theta_K) = \sum_{k=1}^K \pi_{jk} f(x_{ji} \mid \theta_k)$$

Probabilistic Topic Models: Multiple Multinomial Mixtures



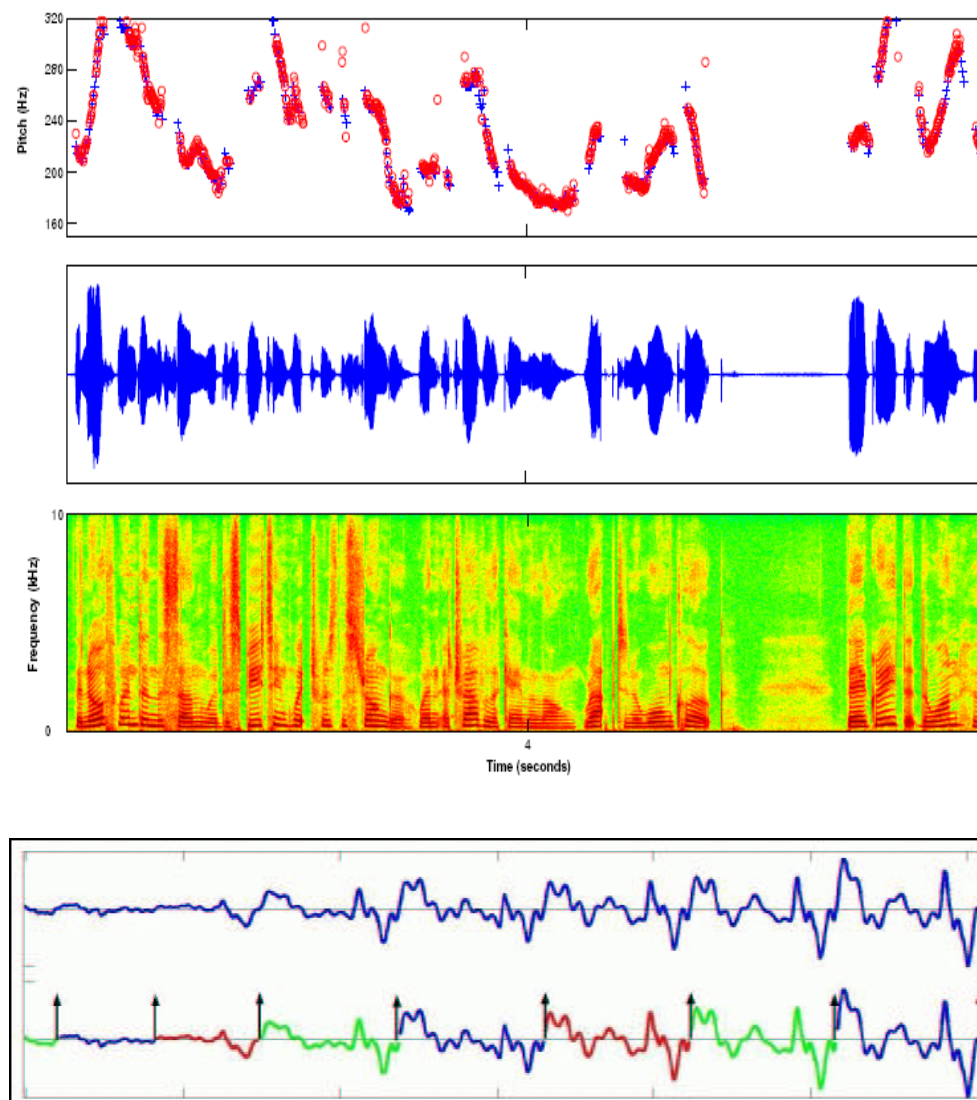
Geometry of Topic Models



- Documents are multinomial distributions over some predefined vocabulary of (tens of thousands) of words
- Topics are multinomial distributions on same vocabulary
- Generative model: Each document is (nearly) a convex combination of the topic distributions

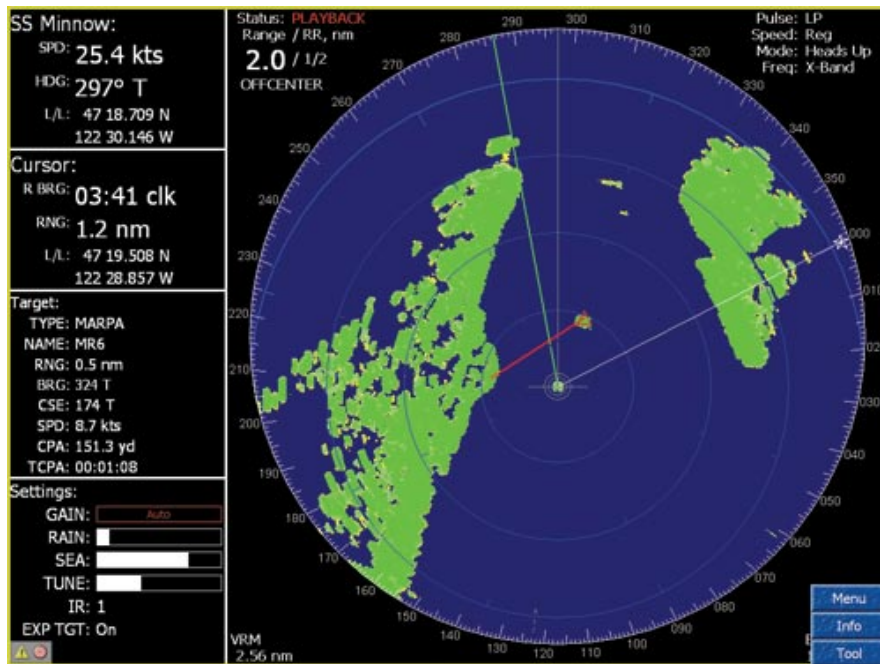
Speech Recognition

- Given an audio waveform, robustly extract & recognize any spoken words
- Statistical models can be used to
 - Provide greater robustness to noise
 - Adapt to accent of different speakers
 - Learn from training

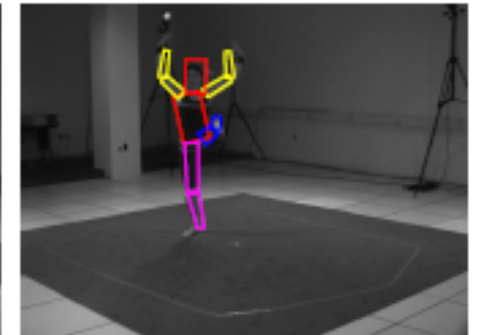
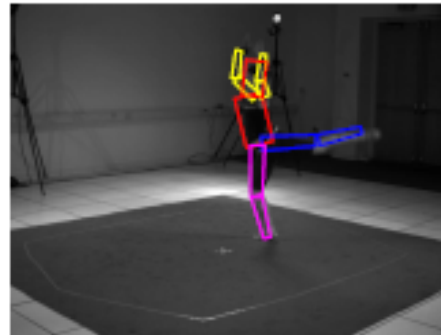
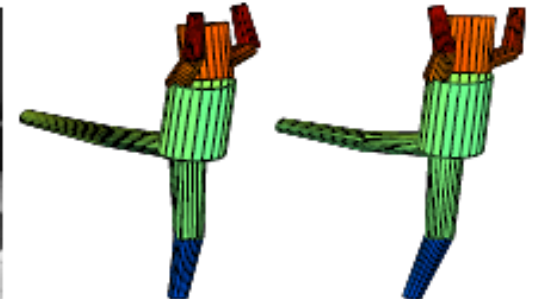
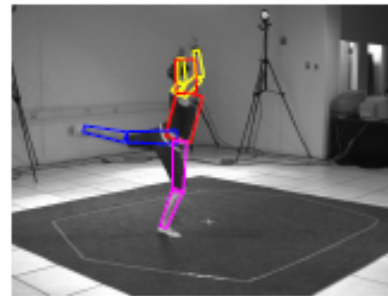


S. Roweis, 2004

Target Tracking



*Radar-based tracking
of multiple targets*

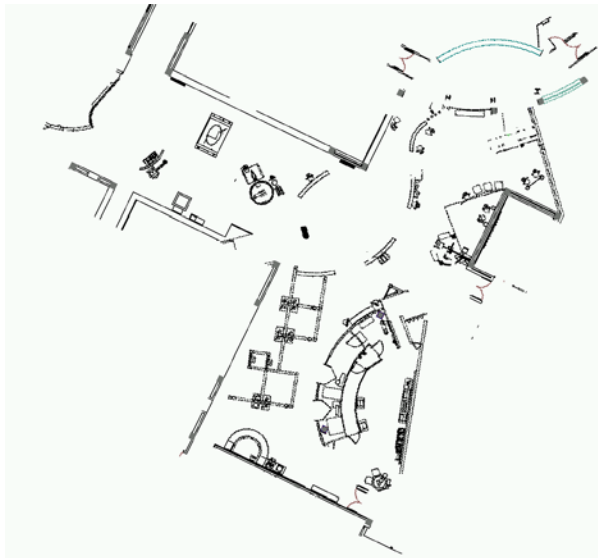


*Visual tracking of
articulated objects*
(L. Sigal et. al., 2006)

- Estimate motion of targets in 3D world from indirect, potentially noisy measurements

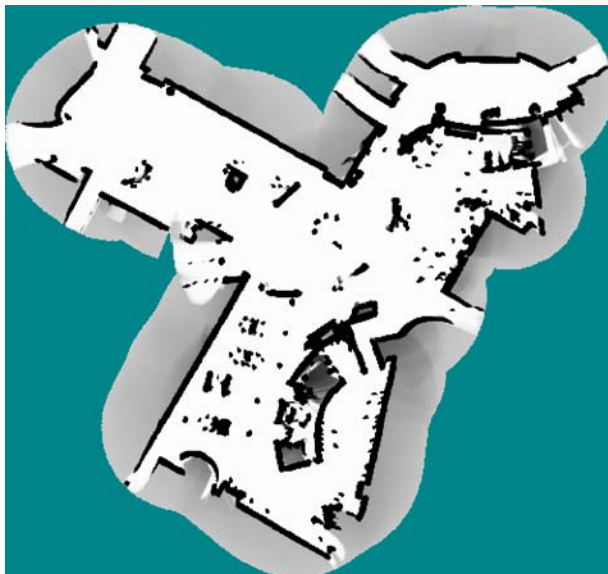
Robot Navigation: *SLAM*

Simultaneous Localization and Mapping



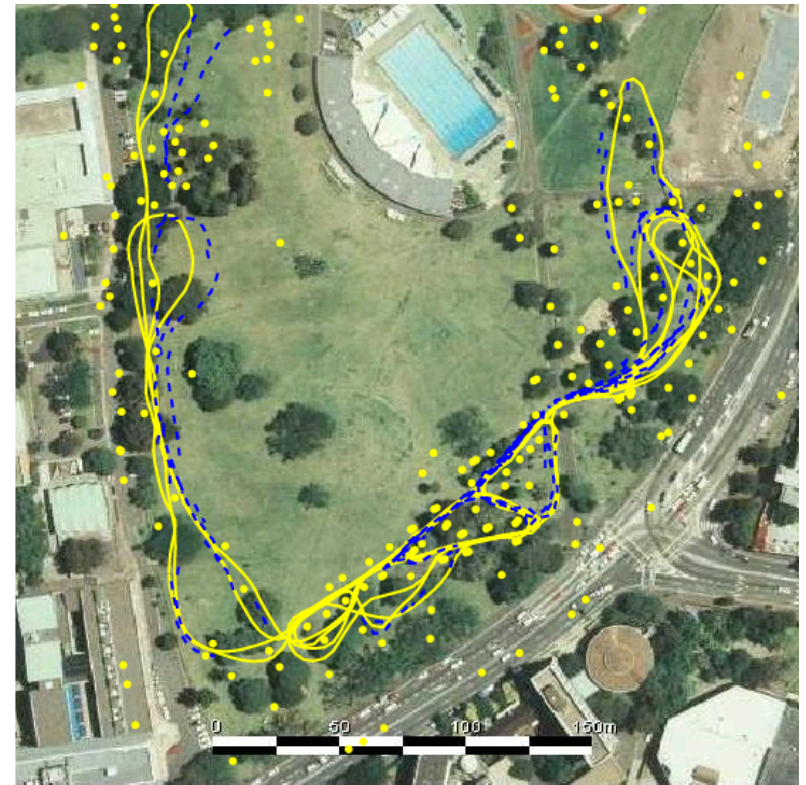
*CAD
Map*

(S. Thrun,
San Jose Tech Museum)



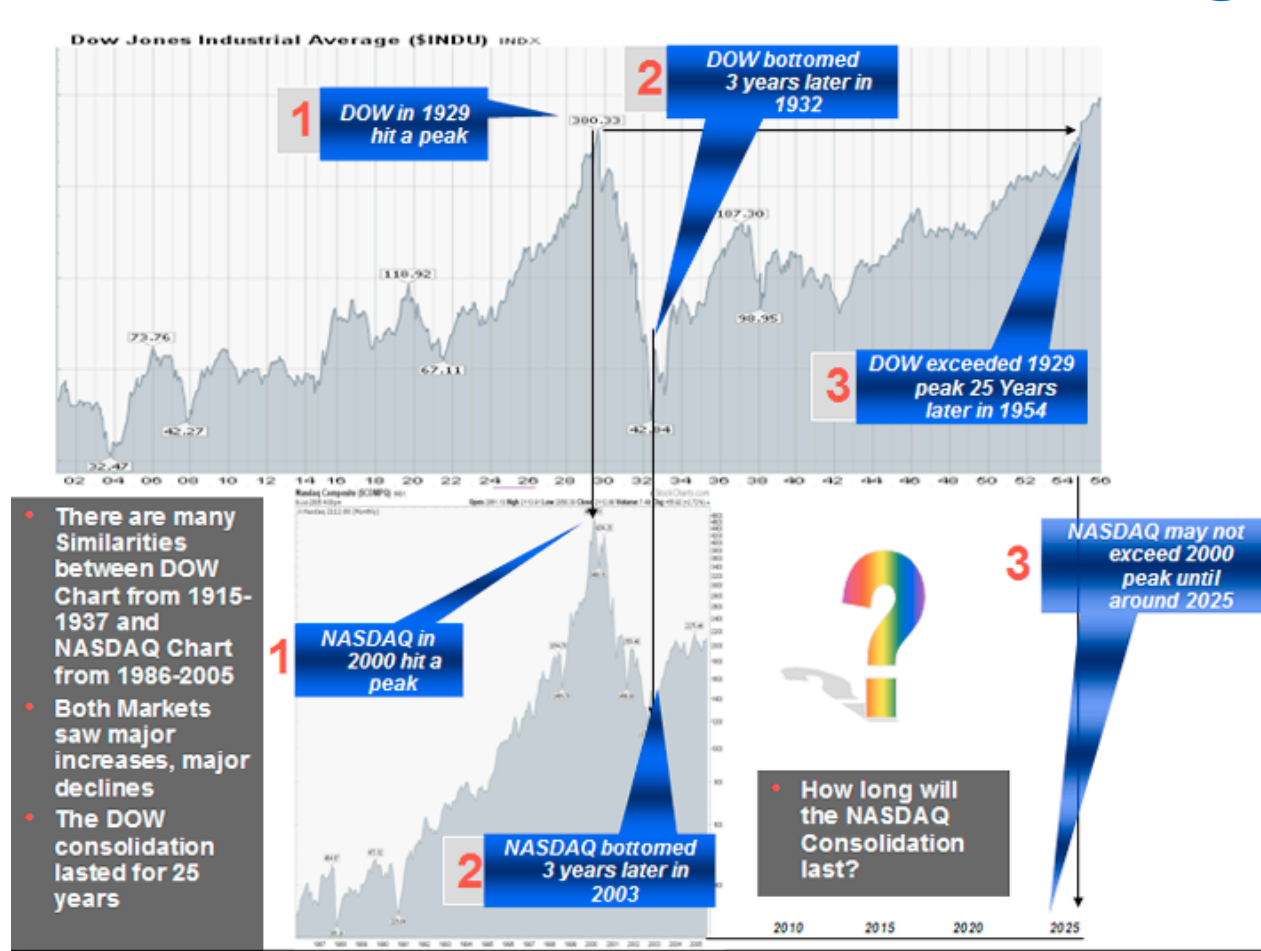
*Estimated
Map*

*Landmark
SLAM
(E. Nebot,
Victoria Park)*



- As robot moves, estimate its pose & world geometry

Financial Forecasting



<http://www.steadfastinvestor.com/>

- Predict future market behavior from historical data, news reports, expert opinions, ...

Special Guest Slides

- David Blei on Latent Dirichlet Allocation
- Mark Johnson on Hidden Markov Models

A gentle introduction to Hidden Markov Models

Mark Johnson
Brown University

November 2009

Outline

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

Sequence labeling

- Input: a sequence $\mathbf{x} = (x_1, \dots, x_n)$ (usually n may vary)
- Output: a *sequence* $\mathbf{y} = (y_1, \dots, y_n)$, where y_i *is a label for* x_i

Part of speech (POS) tagging:

\mathbf{y} : DT JJ NN VBD NNP .
 \mathbf{x} : the big cat bit Sam .

Noun-phrase chunking:

\mathbf{y} : [NP NP NP] - [NP] .
 \mathbf{x} : the big cat bit Sam .

Named entity detection:

\mathbf{y} : [CO CO] - [LOC] - [PER] -
 \mathbf{x} : XYZ Corp. of Boston announced Spade's resignation

Speech recognition: The \mathbf{x} are 100 msec. time slices of acoustic input, and the \mathbf{y} are the corresponding phonemes (i.e., y_i is the phoneme being uttered in time slice x_i)

Sequence labeling with probabilistic models

- General idea: develop a probabilistic model $P(\mathbf{y}, \mathbf{x})$
 - ▶ use this to compute $P(\mathbf{y}|\mathbf{x})$ and $\hat{\mathbf{y}}(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x})$
- A *Hidden Markov Model* (HMM) factors:

$$P(\mathbf{y}, \mathbf{x}) = P(\mathbf{y}) P(\mathbf{x}|\mathbf{y})$$

- $P(\mathbf{y})$ is the probability of sequence \mathbf{y}
 - ▶ In an HMM, $P(\mathbf{y})$ is a *Markov chain*
 - in applications, \mathbf{y} is usually “hidden”
 - ▶ $P(\mathbf{x}|\mathbf{y})$ is the *emission model*.

Outline

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

Markov models of sequences

- Let $\mathbf{y} = (y_1, \dots, y_n)$, where $y_i \in \mathcal{Y}$ and n is fixed

$$\begin{aligned} P(\mathbf{y}) &= P(y_1) P(y_2|y_1) P(y_3|y_1, y_2) \dots P(y_n|y_1, \dots, y_{n-1}) \\ &= P(y_1) \prod_{i=2}^n P(y_i|y_{1:i-1}) \end{aligned}$$

where $y_{i:j}$ is the subsequence (y_i, \dots, y_j) of \mathbf{y}

- In a (first-order) *Markov chain*:

$$P(y_i|y_{1:i-1}) = P(y_i|y_{i-1}) \text{ for } i > 1$$

- In a *homogenous* Markov chain, $P(y_i|y_{i-1})$ does not depend on i , so probability of a sequence is completely determined by:

$$P(y_1) = \iota(y_1) \quad (\textit{start probabilities})$$

$$P(y_i|y_{i-1}) = \sigma(y_{i-1}, y_i) \quad (\textit{transition probabilities}), \text{ so:}$$

$$P(\mathbf{y}) = \iota(y_1) \prod_{i=2}^n \sigma(y_{i-1}, y_i)$$

Markov models of varying-length sequences

- To define a model of sequences $\mathbf{y} = (y_1, \dots, y_n)$ of varying lengths n , need to define $P(n)$
- An easy way to do this:
 - ▶ add a new *end marker* symbol ' \triangleleft ' to \mathcal{Y}
 - ▶ pad all sequences with this end-marker

$$\mathbf{y} = y_1, y_2, \dots, y_n, \triangleleft, \triangleleft, \dots$$

- ▶ set $\sigma(\triangleleft, \triangleleft) = 1$, i.e., \triangleleft is a *sink state*
 - ▶ so $\sigma(y, \triangleleft)$ is probability of sequence ending after state y
- We can also incorporate the start probabilities ι into σ by introducing a new *begin marker* $y_0 = \triangleright$

$$\mathbf{y} = \triangleright, y_1, y_2, \dots, y_n, \triangleleft, \triangleleft, \dots, \text{ so:}$$

$$P(\mathbf{y}) = \prod_{i=1}^{\infty} \sigma(y_{i-1}, y_i) \quad (\text{only compute to } i = n + 1)$$

Outline

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

Hidden Markov models

- An HMM is specified by:
 - ▶ A *state-to-state transition matrix* σ , where $\sigma(y, y')$ is the probability of moving from state y to state y'
 - ▶ A *state-to-observation emission matrix* τ , where $\tau(y, x)$ is the probability of emitting x in state y

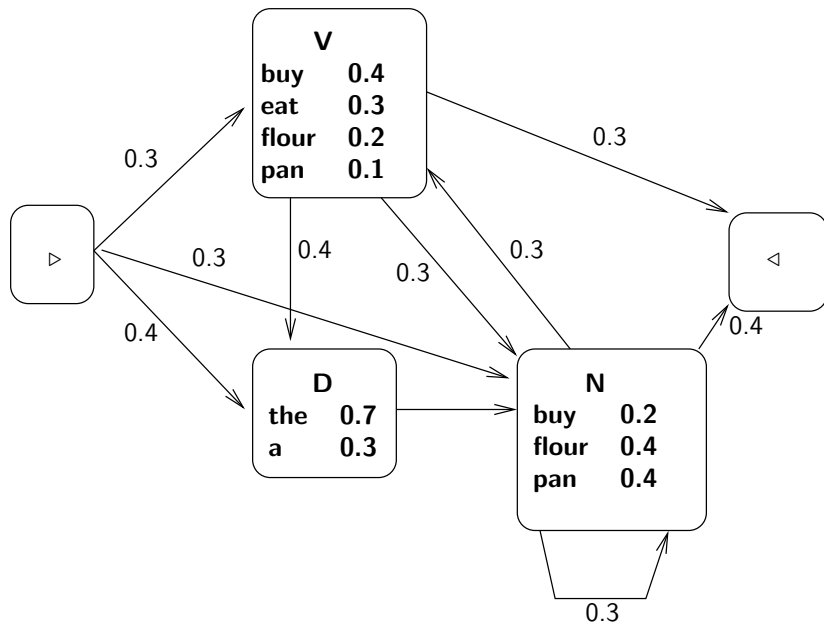
$$P(\mathbf{y}, \mathbf{x}) = P(\mathbf{y}) P(\mathbf{x}|\mathbf{y})$$

$$P(\mathbf{y}) = \prod_{i=1}^{n+1} P(y_i|y_{i-1}) = \prod_{i=1}^{n+1} \sigma(y_{i-1}, y_i)$$

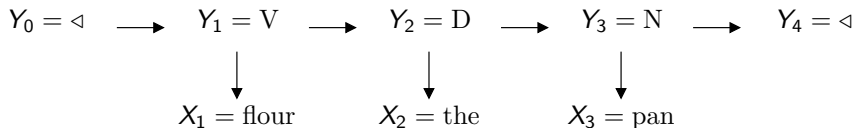
$$P(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^n P(x_i|y_i) = \prod_{i=1}^n \tau(y_i, x_i)$$

- I.e., each x_i is generated conditioned on y_i with probability $\tau(y_i, x_i)$

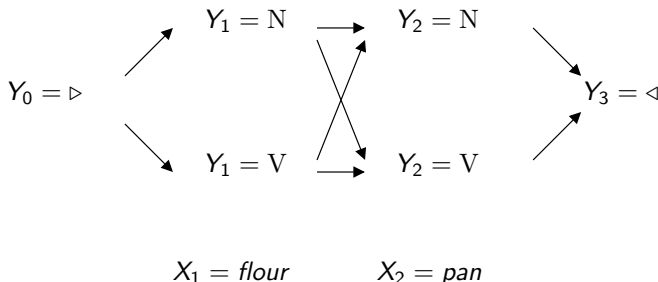
HMMs as stochastic automata



Bayes net representation of an HMM



Trellis representation of HMM



- The trellis representation represents the possible state sequences \mathbf{y} for *a single observation sequence* \mathbf{x}
- A trellis is a DAG with *nodes* (i, y) , where $i \in 0, \dots, n + 1$ and $y \in \mathcal{Y}$ and *edges* from $(i - 1, y')$ to (i, y)
- It plays a crucial role in *dynamic programming algorithms* for HMMs

Outline

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

Two important inference problems

- Given an HMM σ, τ and an observation sequence x :
 - ▶ what is *the most likely state sequence* $\hat{y} = \operatorname{argmax}_y P(x, y)$?
 - ▶ what is the *probability of a sequence* $P(x) = \sum_y P(x, y)$?
- If the sequences are short enough, can be solved by exhaustively enumerating y
- Given a sequence of length n and an HMM with m states, can be solved using dynamic programming in $O(n m^2)$ time

Key observation: $P(\mathbf{x}, \mathbf{y})$ factorizes!

$$P(\mathbf{x}, \mathbf{y}) = \left(\prod_{i=1}^n \sigma(y_{i-1}, y_i) \tau(y_i, x_i) \right) \sigma(y_n, \triangleleft), \text{ so:}$$

$$-\log P(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n -\log(\sigma(y_{i-1}, y_i) \tau(y_i, x_i)) \right) - \log \sigma(y_n, \triangleleft)$$

- Given observation sequence \mathbf{x} , define a trellis with *edge weights*

$$\begin{aligned} w((i-1, y'), (i, y)) &= -\log(\sigma(y_{i-1}, y_i) \tau(y_i, x_i)) \\ w((n, y), (n+1, \triangleleft)) &= -\log(\sigma(y, \triangleleft)) \end{aligned}$$

- Then the $-\log$ probability of any state sequence \mathbf{y} is the sum of the weights of the corresponding path in the trellis

Finding the most likely state sequence \hat{y}

- Goal: find $\hat{y} = \operatorname{argmax}_y P(x, y)$
- Algorithm:
 - ▶ Construct a trellis with edge weights

$$\begin{aligned}w((i-1, y'), (i, y)) &= -\log(\sigma(y_{i-1}, y_i) \tau(y_i, x_i)) \\w((n, y), (n+1, \triangleleft)) &= -\log(\sigma(y, \triangleleft))\end{aligned}$$

- ▶ Run a *shortest path algorithm* on this trellis, which returns a state sequence with lowest $-\log$ probability

The Viterbi algorithm for finding \hat{y}

- Idea: given $\mathbf{x} = (x_1, \dots, x_n)$, let $\mu(k, u)$ be the maximum probability of any state sequence y_1, \dots, y_k ending in $y_k = u \in \mathcal{Y}$.

$$\mu(k, u) = \max_{\mathbf{y}_{1:k} : y_k = u} P(\mathbf{x}_{1:k}, \mathbf{y}_{1:k})$$

- ▶ Base case: $\mu(0, \triangleright) = 1$
- ▶ Recursive case:

$$\mu(k, u) = \max_{u'} \mu(k-1, u') \sigma(u', u) \tau(u, x_k)$$

- ▶ Final case: $\mu(n+1, \triangleleft) = \max_{u'} \mu(n, u') \sigma(u', \triangleleft)$
- To find \hat{y} , for each trellis node (k, u) keep a *back-pointer* to *most likely predecessor*

$$\rho(k, u) = \operatorname{argmax}_{u' \in \mathcal{Y}} \mu(k-1, u') \sigma(u', u) \tau(u, x_k)$$

Why does the Viterbi algorithm work?

$$\begin{aligned}\mu(k, u) &= \max_{\mathbf{y}_{1:k} : y_k = u} P(\mathbf{x}_{1:k}, \mathbf{y}_{1:k}) \\&= \max_{\mathbf{y}_{1:k} : y_k = u} \prod_{i=1}^k \sigma(y_{i-1}, y_i) \tau(y_i, x_i) \\&= \max_{\mathbf{y}_{1:k-1}} \left(\prod_{i=1}^{k-1} \sigma(y_{i-1}, y_i) \tau(y_i, x_i) \right) \sigma(y_{k-1}, u) \tau(u, x_k) \\&= \max_{u'} \left(\max_{\mathbf{y}_{1:k-1} : y_{k-1} = u'} \prod_{i=1}^{k-1} \sigma(y_{i-1}, y_i) \tau(y_i, x_i) \right) \sigma(u', u) \tau(u, x_k) \\&= \max_{u'} \mu(k-1, u') \sigma(u', u) \tau(u, x_k)\end{aligned}$$

The sum-product algorithm for computing $P(\mathbf{x})$

- Goal: compute $P(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{x}, \mathbf{y})$
- Idea: given $\mathbf{x} = (x_1, \dots, x_n)$, let $\alpha(k, u)$ be the *sum of the probabilities* of x_1, \dots, x_k and all state sequences y_1, \dots, y_k ending in $y_k = u$.

$$\alpha(k, u) = \sum_{y_{1:k} : y_k = u} P(\mathbf{x}_{1:k}, \mathbf{y}_{1:k})$$

- ▶ Base case: $\alpha(0, \triangleright) = 1$
- ▶ Recursive case: for $k = 1, \dots, n$:

$$\alpha(k, u) = \sum_{u'} \alpha(k-1, u') \sigma(u', u) \tau(u, x_k)$$

- ▶ Final case: $\alpha(n+1, \triangleleft) = \sum_{u'} \alpha(n, u') \sigma(u', \triangleleft)$
- The α are called *forward probabilities*

Outline

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

Estimating HMMs from *visible data*

- Training data: $D = ((x_1, y_1), \dots, (x_n, y_n))$, where the x_i and y_i are *sequences*
 - ▶ for conceptual simplicity, assume that the data comes as two long sequences $D = (x, y)$
- Then the MLE is:

$$\hat{\sigma}(y', y) = \frac{c(y', y)}{\sum_{y \in \mathcal{Y}} c(y', y)}$$

$$\hat{\tau}(y, x) = \frac{c(y, x)}{\sum_{x \in \mathcal{X}} c(y, x)}, \text{ where:}$$

$c(y', y)$ = the number of times state y' precedes state y in y

$c(y, x)$ = the number of times state y emits x in (x, y)

Estimating HMMs from *hidden data* using EM

- EM iteration at time t :

$$\sigma(y', y)^{(t)} = \frac{\mathbb{E}[c(y', y)]}{\sum_{y \in \mathcal{Y}} \mathbb{E}[c(y', y)]}$$

$$\tau(y, x)^{(t)} = \frac{\mathbb{E}[c(y, x)]}{\sum_{x \in \mathcal{X}} \mathbb{E}[c(y, x)]}, \text{ where:}$$

$$\mathbb{E}[c(y', y)] = \sum_{\mathbf{y} \in \mathcal{Y}^n} c(y', y) P(\mathbf{y} | \mathbf{x}, \boldsymbol{\sigma}^{(t-1)}, \boldsymbol{\tau}^{(t-1)})$$

$$\mathbb{E}[c(y, x)] = \sum_{\mathbf{y} \in \mathcal{Y}^n} c(y, x) P(\mathbf{y} | \mathbf{x}, \boldsymbol{\sigma}^{(t-1)}, \boldsymbol{\tau}^{(t-1)})$$

- Key computational problem: *computing these expectations efficiently*
- Dynamic programming to the rescue!

Expectations require marginal probabilities

- Suppose we could efficiently compute the *marginal probabilities*:

$$P(Y_i = y' | \mathbf{x}) = \frac{\sum_{\mathbf{y} \in \mathcal{Y}^n : y_i = y'} P(\mathbf{y}, \mathbf{x})}{P(\mathbf{x})}$$

- Then we could efficiently compute:

$$E[c(y', \mathbf{x}')] = \sum_{i : x_i = x'} P(Y_i = y' | \mathbf{x})$$

- Efficiently computing $E[c(y', y)]$ requires computation of a different marginal probability

Backward probabilities

- Idea: given $\mathbf{x} = (x_1, \dots, x_n)$, let $\beta(k, v)$ be the *sum of the probabilities* of x_{k+1}, \dots, x_n and all state sequences y_{k+1}, \dots, y_n conditioned on $y_k = v$.

$$\beta(k, v) = \sum_{y_{k:n} : y_k = v} P(\mathbf{x}_{k+1:n}, \mathbf{y}_{k+1:n} | Y_k = v)$$

- ▶ Base case: $\beta(n, v) = \sigma(v, \triangleleft)$ for each $v \in \mathcal{Y}$
- ▶ Recursive case: for $k = n - 1, \dots, 0$:

$$\beta(k, v) = \sum_{v'} \sigma(v, v') \tau(v', x_{k+1}) \beta(k + 1, v')$$

- The β are called *backward probabilities*

Computing marginals from forward and backward probabilities

$$\begin{aligned}P(Y_i = u | \mathbf{x}) &= \frac{1}{P(\mathbf{x})} \sum_{\mathbf{y} \in \mathcal{Y}^n : y_i = u} P(\mathbf{y}, \mathbf{x}) \\&= \frac{1}{P(\mathbf{x})} \sum_{y_{1:k} : y_i = u} \sum_{y_{i:n} : y_i = u} P(\mathbf{x}_{1:i}, \mathbf{y}_{1:i}) P(\mathbf{x}_{i+1:n}, \mathbf{y}_{i:n}) \\&= \frac{1}{P(\mathbf{x})} \left(\sum_{y_{1:k} : y_i = u} P(\mathbf{x}_{1:i}, \mathbf{y}_{1:i}) \right) \left(\sum_{y_{i:n} : y_i = u} P(\mathbf{x}_{i+1:n}, \mathbf{y}_{i:n}) \right) \\&= \frac{1}{P(\mathbf{x})} \alpha(i, u) \beta(i, u)\end{aligned}$$

- So the marginal $P(Y_i | \mathbf{x})$ can be efficiently computed from α and β

Outline

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

Conclusion

- Hidden Markov Models (HMMs) can be used to label a sequence of observations
- There are efficient dynamic programming algorithms to find the most likely label sequence and the marginal probability of a label
- The expectation-maximization algorithm can be used to learn an HMM from unlabeled data
- The expected values required for EM can be efficiently computed