Introduction

• What are class-based language models (CLMs)?

- A large family of language models which use word classes
- "Similar words appear in similar context"
- Most of them maintain a hard-clustering assumption
- One word belongs to only one class
- For theoretical and computational simplicity

• We propose a Bayesian formulation of the CLM

- Naturally supports soft-clustering
- One word belongs to many classes with probabilities
- "Latent Hierarchical Pitman-Yor Allocation for Language Modeling"
- History as "document" in Latent Dirichlet Allocation (LDA) parlance • Cares more about language modeling (perplexity) than topic modeling (word classes)

Class-based Language Models (CLMs)

Definition

$$P(w \mid h) = P(c(w) \mid h)P(w \mid c(w), h)$$

- Finding classes
- Agglomerative clustering (Brown et al, 1992)
- Exchange-based clustering (Martin et al, 1998)
- Random clustering (Emami and Jelinek, 2005)
- Averaging many CLMs with word classes derived from randomly initialized exchange-based clustering gives better results.

$$P(w \mid h) = rac{1}{K} \sum_{k=1}^{K} P_k(c_k(w) \mid h) P_k(w \mid c_k(w))$$

Soft-clustering Class-based Language Models (SCLMs)

Definition

$$P(w \mid h) = \sum_{c \in \mathcal{C}} P(c \mid h) P(w \mid c, h)$$

Random sampling

► Let A denote a class assignment of all words in the training text.

► Given A, class model and word model are just regular LMs.

$$egin{aligned} \mathcal{P}(w\mid h) &= \sum_{A\in\mathcal{A}} \mathcal{P}(A) \sum_{c\in\mathcal{C}} \mathcal{P}_A(c\mid h) \mathcal{P}_A(w\mid c, c) \ &pprox rac{1}{K} \sum_{k=1}^K \sum_{c\in\mathcal{C}} \mathcal{P}_{A_k}(c\mid h) \mathcal{P}_{A_k}(w\mid c, h) \end{aligned}$$

Bayesian Class-based Language Models Yi Su Nuance Communications, Inc.



Soft-clustering Class-based Language Models: Inference

- Model does not subscribe itself to any inference algorithm.
- Collapsed Gibbs sampler comes natural to this model: $P(c_i = j \mid \mathbf{C}_{\neg i}, \mathbf{w}) \propto P(c_i = j \mid \mathbf{C}_{\neg i}, \mathbf{w}_{\neg i}) \cdot P(w_i \mid c_i = j, \mathbf{C}_{\neg i}, \mathbf{w}_{\neg i})$ (5)
- Two terms on the right hand side can be computed with "leave-one-out" versions of class and word models, respectively.
- Taking multiple samples is trivially parallelizable.

Fully Bayesian Formulation

Smoothing is to frequentists as prior is to Bayesians.

- Hierarchical Pitman-Yor (HPY) prior is the Bayesian counterpart of Kneser-Ney (KN)
- Or Kneser-Ney smoothing is a frequentist approximation of HPY.

Soft-clustering Class-based Hierarchical Pitman-Yor LMs (SCHPYLMs)

- Uses HPY LM in place of KN LM in the SCLM
- A fully Bayesian formulation of the CLM
- Plug in a Gibbs sampler of the HPY LM to get a sampler of the whole model

Fully Bayesian Formulation: Generative Process

1 . Choose parameters d_j and θ_j for $j \in \{0, 1, \cdots, n-1\}$:	
$d_j \sim \text{Beta}(1, 1)$ $\theta_j \sim \text{Gamma}(1, 1)$	(6)
2. Choose parameters e_j and μ_j for $j \in \{0, 1, \cdots, n\}$:	
$m{e}_j \sim ext{Beta}(1, 1)$ $\mu_j \sim ext{Gamma}(1, 1)$	(7)
3. For every word w_i and its history $h_i = w_{i-1} \cdots w_{i-n+1}$:	
3.1 Generate <i>c_i</i> :	
$egin{aligned} G_{\phi}(m{c}) &\sim \mathrm{PY}(m{d}_0, heta_0, m{G}_0(m{c})) \ G_{w_{i-1}}(m{c}) &\sim \mathrm{PY}(m{d}_1, heta_1, m{G}_{\phi}(m{c})) \ G_{w_{i-1}w_{i-2}}(m{c}) &\sim \mathrm{PY}(m{d}_2, heta_2, m{G}_{w_{i-1}}(m{c})) \end{aligned}$	
$c_i \mid h_i \sim \operatorname{Mult}(G_{h_i}(c)),$	(8)
3.2 Generate w _i :	
$egin{aligned} & H_{\phi}(m{w}) \sim \mathrm{PY}(m{e}_0, \mu_0, H_0(m{w})) \ & H_{c_i}(m{w}) \sim \mathrm{PY}(m{e}_1, \mu_1, H_{\phi}(m{w})) \ & H_{c_i w_{i-1}}(m{w}) \sim \mathrm{PY}(m{e}_2, \mu_2, H_{c_i}(m{w})) \ & H_{c_i w_{i-1} w_{i-2}}(m{w}) \sim \mathrm{PY}(m{e}_3, \mu_3, H_{c_i w_{i-1}}(m{w})) \end{aligned}$	
$w_i \mid c_i, h_i \sim \operatorname{Mult}(H_{c_i h_i}(w)).$	(9)

Experiments: Perplexity



Figure: Perplexity as function of the number of classes

Experiments: Word Error Rate (WER)

- Setup
- Lattice-rescoring
- IBM 2004 Rich Transo **Conversational Telep** Speech (RT-04 CTS)
- 20M words Fisher dat
- ► 30K vocabulary
- Interpolation with a big LM built from many other sources
- Results
- Follow the same trend as perplexity
- Match best published results on the same setup
- WER reduction significant at 0.001

Conclusions

- We proposed

- 22% perplexity reduction on WSJ ▶ 6% WER reduction on IBM RT-04 CTS
- for stimulating discussions.



NUANCE

Setup

- IM words Wall Street Journal (WSJ)
- ► 10K vocabulary, 3-grams
- ▶ 100 samples per experiment
- 800 iterations of burn-in
- Results
- Perplexity decreases then increases, as expected
- Soft-clustering beats hard-clustering
- HPY beats KN

	Model	w/o Interp.	w/ Interp.
scription	Kneser-Ney	14.4	13.5
	CLM	14.2	13.4
	SCLM	13.7	13.3
system	Hier. Pitman-Yor	14.1	13.4
ita	CHPYLM	13.7	13.2
	SCHPYLM	13.5	13.1

Table: Lattice-rescoring WERs on IBM RT-04 CTS

Soft-clustering Class-based Hierarchical Pitman-Yor Language Models A simple collapsed Gibbs sampler for inference

Great performance in perplexity and word error rate

Many thanks to Max Bisani, Ernie Pusateri, Wilson Tam and Paul Vozila