

Language Models & Smoothing

$$p(W = w) = p(W_1 = w_1, W_2 = w_2, \dots, W_{L+1} = \text{stop})$$

$$= \left(\prod_{\ell=1}^L p(W_\ell = w_\ell \mid \mathbf{w}_{1:\ell-1} = \mathbf{w}_{1:\ell-1}) \right) p(W_{L+1} = \text{stop} \mid \mathbf{w}_{1:L} = \mathbf{w}_{1:L})$$

$$= \left(\prod_{\ell=1}^L p(w_\ell \mid \text{history}_\ell) \right) p(\text{stop} \mid \text{history}_L)$$

$$\text{perplexity}(p(\cdot); \mathbf{w}) = 2^{\left(-\frac{\log_2 p(\mathbf{w})}{|\mathbf{w}|}\right)}$$

$$= p(\mathbf{w})^{-\frac{1}{|\mathbf{w}|}}$$

$$= \sqrt[|\mathbf{w}|]{\frac{1}{p(\mathbf{w})}}$$

$$= \sqrt[|\mathbf{w}|]{\frac{1}{\prod_{i=1}^{|\mathbf{w}|} p(w_i \mid w_{i-N+1}, \dots, w_{i-1})}}$$

Maximum likelihood (“relative frequency”) estimation:

$$\hat{p}_{\text{MLE}}(w_N \mid \langle w_1, \dots, w_{N-1} \rangle) = \frac{\text{count}(\langle w_1, w_2, \dots, w_{N-1}, w_N \rangle)}{\sum_{v \in \Sigma} \text{count}(\langle w_1, w_2, \dots, w_{N-1}, v \rangle)}$$

Add- λ smoothing (commonly, $\lambda = 1$, in which case this is also called Laplace smoothing):

$$\hat{p}_{\text{add-}\lambda}(w_N \mid \langle w_1, \dots, w_{N-1} \rangle) = \frac{\lambda + \text{count}(\langle w_1, w_2, \dots, w_{N-1}, w_N \rangle)}{\lambda |\Sigma| + \sum_{v \in \Sigma} \text{count}(\langle w_1, w_2, \dots, w_{N-1}, v \rangle)}$$

Good-Turing discounted counts:

$$c^* = \frac{(c+1) \times N_{c+1}}{N_c}$$

Linear interpolation (“mixture model”):

$$\hat{p}_{\text{interp}}(w_N \mid \langle w_1, \dots, w_{N-1} \rangle) = \lambda_N \times \hat{p}_{\text{MLE}}(w_N \mid \langle w_1, \dots, w_{N-1} \rangle)$$

$$+ \lambda_{N-1} \times \hat{p}_{\text{MLE}}(w_N \mid \langle w_2, \dots, w_{N-1} \rangle)$$

$$\vdots$$

$$+ \lambda_2 \times \hat{p}_{\text{MLE}}(w_N \mid \langle w_{N-1} \rangle)$$

$$+ \lambda_1 \times \hat{p}_{\text{MLE}}(w_N)$$

$$+ \lambda_0 \times \frac{1}{|\Sigma|}$$

Noisy Channel & Edit Distance

$$D_{0,0} = 0$$

$$D_{i,j} = \min \begin{cases} D_{i-1,j} + \text{inscost}(t_i) \\ D_{i,j-1} + \text{delcost}(s_j) \\ D_{i-1,j-1} + \text{substcost}(t_i, s_j) \end{cases}$$

$$D_{0,0} = 0$$

$$D_{i,j} = \min \begin{cases} D_{i-1,j} + \text{inscost}(t_i) \\ D_{i,j-1} + \text{delcost}(s_j) \\ D_{i-1,j-1} + \text{substcost}(t_i, s_j) \\ D_{i-2,j-2} + \text{transcost}(s_{j-1}, s_j) \text{ if } s_{j-1} = t_i \text{ and } s_j = t_{i-1} \end{cases}$$

+ Noisy Channel

The Noisy Channel Model assumes we want to recover the most likely system output y for an observed input x by modeling the probability distribution $p(\mathbf{y} \mid \mathbf{x})$

$$y^* = \arg \max_y p(y \mid x)$$

$$= \arg \max_y \frac{p(x \mid y) \times p(y)}{p(x)} \quad \text{Bayes' rule}$$

$$= \arg \max_y \frac{p(x \mid y) \times p(y)}{\sum_{y'} p(x \mid y') \times p(y')}$$

$$= \arg \max_y p(x \mid y) \times p(y)$$

Classification

+ Naïve Bayes Classifier

$$\phi_j \leftarrow [\Phi(x)]_j$$

$$\text{return}$$

$$\arg \max_{y'} p(y') \times \prod_j p(\phi_j \mid y')$$

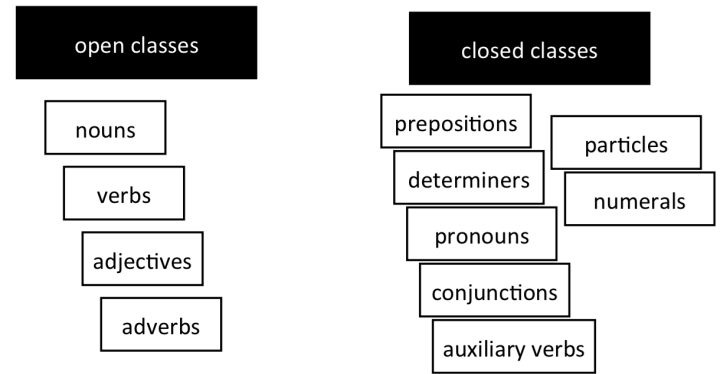
+ Naïve Bayes Learner

$$\mathbf{X} \rightarrow \left\{ \begin{array}{l} \forall y, p(y) \leftarrow \frac{\text{count}(y)}{N} \\ \forall y, \forall j, \forall f, p(\phi_j(x) = f \mid y) \leftarrow \frac{\text{count}(f, y)}{\text{count}(y)} \end{array} \right. \rightarrow p$$

Sentiment Analysis

Social Media (tweets), product reviews, discussion forum posts, and blog posts

Part of Speech Tags



CC	Conj.	PRP	Pronoun, personal
CD	Numeral, Cardinal	RB	Adverb
DT	Determiner	RBR	Adverb, comparative
IN	Prep., Conj.	RP	Particle
JJ	Adj. or Numeral	VB	Verb, base form
NN	Noun, common, s.	VBD	Verb, past tense
NNP	Noun, proper, s.	VBG	V., present participle
NNPS	Noun, proper, pl.	VBN	V., past participle
NNS	Noun, common, pl.	VBP	V., present, not 3rd

Two difficulties:

- 1) What tags? Unknown words?
- 2) Disambiguating words (more than one tag)

Hidden Markov Models

+ Markov Property: the state we move to at time t depends only on the state we were at time t-1

+ Symbols are observed, state sequence is hidden

- q_0 : start state (“silent”)
- q_f : final state (“silent”)
- Q : set of “normal” states (excludes q_0 and final q_f)
- Σ : vocabulary of observable symbols
- γ_{ij} : probability of transitioning to q_j given current state q_i
- $\eta_{i,w}$: probability of emitting $w \in \Sigma$ given current state q_i

+ Viterbi Algorithm

$$y_n^* = \arg \max_{q_i \in Q} p(Y_1 = y_1^*, \dots, Y_{n-1} = y_{n-1}^*, Y_n = q_i \mid x)$$

$$= \arg \max_{q_i \in Q} V[n-1, y_{n-1}^*] \cdot \gamma_{y_{n-1}^*, i} \cdot \eta_{i, x_n} \cdot \gamma_{i, f}$$

$$= \arg \max_{q_i \in Q} \gamma_{y_{n-1}^*, i} \cdot \eta_{i, x_n} \cdot \gamma_{i, f}$$

$$V[0, q_0] = 1$$

$$V[t, q_j] = \max_{q_i \in Q \cup \{q_0\}} V[t-1, q_i] \cdot \gamma_{i, j} \cdot \eta_{j, x_t}$$

$$\text{goal} = \max_{q_i \in Q} V[n, q_i] \cdot \gamma_{i, f}$$

$$V[*, *] \leftarrow 0$$

$$V[0, q_0] \leftarrow 1$$

for $t = 1 \dots n$

foreach q_j

foreach q_i

$$V[t, q_j] \leftarrow \max\{V[t, q_j], V[t-1, q_i] \times \gamma_{i, j} \times \eta_{j, x_t}\}$$

foreach q_i

$$\text{goal} \leftarrow \max\{\text{goal}, V[n, q_i] \times \gamma_{i, f}\}$$

return goal

Syntactic Representation

- Vocabulary of terminal symbols, Σ
- Set of nonterminal symbols (a.k.a. variables), N
- Special start symbol $S \in N$
- Production rules of the form $X \rightarrow \alpha$

where

$X \in N$

$\alpha \in (N \cup \Sigma)^*$

Grammatical: said of a sentence in the language

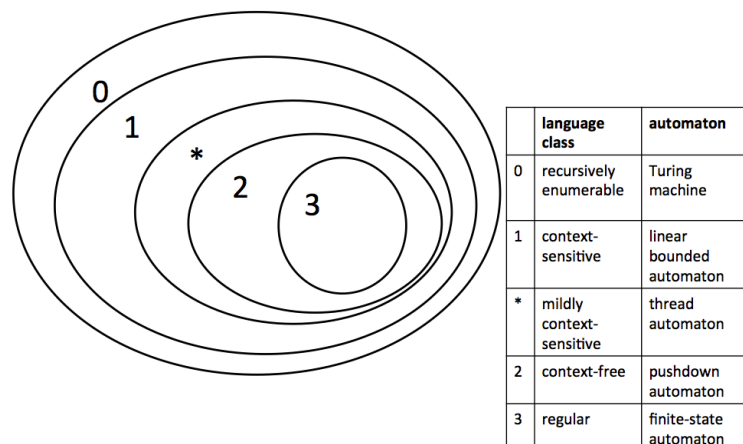
Ungrammatical: said of a sentence **not** in the language

Derivation: sequence of top-down production steps

Parse tree: graphical representation of the derivation

A string is grammatical iff there exists a derivation for it.

Chomsky Hierarchy



+ Pumping Lemma

If L is an infinite regular language, then there are strings x , y , and z such that y is not empty, $xy^n z$ is in L , for all $n \geq 0$

Context-free Recognition & CKY Algorithm

Requires the grammar be in Chomsky normal form

for $i = 1 \dots n$

$C[i-1, i] = \{V \mid V \rightarrow w_i\}$

for $\ell = 2 \dots n$ // width

for $i = 0 \dots n - \ell$ // left boundary

$k = i + \ell$ // right boundary

for $j = i + 1 \dots k - 1$ // midpoint

$C[i, k] = C[i, k] \cup \{V \mid V \rightarrow YZ, Y \in C[i, j], Z \in C[j, k]\}$

$C[j, k]$

return true if $S \in C[0, n]$

$C[i-1, i, w_i] = \text{TRUE}$

$$C[i-1, i, V] = \begin{cases} \text{TRUE} & \text{if } V \rightarrow w_i \\ \text{FALSE} & \text{otherwise} \end{cases}$$

$$C[i, j, V] = \begin{cases} \text{TRUE} & \text{if } \exists j, Y, Z \text{ such that} \\ & V \rightarrow YZ \\ & \text{and } C[i, k, Y] \\ & \text{and } C[k, j, Z] \\ & \text{and } i < k < j \\ \text{FALSE} & \text{otherwise} \end{cases}$$

goal = $C[0, n, S]$

\mathcal{L}_1 Grammar	\mathcal{L}_1 in CNF
$S \rightarrow NP VP$	$S \rightarrow NP VP$
$S \rightarrow Aux NP VP$	$S \rightarrow X1 VP$
	$X1 \rightarrow Aux NP$
$S \rightarrow VP$	$S \rightarrow book \mid include \mid prefer$
	$S \rightarrow Verb NP$
	$S \rightarrow X2 PP$
	$S \rightarrow Verb PP$
	$S \rightarrow VP PP$
$NP \rightarrow Pronoun$	$NP \rightarrow I \mid she \mid me$
$NP \rightarrow Proper-Noun$	$NP \rightarrow TWA \mid Houston$
$NP \rightarrow Det Nominal$	$NP \rightarrow Det Nominal$
$Nominal \rightarrow Noun$	$Nominal \rightarrow book \mid flight \mid meal \mid money$
$Nominal \rightarrow Nominal Noun$	$Nominal \rightarrow Nominal Noun$
$Nominal \rightarrow Nominal PP$	$Nominal \rightarrow Nominal PP$
$VP \rightarrow Verb$	$VP \rightarrow book \mid include \mid prefer$
$VP \rightarrow Verb NP$	$VP \rightarrow Verb NP$
$VP \rightarrow Verb NP PP$	$VP \rightarrow X2 PP$
	$X2 \rightarrow Verb NP$
$VP \rightarrow Verb PP$	$VP \rightarrow Verb PP$
$VP \rightarrow VP PP$	$VP \rightarrow VP PP$
$PP \rightarrow Preposition NP$	$PP \rightarrow Preposition NP$

Figure 13.8 \mathcal{L}_1 Grammar and its conversion to CNF. Note that although they aren't shown here all the original lexical entries from \mathcal{L}_1 carry over unchanged as well.