



Graphical Model Research in Audio, Speech, and Language Processing

Jeff A. Bilmes University of Washington Department of EE, SSLI-Lab



Outline

- I. Graphical Models Review
- **II.** Speech Recognition Overview
- III. Goals for GMs in Speech/Language
 - A. Explicit control
 - **B**. Latent Modeling (audio, speech, language)

GMs in Audio, Speech, and Language

- C. Observation Modeling
- D. Structure Learning
- IV. Toolkits and Inference

Jeff A. Bilmes

<section-header><section-header><section-header><section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item>

Graphical Models (GMs)

GMs give us:

- I. Structure: A method to explore the structure of "natural" phenomena (causal vs. correlated relations, properties of natural signals and scenes)
- **II.** Algorithms: A set of algorithms that provide "efficient" probabilistic inference and statistical decision making
- **III. Language**: A mathematically formal, abstract, visual language with which to efficiently discuss and intuit families of probabilistic models and their properties.

Jeff A. Bilmes





Directed GM (DGMs) (Bayesian Networks)

- When is $X_A \parallel X_B \mid X_C?$
- Only when *C* d-separates *A* from *B*, i.e. if: for all paths from A to B, there is a *v* on the path s.t. either one of the following holds:
 - 1. Either $\rightarrow v \rightarrow$ or $\leftarrow v \rightarrow$ and $v \in C$
 - 2. $\rightarrow v \leftarrow$ and neither v nor any descendants are in C
- Equivalent to "directed local Markov property" (CI of non-descendants given parents), plus others (again see Lauritzen '96)

Jeff A. Bilmes





Undirected GMs

- When is $X_A \coprod X_B \mid X_C?$
- Only when *C* separates *A* from *B*. I.e., if: for all paths from A to B, there is a *v* on the path s.t. *v* ∈ *C*
- Simpler semantics than Bayesian networks.
- Equivalent to "global Markov property", plus others (again see Lauritzen '96)

Jeff A. Bilmes





Jeff A. Bilmes











Ideal case, use Bayes decision rule

$$(K^*, W^*_{1:K}) = \underset{K, W_{1:K}}{\operatorname{argmax}} \Pr(W_{1:K}, K | X_{1:T})$$

$$= \underset{K,W_{1:K}}{\operatorname{argmax}} \Pr(X_{1:T} | W_{1:K}, K) \Pr(W_{1:K})$$

•Bayes Decision Theory (see Duda & Hart 73)

Jeff A. Bilmes

GMs in Audio, Speech, and Language

Generative vs. Discriminative Models

• Ideal case: discriminative model

 $P(W_{1:K} | X_{1:T})$

- Too many classes for a discriminative model
 - 100k words, $K = 10 \Rightarrow (100k)^{10}$ classes
 - (not to mention we didn't consider all other K's)
- Generative model can help: $P(x_{1:T})$
- Use the "natural" hierarchy in speech/language:
 - Sentences are composed of words (W)
 - Words (W) are composed of phones (Q)
 - Phones (Q) are composed of Markov chain states (S)
 - States (S) are composed of acoustic feature vector sequences (X)
 - Acoustic feature vector sequences (X) are composed of noisy (e.g., channel distorted) versions thereof (Y)

Jeff A. Bilmes







- $P(x|q_t)$ where q_t is a tri-phone
- **x**-**y**+**z** notation: phone **y** with
 - left context of x
 - right context of z
- Example transcription of "Beat it" : sil b iy t ih t sil
 - sil
 - sil-b+iy
 - **b-iy-t** (or **b-iy+dx** for Americans)
 - iy-t+ih (or iy-dx+ih)
 - $t-ih+t \qquad (or dx-ih+t)$
 - ih-t+sil
 - sil
- To further increase states: Word internal vs. crossword tri-phones

Jeff A. Bilmes





Solution implemented using search via dynamic programming

 $w^* = \max_{s,q,w} p(x,q,s,w)$

- Need optimized search algorithms
 - Viterbi decoding, time synchronous
 - stack decoding, A* search, time asynchronous
 - Both will heavily prune the search space, thus achieving a form of "approximate" inference













Challenges in Speech Recognition

- > 60k words, exhaustive examination of all words is infeasible since |W|²|Q||S| states.
- Even HMM decoding is a challenge
 - Clearly, large grid approach is infeasible
 - Pruning with a beam: try to discard unlikely partial hypotheses as soon as possible (without increasing error)
 - Explore word sequences in parallel (multiple partial hypotheses are considered at same time)

GMs in Audio, Speech, and Language

Jeff A. Bilmes



The Savior: Parameter Tying Generative model + speech/language hierarchy allows

for massive amounts of parameter tying or sharing.

- Same words in difference sentences or different parts of same sentence are the same
- Same phones (subwords) in different words or in different parts of same word are the same
- Certain states in different phones are merged
 - E.g., p(x|S=i) = p(x|S=j) for the right i and j.
- Certain observation parameters (e.g., means) are shared.
- Various ways to accomplish this:
 - backing off (like in language model)
 - [a-b+c] model backs off to [b+c] or to [a-b] etc.
 - Smoothing, interpolation, and mixing
 - clustering (widely used)
 - Decision tree clustered tri-phones
 - both bottom up and top down clustering procedures.

Jeff A. Bilmes

GMs in Audio, Speech, and Language

Four Main Goals for GMs in Speech/Language

- 1. *Explicit Control*: Derive graph structures that themselves *explicitly* represent control constructs
 - E.g., parameter tying/sharing, state sequencing, smoothing, mixing, backing off, etc.
- 2. Latent Modeling: Use graphs to represent latent *information* in speech/language, not normally represented.
- 3. *Observation Modeling*: represent structure over observations.
- **4.** *Structure learning*: Derive *structure* automatically, ideally to improve error rate while simultaneously minimizing computational cost.

Jeff A. Bilmes





Key Points

- Graph explicitly represents parameter sharing
- Same phone at different parts of the word are the same: phone /aa/ in positions 2, 4, and 6 of the word "yamaha"
- Phone-dependent transition indicator variables yield geometric phone duration distributions for each phone
- Counter variable ensures /aa/'s at different positions move only to correct next phone
- Some edge implementations are deterministic (green) and others are random (red)
- End of word observation, gives zero probability to variable assignments corresponding to incomplete words.

Jeff A. Bilmes



































Skip Bi-gram

- Often there is silence between words

 "fool me once <sil> shame on <sil> shame on you"
- Silence might not be good predictor of next word
- But silence lexemes should be represented since acoustics quite different during silence.
- Goal: allow silence between words, but retain true word predictability skipping silence regions.
- Switching parents can facilitate such a model.

Jeff A. Bilmes



Skip bi-gram with conditional mixtures

 $p_{bigram}(w_t \mid r_{t-1}) = P(\alpha_t = 1)P(w_t) + P(\alpha_t = 2)P(w_t \mid w_{t-1})$











Explicit Smoothing

- Disjoint partition of vocabulary based on training-data counts: = {unk}∪ ∪
- = singletons, = "many-tons", unk=unknown
- ML distribution gives zero probability to unk.
- Goal: Directed GM that represents:

$$p(w) = \begin{cases} 0.5 p_{ml}() & \text{if } w = unk \\ 0.5 p_{ml}(w) & \text{if } w \in \\ p_{ml}(w) & \text{otherwise} \end{cases}$$

• Word variable is <u>like</u> a switching parent of itself (but of course can't be)

Jeff A. Bilmes









Factored Language Models

- Decompose words into smaller morphological or class-based units (e.g., morphological classes, stems, roots, patterns, or other automatically derived units).
- Produce probabilistic models over these units to attempt to improve language modeling accuracy and parameter estimation

GMs in Audio, Speech, and Language

Jeff A. Bilmes

Example with Words, Stems, and Morphological classes $M_{t-3} \quad M_{t-2} \quad M_{t-1} \quad M_{t}$ $S_{t-3} \quad S_{t-2} \quad S_{t-1} \quad S_{t}$ $W_{t-3} \quad W_{t-2} \quad W_{t-1} \quad W_{t}$ $P(w_{t} | s_{t}, m_{t}) \quad P(s_{t} | m_{t}, w_{t-1}, w_{t-2}) \quad P(m_{t} | w_{t-1}, w_{t-2})$























Discriminative structure learning

- Structure is typically learned to optimize marginal likelihood (e.g., statistical predictability)
- When the underlying goal is classification (regression), discriminative structure learning
- Structure is chosen to optimize conditional posterior of class variable (more generally, conditional likelihood) rather than marginal likelihood.
- Can still use generative models
- Structure edges can "switch" depending on current condition

Discriminative vs. Generative Models

- p(X|C) vs. p(C|X)
- Goals for recognition (classification) are different than for generative accuracy (e.g., synthesis)
- Generative models natural for ASR
- Approach: retain generative models but train discriminatively
 - Discriminative parameter training can occur for parameters of generative models also (e.g., maximum mutual information estimation on HMMs)
 - Discriminative structure learning.

Jeff A. Bilmes

GMs in Audio, Speech, and Language

Discriminative Generative Models (DG models)

 $\mathcal{F} = \{ f(x:m) : \operatorname{argmax} p(x \mid m) p(m) = \operatorname{argmax} f(x,m) p(m) \}$ m m

So choose the *f* that satisfies the above, but is as simple (few parameters, easy to compute) as possible.

Model might no longer "generates" samples accurately, but discriminates well.

Jeff A. Bilmes







Jeff A. Bilmes



ge	Com enerat	parison of g tive-discrim	genera inativ	tive vs. ve models.
	CASE	TYPE	WER	PARAMS
*	1 🐗	Generative	32.0%	207k
*	2 🌾	HMM	5.0%	157k
	3 🐗	Discriminative	4.6%	157k
Jeff A. Bilmes			GMs in Audio, Speech, and Language	















GMTK: Graphical Models Toolkit

- A GM-based software system for speech, language, and time-series modeling
- One system Many different underlying statistical models (more than an HMM)
- <u>Complements</u> rather than replaces other ASR and GM systems (e.g., HTK, AT&T, ISIP, BNT, BUGS, Hugin, etc.)
- Ultimately will be open-source, freely available
- Long-term multi-year goal: improve features, computational speed, and portability.

Jeff A. Bilmes



GMTK Features

- Textual Graph Language
- Switching Parent Functionality
- Linear/Non-linear Dependencies on observations
- Arbitrary low-level parameter sharing (EM/GEM training)
- Gaussian Vanishing/Splitting algorithm.
- Decision-Tree-Based implementations of dependencies (deterministic and sparse)
- Full inference, single pass decoding possible on smaller tasks (current version)
- Sampling Methods
- Log space Exact Inference Memory O(logT)

Jeff A. Bilmes

GMTK Structure file for HMM			
frame : 0 {			
variable : state {			
<pre>type : discrete hidden cardinality 4000;</pre>			
switchingparents : nil;			
<pre>conditionalparents : nil using DenseCPT("pi");</pre>			
}			
variable : observation {			
type : continuous observed 0:39;			
switchingparents : nil;			
<pre>conditionalparents : state(0) using mixGaussian mapping("state2obs");</pre>			
}			
}			
frame : 1 {			
variable : state {			
type : discrete hidden cardinality 4000;			
switchingparents : nil;			
<pre>conditionalparents : state(-1) using DenseCPT("transitions");</pre>			
}			
variable : observation {			
type : continuous observed 0:39;			
switchingparents : nil;			
<pre>conditionalparents : state(0) using mixGaussian mapping("state2obs");</pre>			
}			
}			
Jeff A. Bilmes GMs in Audio, Speech, and Language			











Gaussians Represented as Bayesian Networks

- Factor concentration matrix: K = U'DU
 - D = positive diagonal of conditional variances
 - U = unit upper-triangular matrix
- det(U) = 1 so det(U'DU) = det(D)
- Gaussian becomes:

$$f(x) = |2\pi D|^{-1/2} e^{-\frac{1}{2}(x-\mu)'U'DU(x-\mu)}$$

Jeff A. Bilmes







GMTK Splitting/Vanishing Algorithm

- Determines number Gaussian components/state
- Split Gaussian if it's component probability ("responsibility") rises above a number-ofcomponents dependent threshold
- Vanish Gaussian if it's component probability falls below a number-of-components dependent threshold
- Use a splitting/vanishing schedule, one set of thresholds per each EM training iteration.

Jeff A. Bilmes



Linear and Log space exact inference

- Exact inference O(T*S) space and time complexity, S = clique state space size
- Log-space inference O(log(T)*S) space at an extra cost of a factor of log(T) time.
- Can use both linear and log space inference at same time (for optimal tradeoff).
- This is same idea as what has been called the Island Algorithm







The GMTK Triangulation Engine (an anytime algorithm)

- User specifies an amount of time (2mins, 3 hours, 4 days, 5 weeks, etc.) to spend triangulating
- User does not worry about intricacies of graph triangulation
- Uses a "boundary algorithm" to find chunks of DBN to triangulate (UAI'2003)
- Many heuristics implemented: min-fill in, min size, min weight, maximum cardinality search, simulated annealing, exhaustive elimination, and exhaustive triangulation

Jeff A. Bilmes

Current Status

- I. System available at:
 - A. http://ssli.ee.washington.edu/~bilmes/gmtk
 - B. ~100 pages of documentation
 - C. Book chapter on use of graphical models for speech and language
 - D. JHU'2001 Workshop technical report
- II. GMTK Triangulation "Engine" running and ready

Jeff A. Bilmes



Approximate inference in DBNs

- Standard approximate inference methods
 - Pruning as is done performed by modern speech recognition systems
 - Variational and mean-field approaches
 - Loopy belief propagation
 - Sampling, particle filtering, etc.
- All techniques for approximate inference in DBNs are relevant to the speech/language case as well.

Jeff A. Bilmes

GMs in Audio, Speech, and Language

Conclusions

- Many models and many techniques
- We have just scratched the surface, still a relatively young research area.
- Key challenges summary:
 - Explicit Control Structures
 - Structure learning
 - Fast inference techniques
 - Identifying interesting latent variables
 - Structural Discriminability

Jeff A. Bilmes

