

Statistical Language Model Adaptation

Reference:

J. R. Bellegarda, “Statistical language model adaptation: review and perspectives” , 2003

Outline

- Introduction
- Adaptation techniques
- Conclusions

Why adaptation?

- Natural language is highly variable in several aspects
 - Language evolves as does the world it seeks to describe
 - The vocabulary changes dynamically with time
 - “protemics” to the utter demise of “ague”
 - Different domains tend to involve relatively disjoint concepts with markedly different word sequence statistics
 - Subject matter affects the underlying semantic characteristics
 - Interest rate to banking application vs. a general conversation on gaming platforms

Why adaptation? (cont.)

- People naturally adjust their use of the language based on the task at hand
 - Typical syntax or grammatical infrastructure
 - Formal technical papers vs. casual e-mails
- People style of discourse may independently vary due to a variety of factors such as socio-economics, emotional state
 - More pronounced on spoken natural language

Why adaptation? (cont.)

- Lexical, syntactic or semantic characteristics of the discourse in the training and recognition tasks are quite likely differ
- The performance of any statistical approach, i.e. n-gram modeling, always suffers from such mismatch
- SLMs is extremely brittle across domains and even within domain when training and recognition domain involve moderately disjoint time periods
- Adaptation techniques
 - Model adaptation
 - Constraint specification
 - Meta-information extraction

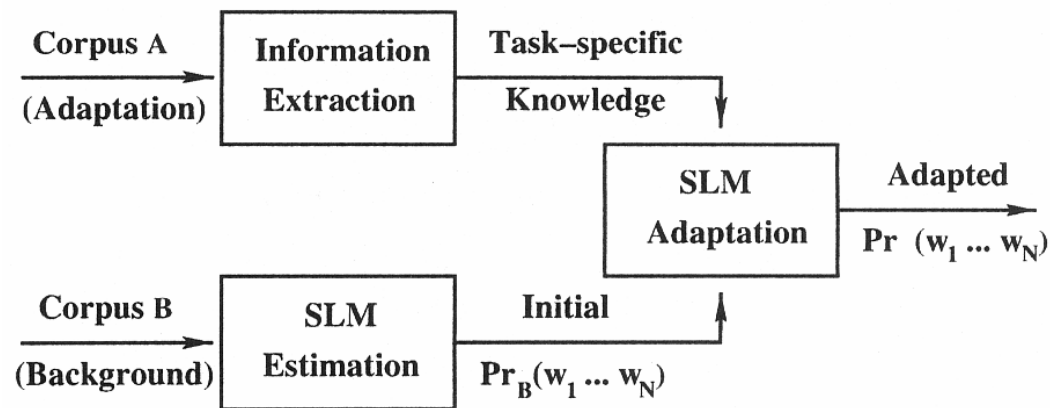
Adaptation framework

- LM estimate:

$$\Pr(w_1, \dots, w_N) = \prod_{q=1}^N \Pr(w_q | h_q)$$

- N-gram model: $h = w_{q-n+1}, \dots, w_{q-1}$

- Framework for SLM adaptation:



- The general idea is to dynamically modify the background SLM estimate on the basis of what can be extracted from A

Adaptation data

- Corpus A may already be available
 - Domain-specific texts for some purpose
 - Wizard of Oz experiments

- Corpus A is not available
 - Use Monte Carlo methods for choosing rules at the branching points of the grammar and then create an artificial corpus
 - Use small amount of adaptation data to weigh the different rules and retain more realistic corpus (bootstrap SLM)

Adaptation data (cont.)

- Use multiple sentence hypotheses weighted by posteriori likelihood from an n-best list
 - The effect of a recognition error is distributed over several competing hypotheses
- sausages
- Once enough adaptation material has been accumulated to suitable characterize the new domain
 - To gather similar data by Information Retrieval
- Data augmentation
 - Adjust the count of corpus

Data augmentation

- D. Janiszek, et al, “DATA AUGMENTATION AND LANGUAGE MODEL ADAPTATION”, 2001
- represent words as vectors after SVD
 - Augmented count

$$a_{ij} = c_{ij} + \sum_{h_k} c_{ik} \times f(d_{jk}^\alpha) \text{ with } h_k \in \Gamma_j^\alpha(\theta), f(d_{jk}^\alpha) = e^{-\frac{d_{jk}^\alpha}{D}}$$

a_{ij} : augmented count for word i with history j

c_{ij} : original count for word i with history j

d_{jk}^α : distance of history j and k

$\Gamma_j^\alpha(\theta)$: set of histories having distance lower than threshold θ with history j in the space S_α

Model interpolation

- Corpus A is used to derive a task-specific (dynamic) SLM which is then combined with the background (static) SLM
 - Provides fertile grounds
- Model interpolation
 - Model merging
 - Dynamic cache models
 - MAP adaptation

Model merging

- Only for certain idiosyncratic word sequences, particularly frequent in the current task, may the dynamic model outperform the initial estimate SLM
 - Take advantage of the new information as appropriate

- Linear interpolation

$$\Pr(w_q | h_q) = (1 - \lambda) \Pr_A(w_q | h_q) + \lambda \Pr_B(w_q | h_q)$$

- Back-off

$$\Pr(w_q | h_q) = \begin{cases} \Pr_A(w_q | h_q) & \text{if } C_A(hq, wq) \geq T \\ \beta \Pr_B(w_q | h_q) & \text{otherwise} \end{cases}$$

Dynamic Cache models

- A special case of linear interpolation
- Cache models exploit self-triggering words inside the corpus A to capture short-term shifts in word-use frequencies
 - Unigram case of general model merging strategy
- To propagate the power of the method
 - Class n-gram model

$$\Pr(w_q | h_q) = \sum_{\{c_q\}} \Pr(w_q | c_q) \Pr(c_q | h_q)$$

MAP adaptation

- It is argued that combination should be done at the frequency count level rather than the model level

$$P(W | \theta)P(\theta) \propto \prod_{k=1}^{|\mathcal{h}|} \prod_{i=1}^{|\mathcal{v}|} \phi_{h_k, w_i}^{C_{h_k, w_i} + \nu_{h_k, w_i} - 1} \quad \text{Note: Multinomial * Dirichlet}$$

$$\log P(W | \theta)P(\theta) \propto \sum_{k=1}^{|\mathcal{h}|} \sum_{i=1}^{|\mathcal{v}|} (C_{h_k, w_i} + \nu_{h_k, w_i} - 1) \log \phi_{h_k, w_i} + \sum_{k=1}^{|\mathcal{h}|} l_{h_k} \left(\sum_{i=1}^{|\mathcal{v}|} \phi_{h_k, w_i} - 1 \right)$$

$$\phi_{h_k, w_i} = \frac{C_{h_k, w_i} + \nu_{h_k, w_i} - 1}{\sum_{j=1}^{|\mathcal{v}|} (C_{h_k, w_j} + \nu_{h_k, w_j} - 1)} \quad , \nu_{h_k, w_i} = \frac{1}{\varepsilon} C_B(h_k, w_i) + 1$$

- Count merging

$$\Pr(w_q | h_q) = \frac{\varepsilon C_A(h_q, w_q) + C_B(h_q, w_q)}{\varepsilon C_A(h_q) + C_B(h_q)}$$

Constraint specification

- Extract feature from corpus A that the adapted SLM is constrained to satisfy

- Exponential models
 - MDI adaptation
 - Unigram constraints

MDI adaptation

- The feature extracted from corpus A are considered as important properties
- Minimize KL distance from background distribution

$$\min_{Q(h,w)} \sum_{\{(h,w)\}} Q(h,w) \log \frac{Q(h,w)}{\Pr_B(h,w)}$$

which satisfies

$$\sum_{\{(h,w)\}} I_k(h,w) Q(h,w) = \alpha_A(\hat{h}_k, \hat{w}_k), 1 \leq k \leq \mathbf{K}$$

- Take joint background distribution as prior distribution

$$\Pr(h,w) = \frac{\Pr_B(h,w)}{Z(h,w)} \prod_{k=1}^K \exp\{\lambda_k I_k(h,w)\}$$

Unigram Constraints

- Special case of MDI adaptation
 - Only unigram features are reliably estimated

$$\sum_{\{(h,w)\}} I_k(h,w) Q(h,w) = \alpha_A(\hat{w}_k), 1 \leq k \leq \mathbf{K}$$

- Similar to dynamic cache model
- GIS

$$\Pr(h,w) = \Pr_B(h,w) \frac{\alpha_A(w)}{\Pr_B(w)}$$

- The adapted SLM is simply the background SLM scaled by the scaling factor

Topic Information

- Exploit the information about the underlying subject matter from corpus A
- Mixture models
 - Linear interpolate K n-gram that the resulting mixture best matches the adaptation data A

$$\Pr(w_q | h_q) = \sum_{k=1}^K \lambda_{A,k} \Pr_{B,k}(w_q | h_q)$$

- One of these background model is trained on the entire corpus B
- count merge from mixture model with different weight of topics
- Drawback: inherent fragmentation of training data

Explicit topic models

- Mixture modeling includes topic information indirectly
 - Express the topic contribution more directly
 - No assumption that each history belongs to exactly one topic
 - Conditional independence assumption on word and topic

$$\Pr(w_q | h_q) = \sum_{k=1}^K \Pr(w_q | t_k) \Pr(t_k | h_q)$$

- Topic n-gram is assumed to be unaffected
- The topic assignment is adapted

$$\Pr(w_q | t_k) = (1 - \lambda) \Pr_A(w_q | t_k) + \lambda \Pr_B(w_q | t_k)$$

Trigger

- Rosenfeld, “A maximum entropy approach to adaptive statistical language model”, 1996
- If a word sequence A is significantly correlated with another word sequence B
 - $(A \rightarrow B)$ is considered a “trigger pair”, with A being the trigger and B the triggered sequence
 - If A, B are single words, possible triggers are more than bigram
- Correlation measure
 - Cross product ratio (not enough in determining the utility)
 - Average utility ? (the more common, the higher)
 - Average mutual information

Trigger (cont.)

- Average mutual information

$$I(A_o : B) = P(A_o, B) \log \frac{P(B | A_o)}{P(B)} + P(A_o, \bar{B}) \log \frac{P(\bar{B} | A_o)}{P(\bar{B})} \\ + P(\bar{A}_o, B) \log \frac{P(B | \bar{A}_o)}{P(B)} + P(\bar{A}_o, \bar{B}) \log \frac{P(\bar{B} | \bar{A}_o)}{P(\bar{B})}$$

Or
$$I(A_o : B) = \log \frac{P(A_o, B)}{P(A_o)P(B)}$$

– Implementation :

- set the window size (Church et al., Word Association Norms, Mutual Information, and Lexicography, 1990)

Trigger (cont.)

- Carlos Troncoso et al, “Trigger-Based Language Model Adaptation for Automatic Meeting Transcription”, 2005
- Focus on the transcription of meetings and lectures rather than newspaper articles
 - the subject is more homogeneous in meeting session (document)
 - the topic of newspaper is too general to extract task-specific triggers
- Extract trigger from initial speech recognition results

Trigger (cont.)

- Selection methods
 - TD/IDF

$$v_{ik} = \frac{tf_{ik} \log(N / df_k)}{\sqrt{\sum_{j=1}^T (tf_{ij})^2 [\log(N / df_k)]^2}}$$

- The TF part is calculated from the K-best hypotheses
 - The IDF part is computed from a fraction of a large corpus
- Create all possible word pairs including self-triggers and then use POS-based filter to discard function words

Trigger (cont.)

- Probability of trigger pairs

$$P_{TP}^{IT}(w_2 | w_1) = \frac{N(w_1, w_2)}{\sum_j N(w_1, w_j)}$$

- Proposed trigger-based language model

$$P_{LM}(w_i | H) = \frac{1}{|H|} \sum_{h \in H} P_{LM}(w_i | h)$$

$$P_{LM}(w_i | H) = \begin{cases} P_{NG}(w_i | H), & \text{if } P_{TP}^{IT}(w_j | h) = 0, \forall j \\ \lambda P_{NG}(w_i | H) + (1 - \lambda) P_{TP}^{IT}(w_j | h), & \text{else} \end{cases}$$

- Proposed back-off method

$$P_{LM}(w_i | H) = \lambda P_{NG}(w_i | H) + (1 - \lambda) (\delta P_{TP}^{LC}(w_i | h) + (1 - \delta) P_{TP}^{IT}(w_i | h))$$

Latent Semantic Analysis

- Bellegarda, “Exploiting latent semantic information in statistical language modeling”, 2000
- The paradigm relies on the concept of a document, i.e., a “bag-of-words” entity forming a semantically homogeneous unit
- The resulting semantic knowledge is encapsulated in a continuous vector space

Latent Semantic Analysis

- Step 1: Build a matrix of co-occurrences between words and documents
 - Normalized entropy

$$\varepsilon_i = -\frac{1}{\log N} \sum_{j=1}^N \frac{c_{i,j}}{t_i} \log \frac{c_{i,j}}{t_i}$$

t_i : the total number of times of w_i occurs in Corpus

$c_{i,j}$: the total number of times of w_i occurs in document j

- Cell of Matrix W

$$w_{i,j} = (1 - \varepsilon_i) \frac{c_{i,j}}{n_j}$$

Latent Semantic Analysis (cont.)

- Step 2: Singular Value Decomposition

$$W \approx \hat{W} = USV^T$$

U : left singular matrix, $U^T U = I_R$

V : right singular matrix, $V^T V = I_R$

S : diagonal matrix, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_R \geq 0$

row vector $u_i S$, $v_j S$ characterize

the position of word w_i and document d_j in the space

Latent Semantic Analysis (cont.)

- Step 3: Create pseudo document vector

- SVD expansion: $\tilde{d}_p = US\tilde{v}_p^T$

- pseudo document vector

$$\tilde{\tilde{v}}_p = \tilde{v}_p S = \tilde{d}_p^T U$$

- If the new document contains language patterns which are inconsistent with those extracted from W , the SVD expansion will no longer apply

Latent Semantic Analysis (cont.)

- N-gram + LSA Language Modeling
 - The closeness between word and pseudo document vector

$$\Pr_{LSA}(w_q | H_{q-1}, S) = \Pr_{LSA}(w_q | \tilde{d}_{q-1})$$

- \tilde{d}_{q-1} and \tilde{d}_q differ only in one coordinate and the entropy does not change appreciably

$$\tilde{d}_q = \frac{n_q - 1}{n_q} \tilde{d}_{q-1} + \frac{1 - \varepsilon_i}{n_q} [0 \dots 1 \dots 0]^T$$

- Pseudo document vector can be efficiently updated directly in the LSA space

$$\tilde{v}_q = \tilde{v}_q S = \frac{1}{n_q} [(n_q - 1) \tilde{v}_{q-1} + (1 - \varepsilon_i) u_i]$$

Latent Semantic Analysis (cont.)

- LSA probability

- Closeness

$$K(w_q, \tilde{d}_{q-1}) = \cos(u_q S^{1/2}, \tilde{v}_{q-1} S^{1/2}) = \frac{u_q S \tilde{v}_{q-1}^T}{\|u_q S^{1/2}\| \|\tilde{v}_{q-1} S^{1/2}\|}$$

- Probability

$$\Pr_{LSA}(w_q | \tilde{d}_{q-1}) = \frac{\pi - \mathbf{acos}(K)}{\pi}$$

- Integration with N-grams

$$\Pr(w_q | H_{q-1}^{(n+l)}) = \frac{\Pr_{N\text{-gram}}(w_q | w_{q-n+1}, \dots, w_{q-1}) \frac{\Pr_{LSA}(w_q | d_{q-1})}{\Pr_{N\text{-gram}}(w_q)}}{\sum_{w_i \in V} \Pr_{N\text{-gram}}(w_i | w_{i-n+1}, \dots, w_{i-1}) \frac{\Pr_{LSA}(w_i | d_{q-1})}{\Pr_{N\text{-gram}}(w_i)}}$$

Latent Semantic Analysis (cont.)

- Adaptation:

$$\frac{\Pr_{LSA}(w_q | d_{q-1})}{\Pr_{N-gram}(w_q)} = (1 - \lambda) \frac{\Pr_{A-LSA}(w_q | d_{q-1})}{\Pr_{A-N-gram}(w_q)} + \lambda \frac{\Pr_{B-LSA}(w_q | d_{q-1})}{\Pr_{B-N-gram}(w_q)}$$

- Other issues

- History size or importance of history

$$\tilde{v}_q = \tilde{v}_q S = \frac{1}{n_q} \left[(n_q - 1) \tilde{v}_{q-1} + (1 - \varepsilon_i) u_i \right]$$

- Word Smoothing

$$\Pr_{LSA}(w_q | \tilde{d}_{q-1}) = \sum_{k=1}^K \Pr(w_q | C_k) \Pr(C_k | \tilde{d}_{q-1})$$

- Document Smoothing

$$\Pr_{LSA}(w_q | \tilde{d}_{q-1}) = \sum_{l=1}^L \Pr(w_q | D_l) \Pr(D_l | \tilde{d}_{q-1})$$

Latent Dirichlet Allocation

- Yik-Cheung Tam et al. ,”Dynamic Language Model Adaptation using Variational Bayes Inference”, 2005
- In broadcast news, a document usually refers to a piece of news story within which the latent topics are consistent
- LDA is a Bayesian extension of a mixture of unigram models where the topic mixture weight θ is drawn from a prior Dirichlet distribution

$$f(\theta; \alpha) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

- LDA makes an assumption in which the words in a document are conditionally independent given the latent topic sequence

$$f(w_1^n | \theta) = \prod_{i=1}^n \sum_{k=1}^K f(w_i | z_i = k) \cdot f(z_i = k | \theta) = \prod_{i=1}^n \sum_{k=1}^K \beta_{w_i k} \cdot \theta_k$$

Latent Dirichlet Allocation (cont.)

- A document is generated by firstly sampling a mixture weight θ from its prior distribution and then repeatedly sample a topic k from θ and a word from the k -th latent unigram until all words in the document are generated

$$f(w_1^n) = \int_{\theta} f(w_1^n | \theta) \cdot f(\theta; \alpha) d\theta$$

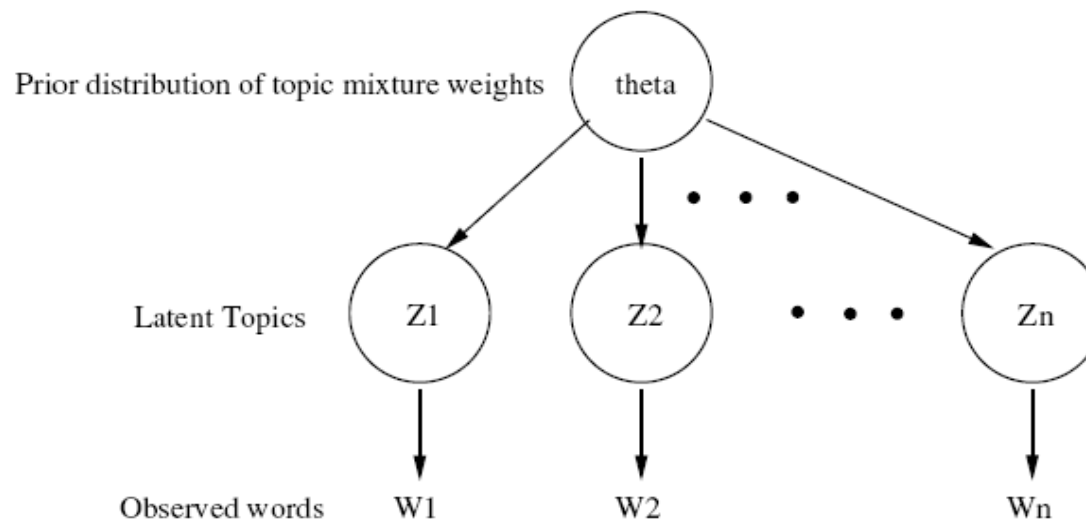


Figure 1: Graphical representation of Latent Dirichlet Allocation.

Variational Bayes approach (cont.)

- Jensen's inequality

$$\log \sum_i q_i \cdot \frac{f_i}{q_i} \geq \sum_i q_i \cdot \log \frac{f_i}{q_i} = E_q \left[\log \frac{f(\cdot)}{q(\cdot)} \right] \text{ where } \sum_i q_i = 1$$

- Optimizing the exact likelihood is computationally intractable
 - optimize the lower bound

$$\begin{aligned} \log f(w_1^n) &= \log \sum_{\theta, z_1^n} q(\theta, z_1^n; \Gamma) \frac{f(\theta, w_1^n, z_1^n; \Lambda)}{q(\theta, z_1^n; \Gamma)} \\ &\geq \sum_{\theta, z_1^n} q(\theta, z_1^n; \Gamma) \log \frac{f(\theta, w_1^n, z_1^n; \Lambda)}{q(\theta, z_1^n; \Gamma)} = E_q \left[\log \frac{f(\theta, w_1^n, z_1^n; \Lambda)}{q(\theta, z_1^n; \Gamma)} \right] = Q(\Lambda, \Gamma) \end{aligned}$$

- To choose the tractable variational distribution that the integration can be done in a tractable manner (interpret q an approximation to the posteriori distribution over f)

$$\Lambda = \{ \{ \alpha_k \}, \{ \beta_{vk} \} \}, \quad q(\theta, z_1^n; \Gamma) = q(\theta) \cdot \prod_{i=1}^n q(z_i)$$

Dirichlet * multinomial

Variational Bayes approach (cont.)

- Parameters for a single document

$$\gamma_k = \alpha_k + \sum_{i=1}^n q(z_i = k)$$

$$q(z_i = k) \propto \beta_{w_i k} \cdot e^{E_q[\log \theta_k]}$$

$$\beta_{vk} \propto \sum_{i=1}^n q(z_i = k) \delta(w_i, v)$$

- Adaptation

- Compute the Maximum A-posterior likelihood

$$f(w|h) \approx \int_{\theta} \sum_{k=1}^K f(w|z=k) f(z=k|\theta) q(\theta|h) \approx \sum_{k=1}^K \beta_{wk} \cdot \hat{\theta}_k, \hat{\theta}_k = \frac{\gamma_k}{\sum_{k=1}^K \gamma_k}$$

- $\hat{\theta}$ is the mode of the variational posteriori

Syntactic infrastructure

- Background and recognition tasks share a common grammatical infrastructure
 - Grammatical constraints are portable from corpus B to corpus A
- Structured language models
 - Hierarchical framework by using syntactic information to determine equivalence classes on the n-gram history

$$\Pr(w_q | \bar{h}_q) = \frac{1}{Z(\bar{h}_q)} \sum_{\{\pi_q\}} \Pr(w_q | \bar{h}_q, \pi_q) \Pr(\bar{h}_q, \pi_q)$$

- Each parse gives rise to a unique headword history, and it is expedient to simplify the model by limit to n-1 headwords
- Main caveat in structured language modeling remain the reliance on the parser, and particularly the implicit assumption that the correct parse will be assigned a high probability

Syntactic infrastructure (cont.)

- Syntactic triggers
 - Exploit syntactic structure contained in previous sentences
 - Full parse or the syntactic/semantic tag

- Although not yet implemented in an adaptation context, this concept may ultimately provide the necessary framework to extend the benefits of structured language modeling to a span greater than that of a sentence

Multiple sources

- Whole sentence models
 - It has been argued hindrance to modeling linguistic supra-structure such as person and number agreement, semantic coherence, and even length
 - Adopt a “bag-of-features” approach to each sentence

$$\Pr(\sigma) = \frac{\Pr_0(\sigma)}{Z} \prod_{k=1}^K \exp\{\lambda_k I_k(\sigma)\}$$

- Normalization is infeasible since it involves summation

Conclusions

- Language model adaptation refers to the process of exploiting specific, albeit limited, knowledge about recognition task to compensate for any mismatch between and recognition
- Dynamically modifying the language model statistics according to this information
- Several techniques are surveyed