

Unsupervised Natural Language Parsing

Kewei Tu, ShanghaiTech University

Yong Jiang, Alibaba DAMO Academy

Wenjuan Han, National University of Singapore

Yanpeng Zhao, University of Edinburgh

March 29, 2021

2. Generative Approaches: Parameter Learning

Motivation

Structure learning is hard!

Motivation

Structure learning is hard!

Another solution is to directly perform parameter learning with a fixed structure (e.g., enumerating all possible rules).

- ▶ Most research on parameter learning focuses on DMV, a form of generative dependency grammars.
- ▶ These parameter learning methods are generally applicable for PCFGs.

Motivation

Structure learning is hard!

Another solution is to directly perform parameter learning with a fixed structure (e.g., enumerating all possible rules).

- ▶ Most research on parameter learning focuses on DMV, a form of generative dependency grammars.
- ▶ These parameter learning methods are generally applicable for PCFGs.
- ▶ Before introducing the DMV model, let us review a preliminary model: HMM.

Unsupervised POS Tagging with HMM

- ▶ Goal: find syntax clusters for each word in a sentence.

Unsupervised POS Tagging with HMM

- ▶ Goal: find syntax clusters for each word in a sentence.
- ▶ a.k.a. POS induction: give each word x_i a label z_i for sentence x .

Unsupervised POS Tagging with HMM

- ▶ Goal: find syntax clusters for each word in a sentence.
- ▶ a.k.a. POS induction: give each word x_i a label z_i for sentence x .

General approach for POS induction: Hidden Markov models (HMMs)

- ▶ Start probabilities: $P(z_1|z_0 = \textit{START})$
- ▶ Transition probabilities: $P(z_i|z_{i-1})$
- ▶ Emission probabilities: $P(x_i|z_i)$

Unsupervised POS Tagging with HMM

- ▶ Goal: find syntax clusters for each word in a sentence.
- ▶ a.k.a. POS induction: give each word x_i a label z_i for sentence x .

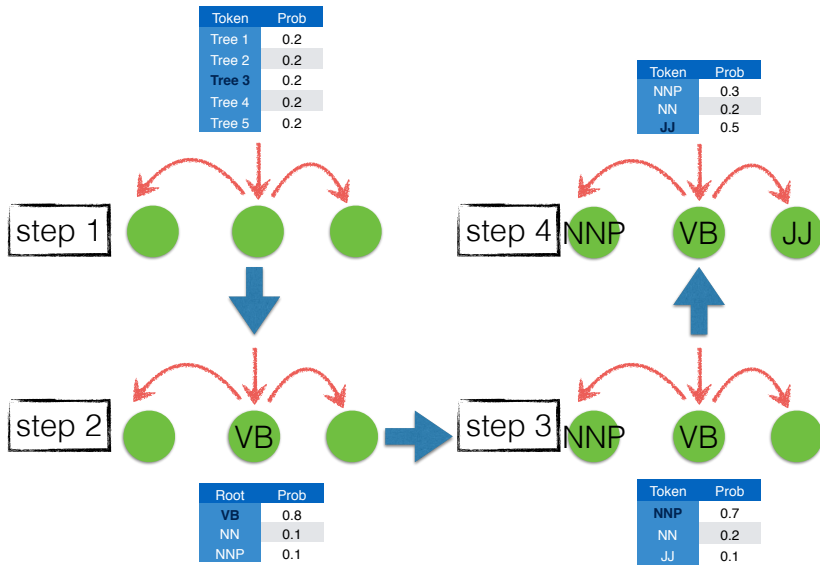
General approach for POS induction: Hidden Markov models (HMMs)

- ▶ Start probabilities: $P(z_1|z_0 = START)$
- ▶ Transition probabilities: $P(z_i|z_{i-1})$
- ▶ Emission probabilities: $P(x_i|z_i)$

Example: To generate a sentence I SWAM with POS sequence PRONOUN VERB.

$$\begin{aligned} P(x, z) &= P(z_1 = \textit{Pronoun} | z_0 = \textit{START}) \\ &\cdot P(z_2 = \textit{Verb} | z_1 = \textit{Pronoun}) \cdot P(x_1 = \textit{I} | z_1 = \textit{Pronoun}) \\ &\cdot P(x_2 = \textit{swam} | z_2 = \textit{Verb}) \end{aligned}$$

Old Dependency Models



[Paskin (2002) and Carroll & Charniak (1992)]

DMV with an example

step 1

↓
VB

Token	Prob
VB	0.8
NN	0.1
NNP	0.1

step 5

NNP VB

Decision	Prob
CONT	0.9
STOP	0.1

step 2

↖ ↘
↓
VB

Decision	Prob
CONT	0.9
STOP	0.1

step 6/7

NNP VB JJ

Token	Prob
NNP	0.3
NN	0.2
JJ	0.5
Decision	Prob
CONT	0.2
STOP	0.8

step 3

NNP VB

Token	Prob
NNP	0.7
NN	0.2
JJ	0.1

step 8/9

NNP VB JJ

Decision	Prob
CONT	0.1
STOP	0.9
Decision	Prob
CONT	0.2
STOP	0.8

step 4

↖ ↘
↖ ↘
NNP VB

Decision	Prob
CONT	0.3
STOP	0.7

step 10/11

NNP VB JJ

Decision	Prob
CONT	0.4
STOP	0.6
Decision	Prob
CONT	0.3
STOP	0.7

DMV Model Representation (Klein & Manning, ACL 2004)

Formal Definition (Dependency Model with Valence):

- ▶ Sentence x , parse tree z , model joint probability $P(x, z; \Theta)$
- ▶ three kinds of grammars rules: root, attach and decision.
- ▶ $dir(h, c)$: dependency direction from parent token h to child token c .
- ▶ $val(h)$ indicates valency: whether h has already generated a child.

DMV Model Representation (Klein & Manning, ACL 2004)

Formal Definition (Dependency Model with Valence):

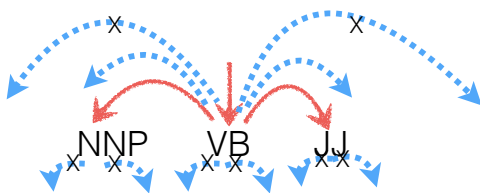
- ▶ Sentence x , parse tree z , model joint probability $P(x, z; \Theta)$
- ▶ three kinds of grammars rules: root, attach and decision.
- ▶ $dir(h, c)$: dependency direction from parent token h to child token c .
- ▶ $val(h)$ indicates valency: whether h has already generated a child.

Rule Schema:

- ▶ Root: $p_{root}(c)$
- ▶ Attach: $p_{attach}(c|h, dir)$
- ▶ Decision:
 $p_{decision}(CONT|h, dir, val), p_{decision}(STOP|h, dir, val)$

DMV for Computing Sentence & Parse Probability

$$P(x, z; \Theta) = p_{root}(r(x, z)) \times \prod_{(h,c) \in z} (p_{attach}(c|h, dir) p_{decision}(STOP|h, dir, val)) \times \prod_{dir \in \{\leftarrow, \rightarrow\}} p_{decision}(CONT|h, dir, val))$$



Extensions of DMV

- ▶ Headden III et al. (2009) introduced the valence into the condition of attach sampling.

$$p_{attach}(c|h, dir) \Rightarrow p_{attach}(c|h, dir, val)$$

Extensions of DMV

- ▶ Headden III et al. (2009) introduced the valence into the condition of attach sampling.

$$p_{attach}(c|h, dir) \Rightarrow p_{attach}(c|h, dir, val)$$

- ▶ Spitkovsky et al. (2012) conditioned decision and child token generation on sibling words, sentence completeness, and punctuation context.

Extensions of DMV

- ▶ Headden III et al. (2009) introduced the valence into the condition of attach sampling.

$$p_{attach}(c|h, dir) \Rightarrow p_{attach}(c|h, dir, val)$$

- ▶ Spitkovsky et al. (2012) conditioned decision and child token generation on sibling words, sentence completeness, and punctuation context.
- ▶ Yang et al. (2020) proposed a second-order extension of DMV that incorporates grandparent-child or sibling information (p here).

$$p_{attach}(c|h, dir) \Rightarrow p_{attach}(c|h, p, dir, val)$$

Three Different Models of Representing Rule Probabilities

Methods	Representation (Parameters: Θ)
Table-based	$p(c h, dir)$
Feature-based [B-K et al. 2010]	$p(c h, dir) \propto w^T f(h, c, dir)$
Neural-based [Jiang et al. 2016]	$p(c h, dir) = \text{softmax}_c(f(h, dir))$

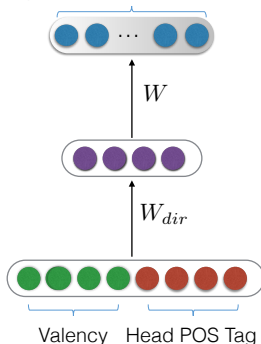
Token	Prob
NNP	0.7
NN	0.2
JJ	0.1

Table-based model

$$p \propto \begin{bmatrix} 0.8, -0.2, 0.3, 0.4 \\ 0, 1, 0, 1 \end{bmatrix}$$

Feature-based model

Outputs (CHILD or DECISION)



Neural-based model

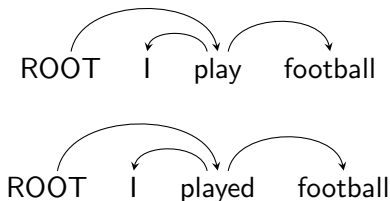
Differences

Drawbacks of Table-based methods: Symbols are independent with each other. However, some words behave alike in parsing.



Differences

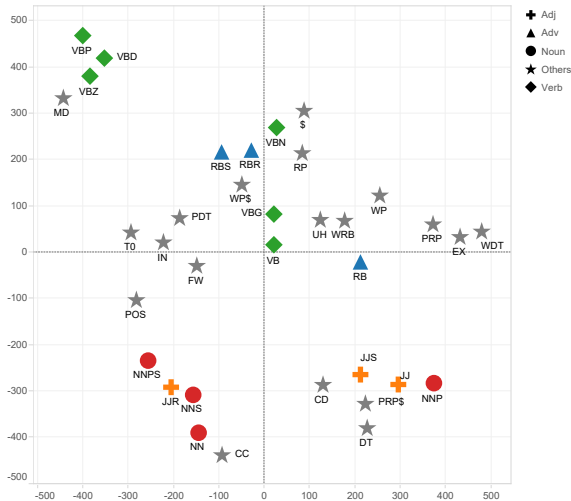
Drawbacks of Table-based methods: Symbols are independent with each other. However, some words behave alike in parsing.



The feature-based and neural-based methods can tackle this problem:

- ▶ Utilizing hand-crafted sparse features (log-linear model).
- ▶ Neuralize the grammar.

Learned Correlations of POS tags for Neural-based Models



[Jiang et al. 2016]

Comparisons

Methods	Pros	Cons
Table-based	Simple parameter learning	Independent symbols.
Feature-based	Modeling symbol similarity	Need manual-designs
Neural-based	Automatic learned similarity	-

Supervised Parameter Estimation

Given a set of annotated sentences $\mathcal{X} = x^1, x^2, \dots, x^N$ and parses $\mathcal{Z} = z^1, z^2, \dots, z^N$, how to learn parameters Θ of generative grammars.

$$J(\Theta) = -\frac{1}{N} \sum_{d=1}^N \log \underbrace{P(x^d, z^d)}_{\text{joint probability}} = -\frac{1}{N} \sum_{d=1}^N \log \underbrace{\sum_{r \in \mathcal{R}(x^d, z^d)} p(r)}_{\text{rule factorization}}$$

where $p(r)$ is normalized.

Supervised Parameter Estimation

Given a set of annotated sentences $\mathcal{X} = x^1, x^2, \dots, x^N$ and parses $\mathcal{Z} = z^1, z^2, \dots, z^N$, how to learn parameters Θ of generative grammars.

$$J(\Theta) = -\frac{1}{N} \sum_{d=1}^N \log \underbrace{P(x^d, z^d)}_{\text{joint probability}} = -\frac{1}{N} \sum_{d=1}^N \log \underbrace{\sum_{r \in \mathcal{R}(x^d, z^d)} p(r)}_{\text{rule factorization}}$$

where $p(r)$ is normalized.

- ▶ Table-based: parameter estimation based on cooccurrence counts. The following is an Attach rule example:

$$p(c|h, \rightarrow) = \frac{\text{count}(c, h, \rightarrow)}{\text{count}(h, \rightarrow)}$$

- ▶ Feature-based & Neural-based: parameter estimation based on gradient-based algorithms.

Unsupervised Learning

Given a set of unannotated sentences $\mathcal{X} = x^1, x^2, \dots, x^N$, how to learn parameters Θ .

- Maximum likelihood estimation:

$$\begin{aligned} J(\Theta) &= -\frac{1}{N} \sum_{d=1}^N \log P(x^d) = -\frac{1}{N} \sum_{d=1}^N \log \underbrace{\sum_{z^d} P(x^d, z^d)}_{\text{marginalized likelihood}} \\ &= -\frac{1}{N} \sum_{d=1}^N \log \sum_{z^d} \prod_{r \in \mathcal{R}(x^d, z^d)} p(r) \end{aligned}$$

Unsupervised Learning

Given a set of unannotated sentences $\mathcal{X} = x^1, x^2, \dots, x^N$, how to learn parameters Θ .

- ▶ Maximum likelihood estimation:

$$\begin{aligned} J(\Theta) &= -\frac{1}{N} \sum_{d=1}^N \log P(x^d) = -\frac{1}{N} \sum_{d=1}^N \log \underbrace{\sum_{z^d} P(x^d, z^d)}_{\text{marginalized likelihood}} \\ &= -\frac{1}{N} \sum_{d=1}^N \log \sum_{z^d} \prod_{r \in \mathcal{R}(x^d, z^d)} p(r) \end{aligned}$$

- ▶ Here $p(r)$ can be parameterized by table-based, feature-based or neural-based methods.
- ▶ EM algorithm to optimize the objective function.

Unsupervised Learning: EM algorithms-Part 1

For unsupervised learning, a guess-and-update approach is utilized.

$$\begin{aligned}\log P(\mathbf{x}; \Theta) &= \sum_z q(z) \log P(\mathbf{x}; \Theta) \\&= \sum_z q(z) \log \frac{P(\mathbf{x}, z; \Theta)}{P(z|\mathbf{x}; \Theta)} \\&= \sum_z q(z) \log \frac{P(\mathbf{x}, z; \Theta)}{q(z)} \frac{q(z)}{P(z|\mathbf{x}; \Theta)} \\&= \sum_z q(z) \log \frac{P(\mathbf{x}, z; \Theta)}{q(z)} + \text{KL}(q(z) || P(z|\mathbf{x}; \Theta)) \\&\geq \sum_z q(z) \log \frac{P(\mathbf{x}, z; \Theta)}{q(z)}\end{aligned}$$

Unsupervised Learning: EM algorithms-Part 2

We can obtain a new objective function:

$$\begin{aligned} J'(\Theta, Q(z)) &= - \sum_{d=1}^N \left(\sum_{z^d} q(z^d) \log \frac{P(x^d, z^d; \Theta)}{q(z^d)} \right) \\ &= J(\Theta) + \sum_{d=1}^N \text{KL}(q(z^d) || P(z^d | x^d; \Theta)) \end{aligned}$$

Unsupervised Learning: EM algorithms-Part 2

We can obtain a new objective function:

$$\begin{aligned} J'(\Theta, Q(z)) &= - \sum_{d=1}^N \left(\sum_{z^d} q(z^d) \log \frac{P(x^d, z^d; \Theta)}{q(z^d)} \right) \\ &= J(\Theta) + \sum_{d=1}^N \text{KL}(q(z^d) || P(z^d | x^d; \Theta)) \end{aligned}$$

► E-step, fix Θ , optimize $q(z^d)$:

$$\begin{aligned} \arg \min_{Q(z)} J'(\Theta, Q(z)) &= \arg \min_{Q(z)} \sum_{d=1}^N \text{KL}(q(z^d) || P(z^d | x^d; \Theta)) \\ \Rightarrow q(z^d) &= P(z^d | x^d; \Theta) \end{aligned}$$

► M-step, fix $q(z^d)$, optimize Θ :

$$J'(\Theta, Q(z)) = - \sum_{d=1}^N \left(\sum_{z^d} q(z^d) \log \frac{P(x^d, z^d; \Theta)}{q(z^d)} \right)$$

Unsupervised Learning: EM algorithms-Part 3

For table-based probabilistic grammars:

- ▶ E-step: utilize dynamic programming (the inside-outside algorithm) to compute a vector of expected frequencies.

$$e(r, x) = E_{q(z)} c(r, x, z)$$

- ▶ M-step: update Θ using expected frequencies.

$$p(r) \propto \sum_x e(r, x)$$

Unsupervised Learning: Online EM algorithms

Cons of EM algorithm:

- ▶ EM algorithm is a batch-style algorithm, suffers from slow convergence.
- ▶ If we have large-scale unlabeled data, performing EM algorithm is very time-consuming.

Online EM algorithm provides significant speed-up.

Unsupervised Learning: Online EM algorithms

Cons of EM algorithm:

- ▶ EM algorithm is a batch-style algorithm, suffers from slow convergence.
- ▶ If we have large-scale unlabeled data, performing EM algorithm is very time-consuming.

Online EM algorithm provides significant speed-up.

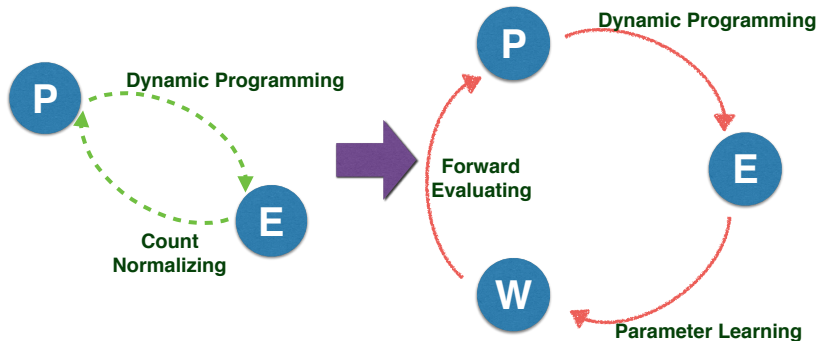
Key ideas

- ▶ Updating parameters (in M-step) after running E-step on a mini-batch of samples rather than the entire corpus.
- ▶ During E-step, interpolating the q distribution with distributions from previous steps.

[Liang et al. 2009]

Unsupervised Learning: Modified EM algorithms

EM algorithm \rightarrow Modified EM algorithm (for feature/neural-based methods).



[B-K et al. 2010, Jiang et al. 2016]

Unsupervised Learning: Direct Gradient Descent

Another approach is directly computing the gradient:

$$\begin{aligned}\nabla_{\Theta} \log P(x) &= \sum_z P(z|x; \Theta) \nabla_{\Theta} \log P(x, z; \Theta) \\ &= \sum_z P(z|x; \Theta) \sum_{r \in \mathcal{R}(x^d, z^d)} c(r, x, z) \nabla_{\Theta} \log p(r) \\ &= \sum_{r \in \mathcal{R}(x^d, z^d)} e(r, x) \nabla_{\Theta} \log p(r)\end{aligned}$$

A trick from [Eisner. 2016]: we can use back-propagation to calculate the expected frequencies $e(r, x)$.

$$e(r, x) = \frac{\partial \log P(x; \Theta)}{\partial \log p(r; \Theta)}$$

No need for the outside algorithm.

Problems

- ▶ MLE objective only aims to explain the training data, which lacks of inductive bias.
- ▶ Local optima problem.

Improvements for MLE: Maximum A Posteriori

$$J(\Theta) = -\log P(\Theta|X) \propto -\sum_{d=1}^N \log \underbrace{\sum_{z^d} P(x^d, z^d|\Theta)}_{\text{marginalized likelihood}} - \underbrace{\log P(\Theta)}_{\text{prior}}$$

Improvements for MLE: Maximum A Posteriori

$$J(\Theta) = -\log P(\Theta|X) \propto -\sum_{d=1}^N \log \underbrace{\sum_{z^d} P(x^d, z^d|\Theta)}_{\text{marginalized likelihood}} - \underbrace{\log P(\Theta)}_{\text{prior}}$$

- Probabilistic grammars are built out of multinomial distributions, so the Dirichlet distributions are natural candidates as priors which encourage smoothness or sparsity [Kurihara and Sato, 2004; Johnson et al. 2007; Tu et al. 2016].

Improvements for MLE: Maximum A Posteriori

$$J(\Theta) = -\log P(\Theta|X) \propto -\sum_{d=1}^N \log \underbrace{\sum_{z^d} P(x^d, z^d|\Theta)}_{\text{marginalized likelihood}} - \underbrace{\log P(\Theta)}_{\text{prior}}$$

- ▶ Probabilistic grammars are built out of multinomial distributions, so the Dirichlet distributions are natural candidates as priors which encourage smoothness or sparsity [Kurihara and Sato, 2004; Johnson et al. 2007; Tu et al. 2016].
- ▶ Cohen and Smith (2008; 2009) leverage logistic-normal prior distributions to encourage symbol correlation.

Improvements for MLE: Maximum A Posteriori

$$J(\Theta) = -\log P(\Theta|X) \propto -\sum_{d=1}^N \log \underbrace{\sum_{z^d} P(x^d, z^d|\Theta)}_{\text{marginalized likelihood}} - \underbrace{\log P(\Theta)}_{\text{prior}}$$

- ▶ Probabilistic grammars are built out of multinomial distributions, so the Dirichlet distributions are natural candidates as priors which encourage smoothness or sparsity [Kurihara and Sato, 2004; Johnson et al. 2007; Tu et al. 2016].
- ▶ Cohen and Smith (2008; 2009) leverage logistic-normal prior distributions to encourage symbol correlation.
- ▶ EM is sometimes not useable in MAP inference, so variational inference and MCMC are often used.

Improvements for MLE: Viterbi Likelihood

$$\begin{aligned} J(\Theta) &= -\frac{1}{N} \log \sum_{d=1}^N \max_{z^d} P(x^d, z^d) \\ &= -\frac{1}{N} \log \sum_{d=1}^N \max_{z^d} \prod_{r \in \mathcal{R}(x^d, z^d)} P(r) \end{aligned}$$

Improvements for MLE: Viterbi Likelihood

$$\begin{aligned} J(\Theta) &= -\frac{1}{N} \log \sum_{d=1}^N \max_{z^d} P(x^d, z^d) \\ &= -\frac{1}{N} \log \sum_{d=1}^N \max_{z^d} \prod_{r \in \mathcal{R}(x^d, z^d)} P(r) \end{aligned}$$

- ▶ Optimized with the hard EM algorithm.
- ▶ Stronger performance in unsupervised parsing.
- ▶ Seen as a special case of entropy regularized model learning.

[Spitkovsky et al. 2010, Tu et al. 2012]

Improvements for MLE: Contrastive Estimation

Contrastive Estimation [Smith and Eisner, 2005a,b]

$$J(\Theta) = -\frac{1}{N} \log \sum_{d=1}^N \frac{\sum_{z^d} s(x, z)}{\sum_{x' \in N(x)} \sum_{z^d} s(x', z)}$$

where $s(x, z)$ is the score of x and z .

Improvements for MLE: Contrastive Estimation

Contrastive Estimation [Smith and Eisner, 2005a,b]

$$J(\Theta) = -\frac{1}{N} \log \sum_{d=1}^N \frac{\sum_{z^d} s(x, z)}{\sum_{x' \in N(x)} \sum_{z^d} s(x', z)}$$

where $s(x, z)$ is the score of x and z .

- ▶ The intuition is to assign higher weight to appeared samples and decrease the weight for neighborhood samples.
- ▶ Choice of neighborhood: linguistic knowledge.
- ▶ Examples: deleting words from x , transposing two words.

Improvements for MLE: Posterior Regularization

Basic idea:

- ▶ Uses constraints on posterior distribution to guide parameter learning.
- ▶ Knowledge of unlikely parses simplify learning.

Improvements for MLE: Posterior Regularization

Basic idea:

- ▶ Uses constraints on posterior distribution to guide parameter learning.
- ▶ Knowledge of unlikely parses simplify learning.

$$J'(\Theta, Q(z)) = J(\Theta) + \sum_{d=1}^N \text{KL}(q(z^d) || P(z^d | x^d; \Theta)) + \sum_{d=1}^N f(q(z^d))$$

- ▶ The new term is only dependent on q .
- ▶ Only the E step is affected and modified. M step remains the same.

[Ganchev et al. 2010]

Different Forms of Posterior Constraints

- ▶ Entropy constraints [Tu and Honavar, 2012]

$$f(q(z)) = - \sum_z q(z) \log q(z)$$

- ▶ Linguistic constraints [Naseem et al., 2010]

$$f(q(z)) = E_{q(z)} \phi(x, z)$$

- ▶ Sparsity constraints [Gillenwater et al., 2010]

$$f(q(z)) = \|E_{q(z)} \phi(x, z)\|_\beta$$

- ▶ Bounding recursion depth [Noji et al. 2016].

$\phi(x, z)$ is a decomposed function, which enables effective learning!

Posterior Regularization #1: Entropy Constraints

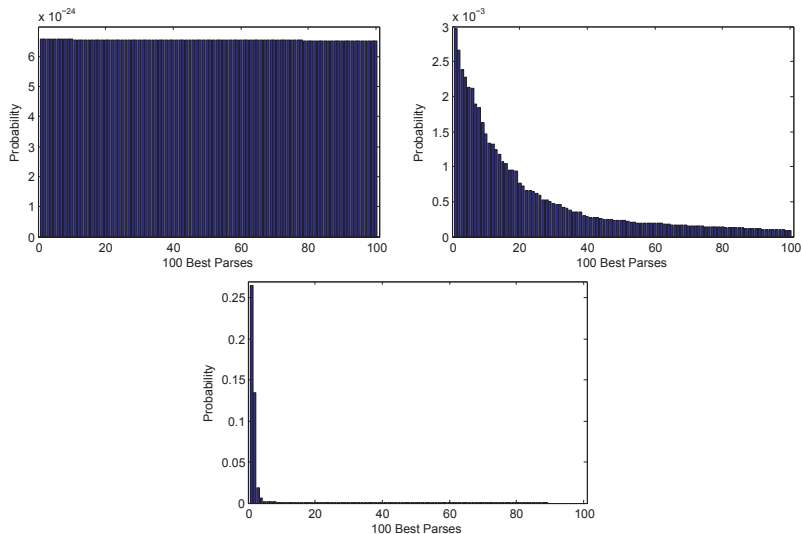
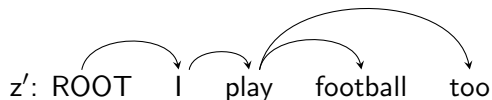
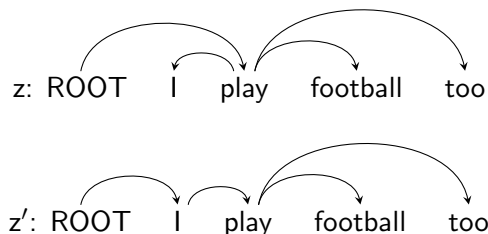


Figure 1: up-left/up-right/down: distribution from random grammar/EM learned model/supervisedly learned grammar, [Tu et al. 2012]

Posterior Regularization #2: Linguistic Constraints



Posterior Regularization #2: Linguistic Constraints



A set of predefined linguistic rules [Naseem et al. 2010], like:

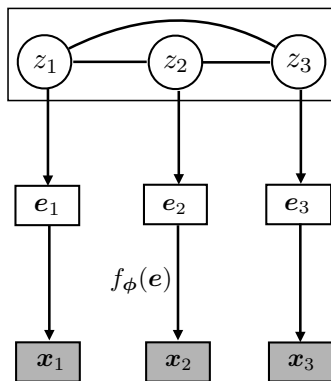
VERB \rightarrow VERB	NOUN \rightarrow NOUN
VERB \rightarrow NOUN	NOUN \rightarrow ADJ
VERB \rightarrow PRON	NOUN \rightarrow DET
VERB \rightarrow ADV	NOUN \rightarrow NUM
VERB \rightarrow ADP	NOUN \rightarrow CONJ
ADJ \rightarrow ADV	ADP \rightarrow NOUN

In this example, $\phi(x, z) = 2$, $\phi(x, z') = 1$

Improvements for Avoiding Local Minimum

- ▶ Deterministic annealing [Smith and Eisner, 2004]: start with a concave objective and gradually move to the actual non-concave objective.
- ▶ Structural annealing [Smith and Eisner, 2006]: gradually decrease structural biases.
- ▶ Curriculum learning [Spitkovsky et al. 2010; Tu et al. 2011]: start learning from short sentences; gradually increase training sentence length limit.
- ▶ Switching between different objectives [Spitkovsky et al. 2013].
- ▶ Treating different learning algorithms and configurations as modules and connecting them to form a network [Spitkovsky et al., 2013].
- ▶ Gibbs sampling [Johnson et al., 2007]: may incorporate constraints and biases, e.g., depth-bound [Jin et al., 2018a,b], subtree reducibility [Mareček and Žabokrtský, 2012; Mareček and Straka, 2013].

Flow-based Models: Improve Syntax with Semantics



- ▶ Jointly learn discrete syntactic structure and continuous word representations (semantic).
- ▶ Latent embedding e , pretrained embeddings x and invertible function $f_\phi(e)$.
- ▶ Training with normalizing flow helps induce better parsers.

[He et al. 2018, Jin et al. 2019]

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.
- ▶ A typical generative model in unsupervised dependency parsing: DMV

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.
- ▶ A typical generative model in unsupervised dependency parsing: DMV
- ▶ Three different kinds of representations: table-based, feature-based and neural based methods.

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.
- ▶ A typical generative model in unsupervised dependency parsing: DMV
- ▶ Three different kinds of representations: table-based, feature-based and neural based methods.
- ▶ EM algorithm is one of standard approach to learn the grammar parameter, which MLE objectives.

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.
- ▶ A typical generative model in unsupervised dependency parsing: DMV
- ▶ Three different kinds of representations: table-based, feature-based and neural based methods.
- ▶ EM algorithm is one of standard approach to learn the grammar parameter, which MLE objectives.
- ▶ Recently, directly gradient descent is more popular.

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.
- ▶ A typical generative model in unsupervised dependency parsing: DMV
- ▶ Three different kinds of representations: table-based, feature-based and neural based methods.
- ▶ EM algorithm is one of standard approach to learn the grammar parameter, which MLE objectives.
- ▶ Recently, directly gradient descent is more popular.
- ▶ Many improvements for MLE: viterbi likelihood, MAP estimation, contrastive estimation.

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.
- ▶ A typical generative model in unsupervised dependency parsing: DMV
- ▶ Three different kinds of representations: table-based, feature-based and neural based methods.
- ▶ EM algorithm is one of standard approach to learn the grammar parameter, which MLE objectives.
- ▶ Recently, directly gradient descent is more popular.
- ▶ Many improvements for MLE: viterbi likelihood, MAP estimation, contrastive estimation.
- ▶ Posterior regularization is a useful approach to encode knowledge.

Summary on Parameter Learning of Generative Approaches

- ▶ Structure learning is hard. Parameter learning is much easier and draws more attention.
- ▶ A typical generative model in unsupervised dependency parsing: DMV
- ▶ Three different kinds of representations: table-based, feature-based and neural based methods.
- ▶ EM algorithm is one of standard approach to learn the grammar parameter, which MLE objectives.
- ▶ Recently, directly gradient descent is more popular.
- ▶ Many improvements for MLE: viterbi likelihood, MAP estimation, contrastive estimation.
- ▶ Posterior regularization is a useful approach to encode knowledge.
- ▶ Many technologies can improve avoiding local minimum.