

آشنایی با زبان‌شناسی رایانشی

برچسب‌زنی اجزای کلام (Part-of-Speech Tagging)

هادی ویسی

h.veisi@ut.ac.ir

دانشگاه تهران - دانشکده علوم و فنون نوین



- معرفی اجزای کلام: Part-of-Speech (POS)
- برچسب‌زنی اجزای کلام (POS Tagging)
 - کاربردها
 - چالش‌ها
- روش‌های برچسب‌زنی اجزای کلام
 - روش‌های مبتنی بر قاعده (Rule-Based POS Tagging)
 - مثال
 - روش‌های آماری (Probabilistic/Stochastic POS Tagging)
 - مدل مخفی مارکوف (HMM) و الگوریتم ویتربی
 - روش مبتنی بر تبدیل (Transformation-Based Tagging)
- ارزیابی سیستم‌های برچسب‌زنی اجزای کلام
- سیستم نمونه برچسب‌زنی (فارسی)



... اجزای کلام: (POS) Part-of-Speech

- بیانگر مقوله نحوی که هر کلمه به آن تعلق دارد
- نام‌های مختلف

• part-of-speech(POS), tags, lexical tags, word classes

○ مثال

- {من، تو، او، ...} یک {کتاب، گوسفند، درخت، ...} را {دیدم، خریدم، فروختم، ...}

○ برچسب‌زنی (Tagging) = POS Tagging

- فرایند انتساب مقوله نحوی به هر کلمه در پیکره متنی

○ ایده قدیمی است

- پیشنهاد ۸ دسته توسط Thrax در ۱۰۰ ق.م.

○ noun, verb, pronoun, preposition, adverb, conjunction, participle, article

- مشابه برچسب‌های مورد استفاده در کتب درسی امروزی

○ noun, verb, adjective, preposition, adverb, conjunction, pronoun, interjection



... اجزای کلام: (POS) Part-of-Speech

○ POS های اصلی در انگلیسی

- اسم (Noun)
 - she, who, my, others, ...
- فعل (Verb)
 - Ali, IBM, Book, ...
- حرف ربط (Conjunctions)
 - and, but, if, ...
- افعال کمکی (Auxiliaries)
 - (can, should, are, ...)
- صفت (Adjective)
 - Good, Beautiful, Young, ...
- حروف (Particles)
 - Slowly, Very, Fortunately, ...
- شماره‌ها (Numerals)
 - on, to, by, from, with, ...
- حرف تعریف (Determiner)
 - a, an, the



اجزای کلام: (POS) Part-of-Speech ...

○ بسته به کاربرد، ممکن است POS های جزئی‌تری نیز به کار روند

- اسم (Noun)
 - فرد یا جمع (Singular و Plural)
 - صفت (Adjective)
 - تفضیلی یا عالی (Comparative یا Superlative)
 - فعل (Verb)
 - اصلی یا وجهی (Main یا Modal)
 - ضمیر (Pronoun)
 - شخصی یا مالکیت (Personal یا Possessive)
 - ...
- **Noun (person, place or thing)**
 - Singular (NN): dog, fork
 - Plural (NNS): dogs, forks
 - Proper (NNP, NNPS): John, Springfields
 - Personal pronoun (PRP): I, you, he, she, it
 - Wh-pronoun (WP): who, what
 - **Verb (actions and processes)**
 - Base, infinitive (VB): eat
 - Past tense (VBD): ate
 - Gerund (VBG): eating
 - Past participle (VBN): eaten
 - Non 3rd person singular present tense (VBP): eat
 - 3rd person singular present tense: (VBZ): eats
 - Modal (MD): should, can
 - To (TO): to (to eat)

○ برای علائم نقطه‌گذاری هم ممکن است برچسب‌های متفاوتی به کار رود



... اجزای کلام: (POS) Part-of-Speech

○ انواع اجزای کلام (POS: Part-Of-Speech)

- واژه‌های محتوایی (Content/Lexical Words) = گروه باز = پذیرش اعضای جدید
 - حامل معنی و استفاده به صورت مستقل
 - در برخی دسته‌ها، اعضای نامتناهی دارند (مانند اسم): تولید اعضای جدید (اسم خاص)

مثال	نقش	توصیف	برچسب
cat	نامیدن موجودات	Noun (اسم)	N
forget	نامیدن رویداد یا شرایط	Verb (فعل)	V
yellow	توصیفی	Adjective (صفت)	Adj
quickly	حالت عملکرد	Adverb (قید)	Adv
oh!	واکنش	Interjection (حرف ندا)	Interj

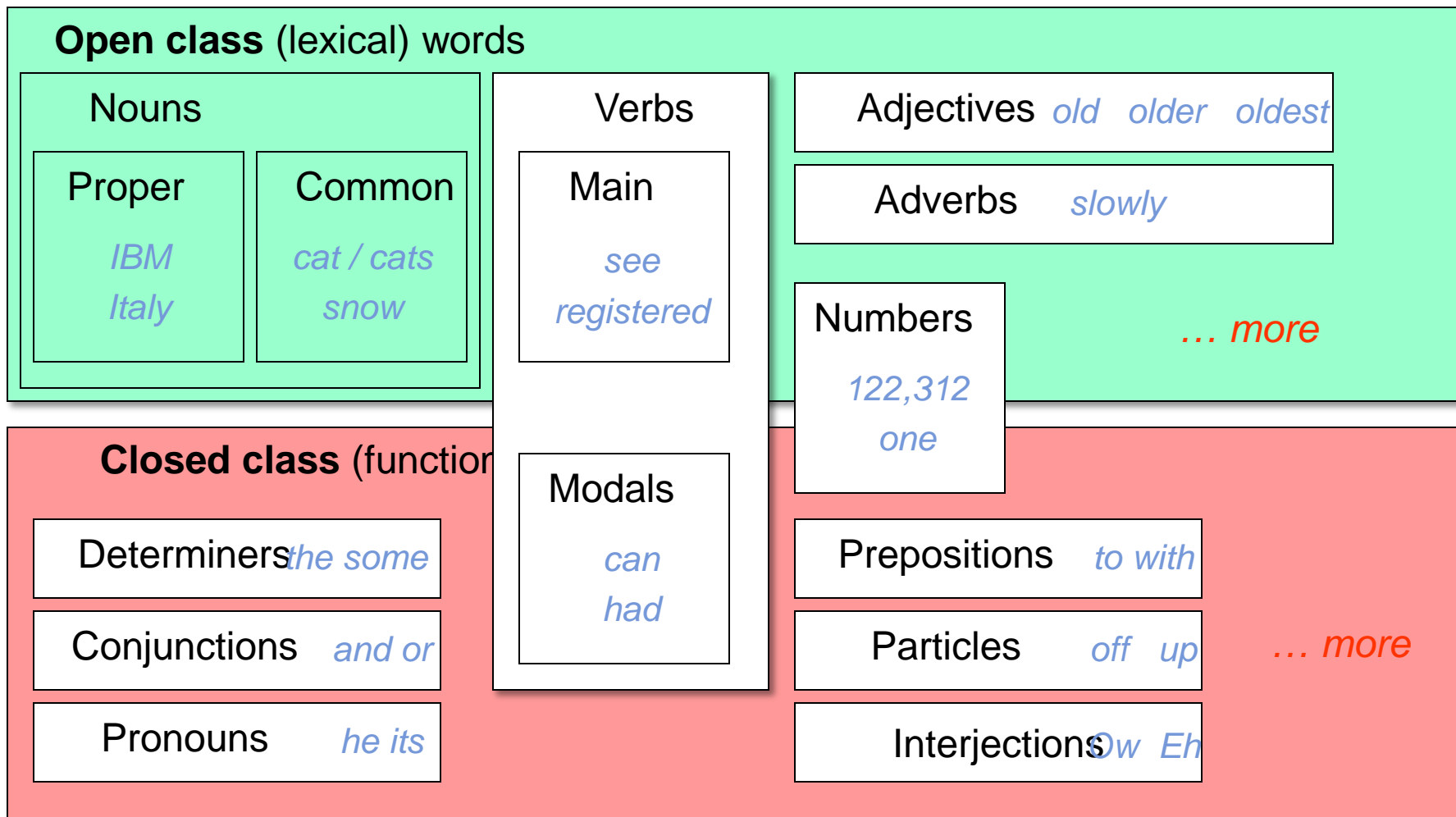
- واژه‌های نقشی (Function Words): گروه بسته = اعضای محدود و ثابت

○ کوتاه، پرتکرار و دارای ساختار نحوی مشخص و استفاده برای بیان رابطه گرامری با سایر کلمات

مثال	نقش	توصیف	برچسب
and	پیوند دهنده‌ی عبارات	Conjunction (حرف ربط)	Conj
the	مشخص‌کننده‌ی شناخته‌بودن	Determiner (حرف تعریف)	Det
from	روابط زمانی، مکانی، جهت‌ی	Preposition (حرف اضافه)	Prep
she	ارجاع ساده شده	Pronoun (ضمیر)	Pron



اجزای کلام: (POS) Part-of-Speech ...





اجزای کلام: (POS) Part-of-Speech ...

Penn Treebank

مثال	توصیف	برچسب
and	حرف ربط پیوند دهنده	CC
two	عدد شمارشی	CD
the	حرف تعریف	DT
there	عبارت وجود داشتن there	EX
omerta	واژه‌ی خارجی	FW
over, but	حرف اضافه، حرف ربط پیرو	IN
yellow	صفت	JJ
better	صفت، مقایسه‌ای	JJR
best	صفت، تفضیلی	JJS
	نشانه‌ی عنصر لیست	LS
might	فعل وجهی	MD
rock, water	اسم، مفرد یا جمعی	NN
rocks	اسم، جمع	NNS
Joe	اسم خاص، مفرد	NNP
Red Guards	اسم خاص، جمع	NNPS
all (all the girls)	حرف تعریف پیشین	PDT
's	نشانه‌ی ملکی	POS
I	ضمیر شخصی	PRP
Mine	ضمیر ملکی	PRP\$
Quickly	قید	RB
higher	قید، مقایسه‌ای	RBR
highest	قید، تفضیلی	RBS
up (take up.)	وند فعلی	RP
to	To	TO
hey!	حرف ندا	UH
choose	بن فعل	VB
chose	فعل، زمان گذشته	VBD
choosing	فعل، استمراری یا صفت فاعلی	VBG
chosen	فعل، صفت مفعولی	VBN
jump	فعل، مفرد غیر سوم شخص زمان حال	VBP
jumps	فعل، سوم شخص مفرد زمان حال	VBZ
which	حرف تعریف Wh	WDT
who	ضمیر Wh	WP
whose	ضمیر ملکی Wh	WP\$
when	قید Wh	WRB

مجموعه برچسب‌های (tag set) انگلیسی

• برچسب‌های Penn Treebank شامل ۴۵ برچسب

○ امکان کاهش تا ۱۲ برچسب رایج

• برچسب‌های Brown Corpus شامل ۸۷ برچسب

○ پیکره یک میلیون کلمه‌ای از ۵۰۰ نویسنده در ژانرهای مختلف

• برچسب‌های C5 شامل ۶۱ برچسب

○ در Lancaster UCREL project's CLAWS

• برچسب‌های C7 شامل ۱۴۶ برچسب

در پیکره متنی زبان فارسی (دکتر بیجن خان)

• حدود ۶۶۰ برچسب

○ کاهش تعداد برچسب‌ها به حدود ۸۰ مورد



... اجزای کلام: (POS) Part-of-Speech

برچسب‌های Penn Treebank (رایج در انگلیسی)

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &</i>
CD	cardinal number	<i>one, two</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>	VBN	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, sing.	<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>			



... اجزای کلام: (POS) Part-of-Speech

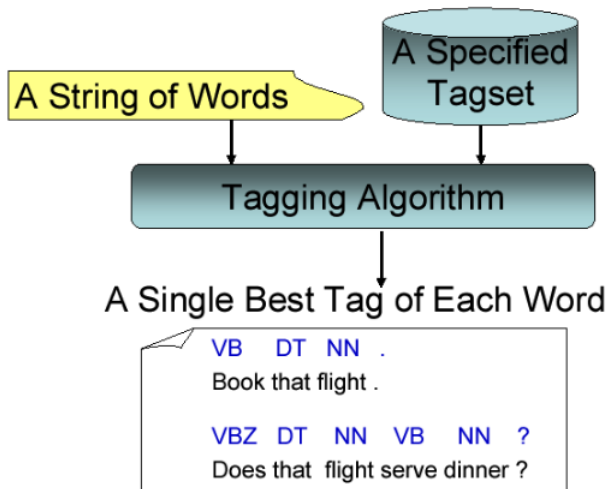
برچسب‌های پیکره بیجن خان (فارسی)

No	Tag	Description	No	Tag	Description
1	ADJ	Adjective, General	21	MQUA	Modifier of Quantifier
2	ADJ_CMPR	Adjective, Comparative	22	MS	Mathematic Symbol
3	ADJ_INO	Past Participle	23	N_PL	Noun, Plural
4	ADJ_ORD	Adjective, Ordinal	24	N_SING	Noun, Singular
5	ADJ_SIM	Adjective, Simple	25	NN	Number
6	ADJ_SUP	Adjective, Superlative	26	NP	Noun Phrase
7	ADV	Adverb, General	27	OH	Oh Interjection (حرف ندا)
8	ADV_EXM	Adverb, Exemplar	28	OHH	Oh noun (منادی)
9	ADV_I	Adverb, Question	29	P	Preposition
10	ADV_NEGG	Adverb, Negation	30	PP	Prepositional Phrase
11	ADV_NI	Adverb, Not Question	31	PRO	Pronoun
12	ADV_TIME	Adverb, Time	32	PS	Pseudo-Sentence
13	AR	Arabic Word	33	QUA	Quantifier
14	CON	Conjunction	34	SPEC	Specifier
15	DEFAULT	Default	35	V_AUX	Verb, Auxiliary
16	DELM	Delimiter	36	V_IMP	Verb, Imperative
17	DET	Determiner	37	V_PA	Verb, Past Tense
18	IF	Conditional	38	V_PRE	Verb, Predicative
19	INT	Interjection	39	V_PRS	Verb, Present Tense
20	MORP	Morpheme	40	V_SUB	Verb, Subjunctive



برچسب‌زنی اجزای کلام (POS Tagging)

○ یافتن برچسب کلمات در متن



Penn Treebank
POS tags

- Input: Plays well with others
- Ambiguity: NNS/VBZ UH/JJ/NN/RB IN NNS
- Output: Plays/VBZ well/RB with/IN others/NNS

○ ابتدا باید متن را نرمال کرد و علائم نگارشی را جدا تفکیک کرد



برچسب‌زنی اجزای کلام: کاربردها

- مدل‌سازی زبانی (در بازشناسی گفتار و ...)
 - استفاده در پیش‌بینی کلمه بعدی
 - مثال: در انگلیسی بعد از ضمایر ملکی، اسم و بعد از ضمایر شخصی، فعل می‌آید.
 - در فارسی معمولاً بعد از حرف اضافه اسم یا ضمیر می‌آید.
- سنتز گفتار: حاوی اطلاعاتی در مورد نحوه تلفظ صحیح یک کلمه
 - مثال: تلفظ کلمه object در انگلیسی به صورت OBject (noun) و obJECT (verb)
 - تلفظ کلمه "مرد" در فارسی به صورت "مَرَد" (اسم)، "مُرد" (فعل)
- بازیابی اطلاعات: کمک به استخراج کلمات مهم در متن
- رفع ابهام معنایی
 - کمک به رفع ابهام معنایی (استفاده در سیستم‌های ترجمه و ...)
 - مثال: کلمه watch در انگلیسی و کلمه "در" در فارسی
- تجزیه نحوی (parsing)
 - در تجزیه پایین به بالا، نیاز به دانستن مقوله نحوی کلمات است



برچسب‌زنی اجزای کلام: چالش ...

○ بعضی از کلمات به چند مقوله نحوی تعلق دارند

- حدود ۴۰٪ از کلمات (Token) پیکره Brown دارای بیش از یک برچسب هستند
- حدود ۱۱.۵٪ از انواع کلمات (Word Type) به کار رفته در پیکره را تشکیل می‌دهند

	Original 87-tag corpus	Treebank 45-tag corpus
Unambiguous (1 tag)	44,019	38,857
Ambiguous (2-7 tags)	5,490	8844
Details:		
2 tags	4,967	6,731
3 tags	411	1621
4 tags	91	357
5 tags	17	90
6 tags	2 (<i>well, beat</i>)	32
7 tags	2 (<i>still, down</i>)	6 (<i>well, set, round, open, fit, down</i>)
8 tags		4 (<i>'s, half, back, a</i>)
9 tags		3 (<i>that, more, in</i>)

ابهام انواع
کلمات در پیکره
Brown



برچسب‌زنی اجزای کلام: چالش ...

○ مثال از ابهام در برچسب زدن

- Mrs/NNP Shaefer/NNP never/RB got/VBD
around/RP to/TO joining/VBG
- All/DT we/PRP gola/VBN do/VB is/VBZ go/VB
around/IN the/DT corner/NN
- Chateau/NNP Petrus/NNP costs/VBZ
around/RB 250/CD

- I know **that** he is honest = IN (relativizer)
- Yes, **that** play was nice = DT (determiner)
- You can't go **that** far = RB (adverb)



برچسب‌زنی اجزای کلام: منابع اطلاعات (ویژگی‌ها) ...

○ از چه اطلاعاتی برای برچسب‌زنی اجزای کلام می‌توان استفاده کرد؟

• خود کلمه (به تنهایی)

○ Word the: the → DT

○ حرف بزرگ در اول کلمه (انگلیسی)

○ کلمات با حرف اول بزرگ، معمولاً به عنوان اسم (خاص) برچسب زده می‌شوند اما برای کلمات اول جمله بهتر است حالت حروف کوچک آن هم بررسی شود

○ Lowercased word Importantly: importantly → RB

○ وجود برخی پسوندها و پیشوندها: بیانگر نوع کلمه است

○ Prefixes unfathomable: un- → JJ

○ Suffixes Importantly: -ly → RB

○ حالت حروف بزرگ کلمات

○ Capitalization Meridian: CAP → NNP

○ ساختار کلمه (عبارت)

○ Word shapes 35-year: d-x → JJ



برچسب‌زنی اجزای کلام: منابع اطلاعات (ویژگی‌ها)

○ از چه اطلاعاتی برای برچسب‌زنی اجزای کلام می‌توان استفاده کرد؟

• کلمات همسایه و دنباله کلمات

- Bill saw that man yesterday
- NNP NN DT NN NN
- VB VB(D) IN VB NN

○ در این مثال، ترکیب DT VB (حرف تعریف قبل از فعل) کمتر رخ می‌دهد

• احتمال برچسب‌ها/کلمات

- احتمال اینکه کلمه man فعل (VB) باشد، خیلی کم است
- احتمال پست سر هم آمدن برچسب‌های مختلف پشت سر هم
- در فارسی: احتمال اینکه بعد از «اسم» برچسب «صفت» بیاید زیاد است اما عکس آن کم است



برچسب‌زنی اجزای کلام: روش‌ها

○ روش‌های مبتنی بر قاعده (Rule-Based POS Tagging)

- استفاده از یک واژگان با برچسب(های) اجزای کلام برای هر کلمه و قوانینی برای رفع ابهام
- مثال: EngCG (Constraint Grammar)

○ روش‌های مبتنی بر یادگیری از داده

- آماری (Probabilistic/Stochastic POS Tagging)
 - استفاده از احتمال برای انتخاب برچسب‌ها
 - استفاده از روش‌های یادگیری آماری و یک پیکره دارای برچسب برای محاسبه احتمال‌ها
 - مثال: HMM Tager
- شبکه عصبی مصنوعی (Artificial Neural Network)

○ ترکیبی: مبتنی بر تبدیل (Transformation-Based POS Tagging)

- استفاده از قواعد برای رفع ابهام (روش‌های مبتنی بر قاعده)
- استفاده از یادگیری ماشین برای استخراج خودکار قوانین از پیکره متنی
- Brill's tagger

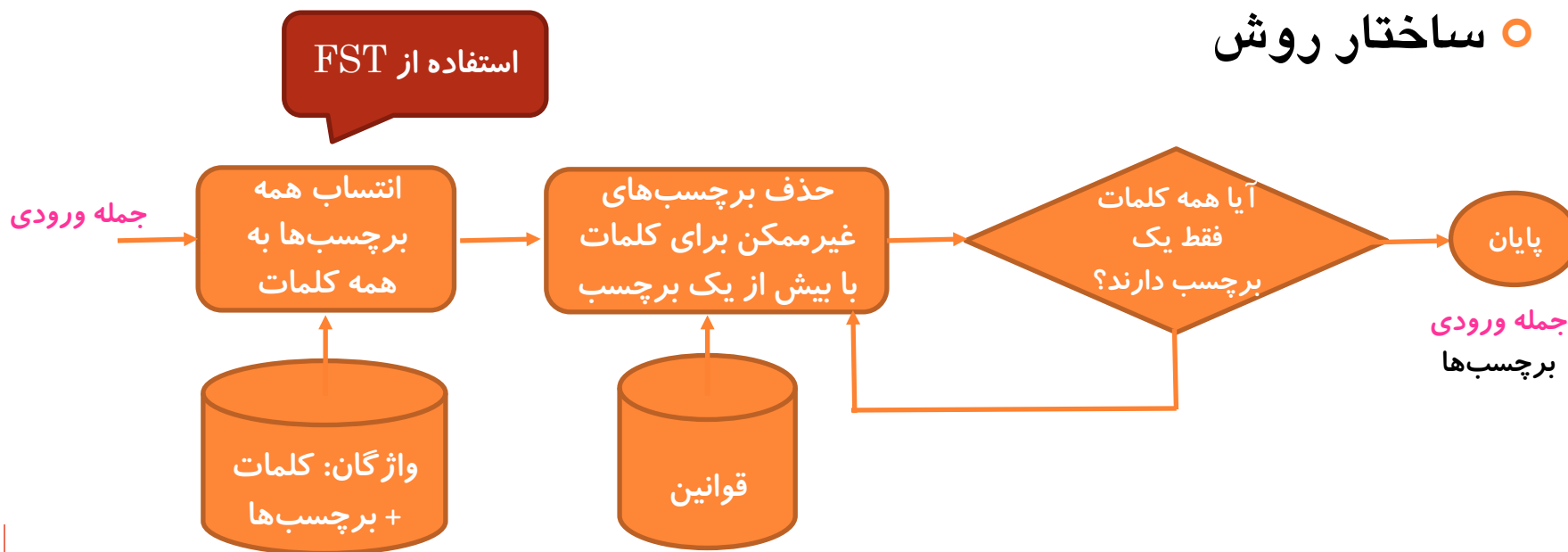


برچسب‌زنی اجزای کلام: روش مبتنی بر قاعده ...

○ نیازمندی‌ها

- واژگانی حاوی کلمات و همه برچسب‌های ممکن برای هر کلمه
- مجموعه قوانینی برای حذف برچسب‌های اضافی انتساب داده شده به یک کلمه
- حذف برچسب‌های غیرممکن تا فقط یک برچسب برای هر کلمه باقی بماند

○ ساختار روش





برچسب‌زنی اجزای کلام: روش مبتنی بر قاعده ...

○ مثال ...

• جمله ورودی: She promised to back the bill

- she: PRP
 - promised: VBN, VBD
 - to: TO
 - back: VB, JJ, RB, NN
 - the: DT
 - bill: NN, VB
- واژگان با برچسب‌ها

• گام ۱: انتساب همه برچسب‌های ممکن به کلمات

			NN			
			RB			
	VBN		JJ			VB
PRP	VBD	TO	VB	DT		NN
She	promised	to	back	the		bill



برچسب‌زنی اجزای کلام: روش مبتنی بر قاعده ...

مثال

گام ۲: اعمال قوانین برای حذف برخی از برچسب‌ها

E.g., *Eliminate VBN if VBD is an option when VBN|VBD follows "<start> PRP"*

			NN			
			RB			
	VBN		JJ			VB
PRP	VBD	TO	VB	DT	NN	
She	promised	to	back	the	bill	



E.g., *Eliminate VBN if VBD is an option when VBN|VBD follows "<start> PRP"*

			NN			
			RB			
			JJ			VB
PRP	VBD	TO	VB	DT	NN	
She	promised	to	back	the	bill	



برچسب‌زنی اجزای کلام: روش مبتنی بر قاعده

Word	POS	Additional POS features
smaller	ADJ	COMPARATIVE
entire	ADJ	ABSOLUTE ATTRIBUTIVE
fast	ADV	SUPERLATIVE
that	DET	CENTRAL DEMONSTRATIVE SG
all	DET	PREDETERMINER SG/PL QUANTIFIER
dog's	N	GENITIVE SG
furniture	N	NOMINATIVE SG NOINDEFDETERMINER
one-third	NUM	SG
she	PRON	PERSONAL FEMININE NOMINATIVE SG3
show	V	PRESENT -SG3 VFIN
show	N	NOMINATIVE SG
shown	PCP2	SVOO SVO SV
occurred	PCP2	SV
occurred	V	PAST VFIN SV

روش EngCG ENGTWOL

- استفاده از یک واژگان با ۵۶۰۰۰ ریشه

- به کارگیری یک Lexicon FST برای انتساب برچسب‌ها به کلمات

Pavlov **PAVLOV N NOM SG PROPER**
 had **HAVE V PAST VFIN SVO**
 HAVE PCP2 SVO
 shown **SHOW PCP2 SVOO SVO SV**
 that **ADV**
 PRON DEM SG
 DET CENTRAL DEM SG
 CS
 salivation **N NOM SG**
 ...

• برای جمله . . . Pavlov had shown that salivation

- استفاده از تعدادی زیادی محدودیت (۳۷۴۴ قانون) برای کاهش برچسب‌ها

ADVERBIAL-THAT RULE

Given input: "that"

if

(+1 A/ADV/QUANT); /* if next word is adj, adverb, or quantifier */

(+2 SENT-LIM); /* and following which is a sentence boundary, */

(NOT -1 SVOC/A); /* and the previous word is not a verb like */

/* 'consider' which allows adjs as object complements */

then eliminate non-ADV tags

else eliminate ADV tag



برچسب‌زنی اجزای کلام: روش آماری ...

ایده

- در نظر گرفتن احتمال وقوع برچسب‌ها (برچسب‌های محتمل) برای کلمات
- با فرض داشتن دنباله کلمات $W=w_1 \dots w_n$ ، دنباله برچسب‌های $T=t_1 \dots t_n$ را طوری پیدا کنید که $P(T | W)$ بیشینه شود

$$\hat{T} = \arg \max_T P(T | W)$$

قانون بیز

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

حذف $P(W)$ تغییری در بیشینه کردن ایجاد نمی‌کند

$$\hat{T} = \arg \max_T P(T | W) = \arg \max_T \frac{P(W | T)P(T)}{P(W)}$$

$$= \arg \max_T \underbrace{P(W | T)}_{\text{Likelihood}} \underbrace{P(T)}_{\text{Prior}} = \arg \max_T P(w_1 w_2 \dots w_n | t_1 t_2 \dots t_n) P(t_1 t_2 \dots t_n)$$

Likelihood Prior

- نحوه محاسبه؟



برچسب‌زنی اجزای کلام: روش آماری ...

○ برای محاسبه نیاز به فرض‌های ساده کننده است

$$\hat{T} = \arg \max_T P(w_1 w_2 \cdots w_n | t_1 t_2 \cdots t_n) P(t_1 t_2 \cdots t_n)$$

- فرض اول: احتمال وقوع یک کلمه فقط به برچسب آن کلمه وابسته است و مستقل از سایر کلمات و برچسب‌های اطراف آن است

$$P(w_1 w_2 \cdots w_n | t_1 t_2 \cdots t_n) \approx \prod_{i=1}^n P(w_i | t_i)$$

- فرض دوم: دیدن یک برچسب فقط به برچسب قبلی آن وابسته است (Bi-Gram)

$$P(t_1 t_2 \cdots t_n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$

- بنابر این داریم:

$$P(w_1 w_2 \cdots w_n | t_1 t_2 \cdots t_n) P(t_1 t_2 \cdots t_n) \approx \prod_{i=1}^n P(w_i | t_i) \cdot \prod_{i=1}^n P(t_i | t_{i-1})$$



برچسب‌زنی اجزای کلام: روش آماری ...

○ مفهوم احتمال $P(t_i | t_{i-1})$ (Tag Transition Probability)

- احتمال آمدن یک برچسب (t_i) بعد از برچسب دیگر (t_{i-1})
- احتمال آمدن «اسم» (NN) یا «صفت» (JJ) بعد از «حرف تعریف» (DT) بالاست
The beautiful story ○
- در پیکره Brown داریم: $P(NN | DT) = 0.49$

○ مفهوم احتمال $P(w_i | t_i)$ (Word Likelihood)

- اگر دنبال کلمه‌ای با برچسب t_i هستیم، احتمال اینکه آن کلمه w_i باشد
- $P(is | VBZ)$ = احتمال اینکه کلمه با برچسب «VBZ» (فعل حال سوم شخص مفرد)، کلمه «is» باشد
○ در پیکره Brown داریم: $P(is | VBZ) = 0.47$



برچسب‌زنی اجزای کلام: روش آماری ...

○ محاسبه احتمال‌ها

- نیاز به یک پیکره متنی داریم که در آن کلمات دارای برچسب باشند
- محاسبه احتمال $P(t_i | t_{i-1})$ (Tag Transition Probability)

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}t_i)}{C(t_{i-1})}$$

- مثال: برای محاسبه $P(NN | DT)$ (در پیکره Brown) - برچسب DT به تعداد 116,454 آمده است که بعد از 56,509 مورد از آنها NN آمده است. پس

$$P(NN | DT) = \frac{56509}{116454} = 0.49$$

- محاسبه احتمال $P(w_i | t_i)$ (Word Likelihood)

$$P(w_i | t_i) = \frac{C(w_i, t_i)}{C(t_i)}$$

- مثال: محاسبه $P(is | VBZ)$ (در پیکره Brown) - برچسب «VBZ» به تعداد 21,627 بار آمده که از میان آنها، تعداد 10,073 کلمه «is» است. پس

$$P(is | VBZ) = \frac{10073}{21627} = 0.47$$

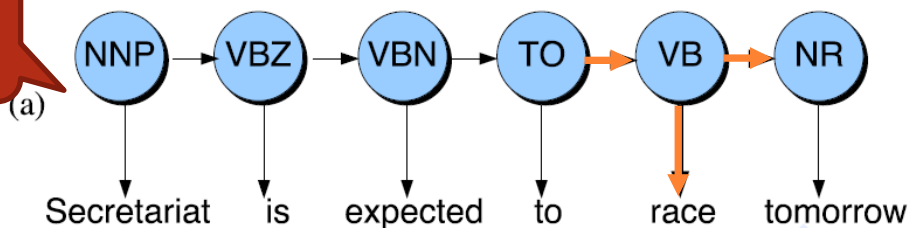


برچسب‌زنی اجزای کلام: روش آماری ...

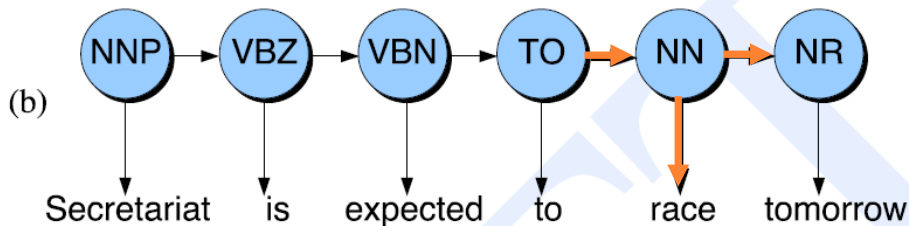
مثال: محاسبه احتمال دنباله برای تعیین برچسب درست

- کلمه race هم می‌تواند فعل (VB) باشد و هم اسم (NN)
- Secretariat/NNP is/BEZ expected/VBN to/TO race/VB tomorrow/NR
- در نظر گرفتن دو حالت از دنباله حالات برای تعیین برچسب race در جمله اول

هر فلش بیانگر یک مقدار احتمال است



تفاوت دو حالت در ۳ مقدار احتمال

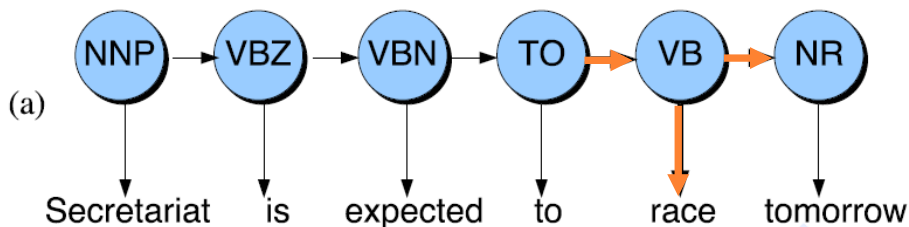




برچسب‌زنی اجزای کلام: روش آماری ...

مثال: محاسبه احتمال دنباله برای تعیین برچسب درست

Secretariat/NNP is/BEZ expected/VBN to/TO race/VB tomorrow/NR

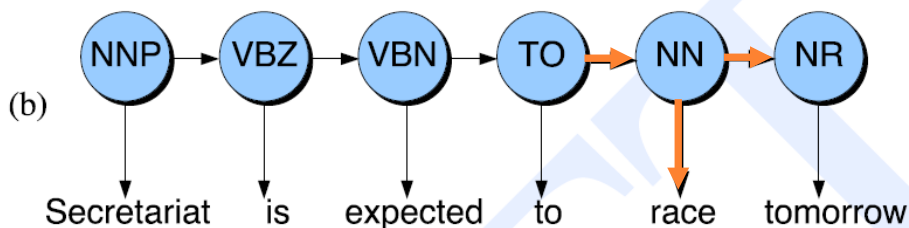


$$P(NN|TO) = .00047$$

$$P(VB|TO) = .83$$

$$P(NR|VB) = .0027$$

$$P(NR|NN) = .0012$$



$$P(\text{race}|NN) = .00057$$

$$P(\text{race}|VB) = .00012$$

مقدار احتمال برای برچسب VB بیشتر است

$$P(VB|TO)P(NR|VB)P(\text{race}|VB) = .00000027$$

$$P(NN|TO)P(NR|NN)P(\text{race}|NN) = .00000000032$$

موارد ابهام دیگر

excepted می‌تواند صفت (J)، فعل گذشته (VBD) یا اسم مفعول (VBN) باشد



برچسب‌زنی اجزای کلام: روش آماری ...

○ محاسبه محتمل‌ترین دنباله از برچسب‌ها

- ساده‌ترین روش: در نظر گرفتن تمام دنباله‌های محتمل و محاسبه احتمال هر یک به روش بیان شده (Brute Force Search)

- با فرض داشتن N برچسب و T کلمه، حداکثر N^T دنباله از برچسب‌ها تولید می‌شود.
 - محاسبات بسیار زیاد

• روش‌های رایج

- مدل مخفی مارکوف (HMM: Hidden Markov Model)
- میدان تصادفی شرطی (CRF: Conditional Random Field)



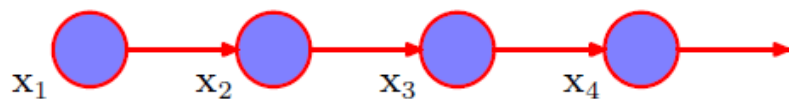
برچسب‌زنی اجزای کلام: روش آماری ...

○ زنجیره مارکوف (Markov Chain)

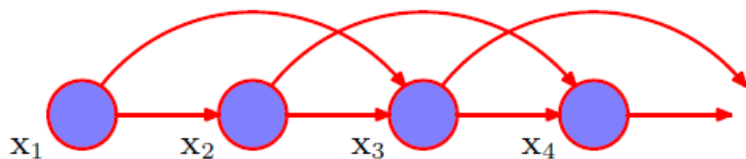
- نوع خاصی از ماشین حالت متناهی (FSA) که در آن به حرکت بین حالت‌ها یک احتمال نسبت داده می‌شود.

- در نظر گرفتن وابستگی بین حالت‌ها

○ درجه اول (first-order Markov chain): هر حالت تنها به یک حالت قبل وابسته است



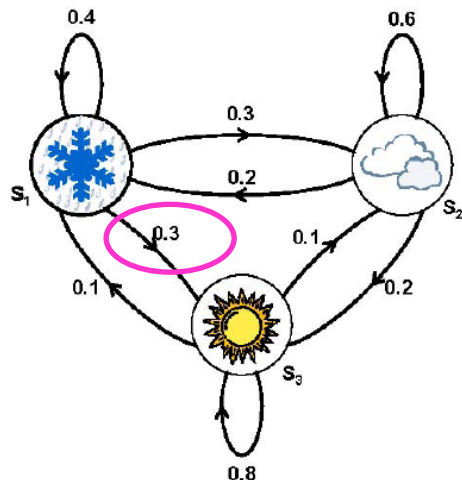
○ درجه دوم (second-order Markov chain): هر حالت به دو حالت قبل وابسته است





برچسب‌زنی اجزای کلام: روش آماری ...

○ زنجیره مارکوف: مثال پیش‌بینی وضعیت هوا ...



- در نظر گرفتن ۳ حالت (state) مختلف

- حالت ۱ (S_1): بارندگی (برف یا باران)

- حالت ۲ (S_2): ابری

- حالت ۳ (S_3): آفتابی

- در نظر گرفتن احتمال انتقال حالت‌ها (state transition probability)

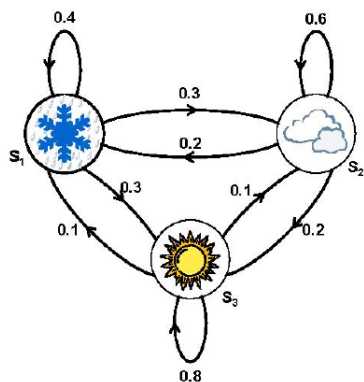
$$A = \{a_{ij}\} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

احتمال رفتن از حالت ۱
(بارندگی) به حالت ۳ (آفتابی)

- سوال ۱: با فرض اینکه امروز آفتابی است، احتمال اینکه هوای ۷ روز آینده به صورت زیر باشد، چقدر است؟ {آفتابی، آفتابی، باران، باران، آفتابی، ابری، آفتابی}



برچسب‌زنی اجزای کلام: روش آماری ...



$$A = \{a_{ij}\} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

○ زنجیره مارکوف: مثال پیش‌بینی وضعیت هوا ...

- سوال ۱: امروز آفتابی است، احتمال اینکه هوای ۷ روز آینده به صورت زیر باشد، چقدر است؟ {آفتابی، باران، آفتابی، باران، آفتابی، ابری، آفتابی}
- استفاده از وابستگی درجه ۱ (وضعیت هر روز به روز قبل)

مشاهده

شامل حالت‌ها و ارتباط بین آنها

$$P(O|Model) = P[S_3, S_3, S_3, S_1, S_1, S_3, S_2, S_3|Model]$$

$$= P[S_3] \cdot P[S_3|S_3] \cdot P[S_3|S_3] \cdot P[S_1|S_3]$$

$$\cdot P[S_1|S_1] \cdot P[S_3|S_1] \cdot P[S_2|S_3] \cdot P[S_3|S_2]$$

$$= \pi_3 \cdot a_{33} \cdot a_{33} \cdot a_{31} \cdot a_{11} \cdot a_{13} \cdot a_{32} \cdot a_{23}$$

$$= 1 \cdot (0.8)(0.8)(0.1)(0.4)(0.3)(0.1)(0.2)$$

$$= 1.536 \times 10^{-4}$$

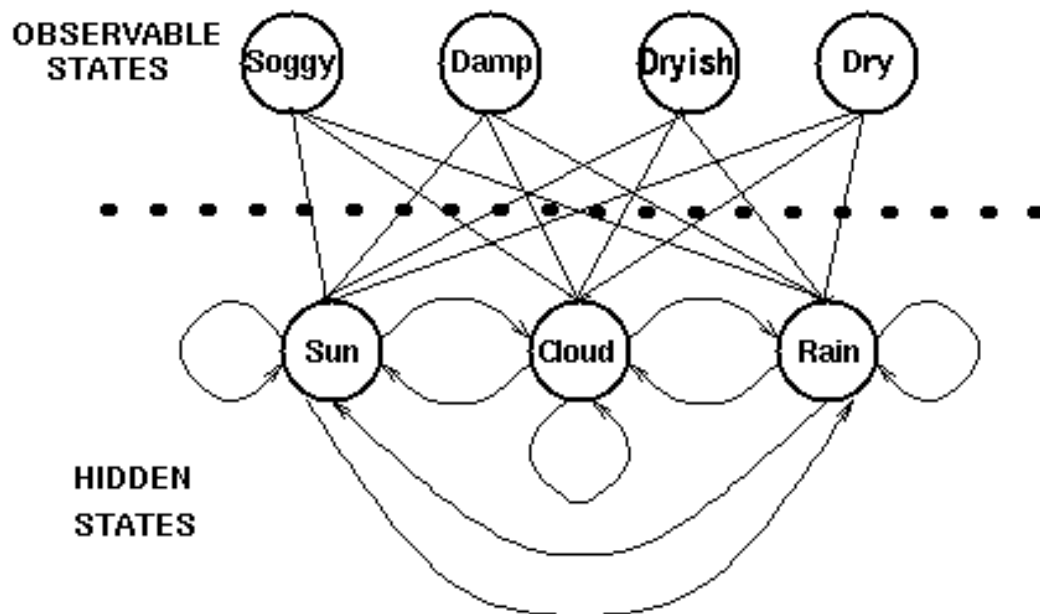
احتمال حالت اولیه
(initial state probability)



برچسب‌زنی اجزای کلام: روش آماری ...

○ زنجیره مارکوف: مثال پیش‌بینی وضعیت هوا

- در آنچه تاکنون بیان شد: **حالت‌ها** با **مشاهده‌ها** یکسان بودند
- در بسیاری از کاربردها، مشاهده‌ها با حالت‌های مساله یکی نیستند
- حالت‌های اصلی **مخفی** هستند و باید مشاهده‌ها را با آنها متناظر کرد



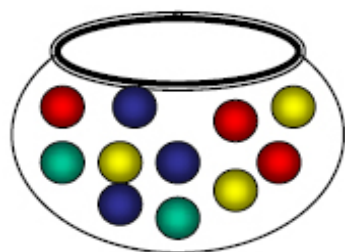
مدل مخفی مارکوف
Hidden Markov Model



برچسب‌زنی اجزای کلام: روش آماری ...

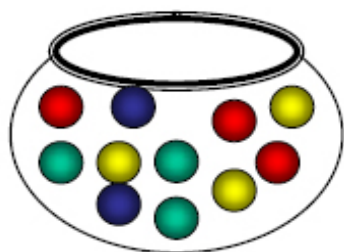
○ مدل مخفی مارکوف: مثال گوی و گلدان ...

- فرض کنید تعداد N گلدان در یک اتاق داریم
- در هر گلدان تعداد زیادی گوی رنگی، شامل M رنگ وجود دارد
- انتخاب هر گوی در هر گلدان متناسب با مقدار احتمال مرتبط است



URN 1

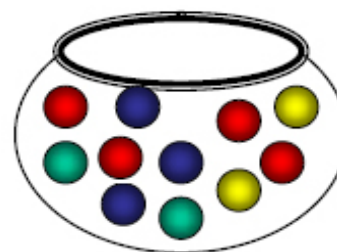
$$\begin{aligned} P(\text{RED}) &= b_1(1) \\ P(\text{BLUE}) &= b_1(2) \\ P(\text{GREEN}) &= b_1(3) \\ P(\text{YELLOW}) &= b_1(4) \\ &\vdots \\ P(\text{ORANGE}) &= b_1(M) \end{aligned}$$



URN 2

$$\begin{aligned} P(\text{RED}) &= b_2(1) \\ P(\text{BLUE}) &= b_2(2) \\ P(\text{GREEN}) &= b_2(3) \\ P(\text{YELLOW}) &= b_2(4) \\ &\vdots \\ P(\text{ORANGE}) &= b_2(M) \end{aligned}$$

...



URN N

$$\begin{aligned} P(\text{RED}) &= b_N(1) \\ P(\text{BLUE}) &= b_N(2) \\ P(\text{GREEN}) &= b_N(3) \\ P(\text{YELLOW}) &= b_N(4) \\ &\vdots \\ P(\text{ORANGE}) &= b_N(M) \end{aligned}$$



برچسب‌زنی اجزای کلام: روش آماری ...

○ مدل مخفی مارکوف: مثال گوی و گلدان

• تعداد N گلدان و M رنگ

• فرآیند

- یک نفر (در اتاقی که ما نمی‌بینیم)، یکی از گلدان‌ها را به صورت تصادفی انتخاب می‌کند
- از داخل گلدان انتخاب شده، یک گوی را بیرون آورده و رنگ آن را اعلام می‌کند
- گوی را به داخل گلدان مربوطه برمی‌گرداند
- بر اساس مقداری تصادفی وابسته به گلدان فعلی، گلدان بعدی انتخاب می‌شود
- مراحل فوق به صورت متوالی تکرار می‌شود

• دنباله مشاهده: دنباله گوی‌ها (رنگ‌ها)

$O = \{GREEN, GREEN, BLUE, RED, YELLOW, RED, \dots, BLUE\}$

• حالت‌ها: گلدان‌ها (از دید مشاهده کننده مخفی است)

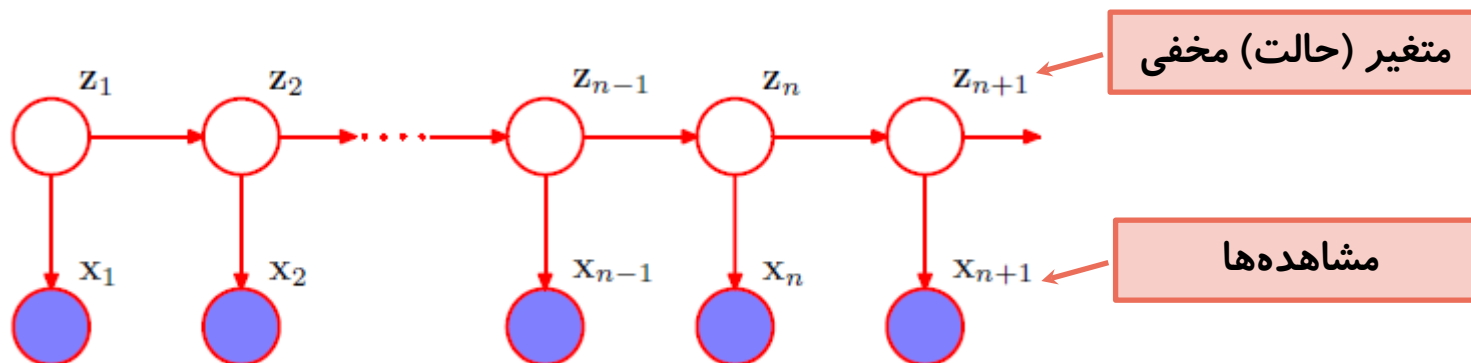
• انتقال حالت‌ها: فرآیند انتخاب گلدان‌ها



برچسب‌زنی اجزای کلام: روش آماری ...

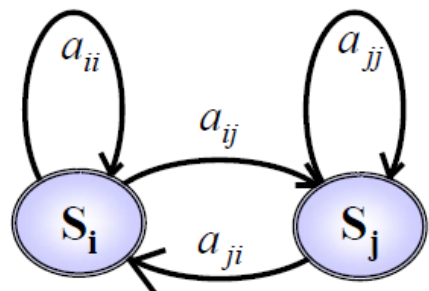
○ مدل مخفی مارکوف: مبانی

- مشاهده‌ها توابع احتمالاتی از حالت‌ها هستند
- دنباله حالت‌ها قابل مشاهده نیستند (مخفی هستند)
- فرض وابستگی درجه اول
- مشاهده‌ها فقط به حالت‌ها وابسته هستند و نه به همدیگر





برچسب‌زنی اجزای کلام: روش آماری ...



$$\begin{matrix} P(v_1 | S_i) \\ P(v_2 | S_i) \\ \vdots \\ P(v_M | S_i) \end{matrix}$$

$$\begin{matrix} P(v_1 | S_j) \\ P(v_2 | S_j) \\ \vdots \\ P(v_M | S_j) \end{matrix}$$

$$S = \{S_1, \dots, S_N\}$$

$$V = \{v_1, v_2, \dots, v_M\}$$

$$a_{ij} = P(q_t = S_j | q_{t-1} = S_i)$$

$$\pi = \{\pi_i\} = P(q_1 = i)$$

در گوی و گلدان: احتمال انتخاب هر کدام از گلدان‌ها در اولین مشاهده (زمان $t=1$)

$$b_j(k) = P(o_t = v_k | q_t = j)$$

در گوی و گلدان: احتمال انتخاب گوی k ام از گلدان j ام

تابع توزیع مشاهده‌ها (مثلاً گاوسی) - احتمال تولید مشاهده $o_t = v_k$ در حالت $q_t = j$

○ مدل مخفی مارکوف: عناصر اصلی

① مجموعه‌ای از N حالت

○ در گوی و گلدان: گلدان‌ها

② مجموعه‌ای از M نماد مشاهده

○ در گوی و گلدان: رنگ‌ها

③ احتمال انتقال حالت‌ها

○ در گوی و گلدان: جابجایی از یک گلدان به گلدان دیگر

④ احتمال اولیه حالت‌ها

○ در گوی و گلدان: احتمال انتخاب هر کدام از گلدان‌ها در اولین مشاهده (زمان $t=1$)

⑤ تابع توزیع برای نماد k ام در حالت j ام

○ در گوی و گلدان: احتمال انتخاب گوی k ام از گلدان j ام

○ تابع توزیع مشاهده‌ها (مثلاً گاوسی) - احتمال تولید مشاهده $o_t = v_k$ در حالت $q_t = j$

○ نمایش یک مدل مخفی مارکوف $\lambda = (A, B, \pi)$



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

○ مدل مخفی مارکوف در برچسب‌زنی اجزای کلام

① مجموعه‌ای از N حالت = هر حالت بیانگر یک برچسب

○ در گوی و گلدان: گلدان‌ها

② مجموعه‌ای از M نماد مشاهده = هر مشاهده بیانگر یک کلمه

○ در گوی و گلدان: رنگ‌ها

③ احتمال انتقال حالت‌ها = احتمال وقوع یک برچسب بعد از دیگری

○ در گوی و گلدان: جابجایی از یک گلدان به گلدان دیگر

④ احتمال اولیه حالت‌ها = احتمال اینکه اولین کلمه چه برچسبی داشته باشد

○ در گوی و گلدان: احتمال انتخاب هر کدام از گلدان‌ها در اولین مشاهده (زمان $t=1$)

⑤ تابع توزیع برای مشاهده k ام در حالت z ام = احتمال اینکه برچسب z ام کلمه k ام باشد

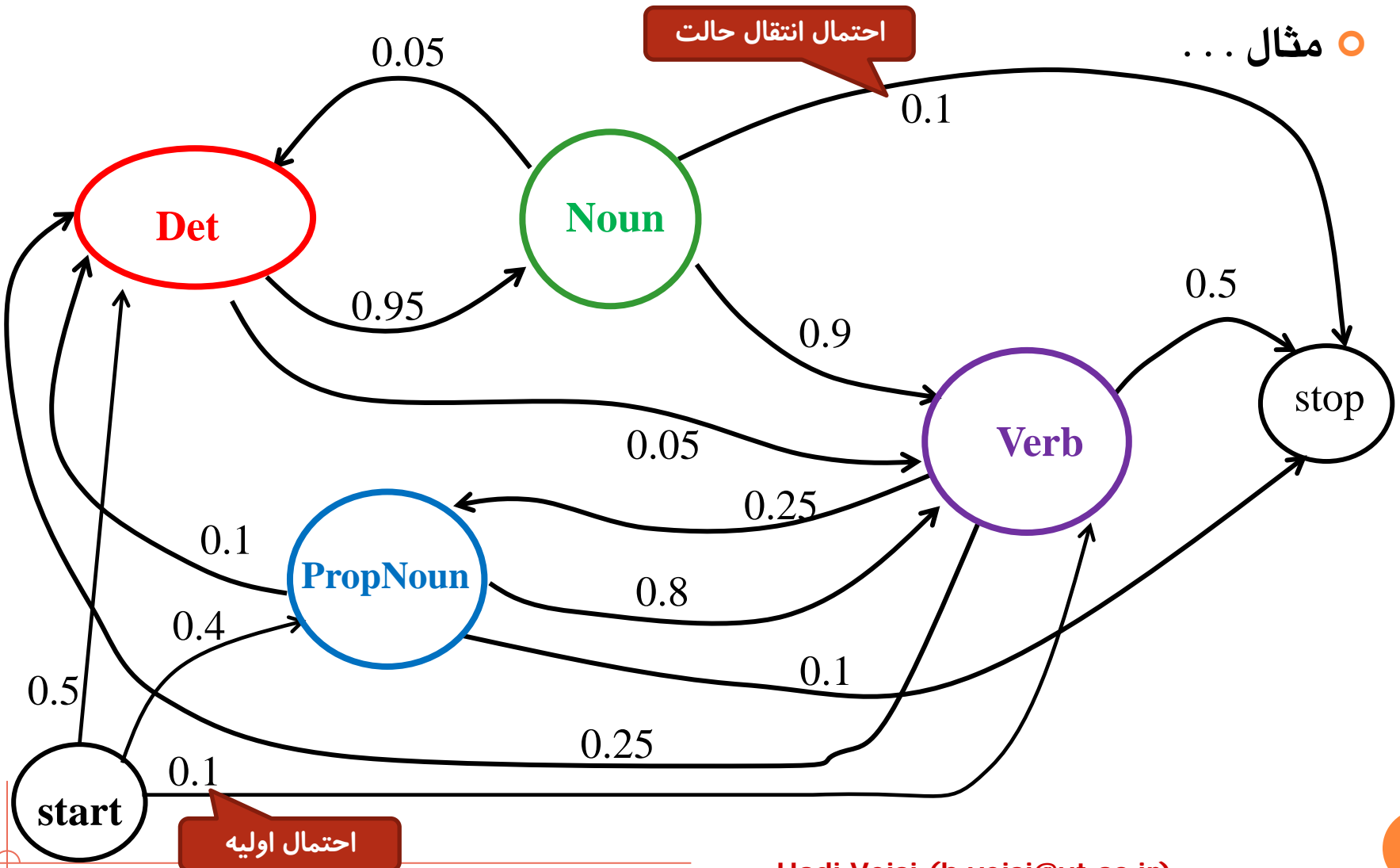
○ در گوی و گلدان: احتمال انتخاب گوی k ام از گلدان z ام

○ تابع توزیع مشاهده‌ها (مثلاً گاوسی) - احتمال تولید مشاهده $O_t = V_k$ در حالت $q_t = z$



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

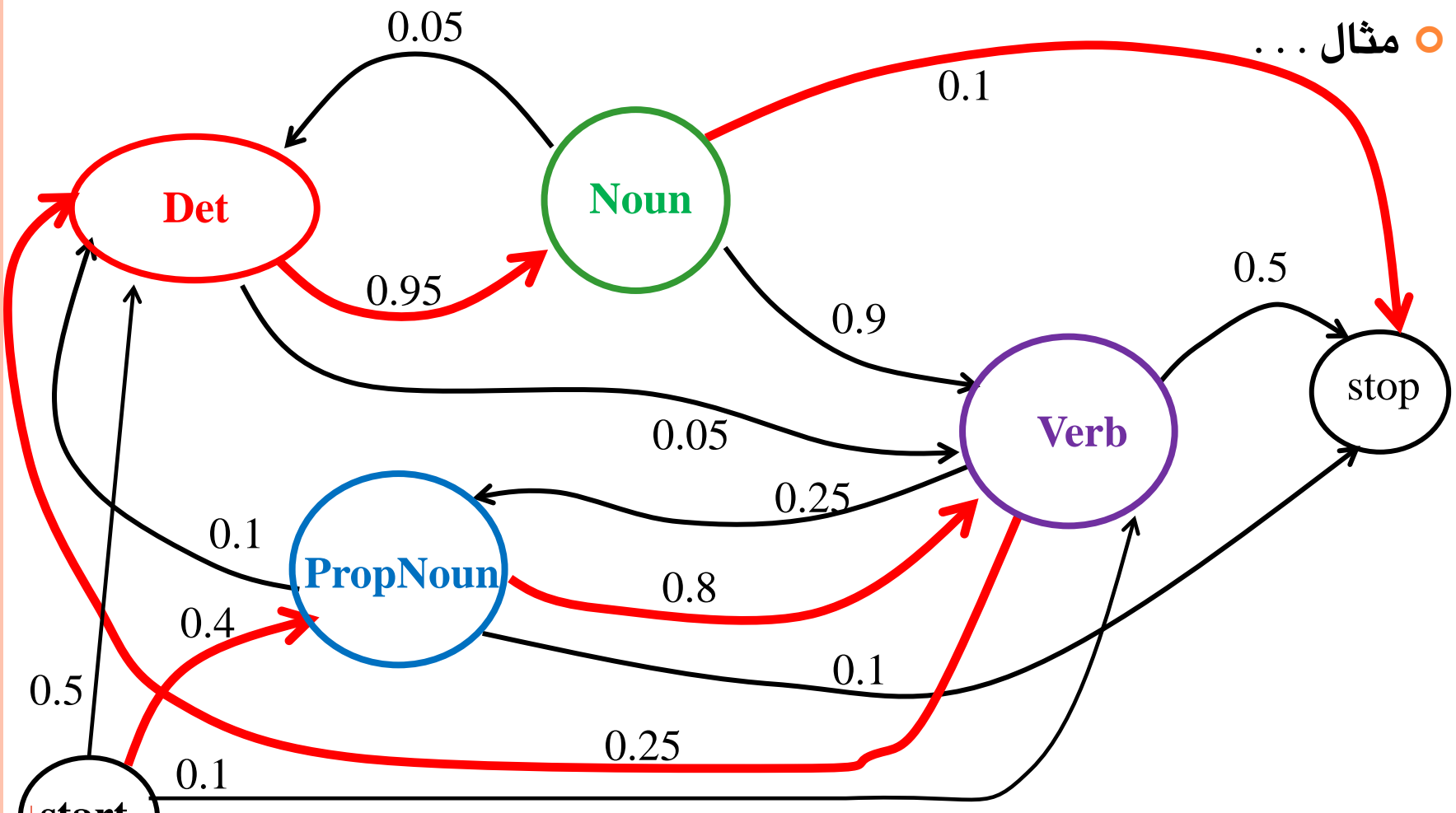
○ مثال ...





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

مثال ...

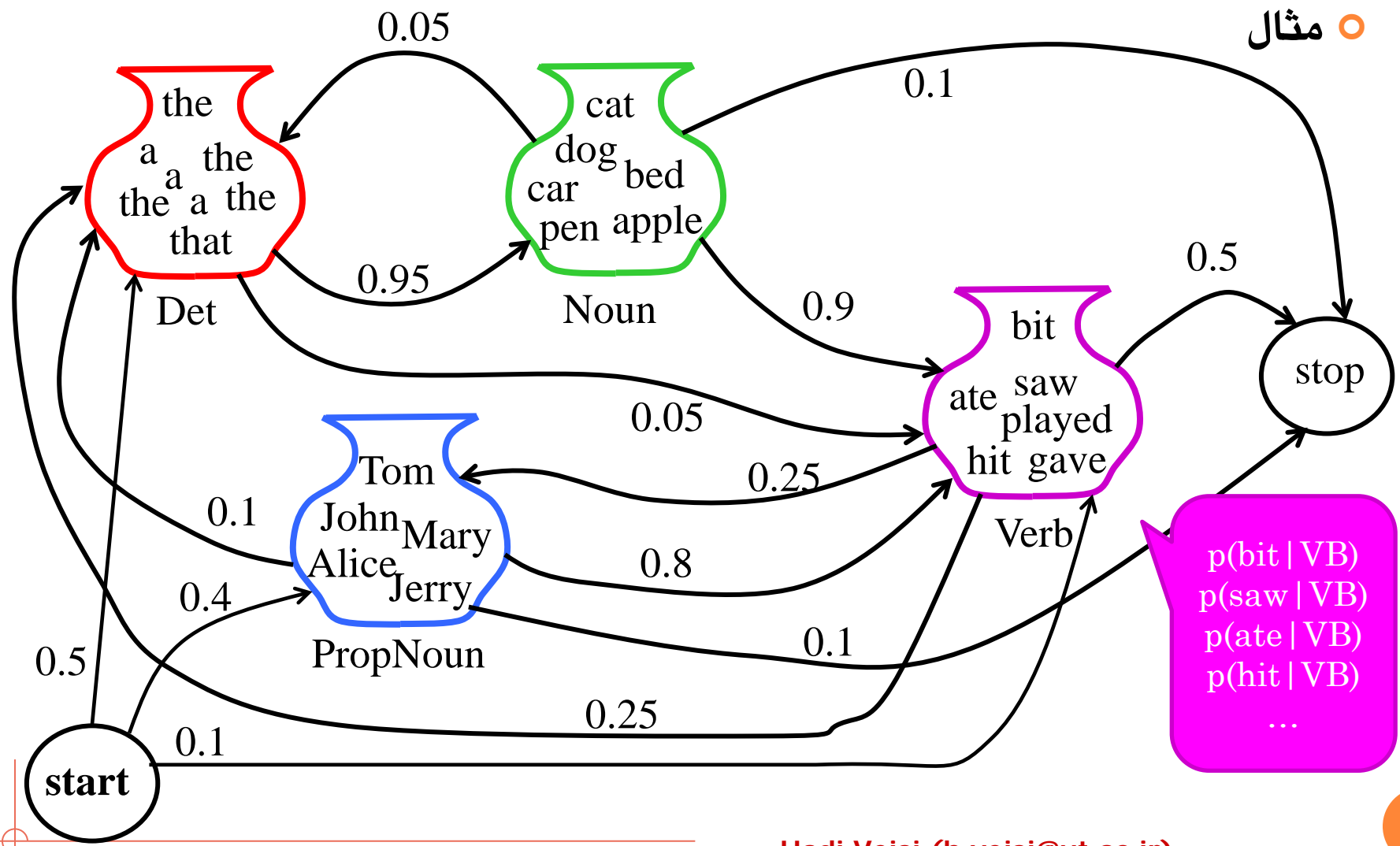


$$P(\text{PropNoun Verb Det Noun}) = 0.4 * 0.8 * 0.25 * 0.95 * 0.1 = 0.0076$$



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

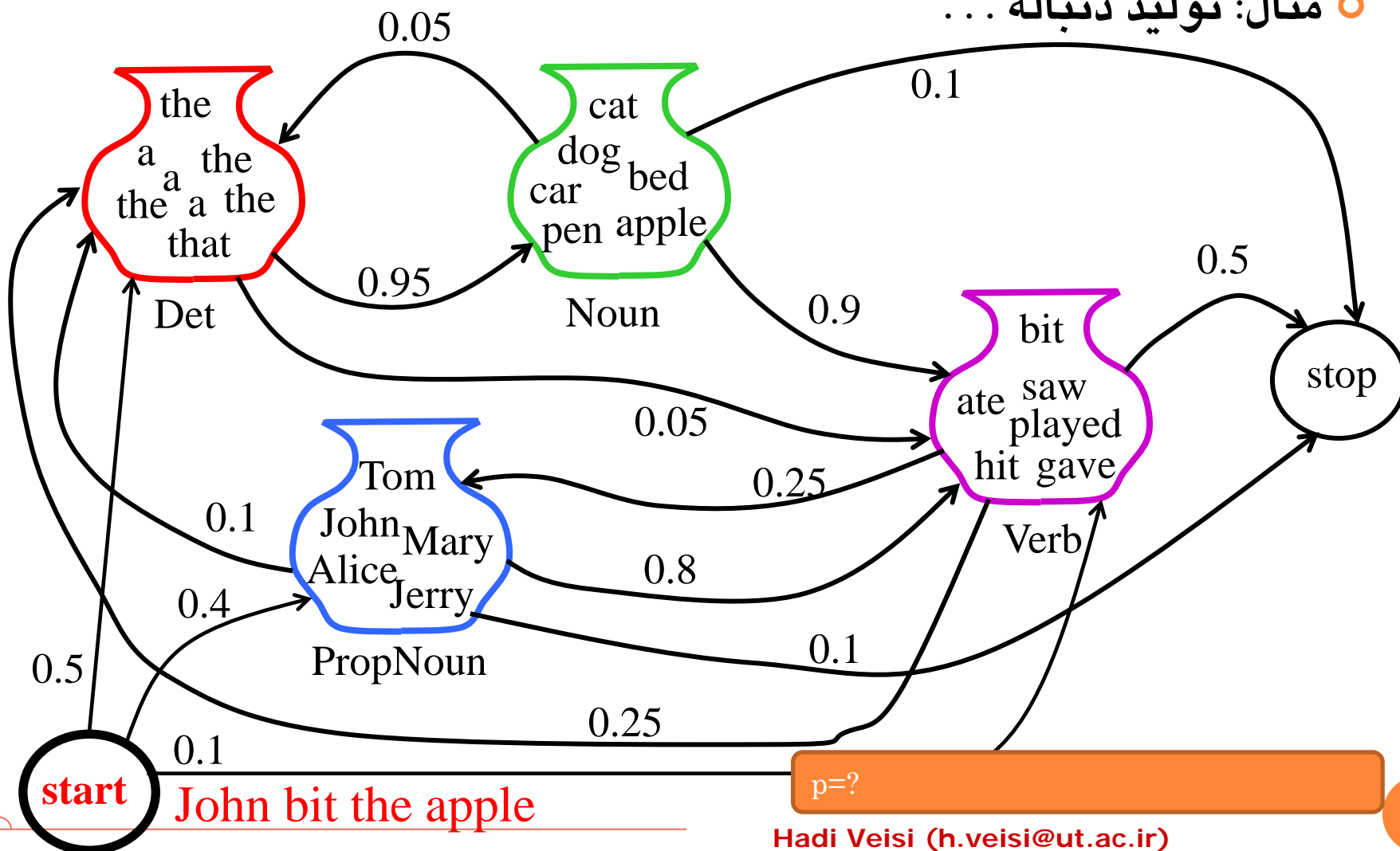
مثال





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

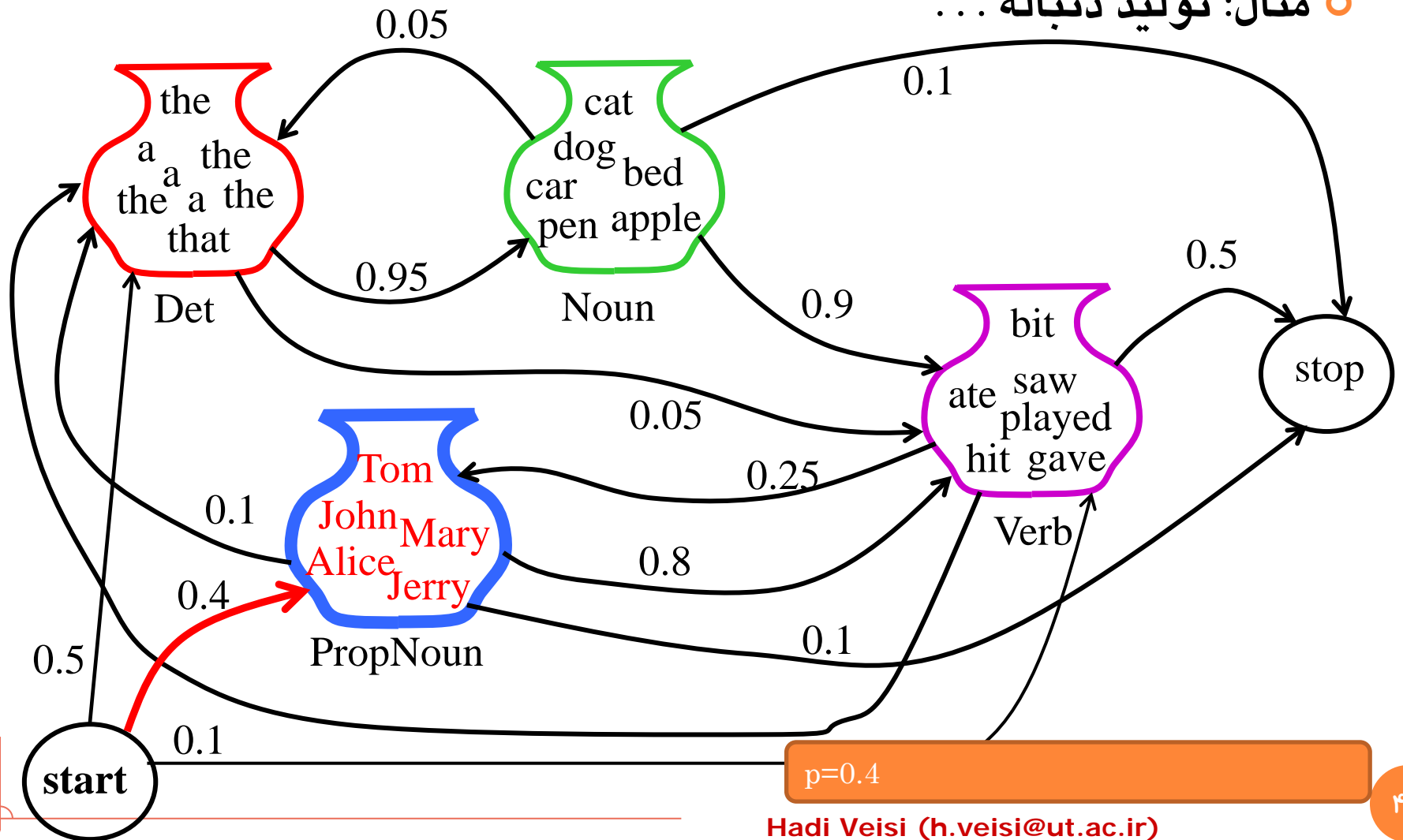
مثال: تولید دنباله ...





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

مثال: تولید دنباله ...

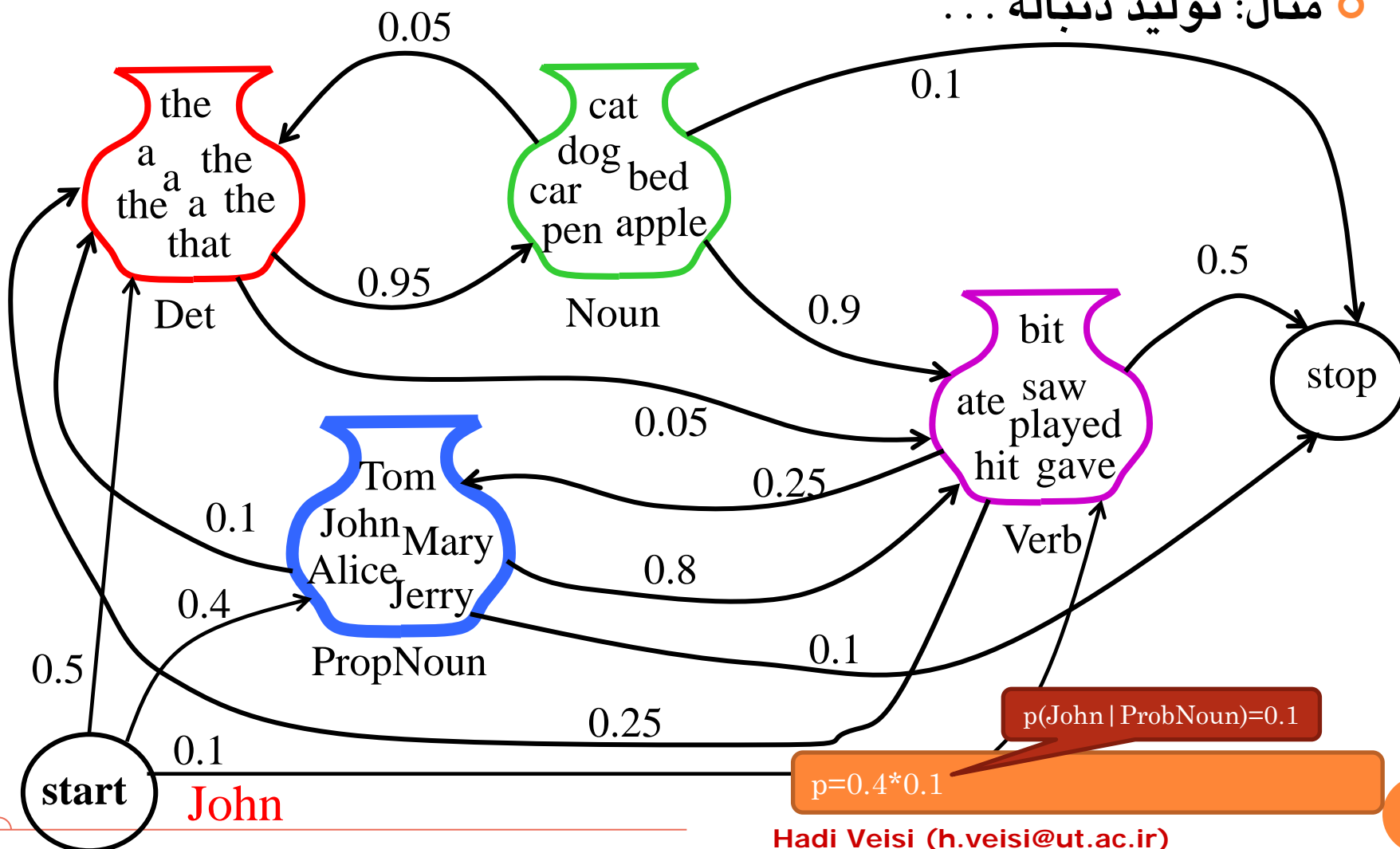


p=0.4



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

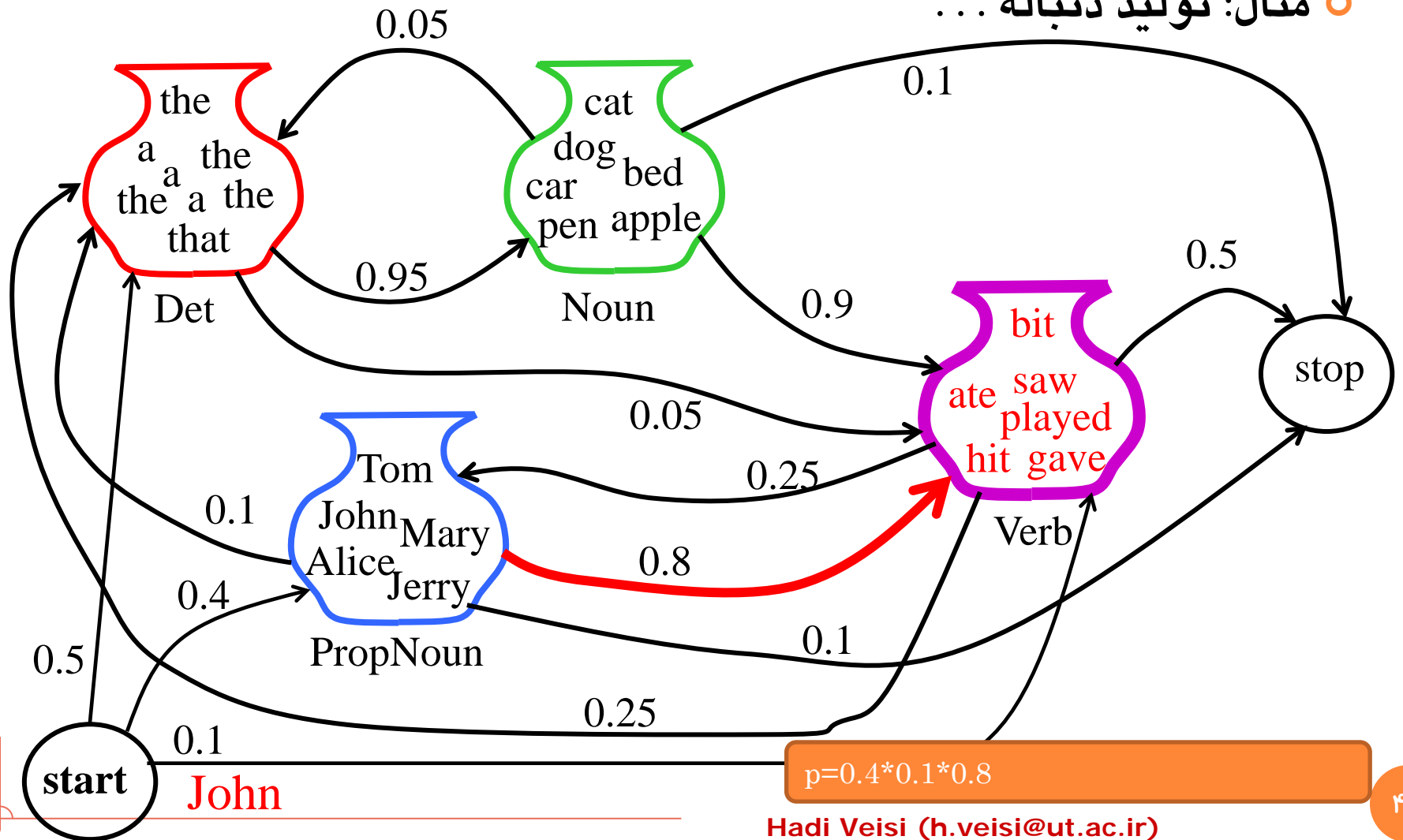
مثال: تولید دنباله ...





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

مثال: تولید دنباله ...

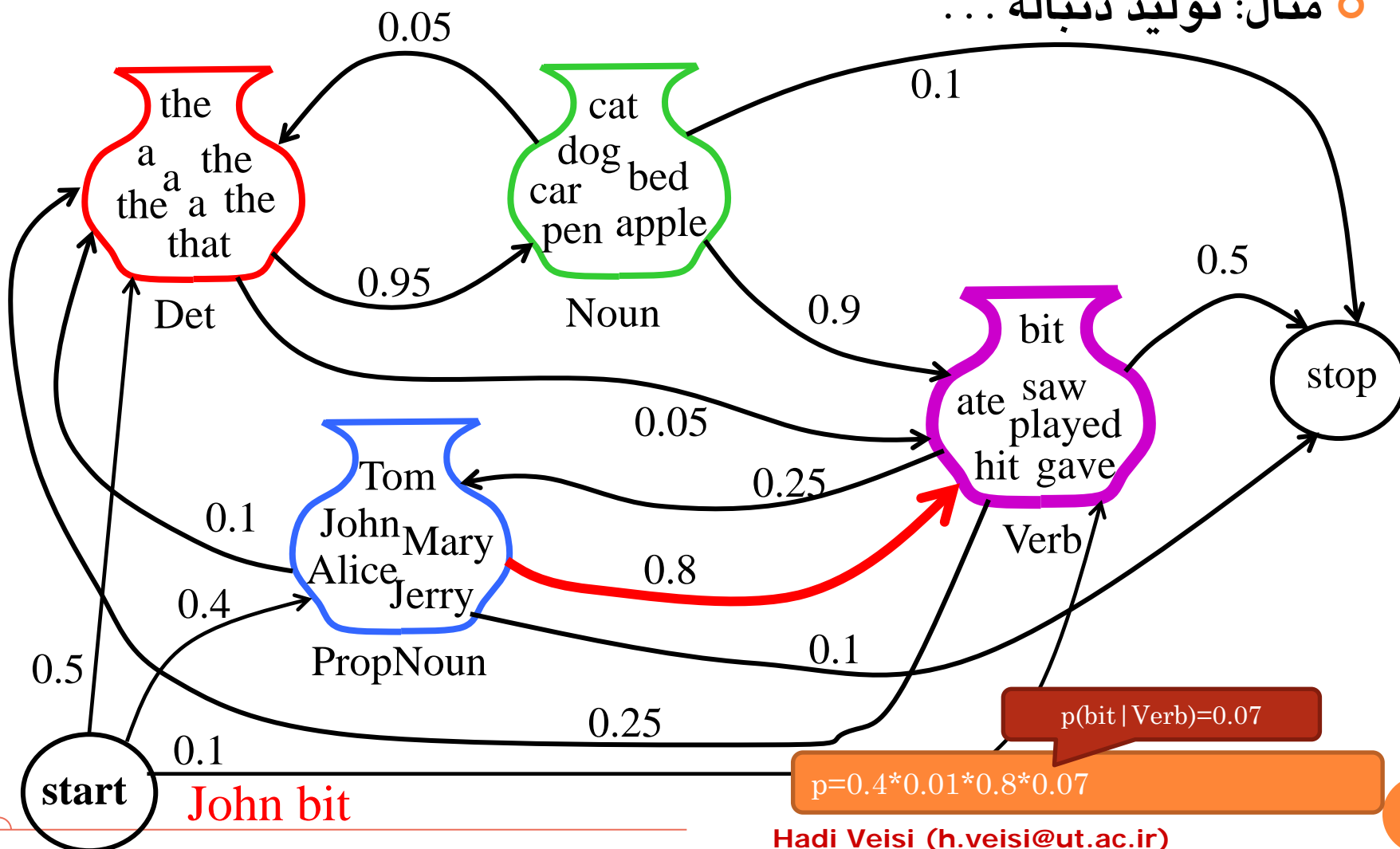


$p = 0.4 * 0.1 * 0.8$



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

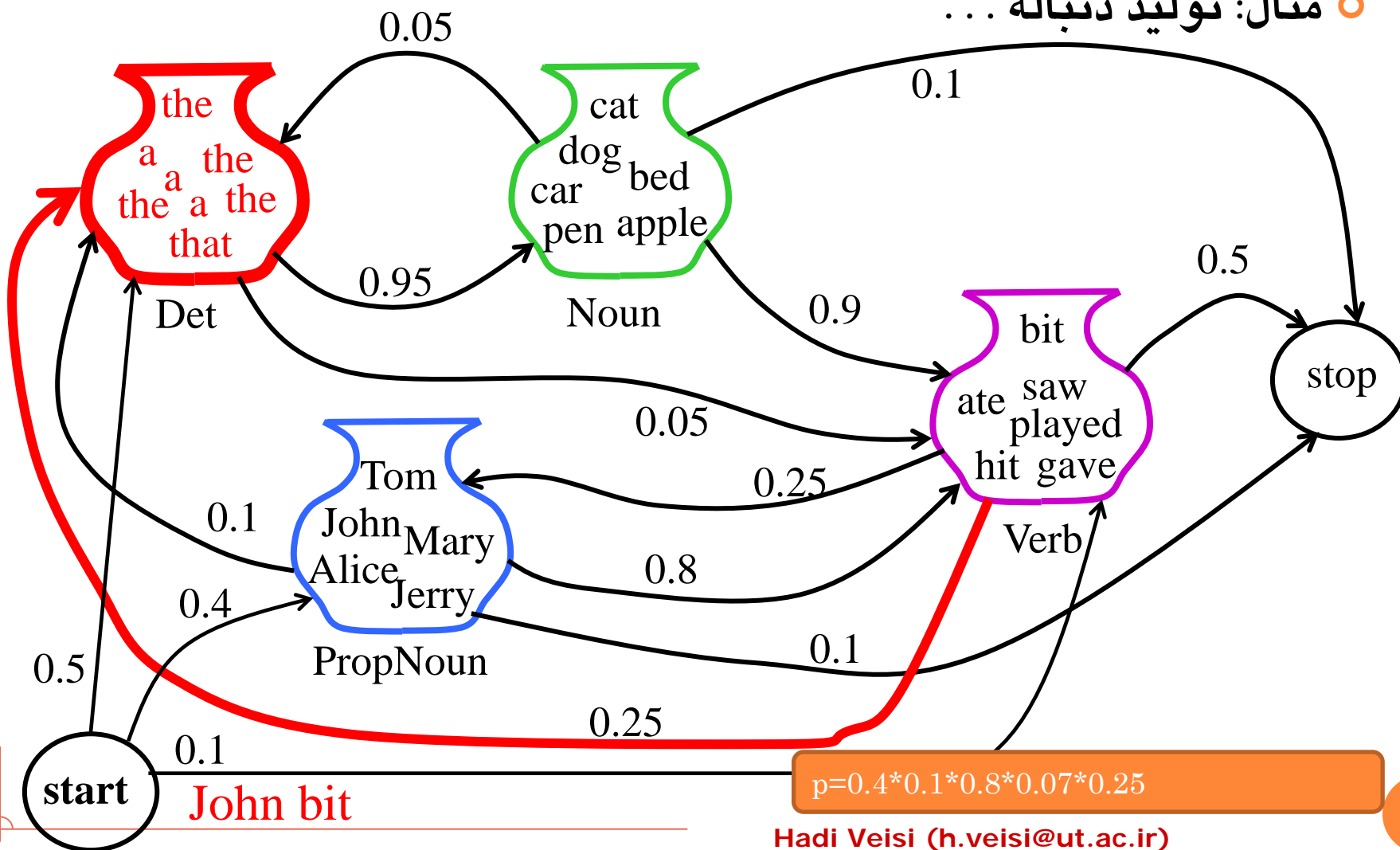
مثال: تولید دنباله ...





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

مثال: تولید دنباله ...

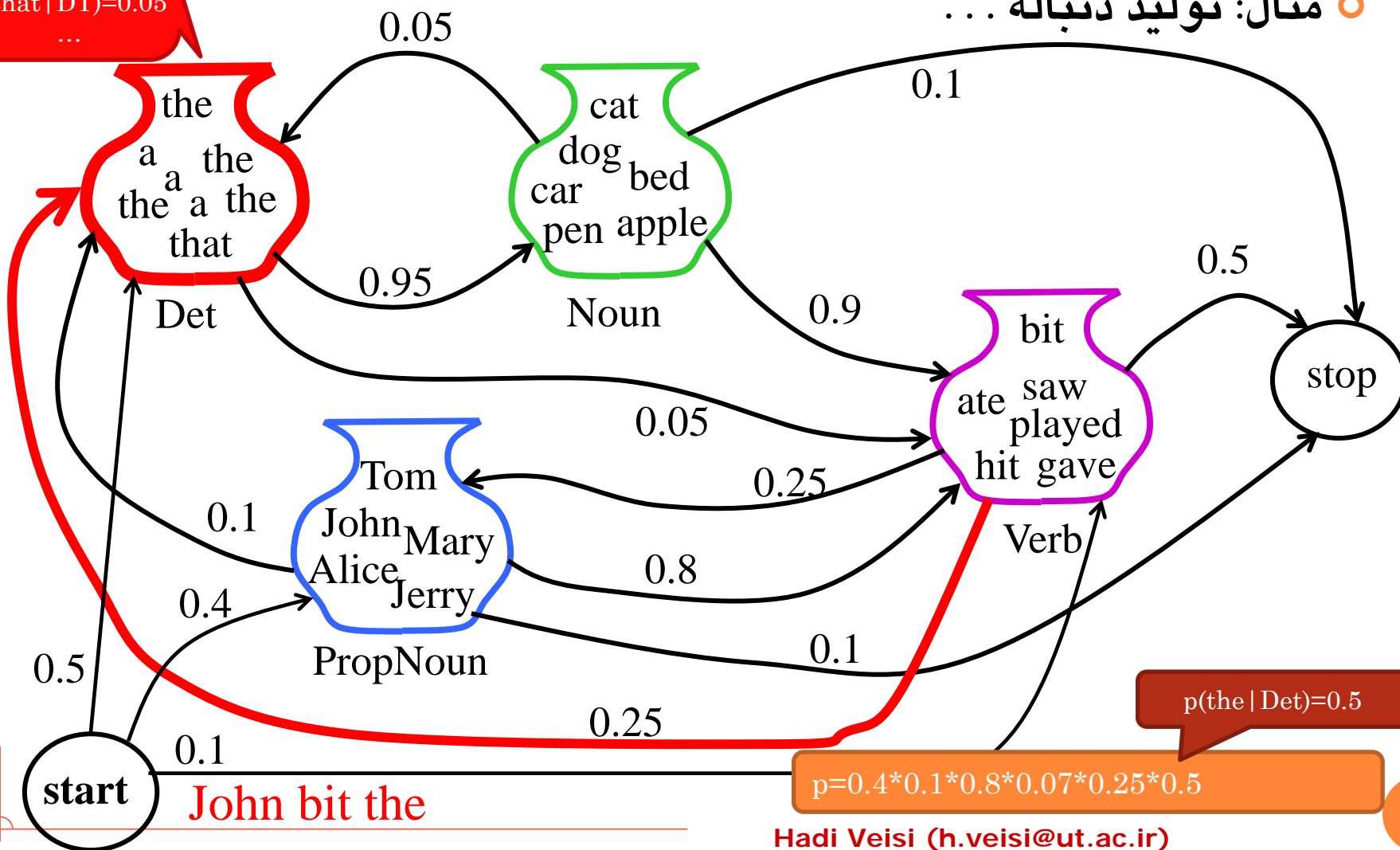




برچسب‌زنی اجزای کلام: روش آماری (HMM)

$p(\text{the} | \text{DT}) = 0.5$
 $p(\text{a} | \text{DT}) = 0.2$
 $p(\text{an} | \text{DT}) = 0.1$
 $p(\text{that} | \text{DT}) = 0.05$
 ...

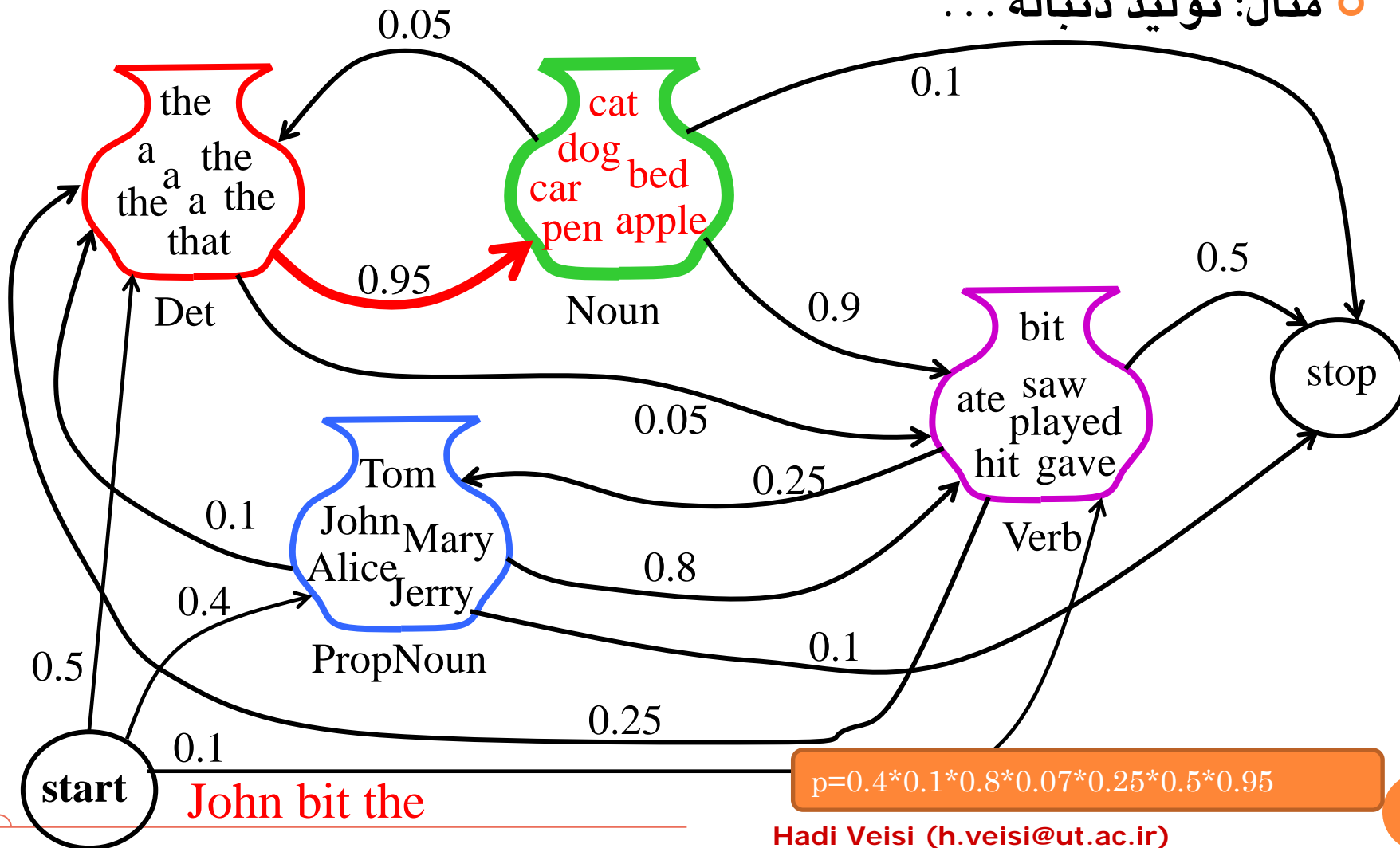
مثال: تولید دنباله ...





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

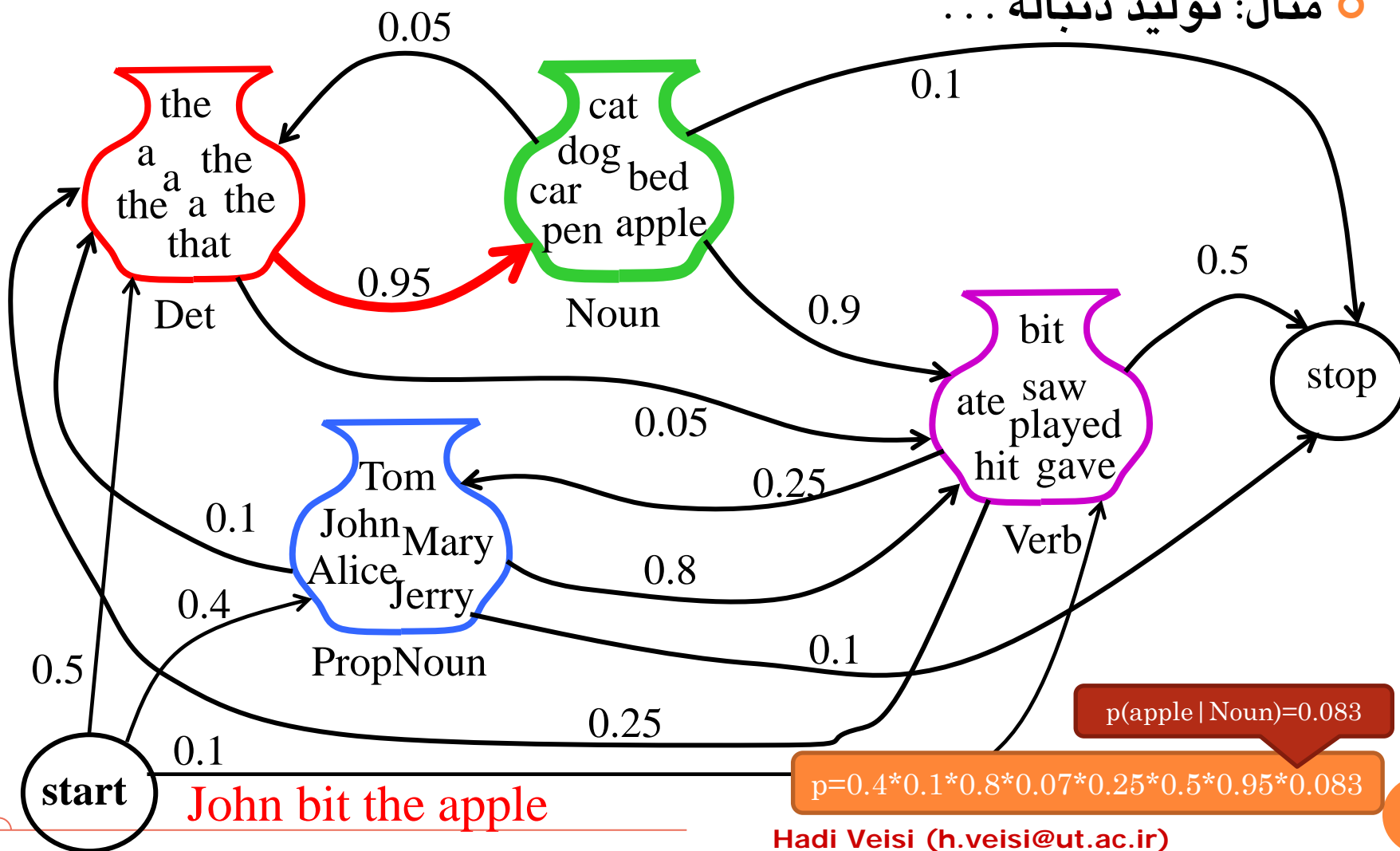
مثال: تولید دنباله ...





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

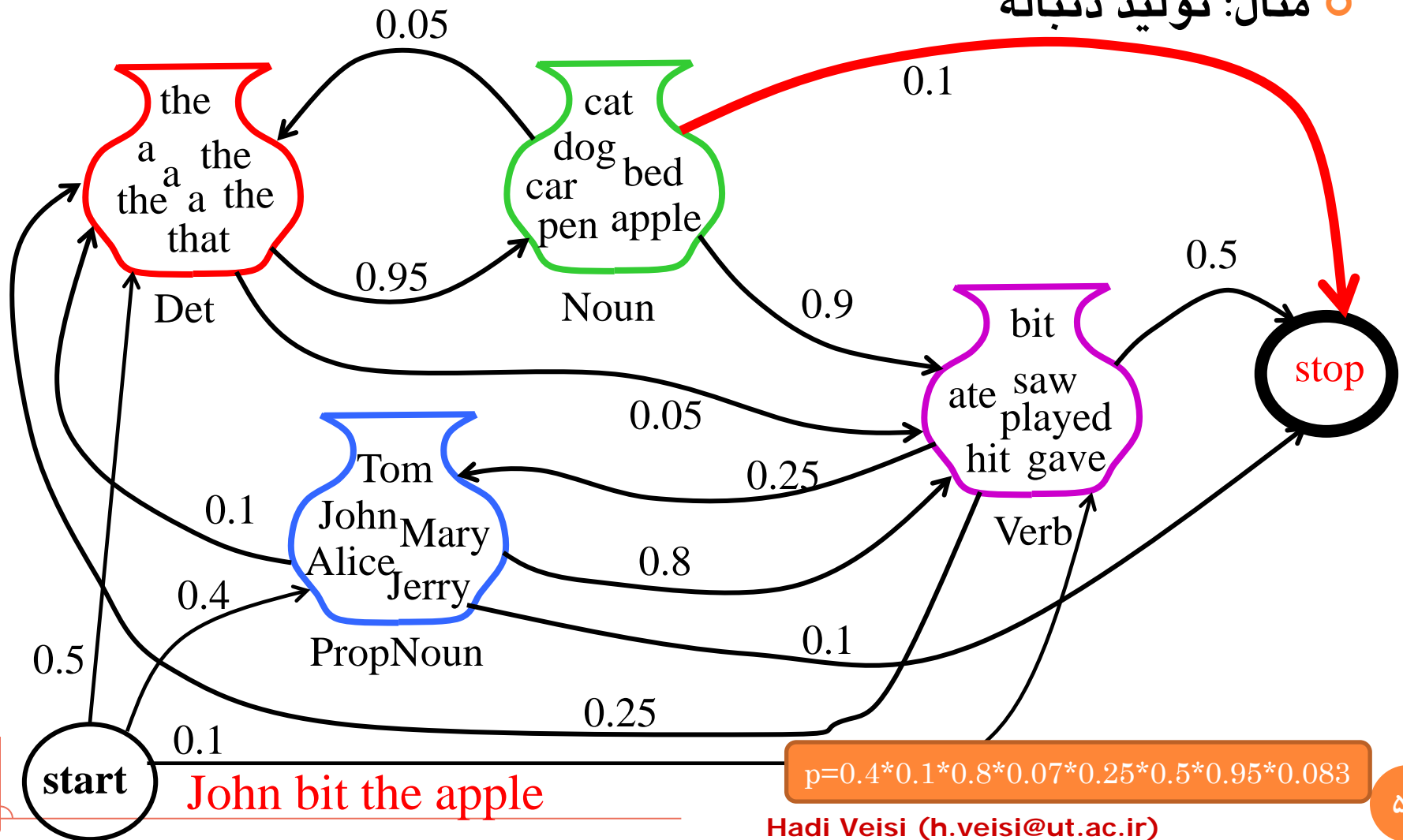
مثال: تولید دنباله ...





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

مثال: تولید دنباله



John bit the apple



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

○ برچسب‌زنی اجزای کلام با روش HMM

• فاز آموزش

- در نظر گرفتن یک واژگان با M کلمه و مجموعه برچسب‌های ممکن (N برچسب)
- در نظر گرفتن یک HMM با تعداد حالت‌های برابر با تعداد برچسب‌ها (N حالت)
- محاسبه احتمال‌های مدل با استفاده از یک پیکره متنی دارای برچسب اجزای کلام برای کلمات
 - احتمال اولیه حالت‌ها (N مقدار: هر حالت یک عدد)
 - احتمال انتقال از یک حالت (برچسب) به حالت دیگر (یک ماتریس $N*N$)
 - احتمال داشتن هر برچسب برای هر کلمه (یک ماتریس $M*N$)

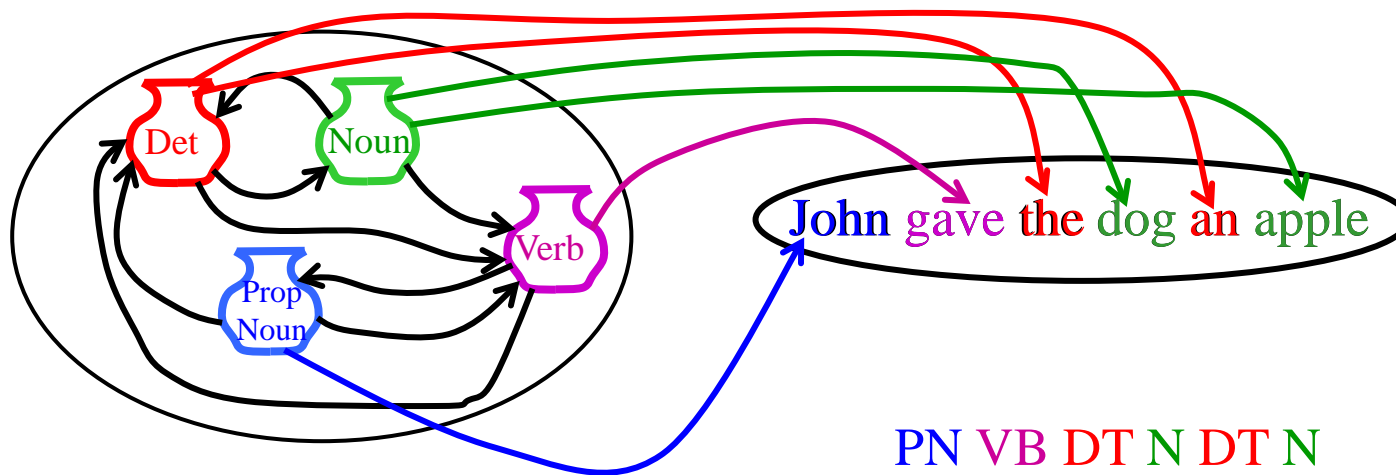
• فاز آزمون (استفاده)

- دریافت یک دنباله از کلمات
- یافتن بهترین برچسب‌های مرتبط = بهترین دنباله حالت در مدل HMM



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

○ برچسب‌زنی اجزای کلام با روش HMM





برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

○ یافتن دنباله حالت‌های (برچسب‌های) بهینه برای یک دنباله از کلمات

- با داشتن دنباله مشاهده (کلمات) $O=O_1O_2...O_T$ و مدل مخفی کارکوف λ ، چگونه می‌توان بهترین دنباله حالت‌های (برچسب‌ها) $Q=q_1q_2...q_T$ که متناسب با مشاهده است، را بدست آورد؟

- مساله دیکدینگ در HMM

- راه‌حل کامل: بررسی تمام دنباله‌های ممکن و انتخاب بهترین آنها

- بسیار زمان‌بر، از مرتبه $O(TN^T)$ که N تعداد حالت‌ها (برچسب‌ها) و T طول دنباله مشاهده‌ها (کلمات) است

○ محاسبه محتمل‌ترین دنباله از برچسب‌ها

- ساده‌ترین روش: در نظر گرفتن تمام دنباله‌های محتمل و محاسبه احتمال هر یک به روش بیان شده (Brute Force Search)

- با فرض داشتن N برچسب و T کلمه، حداکثر N^T دنباله از برچسب‌ها تولید می‌شود. ○ محاسبات بسیار زیاد

• روش‌های رایج

- مدل مخفی مارکوف (HMM: Hidden Markov Model)
- میدان تصادفی شرطی (CRF: Conditional Random Field)

- راه حل بهینه: الگوریتم ویتربی (Viterbi)

- یک روش برنامه نویسی پویا (Dynamic Programming)

- مشابه الگوریتم Minimum Edit Distance



برچسب‌زنی اجزای کلام: روش آماری (HMM) ...

الگوریتم ویتربی

function VITERBI(*observations* of len T , *state-graph* of len N) **returns** *best-path*

create a path probability matrix $viterbi[N+2, T]$

for each state s **from** 1 **to** N **do**

; initialization step

ماتریس احتمالها

$viterbi[s, 1] \leftarrow a_{0,s} * b_s(o_1)$

احتمال اولیه حالت s

ماتریس دنباله حالتها

$backpointer[s, 1] \leftarrow 0$

for each time step t **from** 2 **to** T **do**

; recursion step

for each state s **from** 1 **to** N **do**

$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s',s} * b_s(o_t)$

احتمال اینکه کلمه o_t دارای حالت s باشد

$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s',s}$

$viterbi[q_F, T] \leftarrow \max_{s=1}^N viterbi[s, T] * a_{s,q_F}$

; termination step

$backpointer[q_F, T] \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T] * a_{s,q_F}$

; termination step

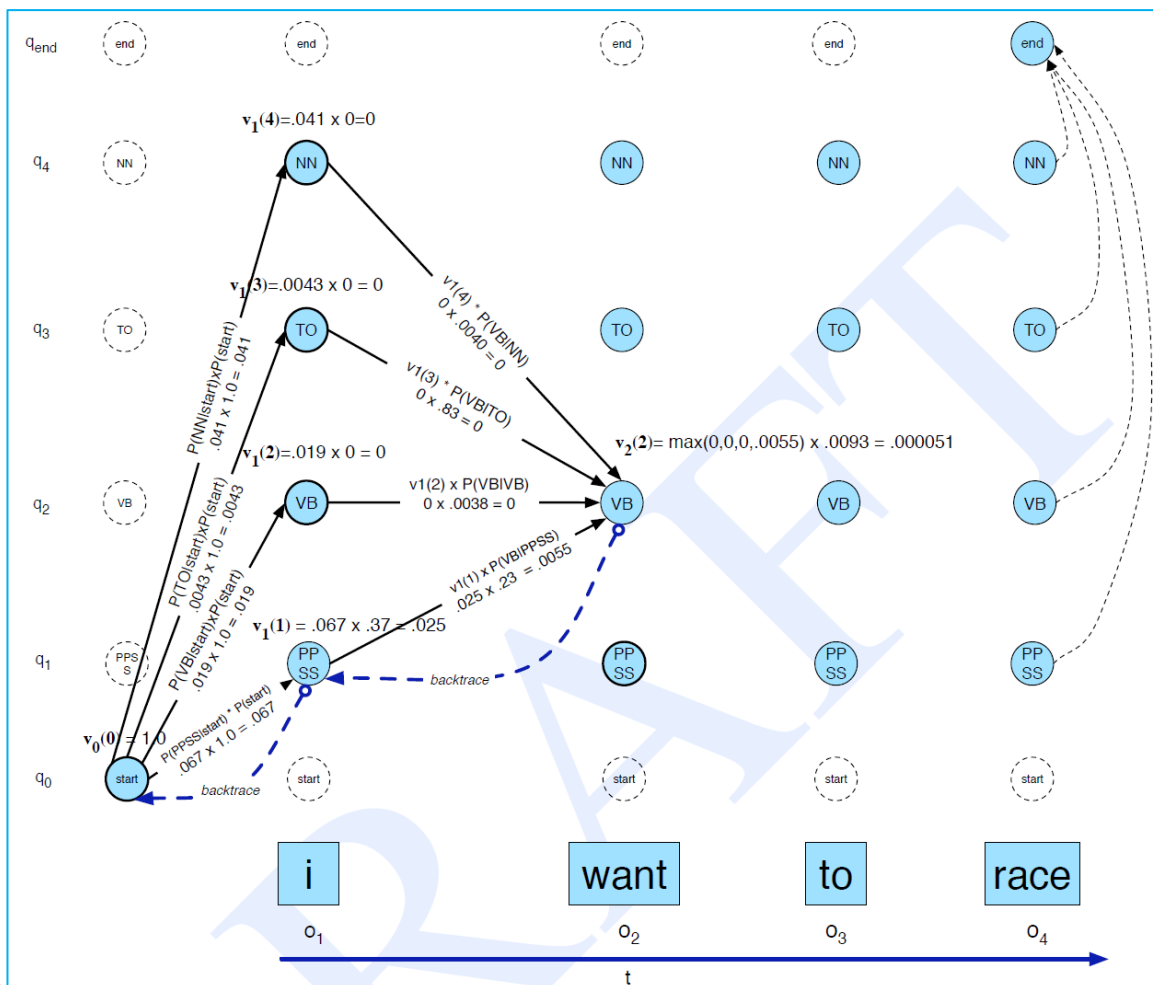
return the backtrace path by following backpointers to states back in time from $backpointer[q_F, T]$



برچسب‌زنی اجزای کلام: روش آماری (HMM)

الگوریتم ویتربی: مثال

• ۴ برچسب و ۴ کلمه



Transition probabilities: $P(t_i | t_{i-1})$

	VB	TO	NN	PPSS
start	0.019	0.0043	0.041	0.067
VB	0.0038	0.0345	0.047	0.070
TO	0.83	0	0.00047	0
NN	0.0040	0.016	0.087	0.0045
PPSS	0.23	0.00079	0.0012	0.00014

Observation likelihoods: $P(w_i | t_i)$

	i	want	to	race
VB	0	0.0093	0	0.00012
TO	0	0	0.99	0
NN	0	0.000054	0	0.00057
PPSS	0.37	0	0	0



برچسب‌زنی اجزای کلام: روش مبتنی بر تبدیل

روش مبتنی بر تبدیل (Transformation-Based Tagging)

- بر اساس ایده یادگیری مبتنی بر تبدیل (TBL: Transformation-Based Learning)
- ترکیب روش‌های مبتنی بر قاعده و روش آماری

- مبتنی بر قاعده: استفاده از قواعد برای شناسایی برچسب‌های نادرست
- آماری: یادگیری قوانین به صورت خودکار از روی داده با روش‌های آماری

الگوریتم

- انتساب محتمل‌ترین برچسب به کلمه (استفاده از واژگان)
- تغییر برچسب بر اساس قوانین
- *“if word-1 is a DT and word is a VRB then change the tag to NN”*
- تکرار الگوریتم (تا رسیدن به شرط توقف)

ساخت قواعد بر اساس قالب‌های (templates) مشخص

The preceding (following) word is tagged **z**.
 The word two before (after) is tagged **z**.
 One of the two preceding (following) words is tagged **z**.
 One of the three preceding (following) words is tagged **z**.
 The preceding word is tagged **z** and the following word is tagged **w**.
 The preceding (following) word is tagged **z** and the word two before (after) is tagged **w**.

Change tags			
#	From	To	Condition
1	NN	VB	Previous tag is TO
2	VBP	VB	One of the previous 3 tags is MD
3	NN	VB	One of the previous 2 tags is MD
4	VB	NN	One of the previous 2 tags is DT
5	VBD	VBN	One of the previous 3 tags is VBZ

Example
 to/TO race/NN → VB
 might/MD vanish/VBP → VB
 might/MD not reply/NN → VB



برچسب‌زنی اجزای کلام: ارزیابی

○ معیار

$$Acc = \frac{\text{تعداد کلمات با برچسب درست}}{\text{تعداد کل کلمات}} \times 100$$

- درصد کلماتی که درست برچسب زده شده است
- محاسبه کارایی روی مجموعه آزمون (Test Set)
 - برابر با ۱۰ تا ۲۰ درصد پیکره دارای برچسب

○ کارایی حدودی (برخی از روش‌ها)

- برچسب زدن بر اساس محتمل‌ترین برچسب
 - دقت کلی: ۹۰٪
 - دقت روی کلمات ناشناس: ۵۰٪
- روش HMM (با احتمال‌های Trigram)
 - دقت کلی: ۹۵٪
 - دقت روی کلمات ناشناس: ۵۵٪
- انسان (کران بالا)
 - حدود ۹۸٪



برچسب‌زنی اجزای کلام: مثال (فارسی) ...

برچسب‌زنی اجزای کلام فارسی

• پیکره: بیجن خان (تعداد برچسب‌ها = ۴۱)

• روش‌ها

○ شبکه عصبی مصنوعی (Artificial Neural Network)

○ پرسپترون چند لایه (MLP)

○ درخت تصمیم (Decision Tree)

○ مدل مخفی مارکوف (HMM) - مبتنی بر bigram

• اطلاعات (ویژگی‌ها) برای ANN و DT

○ کلمه فعلی

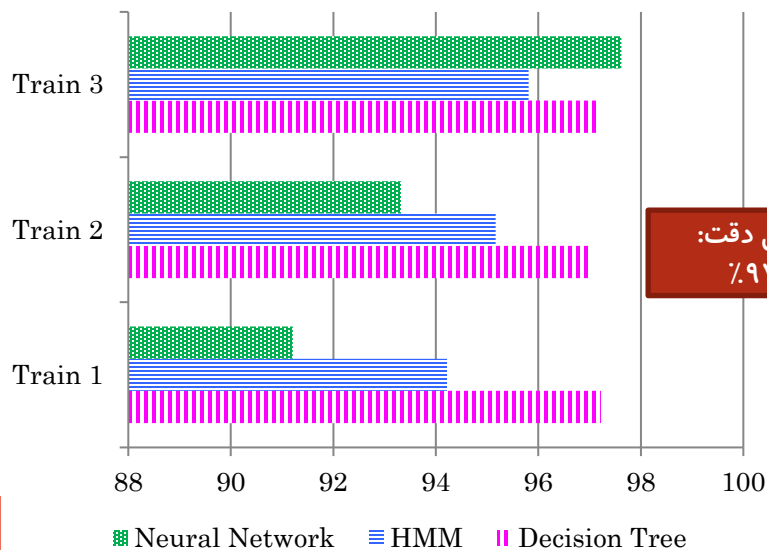
○ کلمه قبلی

○ برچسب کلمه قبلی

○ احتمال همه برچسب‌ها برای کلمه فعلی

○ احتمال همه برچسب‌ها برای کلمه برچسب قبلی

Name	# Samples
Train 1	200K
Train 2	1.0 M
Train 3	2.587 M
Test	10K





برچسب‌زنی اجزای کلام: مثال (فارسی)

○ برچسب‌زنی اجزای کلام فارسی: قدرت تعمیم شبکه عصبی و HMM

- پیکره ۱۰ میلیون کلمه ای بیجن خان با ۴۱ برچسب
- آموزش = ۹۲٪؛ اعتبارسنجی = ۴٪ و آزمون = ۴٪ (حدود ۰.۸٪ کلمات OOV)
- شبکه عصبی LSTM (یک طرفه و دوطرفه) و MLP
- بردار ویژگی Word Vector
- روش HMM دوتایی (Bigram)

MODEL	IV	OOV	Total
HMM	96.30	45.46	95.82
1-layer MLP	94.91	65.0	94.67
2-layer MLP	95.23	66.0	95.00
ULSTM	95.54	50.48	95.16
BLSTM	95.60	50.79	95.23

