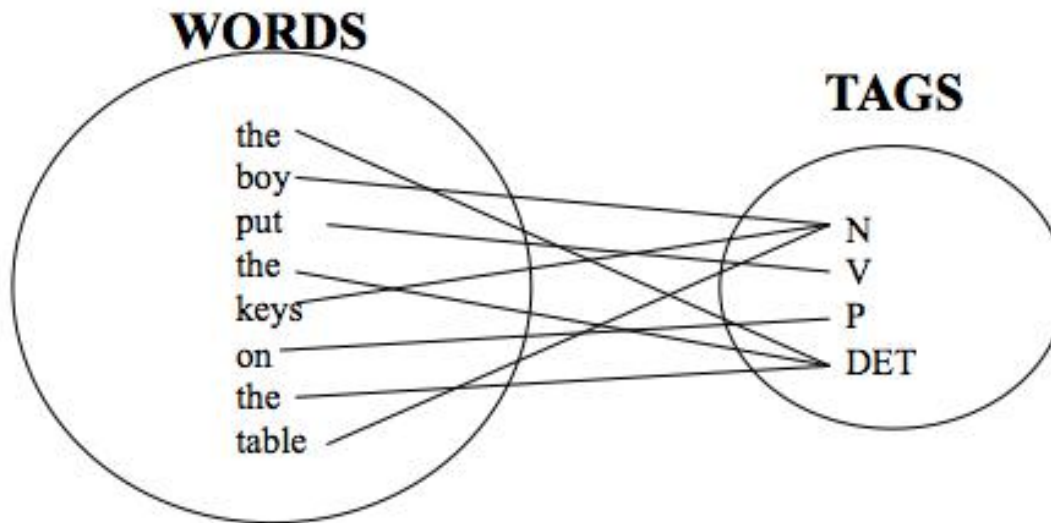


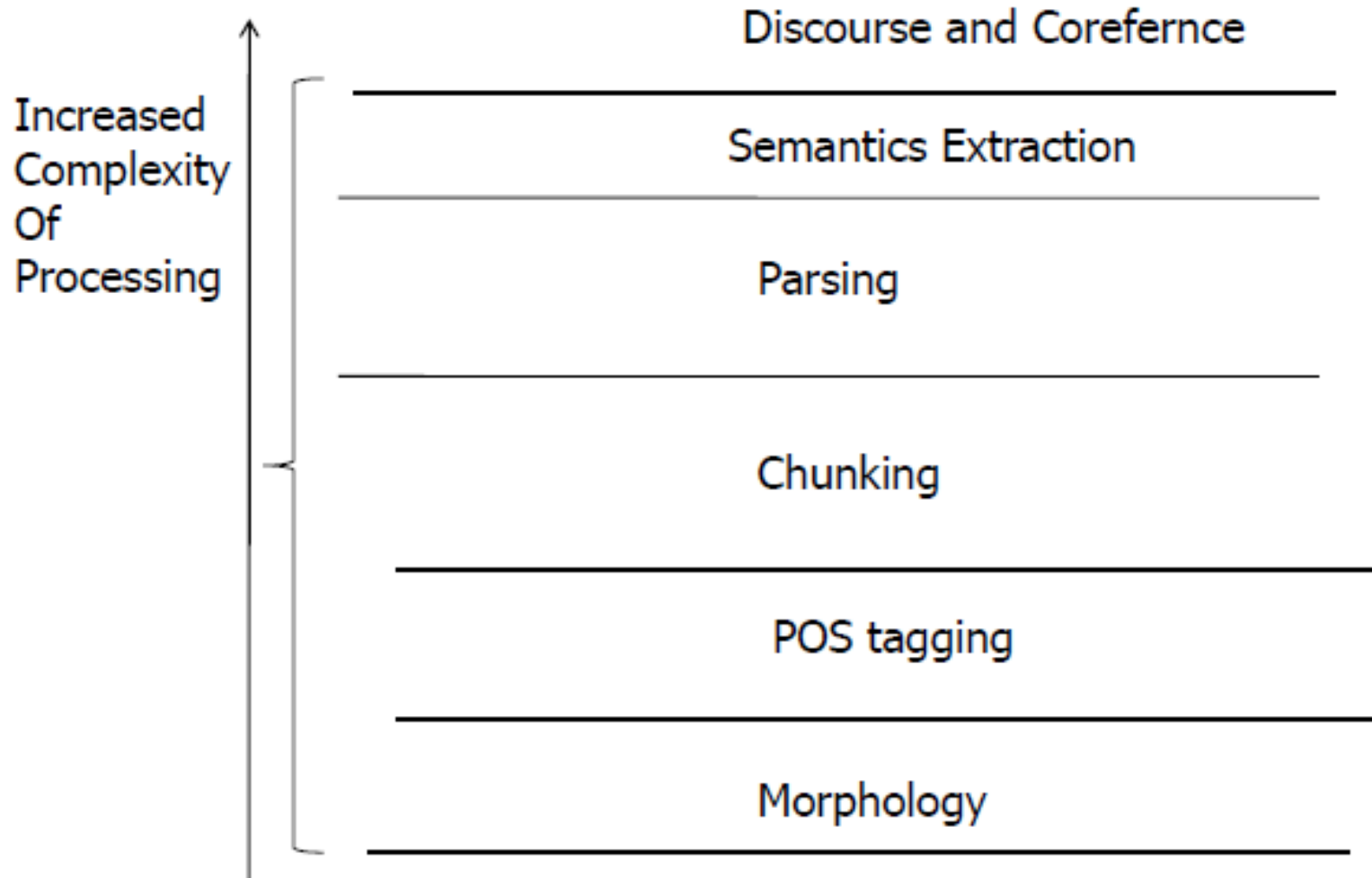
POS tagging

What it is?

- POS Tagging is a process that attaches each word in a sentence with a suitable tag from a given set of tags.



Where does POS tagging fit in



Categories of POS

- Open and closed classes
- Closed classes have a fixed membership of words: determiners, pronouns, prepositions
- Closed class words are usually function word: frequently occurring, grammatically important, often short (e.g. of, it, the, in)
- Open classes: nouns, verbs, adjectives and adverbs

Parts of Speech: How many?

Open class words (content words):

- nouns, verbs, adjectives, adverbs
- mostly content-bearing: they refer to objects, actions, and features in the world
- open class, since new words are added all the time

Parts of Speech: How many?

Closed class words

- pronouns, determiners, prepositions, connectives, ...
- there is a limited number of these
- mostly functional: to tie the concepts of a sentence together

POS examples

- N noun chair, bandwidth, pacing
- V verb study, debate, munch
- ADJ adj purple, tall, ridiculous
- ADV adverb unfortunately, slowly,
- P preposition of, by, to
- PRO pronoun I, me, mine
- DET determiner the, a, that, those

POS tagging: Choosing a tagset

- To do POS tagging, a standard set needs to be chosen
- Could pick very coarse tagsets
N, V, Adj, Adv
- More commonly used set is finer grained,
“UPenn TreeBank tagset”, 45 tags

A Nice Tutorial on POS tags:

<https://sites.google.com/site/partofspeechhelp/>

UPenn TreeBank POS tag set

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VCN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	"	Left quote	<i>(' or ")</i>
POS	Possessive ending	<i>'s</i>	"	Right quote	<i>(' or ")</i>
PRP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([, ({ , <)</i>
PRP\$	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>(] ,) , } , >)</i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>(. ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>(: ; ... - -)</i>
RP	Particle	<i>up, off</i>			

Definition

Example1:

<s> Come in August, and the COEP campus is
abuzz with new and returning students.
</s>

After POS tagging :

<s> Come_VB in_IN August_NNP,_, and_CC
the_DT COEP_NNP campus_NN is_VBZ
abuzz_JJ with_IN new_JJ and_CC
returning_VBG students_NNS.
</s>

POS tagging: Definition

- **Example 2:** “_” The_**DT** guys_**NNS** that_**WDT**
make_**VBP** traditional_**JJ** hardware_**NN**
are_**VBP** really_**RB** being_**VBG**
obsoleted_**VBN** by_**IN**
microprocessorbased_**JJ** machines_**NNS** ,_**,**
”_**”** said_**VBD** Mr._**NNP** Benton_**NNP** ._**.**

Why is POS tagging hard?

- Words often have more than one POS.

Example word: back

- The back door: back/JJ
- On my back: back/NN
- Win the voters back: back/RB
- Promised to back the bill: back/VB

POS tagging problem:

- To determine the POS tag for a particular instance of a word

Brown Corpus: Ambiguous word types

Ambiguity in the Brown corpus:

- 40% of word tokens are ambiguous
- 12% of word types are ambiguous
- Breakdown of ambiguous word types:

Unambiguous (1 tag)	35,340
Ambiguous (2–7 tags)	4,100
2 tags	3,760
3 tags	264
4 tags	61
5 tags	12
6 tags	2
7 tags	1 (“still”)

How bad is the ambiguity problem?

- One tag is usually more likely than the others.
 - In the Brown corpus, race is a noun 98% of the time, and a verb 2% of the time
- A tagger for English that simply chooses the most likely tag for each word can achieve good performance
- Any new approach should be compared against the unigram baseline (assigning each token to its most likely tag)

Deciding the correct POS

Can be difficult even for people:

1. Mrs./NNP Shroff/NNP never/RB got/VBD around/_
to/TO joining/VBG.
2. All/DT we/PRP gotta/VBN do/VB is/VBZ go/VB around/_
the/DT corner/NN.
3. Organic/NNP Onions/NNP costs/VBZ around/_ 250/CD.

Assigning tags:

- Mrs./NNP Shroff/NNP never/RB got/VBD around/RP
to/TO joining/VBG.
- All/DT we/PRP gotta/VBN do/VB is/VBZ go/VB around/IN
the/DT corner/NN.
- Organic/NNP Onions/NNP costs/VBZ around/RB 250/CD.

Relevant knowledge for POS tagging

The word itself:

- Some words may only be nouns, e.g. arrow
- Some words are ambiguous, e.g. flies, like, bank
- Probabilities may help, if one tag is more likely than another

Relevant knowledge for POS tagging

Local context:

- Two determiners rarely follow each other
- Two base form verbs rarely follow each other
- Determiner is almost always followed by adjective or noun

POS tagging: Two approaches

Rule-based Approach:

- Assign each word in the input a list of potential POS tags
- Then reduce down this list to a single tag using hand-written rules

Statistical tagging:

- Get a training corpus of tagged text, learn the transformation rules from the most frequent tags
- Probabilistic: Find the most likely sequence of tags T for a sequence of words W

Probabilistic Tagging: Two different families of models

Problem at hand:

- We have some data $\{(d, c)\}$ of paired observations d and hidden classes c .

Different instances of d and c :

- **Part-of-Speech Tagging**: words are observed and tags are hidden.
- **Text Classification**: sentences/documents are observed and the category is hidden.
- Categories can be positive/negative for sentiments ..
- sports/politics/business for documents ...

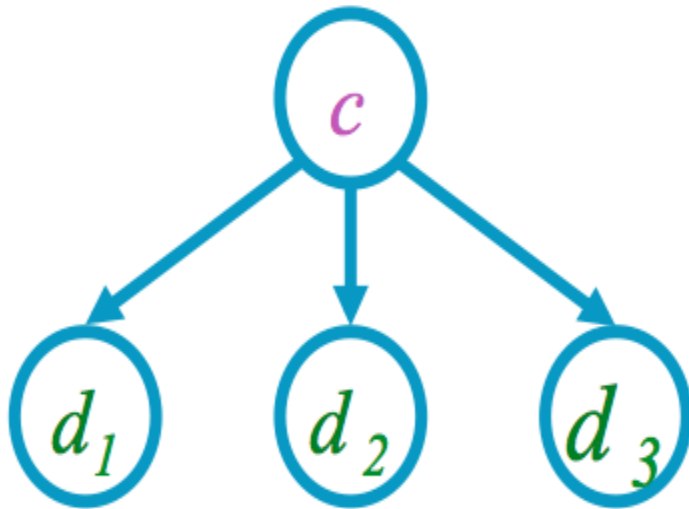
Gives rise to two families?

- Whether they generate the observed data from hidden stuff or the hidden structure given the data?

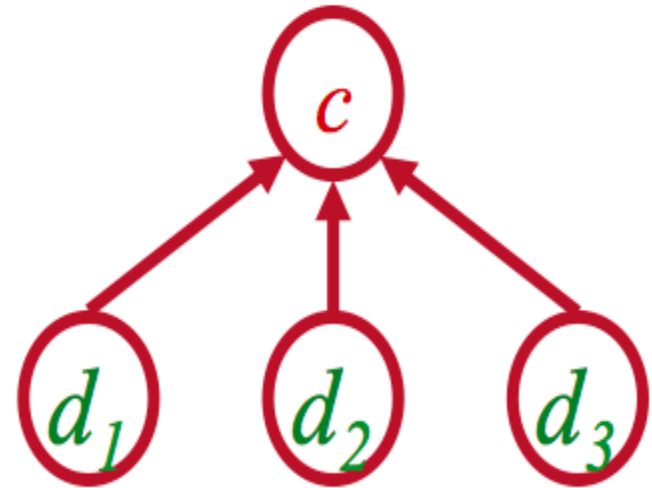
Generative vs. Conditional Models

- **Generative (Joint) Models:**
 - Generate the observed data from hidden stuff, i.e. put a probability over the observations given the class: $P(d, c)$ in terms of $P(d/c)$
 - Egs: Naïve Bayes' classifiers, Hidden Markov Models etc.
- **Discriminative (Conditional) Models:**
 - Take the data as given, and put a probability over hidden structure given the data: $P(c/d)$
 - e.g. Logistic regression, maximum entropy models, conditional random fields

Generative vs. Discriminative Models



Naive Bayes



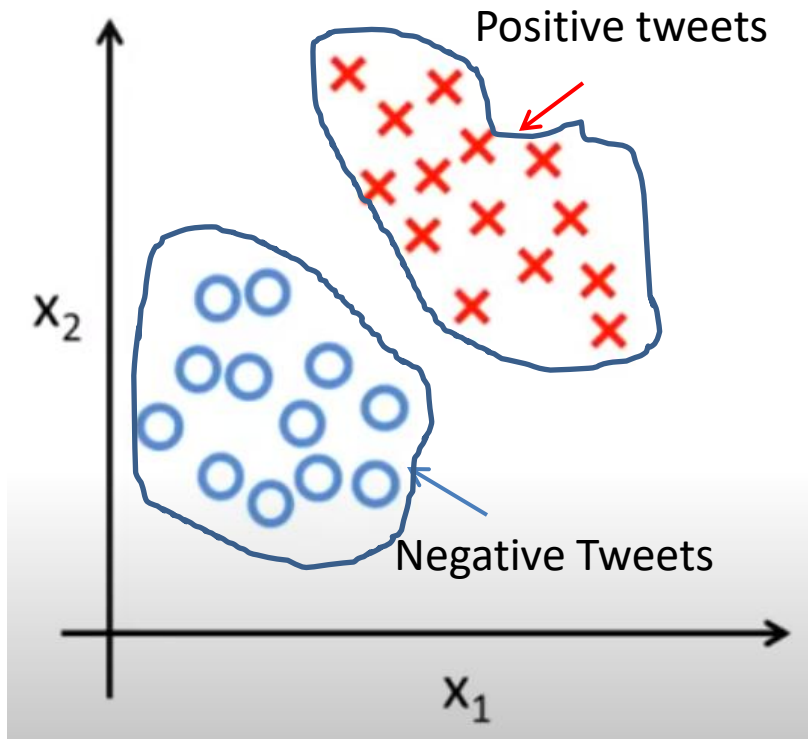
Logistic Regression

Joint vs. conditional likelihood:

- A joint model gives probabilities $P(d/c)$ and tries to maximize this joint likelihood.
- A conditional model gives probabilities $P(c/d)$, taking the data as given and modeling only the conditional probability of the class.

Generative (Joint) Models

Example: Naive Bayes



- In Generative Model, we have to learn $p(x|y)$ and $p(y)$ (class priors) i.e. $p(y=\text{negative tweets})$, $p(y=\text{positive})$
- Generative algorithms try to learn $p(x, y)$ which can be transformed into $p(y|x)$ later to classify the data.

Generative (Joint) Models

- Suppose we have model: $p(x/y)$ and $p(y)$
- Given new x .
- To predict class for x , we need to compute:

$$p(y=1/x) = \frac{p(x/y=1) * p(y=1)}{p(x)} \quad \text{by Bayes rule}$$

Now, we have $p(x/y)$ and $p(y)$ from the model
and

$$p(x) = \sum_y p(x, y) = p(x/y = 1) p(y = 1) + p(x/y = 0) p(y = 0)$$

Example: Generative Model

- Suppose we have trained a generative model, and get a new test example x . Our model tells us that:

$$p(x/y=0)=0.01$$

$$p(x/y=1)=0.03$$

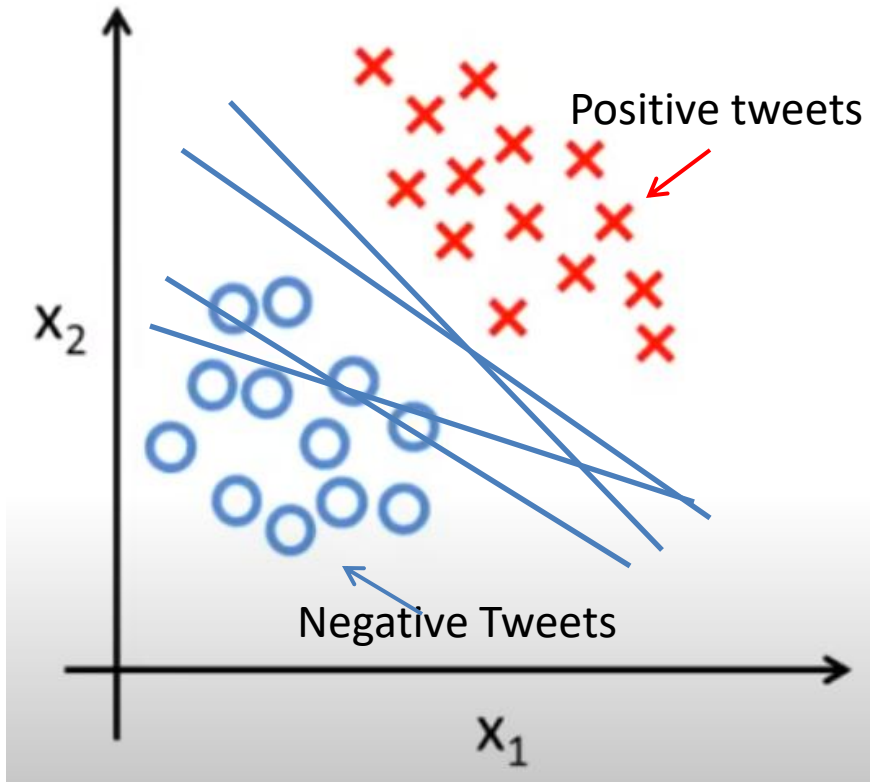
$$p(y=1)=p(y=0)= 0.5$$

What is $p(y=1/x)$?

Solution:

Discriminative (Conditional) Models

Example: Logistic regression



- In Discriminative Model, it directly tries to find a straight line separating the two classes.
- Learns $p(y/x)$ directly

Discriminative (Conditional) Models

- A discriminative algorithm does not care about how the data was generated, it simply categorizes the given data.
- So, discriminative algorithms try to learn $p(y|x)$ directly from the data and then try to classify data.
- Discriminative models do not need to model the distribution of the observed variables.

Mathematics of POS tagging

Argmax Computation

Suppose:

$$x^* = \operatorname{argmax}_x (f(x))$$

Find out value of x which maximizes $f(x)$

Bigram Assumption

Best tag sequence

$$= T^*$$

$$= \operatorname{argmax}_T P(T | W)$$

$$= \operatorname{argmax}_T P(T)P(W | T) \quad (\text{by Bayes Theorem})$$

$$P(T) = P(t_0 = \wedge t_1 t_2 \dots t_{n+1} = .)$$

$$= P(t_0)P(t_1 | t_0)P(t_2 | t_1 t_0)P(t_3 | t_2 t_1 t_0) \dots$$

$$P(t_n | t_{n-1} t_{n-2} \dots t_0)P(t_{n+1} | t_n t_{n-1} \dots t_0)$$

$$= P(t_0)P(t_1 | t_0)P(t_2 | t_1) \dots P(t_n | t_{n-1})P(t_{n+1} | t_n)$$

Bigram Assumption

$$= \prod_{i=0}^{N+1} P(t_i | t_{i-1})$$

Lexical Probability Assumption

$$P(W|T) = P(w_0|t_0-t_{n+1})P(w_1|w_0t_0-t_{n+1})P(w_2|w_1w_0t_0-t_{n+1}) \dots \\ P(w_n|w_0-w_{n-1}t_0-t_{n+1})P(w_{n+1}|w_0-w_nt_0-t_{n+1})$$

Assumption: A word is determined completely by its tag. This is inspired by speech recognition

$$= P(w_0|t_0)P(w_1|t_1) \dots P(w_{n+1}|t_{n+1})$$

$$= \prod_{i=0}^{n+1} P(w_i|t_i)$$

$$= \prod_{i=0}^{n+1} P(w_i|t_i) \quad (\text{Lexical Probability Assumption})$$

Best tag sequence

$$T^* = \operatorname{argmax} P(T)P(W|T)$$

$$P(w_i/t_i)$$

$$= \prod_{i=0}^{N+1} P(t_i|t_{i-1})$$

Process

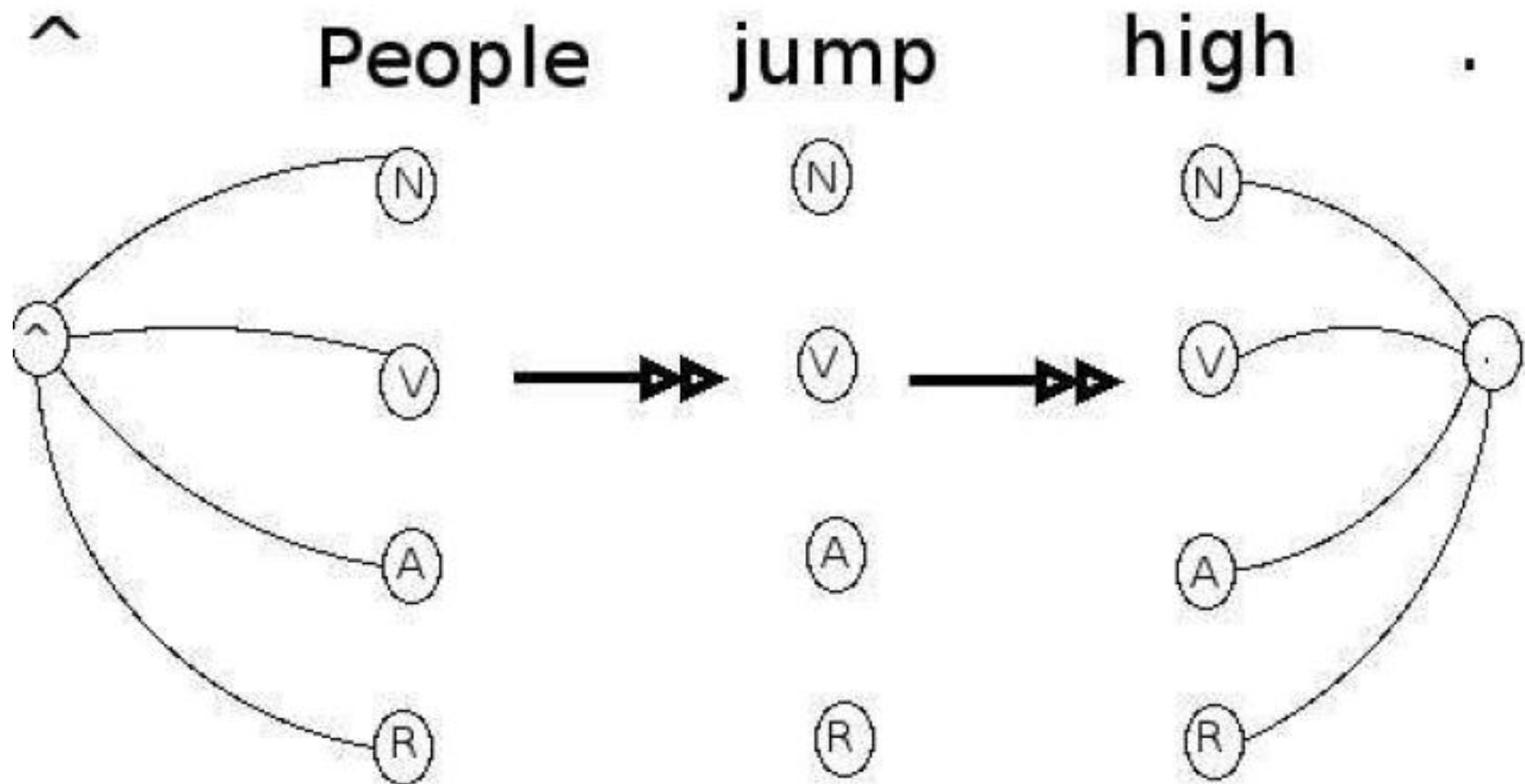
1. List all possible tag for each word in sentence.
2. Choose best suitable tag sequence.

Example

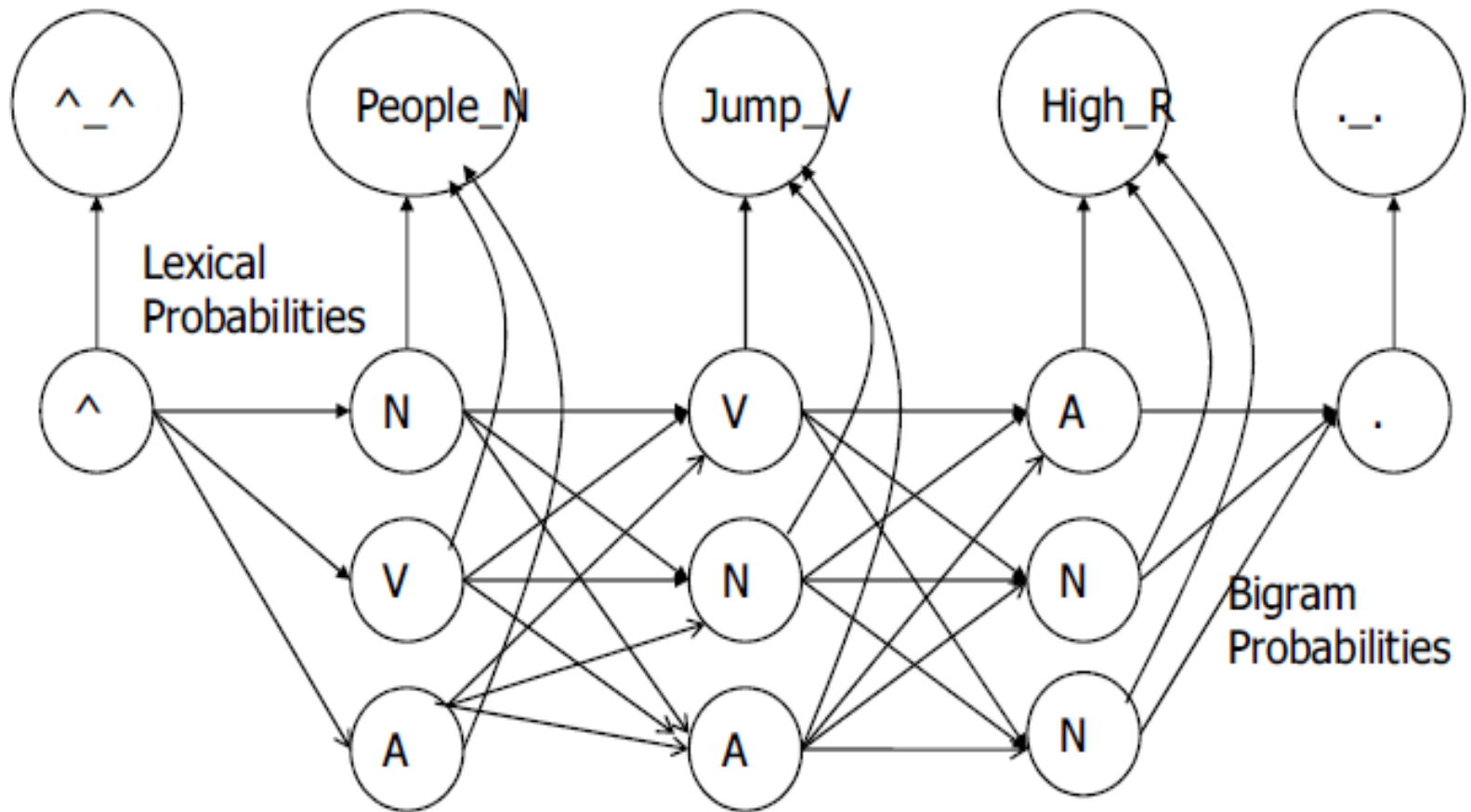
"People jump high".

- People : Noun/Verb/Adjective
- jump : Noun/Verb/Adjective
- high : Noun/Verb/Adjective

Process:



Model



This model is called Generative model.
Here words are observed from tags as states.
This is similar to HMM.

Bigram probabilities from the Corpus

Corpus contains: 300 sentences & 4 categories (N, V, Art, P)

Words: 1998, Nouns: 833, Verbs: 300,

Articles: 558, Prepositions: 307

$$P(\text{ART}/^{\wedge}) = \text{Count} (^{\wedge}, \text{ART}) / \text{Count} (^{\wedge})$$

Category	Count	Pair	Count	Bigram	Prob. Estimate
\wedge	300	\wedge , ART	213	$P(\text{ART}/\wedge)$	0.71
\wedge	300	\wedge , N	87	$P(\text{N}/\wedge)$	0.29
ART	558	ART, N	558	$P(\text{N}/\text{ART})$	1
N	833	N, V	358	$P(\text{V}/\text{N})$	0.43
N	833	N, N	108	$P(\text{N}/\text{N})$	0.13
N	833	N, P	366	$P(\text{P}/\text{N})$	0.44
V	300	V, N	75	$P(\text{N}/\text{V})$	0.35
V	300	V, ART	194	$P(\text{ART}/\text{V})$	0.65
P	307	P, ART	226	$P(\text{ART}/\text{P})$	0.74
P	307	P, N	81	$P(\text{N}/\text{P})$	0.26

Summary of word count in corpus

	N	V	ART	P	TOTAL
flies	21	23	0	0	44
fruit	49	5	1	0	55
like	10	30	0	21	61
a	1	0	201	0	202
the	1	0	300	2	303
flower	53	15	0	0	68
flowers	42	16	0	0	58
birds	64	1	0	0	65
others	592	210	56	284	1142
Total	833	300	558	307	1998

Lexical generation Probabilities: $P(\text{the}/\text{ART}) = \frac{\text{Count}(\text{the as ART})}{\text{Count}(\text{ART})}$

P(the/ART)	0.54	P(like/P)	0.068	P(flower/N)	0.063
P(flies/N)	0.025	P(like/N)	0.012	P(flowers/V)	0.05
P(flies/V)	0.076	P(a/ART)	0.360	P(birds/N)	0.076
P(like/V)	1	P(a/N)	0.001	P(fruit/N)	0.06

Calculation from actual data

- Corpus
 - ^ People Jump High .

Bigram probabilities

	N	V	A
N	0.2	0.7	0.1
V	0.6	0.2	0.2
A	0.5	0.2	0.3

Lexical Probability

	Noun	Verb	Adjective
People	10-5	0.4 X 10-3	10-7
Jump	10-7	10-2	10-7
high	0	0	10-1

- values in cell are P(row-heading/ col-heading)

Observations leading to why probability is needed

1. Many tasks are sequence labeling tasks
2. Tasks carried out in layers
3. Within a layer, there are limited windows of information
4. This naturally calls for strategies for dealing with uncertainty
5. Probability and Markov process give a way for dealing with uncertainty.