



6.863J Natural Language Processing


Lecture 5: What's my line?

Instructor: Robert C. Berwick
berwick@ai.mit.edu



The Menu Bar

- Administrivia
 - Lecture 3 posted; Lab 1a (aka "component II") due yesterday; Lab 1b, due next Monday
- Postmortem: Complexity of Kimmo/fst's – too weak? Too strong? What makes a good computational linguistics representation? A good linguistic representation? A good algorithm?
- Alternatives: morphology w/o a dictionary
- What's my line: take a chance




Is Kimmo necessary?

- Does it explain why many non-human systems never occur (ruling them out)
- Or does it overshoot?

- Ans: it seems to overshoot, in at least 2 ways
- Overshoots detected by computational analysis

6.863J/9.611J SP04 Lecture 5



Overshoot #1: too powerful with dependencies

- More powerful than well-known grammars in linguistics (and computational linguistics)
- We can use kimmo to 'count' – but natural languages don't (or cannot) do this...
- (Recall: we can use Kimmo to output a language with one counting relation: $a^n b^n$ – not a finite-state language)
- But we can do more... nothing stops us from producing a language with m counting relations, e.g, for any n , $\{(x, (cx)^n) \mid x \in \{a^* b^*\}\}$, e.g., for $n=3$, cababcababcabab, cbbbcbbbcbbb...

6.863J/9.611J SP04 Lecture 5

Not captured by context-free language

- (Familiar): $a^n b^n c^n$
- Intuition: use of pushdown stack – can catch one such pairing, but not more

6.863J/9.611J SP04 Lecture 5

So: Kimmo admits more than context-free languages!

- So Kimmo is more powerful than this!
But how powerful is it? We can still parse context-free languages in cubic time (in the length of sentences)
- We shall see that Kimmo is more complex than this!
- Conjecture: all the context-sensitive languages

6.863J/9.611J SP04 Lecture 5

Complexity of Kimmo word recognition

- All these finite-state devices, working in parallel
- There is backup
- Is it intrinsic to the system? Or eradicable?
Or, doesn't matter in practice?

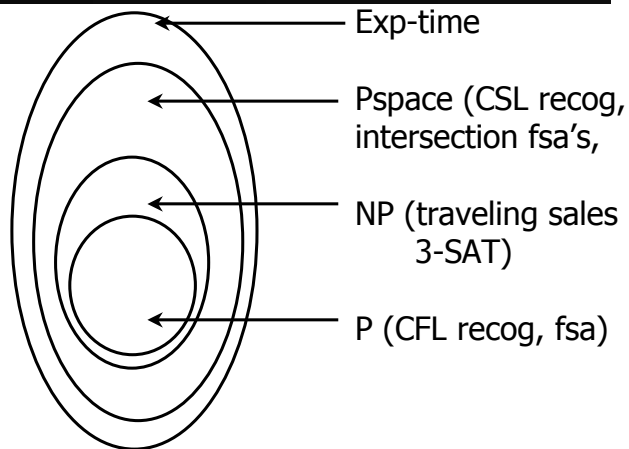
6.863J/9.611J SP04 Lecture 5

Litmus test #2 – computational complexity of Kimmo – word parsing is intractable!

- Kimmo Recognition Problem (KRP):
Given a language defined by an arbitrary (finite) Kimmo dictionary (lexical automata) and a finite set of Kimmo rules, how long in the worst case will it take to recognize whether a form is or is not in the language?
- Kimmo recognition problem is NP-hard
- As hard as any other problem solvable by a nondeterministic Turing machine in polynomial time
- No known det polytime (eg, cubic) algorithm for NP-hard problems...

6.863J/9.611J SP04 Lecture 5

Complexity hierarchy



6.863J/9.611J SP04 Lecture 5

Parsing words with Kimmo is computationally intractable

- Intuition: what if the characters on the surface don't give any clues as to what 'features' they ought to have underlyingly? (e.g., whether a Noun or a Verb, as in *police police police*)
- This seems awfully close to the famous 3-SAT problem: is there an assignment of T(rue), F(alse) to the literals of an arbitrary Boolean formula in 3-conjunctive normal form s.t. the formula evaluates to *true*?
- In fact, we can simulate this problem using Kimmo


6.863J/9.611J SP04 Lecture 5



3-Sat (3-satisfiability) is NP-complete

- *Given* (arb) 3-Sat formula, e.g.,
- There is no known deterministic Turing machine that can figure out quickly (in polynomial time) whether there is an assignment of *true* or *false* to literals x, y, z in order to make the formula evaluates to true just by inspecting the local surface string
- We could guess this in polynomial time – i.e., Nondeterministic Polynomial, or NP time (time measured in length of the formula)

6.863J/9.611J SP04 Lecture 5



Reduction of 3-Sat to Kimmo recognition problem

- For every 3-Sat problem, we can find, in polynomial time, a corresponding Kimmo word recognition problem where there's a valid word if the 3-Sat problem was satisfiable
- If Kimmo recognition could be done in deterministic polynomial time (P) then so could 3-SAT

6.863J/9.611J SP04 Lecture 5

Reduction

Any 3-Sat problem

← Efficient (polynomial time) transformation

Equivalent
Kimmo recognition problem

Answer to original SAT problem

6.863J/9.611J SP04 Lecture 5

The reduction:

Given: arbitrary 3-SAT problem instance, e.g.,

$(x \vee \neg y \vee z) (\neg x \vee \neg z) (x \vee y)$

Fast
(polytime)
transformation

(fixed)
Lexicon, L

Fst's, 1
per variable

$word \in L$ if Sat instance satisfiable

*If we could solve Kimmo recognition easily,
Then we could solve 3-Sat easily*

6.863J/9.611J SP04 Lecture 5



Two components to 3-Sat

- The fact that an x that has a truth assignment in one place, must have the same truth assignment everywhere - what morphological process is that like?
- The fact that every triple must have at least 1 'T' underlyingly (so that the triple is true) - what morphological process is that like?

6.863J/9.611J SP04 Lecture 5



How the reduction works

- Given arbitrary 3-sat formula ϕ , e.g.,
 $(x \vee \neg y \vee z) (\neg x \vee \neg z) (x \vee y)$
- Represent in the form, a 'word':
 $x\text{-}yz,\text{-}xz,xy$
- For each variable x , we have an 'assignment machine' that ensures that x is mapped to T or F throughout the whole formula
- We have one machine (and a fixed dictionary) to checks each disjunction to make sure that at least one disjunct is true in every conjunct

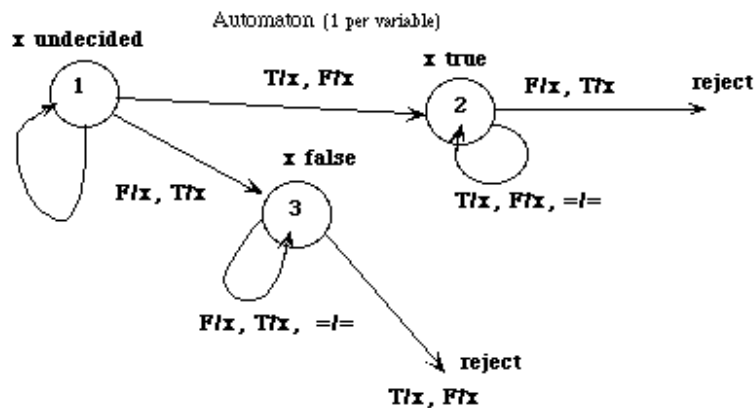
6.863J/9.611J SP04 Lecture 5

Two components

- Agreement: vowel harmony (if round at some point, round everywhere)
- Ambiguity: we can't tell what the underlying value of x is from the surface, but if there's at least one "t" per 'part of word', then we can spell out this constraint in dictionary
- Note that words (like Sat formulas) must be arbitrarily long... (pas de probleme)
- Dictionary is fixed...
- # of Vowel harmony processes corresponds to # of distinct literals

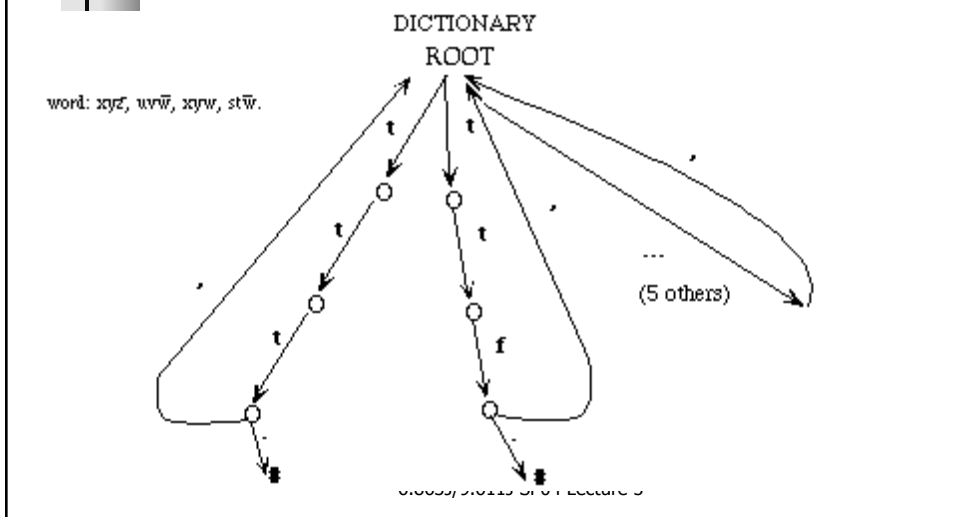
6.863J/9.611J SP04 Lecture 5

Reduce until done: assignment consistency



6.863J/9.611J SP04 Lecture 5

Reduce until done – formula must
eval to *true*



What are the implications?

- FTNs inherently require backup if simulated (in the worst case) – Kimmo at least NP-hard (proof later on)
- Empty elements cause computational complexity (unless restricted – equal length condition) – true in all areas of linguistics
- Composition can save us, but then rule ordering must be watched carefully



Implications

- Do we need a machine powerful enough to represent intractable problems?
- No evidence for unbounded # of counting dependencies or harmony processes...
- Performance? Or do we need something this powerful??

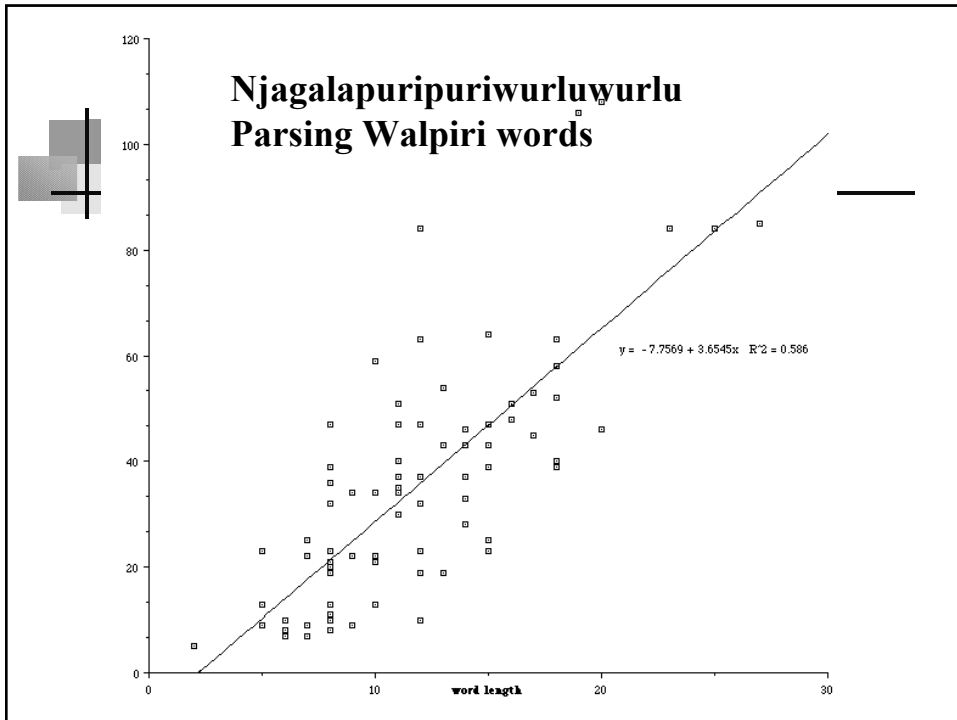
6.863J/9.611J SP04 Lecture 5



Why should we care?

- This is *typical* of a combination of 'agreement and ambiguity' that trickles through all of natural language
- The agreement part – like Turkish vowel harmony
- The ambiguity part – like the *police police police* example
- Suggests that speed won't come from the formalism *all by itself*

6.863J/9.611J SP04 Lecture 5



What if we don't have a dictionary?

- Don't use one
- Learn one from data

Method 1: don't use a dictionary

- Best known method – Porter stemming (Porter, 1980)
- <http://www.tartarus.org/~martin/PorterStemmer/>
<http://snowball.tartarus.org/>
 - For English
 - Most widely used system
 - Manually written rules
 - 5 stage approach to extracting roots
 - Considers suffixes only
 - May produce non-word roots

6.863J/9.611J SP04 Lecture 5

Porter output

Sample Output (English):

consigned	consign	knack	knack
consignment	consign	knackereries	knackeri
consolation	consol	knaves	knavish
consolatory	consolatori	knavish	knavish
consolidate	consolid	knif	knif
consolidating	consolid	knife	knife
consoling	consol	knew	knew

6.863J/9.611J SP04 Lecture 5

Why?

Algorithmic stemmers can be fast (and lean):

E.g.: 1 Million words in 6 seconds on 500 MHz PC

- It is more efficient not to use a dictionary (don't have to maintain it if things change).
- It is better to ignore irregular forms (exceptions) than to complicate the algorithm (not much lost in practice).

6.863J/9.611J SP04 Lecture 5

Output - German

aufeinander	aufeinand	kategorie	kategori
auferlegen	auferleg	kategorien	kategori
auferlegt	auferlegt	kater	kat
auferlegten	auferlegt	katers	kat
auferstanden	auferstand	katze	katz
auferstehen	auferstand	katzen	katz
aufersteht	aufersteht	kätzchen	katzch

6.863J/9.611J SP04 Lecture 5

Method

Porter Stemmers use simple algorithms to determine which affixes to strip in which order and when to apply repair strategies.

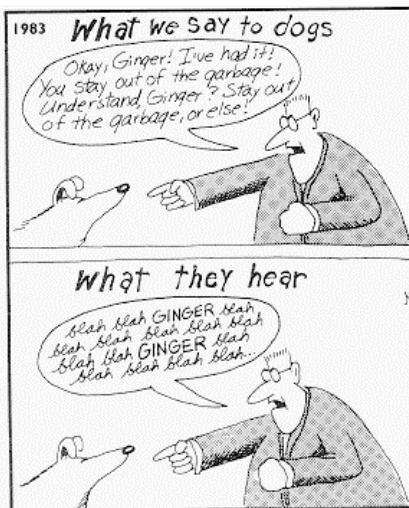
Input	Strip -ed Affix	Repair
hoped	hop	hope (add -e if word is short)
hopped	hopp	hop (delete one if doubled)

Samples of the algorithms are accessible via the Web and can be programmed in any language.

Advantage: easy to see understand, easy to implement.

6.863J/9.611J SP04 Lecture 5

Words are fine – but we need more



6.863J/9.611J SP04 Lecture 5



Paradigmatic example for NLP

- Morphophonemic parsing
- Given surface form, recover underlying form:

6.863J/9.611J SP04 Lecture 5

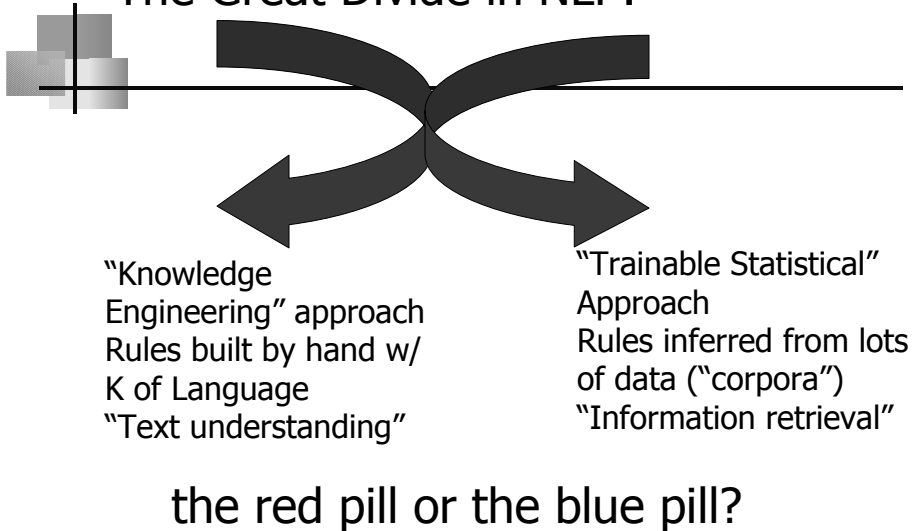


Two ways

- Generative model – concatenate then fix up joints
 - stop + -ing = stopping, fly + s = flies
 - Use a cascade of transducers to handle all the fixups
- Probabilistic model - some constraints on morpheme sequences using prob of one character appearing before/after another
prob(ing | stop) vs. prob(ly| stop)

6.863J/9.611J SP04 Lecture 5

The Great Divide in NLP:



6.863J/9.611J SP04 Lecture 5

Another example: generating language

- Variations in style, syntactic form...
- Where do these come from?

- Jane Austin writes differently from Charlotte Brontë

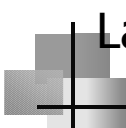
6.863J/9.611J SP04 Lecture 5



The Red Pill

- How to fold statistical information in to the symbolic models?
- Let's try a classic and simple way to define a statistical model for language...

6.863J/9.611J SP04 Lecture 5



Language ID

- "Rabbit and Lukasiewicz are on the menu"
- Is this English or Polish or what?
- Is it "good" (= likely) English?
- Is it "good" (= likely) Polish?

6.863J/9.611J SP04 Lecture 5



Text Categorization

- Automatic Yahoo classification, etc.
- Similar to language ID ...
 - Topic 1 sample: In the beginning God created ...
 - Topic 2 sample: The history of all hitherto existing society is the history of class struggles. ...
- Input text: Matt's Communist Homepage. Capitalism is unfair and has been ruining the lives of millions of people around the world. The profits from the workers' labor ...
- Input text: And they have beat their swords to ploughshares, And their spears to pruning-hooks. Nation doth not lift up sword unto nation, neither do they learn war any more. ...

6.863J/9.611J SP04 Lecture 5



Contextual Spelling Correction

- Which is most probable?
 - ... I think they're okay ...
 - ... I think there okay ...
 - ... I think their okay ...
- Which is most probable?
 - ... by the way, are they're likely to ...
 - ... by the way, are there likely to ...
 - ... by the way, are their likely to ...

6.863J/9.611J SP04 Lecture 5



Topic Segmentation

- Break big document or media stream into indexable chunks
- From NPR's *All Things Considered*:

The U. N. says its observers will stay in Liberia only as long as West African peacekeepers do, but West African states are threatening to pull out of the force unless Liberia's militia leaders stop violating last year's peace accord after 7 weeks of chaos in the capital, Monrovia ... Human rights groups cite peace troops as among those smuggling the arms. I'm Jennifer Ludden, reporting. Whitewater prosecution witness David Hale began serving a 28-month prison sentence today. The Arkansas judge and banker pleaded guilty two years ago to defrauding the Small Business Administration. Hale was the main witness in the Whitewater-related trial that led to the convictions ...

6.863J/9.611J SP04 Lecture 5



Contextual Spelling Correction

- Which is most probable?
 - ... I think they're okay ...
 - ... I think there okay ...
 - ... I think their okay ...
- Which is most probable?
 - ... by the way, are they're likely to ...
 - ... by the way, are there likely to ...
 - ... by the way, are their likely to ...



Speech Recognition

- How do you wreck a nice beach?
- How do you recognize speech?

- Put the file in the folder
- Put the file and the folder

6.863J/9.611J SP04 Lecture 5



Language generation

- Choose randomly among outputs:
 - Visitant which came into the place where it will be Japanese has admired that there was Mount Fuji.
- Top 10 outputs according to bigram probabilities:
 - Visitors who came in Japan admire Mount Fuji.
 - Visitors who came in Japan admires Mount Fuji.
 - Visitors who arrived in Japan admire Mount Fuji.
 - Visitors who arrived in Japan admires Mount Fuji.
 - Visitors who came to Japan admire Mount Fuji.
 - A visitor who came in Japan admire Mount Fuji.
 - The visitor who came in Japan admire Mount Fuji.
 - Visitors who came in Japan admire Mount Fuji.
 - The visitor who came in Japan admires Mount Fuji.
 - Mount Fuji is admired by a visitor who came in Japan.

6.863J/9.611J SP04 Lecture 5

Basic Morphology

Basic Affix Typology (don't seem to need more):

- i-suffix: inflectional suffix

English: *cheer+ed = cheered, fit+ed = fitted, love+ed = loved*

- d-suffix: derivational suffix, changes word type

English: *walk(V)+er = walker(N), happy(A)+ness=happiness(N)*

- a-suffix: attached suffix (enclitics).

Italian *mandargli* = *mandare* + *gli* to send + to him

Algorithmic Method

General Strategy:

- Normal order of suffixes seems to be *d, i, a*.
- Remove from right in order *a, i, d*.
- Generally remove all the *a* and *i* suffixes, sometimes leave the *d* one.

Types of Errors

- Conflation: reply, rep. rep

- Overstemming: wander wand
 news new

- Misstemming: relativity relative

- Understemming: knavish knavish

6.863J/9.611J SP04 Lecture 5

Algorithmic Method

Strategy for German:

- Leave prefixes alone because they can change meaning.
- Put everything in small caps.
- Get rid of *ge-*.
- Get rid of *i* type: *e, em, en, ern, er, es, s, est,*
(e.g, *armes > arm*)
- Get rid of *d* type: *end, ung, ig, ik, isch, lich, heit,*
keit

6.863J/9.611J SP04 Lecture 5



Information Retrieval

Does stemming indeed improve IR?

- No: Harman (1991), Krovetz (1993)
- Possibly: Krovetz (1995)
 - Depends on type of text, and the assumption is that once one moves beyond English, the difference will prove significant.

6.863J/9.611J SP04 Lecture 5



Crosslinguistic Applicability

- Can this type of stemming be applied to all languages?
 - Not to Chinese, for example (doesn't need it).
- Do all languages have the same kind of morphology?
 - No. Stemming assumes basically agglutinative morphology. This is not true crosslinguistically (but the algorithms seem to work pretty well within Indo-

6.863J/9.611J SP04 Lecture 5



Stemming: Methods

- Dictionary approach not enough
 - Example: (Porter, 1991)
 - routed -> route/rout
 - At Waterloo, Napoleon's forces were routed
 - The cars were routed off the highway
 - Here, the (inflected) verb form is polysemous


6.863J/9.611J SP04 Lecture 5



Stemming: Errors

- Understemming: failure to merge
 - Adhere/adhesion
- Overstemming: incorrect merge
 - Probe/probable
 - Claim: *-able* irregular suffix, root: *probare* (Lat.)
- Mis-stemming: removing a non-suffix (Porter, 1991)
 - reply -> rep

6.863J/9.611J SP04 Lecture 5



Stemming: Interaction

- Interacts with noun compounding:
 - Example:
 - operating systems
 - negative polarity items
 - For IR, compounds need to be identified first...

6.863J/9.611J SP04 Lecture 5



Stemming: Porter Algorithm

- Rule format:
 - (condition on stem) suffix₁ -> suffix₂
 - In case of conflict, prefer longest suffix match
- “Measure” of a word is m in:
 - (C) (VC) ^{m} (V)
 - C = sequence of one or more consonants
 - V = sequence of one or more vowels
 - Examples:
 - *tree* C(VC)⁰V
 - *troubles* C(VC)²

6.863J/9.611J SP04 Lecture 5



Stemming: Porter Algorithm

- Step 1a: remove plural suffixation
 - SSES -> SS (caresses)
 - IES -> I (ponies)
 - SS -> SS (caress)
 - S -> (cats)
- Step 1b: remove verbal inflection
 - (m>0) EED -> EE (agreed, feed)
 - (*v*) ED -> (plastered, bled)
 - (*v*) ING -> (motoring, sing)

6.863J/9.611J SP04 Lecture 5



Stemming: Porter Algorithm

- Step 1b: (contd. for *-ed* and *-ing* rules)
 - AT -> ATE (conflated)
 - BL -> BLE (troubled)
 - IZ -> IZE (sized)
 - (*doubled c & ¬(*L v *S v *Z)) -> single c (hopping, hissing, falling, fizzing)
 - (m=1 & *cvc) -> E (filing, failing, slowing)
- Step 1c: Y and I
 - (*v*) Y -> I (happy, sky)

6.863J/9.611J SP04 Lecture 5



Stemming: Porter Algorithm

- Step 2: Peel one suffix off for multiple suffixes
 - (m>0) ATIONAL -> ATE (relational)
 - (m>0) TIONAL -> TION (conditional, rational)
 - (m>0) ENCI -> ENCE (valenci)
 - (m>0) ANCI -> ANCE (hesitanci)
 - (m>0) IZER -> IZE (digitizer)
 - (m>0) ABLI -> ABLE (conformabli) - *able* (step 4)
 - ...
 - (m>0) IZATION -> IZE (vietnamization)
 - (m>0) ATION -> ATE (predication)
 - ...
 - (m>0) IVITI -> IVE (sensitiviti)

6.863J/9.611J SP04 Lecture 5



Stemming: Porter Algorithm

- Step 3
 - (m>0) ICATE -> IC (triplicate)
 - (m>0) ATIVE -> (formative)
 - (m>0) ALIZE -> AL (formalize)
 - (m>0) ICITI -> IC (electriciti)
 - (m>0) ICAL -> IC (electrical, chemical)
 - (m>0) FUL -> (hopeful)
 - (m>0) NESS -> (goodness)

6.863J/9.611J SP04 Lecture 5



Stemming: Porter Algorithm

- Step 4: Delete last suffix
 - (m>1) AL -> (revival) - *revive*, see step 5
 - (m>1) ANCE -> (allowance, dance)
 - (m>1) ENCE -> (inference, fence)
 - (m>1) ER -> (airliner, employer)
 - (m>1) IC -> (gyroscopic, electric)
 - (m>1) ABLE -> (adjustable, mov(e)able)
 - (m>1) IBLE -> (defensible, bible)
 - (m>1) ANT -> (irritant, ant)
 - (m>1) EMENT -> (replacement)
 - (m>1) MENT -> (adjustment)
 - ...

6.863J/9.611J SP04 Lecture 5



Stemming: Porter Algorithm

- Step 5a: remove *e*
 - (m>1) E -> (probate, rate)
 - (m>1 & \neg *cvc) E -> (cease)
- Step 5b: //reduction
 - (m>1 & *LL) -> L (controller, roll)

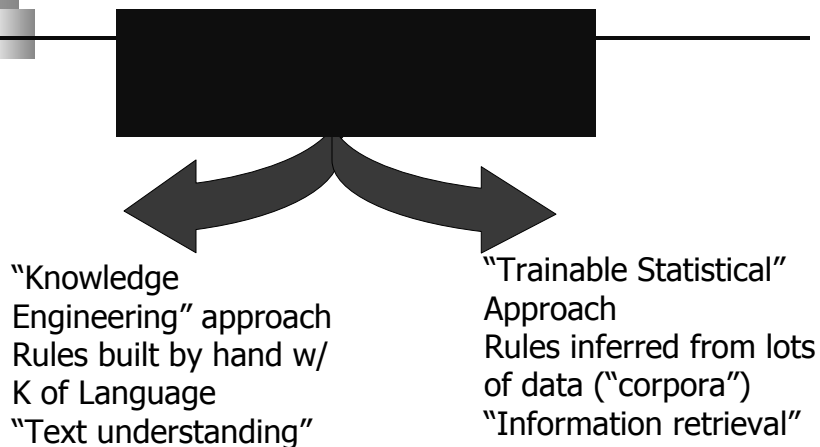
6.863J/9.611J SP04 Lecture 5

Stemming: Porter Algorithm

- Misses (understemming)
 - Unaffected:
 - *agreement* (VC)¹VCC - step 4 (m>1)
 - *adhesion*
 - Irregular morphology:
 - drove, geese
- Overstemming
 - *relativity* - step 2
- Mis-stemming
 - *wander* C(VC)¹VC

6.863J/9.611J SP04 Lecture 5

The Great Divide in NLP: the red pill or the blue pill?



6.863J/9.611J SP04 Lecture 5

A simple example

- We consider a sentence as a sequence of n words; we want to find $P(w_1, \dots, w_n)$
- We can use this to model all the 'noise' that gets into language (the leaks)
- So one idea is to combine the symbolic models (like kimmo) with the 'noise' components, to do better (eg, like econometrics...)

6.863J/9.611J SP04 Lecture 5

Language models, probability & info

- Given a string w , a language model gives us the probability of the string $P(w)$, e.g.,
 - $P(\text{the big dog}) > (\text{dog big the}) > (\text{dgo gib eth})$
 - Easy for humans; difficult for machines
 - Let $P(w)$ be called a language model $L(M)$
- "I have a gub" (Woody Allen)

6.863J/9.611J SP04 Lecture 5

A simple example – which is 'right?'

physical Brainpower not plant is chief , now a 's asset , . firm

a Brainpower not now chief asset firm 's is . plant physical ,

chief a physical , . firm not , Brainpower plant is asset 's now

now not plant Brainpower now physical 's . a chief , asset
firm , is

Brainpower , not physical plant , is now a firm 's chief asset .

Each sentence is a sequence w_1, \dots, w_n .

Task is to find $P(w_1, \dots, w_n)$.

6.863J/9.611J SP04 Lecture 5

N-grams

A simple model of language

- Computes a probability for observed input
- Probability is likelihood of observation being generated by the same source as training data
- Such a model is often called a *language model* (LM)

6.863J/9.611J SP04 Lecture 5



How can we compute this?

- Each of this pr's can be estimated (using frequency counts) from training data
- Can this work directly?
- No – not in practice...why?

6.863J/9.611J SP04 Lecture 5



What's this good for?

- What's my line?

6.863J/9.611J SP04 Lecture 5

Language ID

- "Rabbit and Lukasiewicz are on the menu"
- Is this English or Polish or what?
- We had some notion of using n-gram models ...
- Is it "good" (= likely) English?
- Is it "good" (= likely) Polish?
- Space of events will be not races but character sequences (x_1, x_2, x_3, \dots) where $x_n = \text{EOS}$ (nb, "BOS")

6.863J/9.611J SP04 Lecture 5

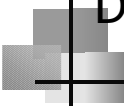
Language ID?

- Let $p(X)$ = probability of text X in English
- Let $q(X)$ = probability of text X in Polish
- Which probability is higher?

"Rabbit and Lukasiewicz are on the menu"

$p(x_1=r, x_2=a, x_3=b, x_4=b, x_5=i, x_6=t, \dots)$

6.863J/9.611J SP04 Lecture 5



Data needs

- How many possible distinct probabilities will be needed?, i.e. parameter values
- Total number of word tokens in our training data
- Total number of unique words: word types is our vocabulary size

6.863J/9.611J SP04 Lecture 5



How can we compute this?

- Each of these p 's can be estimated (using frequency counts) from training data
- Can this work?
- No – not in practice...why?


6.863J/9.611J SP04 Lecture 5



What are the tools we need?

- Crash....

6.863J/9.611J SP04 Lecture 5



Probability you should know..

- Probability notation like $p(X | Y)$:
 - What does this expression mean?
 - How can I manipulate it?
 - How can I estimate its value in practice?
- Probability models:
 - What is one?
 - Can we build one for language?
 - How do I know if my model is any good?

6.863J/9.611J SP04 Lecture 5

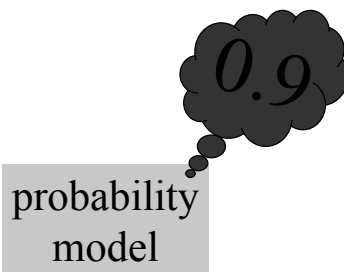
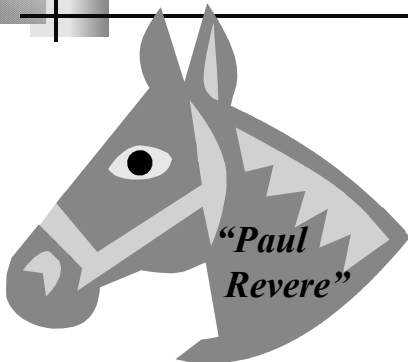
3 Kinds of Statistics

- descriptive: mean MIT SAT (or median)
- confirmatory: statistically significant?
- predictive: wanna bet?

this course – why?

6.863J/9.611J SP04 Lecture 5

Notation for Greenhorns



$$p(\text{Paul Revere wins} \mid \text{weather's clear}) = 0.9$$

6.863J/9.611J SP04 Lecture 5

What does that *really* mean?

$$p(\text{Paul Revere wins} \mid \text{weather's clear}) = 0.9$$

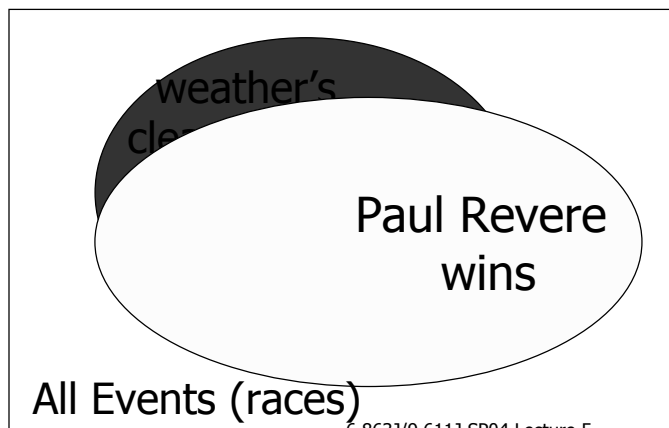
- Past performance?
 - Revere's won 90% of races with clear weather
- Hypothetical performance?
 - If he ran the race in many parallel universes ...
- Subjective strength of belief?
 - Would pay up to 90 cents for chance to win \$1
- Output of some computable formula?
 - Ok, but then which formulas should we trust?

$p(X \mid Y)$ versus $q(X \mid Y)$

6.863J/9.611J SP04 Lecture 5

p is a function on event sets

$$p(\text{win} \mid \text{clear}) \equiv p(\text{win, clear}) / p(\text{clear})$$

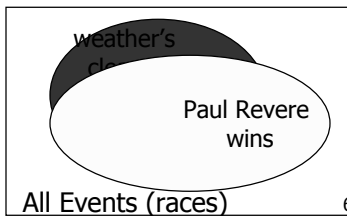


6.863J/9.611J SP04 Lecture 5

p is a function on event sets

$$p(\text{win} \mid \text{clear}) \equiv \underbrace{p(\text{win}, \text{clear})}_{\text{logical conjunction of predicates}} / \underbrace{p(\text{clear})}_{\text{predicate selecting races where weather's clear}}$$

syntactic sugar



p measures total probability of a set of events.

6.863J/9.611J SP04 Lecture 5

The Chain rule – factoring joint events

- $P(\text{GC in Hawaii}, \text{GC alone}, \text{GC low in polls} \mid \text{GC drives drunk}) =$
 $P(\text{GC in Hawaii} \mid \text{GC alone}, \text{GC low in polls}, \text{GC drives drunk}) \times$
 $P(\text{GC alone} \mid \text{GC low in polls}, \text{GC drives drunk}) \times$
 $P(\text{GC low in polls} \mid \text{GC drives drunk})$

Why does this work?

6.863J/9.611J SP04 Lecture 5

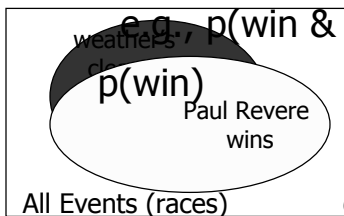
Chain rule – why does it work?

- Remember: $P(X|Y) = P(X,Y)/P(Y)$
- $\frac{HALD}{D} = \frac{HALD}{ALD} \times \frac{ALD}{LD} \times \frac{LD}{D}$
- Simply cancel out the matching terms

6.863J/9.611J SP04 Lecture 5

Required Properties of p (axioms) *most of the*

- $p(\emptyset) = 0$ $p(\text{all events}) = 1$
- $p(X) \leq p(Y)$ for any $X \subseteq Y$
- $p(X) + p(Y) = p(X \cup Y)$ provided $X \cap Y = \emptyset$



p measures total probability of a set of events

6.863J/9.611J SP04 Lecture 5

Commas denote conjunction

~~p(Paul Revere wins, Valentine places, Epitaph shows | weather's clear)~~

what happens as we add conjuncts to left of bar ?

- probability can only decrease
- numerator of historical estimate likely to go to zero:

$$\frac{\# \text{ times Revere wins AND Val places... AND weather's clear}}{\# \text{ times weather's clear}}$$

6.863J/9.611J SP04 Lecture 5

Commas denote conjunction

~~p(Paul Revere wins, Valentine places, Epitaph shows | weather's clear)~~

p(Paul Revere wins | weather's clear, ground is dry, jockey getting over sprain, Epitaph also in race, Epitaph was recently bought by Gonzalez, race is on May 17, ...)

what happens as we add conjuncts to right of bar ?

- probability could increase or decrease
- probability gets more relevant to our case (less *bias*)
- probability *estimate* gets less reliable (more *variance*)

$$\frac{\# \text{ times Revere wins AND weather clear AND ... it's May 17}}{\# \text{ times weather clear AND ... it's May 17}}$$

6.863J/9.611J SP04 Lecture 5

Simplifying Right Side: Backing Off

$p(\text{Paul Revere wins} \mid \text{weather's clear, ground is dry, jockey getting over sprain, Epitaph also in race, Epitaph was recently bought by Gonzalez, race is on May 17, ...})$

not exactly what we want but at least we can get a reasonable estimate of it!

(i.e., more bias but less variance)

try to *keep* the conditions that we suspect will have the most influence on whether Paul Revere wins

6.863J/9.611J SP04 Lecture 5

Simplifying Right Side: Backing Off

$p(\text{Paul Revere wins, Valentine places, Epitaph shows} \mid \text{weather's clear})$

NOT ALLOWED!

but we can do something similar to help ...

6.863J/9.611J SP04 Lecture 5

Factoring Left Side: The Chain Rule

$$\begin{aligned}
 & \cancel{p(\text{Revere, Valentine, Epitaph} \mid \text{weather's clear})} && \cancel{RVEW/W} \\
 = & \cancel{p(\text{Revere} \mid \text{Valentine, Epitaph}, \text{weather's clear})} && = RVEW/VEW \\
 & * p(\text{Valentine} \mid \text{Epitaph}, \text{weather's clear}) && * VEW/EW \\
 & * p(\text{Epitaph} \mid \text{weather's clear}) && * EW/W
 \end{aligned}$$

True because numerators cancel against denominators

Makes perfect sense when read from bottom to top

Moves material to right of bar so it can be ignored

If this prob is unchanged by backoff, we say Revere was **CONDITIONALLY INDEPENDENT** of Valentine and Epitaph (conditioned on the weather's being clear). Often we just **ASSUME** conditional independence to get the nice product above.

6.863J/9.611J SP04 Lecture 5

Language ID?

- Let $p(X)$ = probability of text X in English
- Let $q(X)$ = probability of text X in Polish
- Which probability is higher?

"Rabbit and Lukasiewicz are on the menu"

$p(x_1=r, x_2=a, x_3=b, x_4=b, x_5=i, x_6=t, \dots)$

6.863J/9.611J SP04 Lecture 5

How do we calculate this?

- Use the chain rule in probability...
- But there's a hitch...

6.863J/9.611J SP04 Lecture 5

Apply the Chain Rule

$$\begin{aligned} & p(\mathbf{x}_1=\mathbf{r}, \mathbf{x}_2=\mathbf{a}, \mathbf{x}_3=\mathbf{b}, \mathbf{x}_4=\mathbf{b}, \mathbf{x}_5=\mathbf{i}, \mathbf{x}_6=\mathbf{t}, \dots) \\ &= p(\mathbf{x}_1=\mathbf{r}) && 4470/52108 \\ & * p(\mathbf{x}_2=\mathbf{a} \mid \mathbf{x}_1=\mathbf{r}) && 395/4470 \\ & * p(\mathbf{x}_3=\mathbf{b} \mid \mathbf{x}_1=\mathbf{r}, \mathbf{x}_2=\mathbf{a}) && 5/395 \\ & * p(\mathbf{x}_4=\mathbf{b} \mid \mathbf{x}_1=\mathbf{r}, \mathbf{x}_2=\mathbf{a}, \mathbf{x}_3=\mathbf{b}) && 3/5 \\ & * p(\mathbf{x}_5=\mathbf{i} \mid \mathbf{x}_1=\mathbf{r}, \mathbf{x}_2=\mathbf{a}, \mathbf{x}_3=\mathbf{b}, \mathbf{x}_4=\mathbf{b}) && 3/3 \\ & * p(\mathbf{x}_6=\mathbf{t} \mid \mathbf{x}_1=\mathbf{r}, \mathbf{x}_2=\mathbf{a}, \mathbf{x}_3=\mathbf{b}, \mathbf{x}_4=\mathbf{b}, \mathbf{x}_5=\mathbf{i}) && 0/3 \\ & * \dots = 0 \end{aligned}$$

counts from
Brown corpus

6.863J/9.611J SP04 Lecture 5

How can we compute this?

- Each of this pr's can be estimated (using frequency counts) from training data
- Can this work?
- No – not in practice...why?

6.863J/9.611J SP04 Lecture 5

Forming classes

- "*n-gram*" = sequence of n "words"
 - unigram
 - bigram
 - trigram
 - four-gram...
- In language modeling, the conditioning variables are sometimes called the "history" or the "context."
- The Markov assumption says that the prediction is conditionally independent of ancient history, given recent history.
- I.e., we divide all possible histories into equivalence classes based on the recent history

6.863J/9.611J SP04 Lecture 5

The Markov assumption

- We *approximate* $p(\text{word} \mid \text{all previous words})$
Instead of
 $p(\text{rabbit} \mid \text{Follow the white...})$ we use:
 $P(\text{rabbit} \mid \text{white})$
- This is a *Markov assumption* where past memory is limited to immediately previous state – just 1 state corresponding to the previous word or tag

6.863J/9.611J SP04 Lecture 5

N-grams: limiting history – the Markov assumption

- 0th order Markov model: $P(w_i)$ called a unigram model
- 1st order Markov model: $P(w_i \mid w_{i-1})$ called a bigram model
- 2nd order Markov model: $P(w_i \mid w_{i-2}, w_{i-1})$ called a trigram model

6.863J/9.611J SP04 Lecture 5

Calculation

$$\begin{aligned}
 & p(x_1=r, x_2=a, x_3=b, x_4=b, x_5=i, x_6=t, \dots) \\
 & \approx p(x_1=r) \quad 4470/52108 \\
 & * p(x_2=a \mid x_1=r) \quad 395/4470 \\
 & * p(x_3=b \mid x_1=r, x_2=a) \quad 5/395 \\
 & * p(x_4=b \mid x_2=a, x_3=b) \quad 12/919 \\
 & * p(x_5=i \mid x_3=b, x_4=b) \quad 12/126 \\
 & * p(x_6=t \mid x_4=b, x_5=i) \quad 485 \\
 & * \dots = 7.3e-10 * \dots
 \end{aligned}$$

counts from
Brown corpus

6.863J/9.611J SP04 Lecture 5

Another Independence Assumption

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx p(x_1=h) \quad 4470/52108 \\
 & * p(x_2=o \mid x_1=h) \quad 395/4470 \\
 & * p(x_i=r \mid x_{i-2}=h, x_{i-1}=o) \quad 1417/14765 \\
 & * p(x_i=s \mid x_{i-2}=o, x_{i-1}=r) \quad 1573/26412 \\
 & * p(x_i=e \mid x_{i-2}=r, x_{i-1}=s) \quad 1610/12253 \\
 & * p(x_i=s \mid x_{i-2}=s, x_{i-1}=e) \quad 2044/21250 \\
 & * \dots = 5.4e-7 * \dots
 \end{aligned}$$

counts from
Brown corpus

6.863J/9.611J SP04 Lecture 5

Simplify the Notation

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx p(x_1=h) && 4470/52108 \\
 & * p(x_2=o \mid x_1=h) && 395/ 4470 \\
 & * p(r \mid h, o) && 1417/14765 \\
 & * p(s \mid o, r) && 1573/26412 \\
 & * p(e \mid r, s) && 1610/12253 \\
 & * p(s \mid s, e) && 2044/21250 \\
 & * \dots
 \end{aligned}$$

counts from
Brown corpus

6.863J/9.611J SP04 Lecture 5

Simplify the Notation

$$\begin{aligned}
 & p(x_1=r, x_2=a, x_3=b, x_4=b, x_5=i, x_6=t, \dots) \\
 & \approx p(r \mid \mathbf{BOS}, \mathbf{BOS}) \\
 & * p(a \mid \mathbf{BOS}, r) \\
 & * p(b \mid r, a) \\
 & * p(b \mid a, b) \\
 & * p(i \mid b, b) \\
 & * p(t \mid i, b)
 \end{aligned}$$

the parameters
of a
trigram generator!
Same assumptions
about language.

values of
those
parameters,
as naively
estimated
from Brown
corpus.

4470/ 52108
395/ 4470
1417/14765
1573/26412
1610/12253
2044/21250

* ... These basic probabilities
are used to define p(rabbit)

counts from
Brown corpus

6.863J/9.611J SP04 Lecture 5

Simplify the Notation

$$p(x_1=r, x_2=a, x_3=b, x_4=b, x_5=i, x_6=t, \dots)$$

the parameters of trigram generator!
Same assumptions about language

values of those parameters, as naively estimated from Brown corpus.

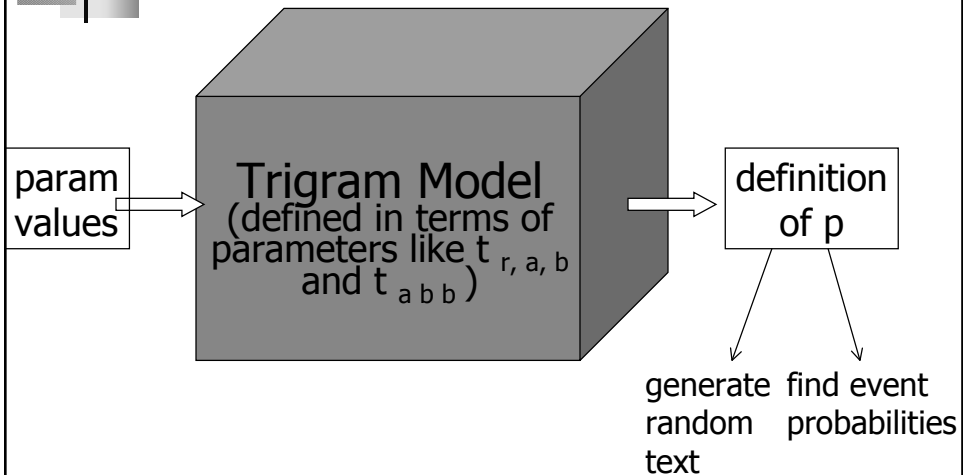
counts from Brown corpus

- $\approx t_{\text{BOS, BOS, r}}$
- $* t_{\text{BOS, r, a}}$
- $* t_{\text{a, BOS, r}}$
- $* t_{\text{b, r, s}}$
- $* t_{\text{r, s, e}}$
- $* t_{\text{t, b, i}}$
- $* \dots$

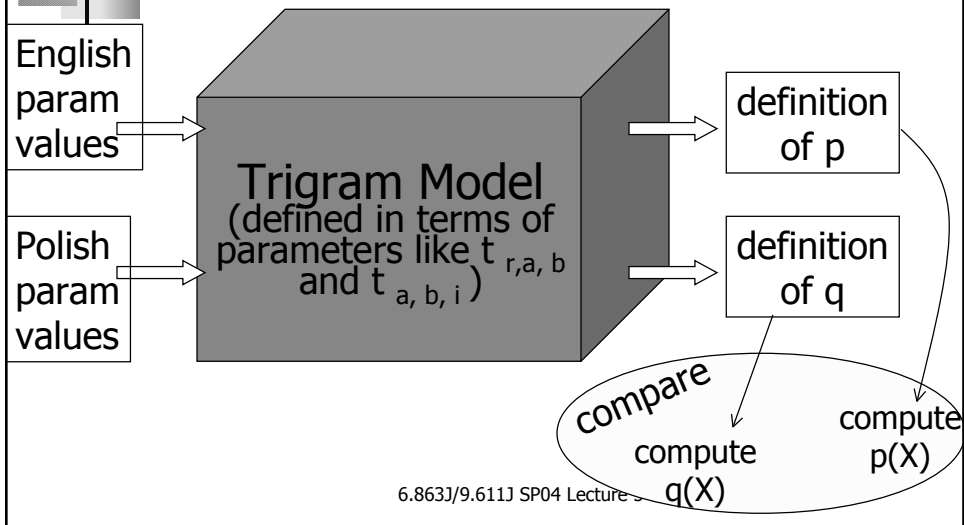
... This notation emphasizes that they're just real variables whose value must be estimated.

4470/52108
395/4470
1417/14765
1573/26412
1610/12253
2044/21250

Definition: Probability Model

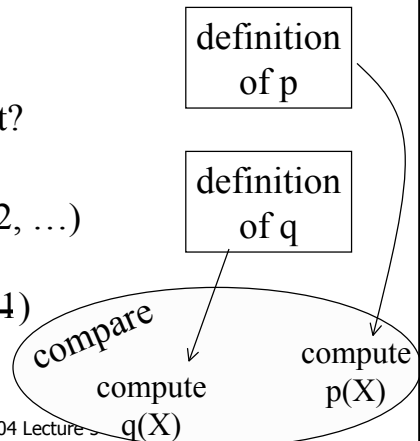


English vs. Polish



What is "X" in $p(X)$?

- Element of some implicit "event space"
 - e.g., race
 - e.g., sentence
- What if event is a whole text?
 - $p(\text{text})$
 - = $p(\text{sentence 1, sentence 2, ...})$
 - = $p(\text{sentence 1})$
 - * $p(\text{sentence 2} | \text{sentence 1})$
 - * ...



Writing like Jane Austen



- Three novels by Jane Austen: *Emma*, *Sense and Sensibility*, *Pride and Prejudice*
- Remove punctuation, keep paragraphs
- Train trigram model on this text

6.863J/9.611J SP04 Lecture 5

Writing like Jane Austen

$f(3\text{gram})$	$f(2\text{gram})$	$f(1\text{gram})$	w_0	w_1	w_2
378	518	10381	I	do	not
366	1366	10381	I	am	sure
214	1917	9182	in	the	world
202	572	6917	she	could	not
189	462	2751	would	have	been
174	184	10381	I	dare	say
173	179	5758	as	soon	as
173	357	11135	a	great	deal
171	332	7573	it	would	be
155	945	3017	could	not	be

6.863J/9.611J SP04 Lecture 5

Conversion

3gram $\frac{f(w_0, w_1, w_3)}{f(w_0, w_1)}$	2gram $\frac{f(w_0, w_1)}{f(w_0)}$	1gram $\frac{f(w_0)}{N}$	w_0	w_1	w_2
0.72	0.04	0.016	I	do	not
0.26	0.13	0.016	I	am	sure
0.11	0.20	0.014	in	the	world
0.35	0.08	0.011	she	could	not
0.40	0.16	0.004	would	have	been
0.94	0.01	0.016	I	dare	say
0.96	0.03	0.009	as	soon	as
0.48	0.03	0.018	a	great	deal
0.51	0.04	0.012	it	would	be
0.16	0.31	0.004	could	not	be

6.863J/9.611J SP04 Lecture 5

3-gram

[Gennethyesse orils of Ted you doorder [6], the Grily
 Capiduatent pildred and For thy werarme: nomiterst halt
 i, what production the Covers, in calt cations on
 wile ars, was name conch rom the exce of the man,
 Winetwentagaint up, and and All. And of Ther so i
 hundal panite days th the res of th rand ung into
 the forD six es, wheralf the hie
 soulsee, frelatche rigat. And the LOperact camen
 unismelight fammedied: and nople,

6.863J/9.611J SP04 Lecture 5

4-gram

[1] By the returall benefit han every familitant of all
Thou go? And At the eld to parises of the nursed by thy
way of all histantly be the ~aciedfag . to the narre
gread abraza of thing, and vas these conwuning clann
com to one language; all Lah, which for the greath
othey die. -

6.863J/9.611J SP04 Lecture 5

5-gram

[Gen 3:1] In the called up history of its opposition of
bourgeois AND Adam to rest, that the existing of
heaven; and land the bourgeois ANger anything but
concealed, the land whethere had doth know ther:
bury thy didst of Terature their faces which went
masses the old society [2] is the breaks out
of oppressor of all which, the proLETARIat goest,
unto German pleast twelves applied in manner with t
hese, first of this polities have all

6.863J/9.611J SP04 Lecture 5



3-word-gram

[Gen 4:25] And Adam gave names to all feudal, patriarchal, idyllic relations. It has but re-established new classes, new conditions of oppression, new forms of struggle in place of the West? The bourgeoisie keeps more and more splitting up into two great lights; the greater light to rule the day of my house is this Eliezer of Damascus.

How far can we go??