### Language Modeling for Information Retrieval

#### Ananthakrishnan R (anand@cse)

### **Outline of the Seminar**

- What is Language Modeling?
- Language Modeling in NLP
  - Speech Recognition, Statistical MT, Classification, IR ...
- The Language Modeling approach to IR
  - the conceptual model
  - smoothing
  - semantic smoothing
    - IR as Statistical Machine Translation

#### Language Modeling

- A *language model* is a probabilistic mechanism for generating text
- Language models estimate the probability distribution of various natural language phenomena

- sentences, utterances, queries ...

## Language Modeling Techniques

- N-grams
- Class-based N-grams
- Probabilistic CFGs
- Decision Trees

**Claude Shannon** was probably the first language modeler "How well can we predict or compress natural language text through simple n-gram models?" Language Modeling for Other NLP Tasks

#### **Applications**

- Speech Recognition
- Statistical Machine Translation
- Classification
- Spell Checking

Play the role of either the prior or the likelihood

# Language Models in Speech Recognition



## Language Models in Statistical Machine Translation



#### Language Models in Text Classification



a language model for each class

Bayesian Classifier

#### The Language Modeling Approach to Information Retrieval

#### The Conceptual Model of IR

• The user has an information need  $\mathcal{I}$ 

- From this need he generates an ideal document fragment d
- He selects a set of key terms from d<sub>y</sub>, and generates a query q from this set





#### The Language Modeling Approach

Using Bayes' law,

$$p(\mathbf{d} \mid \mathbf{q}, \mathcal{U}) = \frac{p(\mathbf{q} \mid \mathbf{d}, \mathcal{U}) p(\mathbf{d} \mid \mathcal{U})}{p(\mathbf{q} \mid \mathcal{U})}$$

For the purpose of ranking,

 $p_{\mathbf{q}}(\mathbf{d}) = p(\mathbf{q} | \mathbf{d}, \mathcal{U}) p(\mathbf{d} | \mathcal{U})$ 

$$p_q(d) = p(q | d) p(d)$$

(user independent)

Language Modeling for IR: Using document language models to assign likelihood scores to queries (Ponte and Croft, 1998)





#### The Language Model

 For each document d in the collection C, for each term w, estimate:



$$p(\mathbf{q} \mid \mathbf{d}) = \prod_{i=1}^{m} p(q_i \mid \mathbf{d})$$



#### The Need for Smoothing

To eliminate "zero" probabilities

 Correct the error in the MLE due to the problem of *data sparsity*



## Smoothing

 Jelinek-Mercer Method: Linear interpolation of the maximum likelihood model with the collection model

$$p_{\lambda}(\mathbf{w} \mid \mathbf{d}) = (1 - \lambda) p(\mathbf{w} \mid \mathbf{d}) + \lambda p(\mathbf{w} \mid \mathcal{C})$$

 (Zhai & Lafferty, 2003) compares Jelinek-Mercer, Dirichlet and the absolute discounting methods for Language Models applied to IR.

#### The Need for Semantic Smoothing



#### **Document-Query Translation Model**



## Semantic Smoothing

Estimate translation probabilities  $t(w_q|w_d)$  for mapping a document term  $w_d$  to a query term  $w_q$ 

$$p(\mathbf{q} \mid \mathbf{d}) = \prod_{i=1}^{m} \sum_{\mathbf{w}} t(q_i \mid \mathbf{w}) p(\mathbf{w} \mid \mathbf{d})$$

(Berger and Lafferty, 1999) (Lafferty and Zhai, 2001)



*p*(*analysis* | review)*p*(*review* | d)

## Compare with The Statistical MT model



### Language Modeling *vis-à-vis the* 'traditional probabilistic model'

$$p(r | q, d) = ?$$

The Probability Ranking Principle (Robertson, 1977)

$$\log \frac{p(d | q, r)}{p(d | q, \overline{r})}$$

The Robertson-Sparck Jones Model (Sparck Jones et al., 2000)

Language Modeling approach:

- Where's the relevance?
- Is the ranking optimal?

(Lafferty & Zhai, 2002) (Laverenko and Croft, 2001)

# Discussion: Advantages of the Language Modeling Approach

- Conditioning on d provides a larger 'foothold' for estimation
- p(d): an explicit notion of the importance of a document
- Document normalization is not an issue

"When designing a statistical model ... the most natural route is to apply a generative model which builds up the output step-by-step. The source channel perspective suggests a different approach: turn the search problem around to predict the input. In speech recognition, NLP, and MT, researchers have time and again found that predicting what is known from competing hypothesis can be easier than directly predicting all of the hypothesis"

# Discussion: Disadvantage of the Language Modeling Approach

 The Robertson-Sparck Jones model can directly use relevance judgements to improve the estimation of p(A<sub>i</sub> | q, r)

This is not possible in the Language Modeling approach

#### References

- (Robertson, 1977) The probability ranking principle in IR
- (Sparck Jones et al., 2000) A probabilistic model of information retrieval: development and comparative experiments.
- (Ponte and Croft, 1998) A language modeling approach to information retrieval
- (Zhai and Lafferty, 2001) A study of smoothing methods for language models applied to ad hoc information retrieval

#### References ...

- (Berger and Lafferty, 1999) Information retrieval as statistical translation
- (Lafferty and Zhai, 2001) Risk minimization and language modeling in information retrieval
- (Lavrenko and Croft, 2001) Relevance-based language models
- (Lafferty and Zhai, 2001) Probabilistic IR models based on document and query generation