

شماره مستند: ۱۹۰/۲۵۳۷/۱/۸



جمهوری اسلامی ایران
دبیرخانه شورای عالی اطلاع رسانی

تحلیل نیازمندی‌های سیستم تشخیص کلمات فینگلیش در لابه‌لای متون زبان

فارسی

نسخه ۱.۰

دانشگاه علم و صنعت ایران

فروردین ۸۸

فهرست

۳	۱مقدمه
۴	۲ساختار متون پینگلیش
۴	۲.۱نگاشت حروف فارسی و انگلیسی
۹	۲.۲تکرار حروف
۹	۲.۳استفاده از کلمات شکسته
۹	۲.۴استفاده از کلمات انگلیسی
۱۱	۲.۵استفاده از حروف و کاراکترهای ویژه در کلمات
۱۲	۳مبدل پینگلیش
۱۲	۳.۱کشف و یادگیری الگوهای تبدیل
۱۴	۳.۲استفاده از الگوهای یادگرفته شده
۱۴	۳.۳الگوریتم کلی تبدیل کلمات پینگلیش
۱۴	۳.۴دقت الگوریتم

۱ مقدمه

پیش از همه‌گیر شدن استفاده از چیدمان فارسی صفحه‌کلید توسط فارسی‌زبانان، رسم‌الخطی به نام پینگلیش بوجود آمده بود که در آن از کاراکترهای انگلیسی برای نوشتن کلمات فارسی استفاده می‌شد. استفاده از این رسم‌الخط در سال‌های اولیه‌ی ورود اینترنت به ایران رواج زیادی داشت، چرا که هنوز استاندارد مشخصی برای تبادل کاراکترها و متون فارسی که اولاً موردقبول اکثریت بوده و ثانیاً در همه‌ی سیستم‌های کامپیوتری قابل دسترس و استفاده باشد، وجود نداشت. امروزه استاندارد یونیکد این مشکل را تا حدود بسیاری مرتفع کرده است اما با این وجود هنوز هم استفاده از رسم‌الخط پینگلیش میان کاربران رایج است. مهم‌ترین حوزه‌ای که امروزه استفاده از پینگلیش در آن رواج دارد، دستگاه‌های قابل حمل ارتباطی مانند PDA و موبایل است. در اینترنت نیز هنوز برخی از کاربران بدلیل عدم آشنایی کامل با چیدمان فارسی صفحه‌کلید، متون خود را پینگلیش تایپ می‌کنند.

حجم قابل توجهی از متون برخی سایت‌ها نیز (مانند فروم‌ها) پینگلیش می‌باشد. برخلاف نوشتن پینگلیش که بنظر ساده می‌آید، خواندن متون پینگلیش (بویژه متونی که بیش از چند پاراگراف باشند) امری خسته‌کننده و وقت‌گیر است. با در دست داشتن ابزاری برای تبدیل متون پینگلیش به فارسی، می‌توان این نقیصه را برطرف نمود.

۲ ساختار متون پینگلیش

بررسی متون پینگلیش نشان می‌دهد که قواعد و قوانین مشخصی بر آن حاکم نمی‌باشد. دلیل این امر نیز تاحدودی واضح است: پینگلیش توسط عامه‌ی کاربران فارسی‌زبان ایجاد شده و مرکز یا مبدا مشخصی برای تدوین قواعد پینگلیش نویسی وجود نداشته است، لذا تنوع زیادی در قواعد و الگوهای پینگلیش نویسی وجود دارد.

بعنوان مثال، کلمه‌ی «اعتماد» در پینگلیش ممکن است به شکل‌های زیر نوشته شود:

etemad, etemaad , e'temad, e'temaad, eatemad, ...

همانگونه که مشاهده می‌شود، گاهی برای ۱ کلمه‌ی فارسی ممکن است بیش از چندین نوع نحوه‌ی نگارش پینگلیش وجود داشته باشد. یکی از دلایل این چندگانگی وجود داشتن چندین کاراکتر انگلیسی برای برخی از حروف فارسی است.

یکی دیگر از ویژگی‌های متون پینگلیش این است که شامل کلمات و اصطلاحات غیررسمی بوده و کلمات در این متون عموماً بصورت شکسته بیان می‌شوند. دخیل کردن حالات و احساسات در بیان کلمات نیز امری متداول در پینگلیش بشمار می‌رود.

۲.۱ نگاشت حروف فارسی و انگلیسی

در جدول ۱، حروف فارسی و معادل‌های رایج آن‌ها در پینگلیش به تفکیک آمده‌اند. در برخی خانه‌های ستون دوم جدول (معادل‌های رایج فینگلیش) کاراکتر \$ مشاهده می‌شود؛ این کاراکتر برای نمایش کاراکتر «خالی» مورد استفاده قرار گرفته است. بعنوان مثال برای حرف ع، به نگاشت زیر در کلمه‌ی «اعتماد» دقت کنید:

اعتماد	etemad
ا	e
ع	\$
ت	t
- -	e
م	m
ا	a
د	d

جدول ۱ - حروف فارسی و معادل‌های آن در پینگلیش

حرف فارسی	معادل‌های رایج فینگلیش	متداولترین معادل	مثال
حروف بیصدا			
ب	b	-	baché
پ	p	-	pedar
ت	t	-	tatilat
ث	s, c	-	mosbat
ج	j, g	j	jahat
چ	ch	-	cheghadr
ح	h	-	mohit
خ	x, kh	kh	khanevadeh
د	d	-	dar
ذ	z	-	begzarim
ر	r	-	rah
ز	z	-	zir
ژ	zh, j	j	mojdeh
س	s, c	s	saket
ش	sh	-	shabih
ص	s, c	s	saboor
ض	z	-	zaroori
ط	t	-	tathir
ظ	z	-	zaher
ع	a, a', ee, ', \$		moeen, mo'een
غ	gh, q	-	ghalat, qalat
ف	f, ph	f	farda
ق	gh, q	-	ghader
ک	k, c	k	kah
گ	g	-	gom
ل	l	-	lazem
م	m	-	mamnoon
ن	n	-	naan
ه	h, \$	h	hadyeh, had
حروف باصدا			
آ	a, aa, \$	a	aab
ا	a, aa, e, \$	a	
ی	i, y, ei, ie, ee, e, iy, ey,	y, i	yeki, ieki, yeky

		ye, yi	
mohit	o	o, \$	ُ (ضمه)
zood, va'de, wared	o,oo,v	o, oo, ou, uo, v, u, w	و
mehman	e	e, \$	ِ (کسره)
mashroot	a	a, \$	َ (فتحه)
			حروف و کاراکترهای غیررسمی
merC, mer3o	-	C, 3o	سی
madre3	-	۳	سه
sa@	-	@	عت
	-	@	ت
raft	-	T	تی
charB	-	B	بی
boD hala	-	D	دی
Zafati (ضیافتی)	-	Z	زی
	-	U (در حالتیکه بعنوان یک کلمه بکار رود)	شما
	-	I (در حالتیکه بعنوان یک کلمه بکار رود)	من
Smaeel	-	S	اس
Geegar	-	G	جی
	-	X	
Nhedam (انهدام)	-	N	ان
Mkanat (امکانات)	-	M	ام

همانگونه که در جدول ۱ نیز مشاهده می‌شود، گوناگونی و تنوع بسیاری در نحوه‌ی تبدیل کلمات فارسی به معادل پینگلیش آن‌ها وجود دارد. در جدول فوق تلاش شده است تا پرکاربردترین این قواعد آورده شوند. با استفاده از جدول ۱، می‌توان جدول دیگری برای نگاشت از حروف انگلیسی به حروف فارسی ایجاد کرد:

جدول ۲- نگاشت حروف انگلیسی به حروف فارسی

متداولترین معادل	معادل محتمل	حرف انگلیسی
		حروف با یک معادل
-	ب	b
-	د	d
-	ر	r
-	ف	f

-	ل	l
-	م	m
-	ن	n
-	و	v
-	و	w
	ی	y
	ع	'
		حروف با چندین معادل
ه	ه	h
	ح	
خ	خ، کس،	x
ت	ت	t
	ط	
ت	تمامی حالات t	T
	+	
	تی	
س	ث	s
	س	
	ص	
	ش = sh	
س	تمامی حالات s	S
	+	
	اس	
س	ث	c
	س	
	ص	
	ک	
	چ = ch	
س	تمامی حالات c	C
	+	
	سی	
پ	ف = ph	p
	پ = p	
پ	تمامی حالات p	P
	+	
	پی	
ژ	ژ	j
	ج	
ج	ج	g

		gh = ق gh = غ گ	
ز		ز ذ ض ظ zh = ژ	z
		(فتحه) ا آ أ عا an - أ ع ی ه وا aa = آ	a
(کسره) ی		(کسره) ی ا eh = ه ee = ع ei = ی ee = ی ey = ی اع	e
ی		ی ای ie = ی iy = ی	i
ک		ک kh = خ	k
و		(ضمه) و او oo = و ou = و	o

و	و او و = uo شما (در حالتیکه یک کلمه‌ی یک حرفی مشاهده کنیم)	u
ی	ق غ	q

همانگونه که در جدول ۲ براحتی قابل مشاهده است، فقط تعداد بسیار کمی از حروف انگلیسی، دارای یک معادل در فارسی هستند و بسیاری از حروف انگلیسی، دارای ۲ یا چند معادل در فارسی هستند. همین امر مشکلات بسیاری را در توسعه‌ی یک مبدل فینگلیش ایجاد می‌کند.

۲.۲ تکرار حروف

در متون فینگلیش، تکرار حروف معمولاً به ۲ دلیل رخ می‌دهد:

۱. برای بیان حرفی مانند «و»، بعنوان مثال: **mashroot**

۲. برای بیان احساسات و ابراز هیجان: **mercccccc**

علاوه بر ۲ علت فوق؛ ممکن است برای بیان «تشدید» نیز از تکرار حروف استفاده شود، اما بررسی متون فینگلیش نشان می‌دهد که کاربران تقریباً در اکثر موارد، تشدید را در نحوه‌ی نگارش کلمات دخالت نداده‌اند. بعنوان مثال املا‌ی کلمه‌ی مفرح بصورت «mofarah» رایج‌تر از «mofarrah» می‌باشد.

۲.۳ استفاده از کلمات شکسته

متون فینگلیش عمدتاً بصورت محاوره‌ای نوشته شده‌اند و کلمات در آن‌ها بصورت شکسته بکار رفته‌اند. بعنوان مثال، کلمه‌ی «خانه» معمولاً بصورت **khoone** نوشته می‌شود و نه **khane**.

۲.۴ استفاده از کلمات انگلیسی

الگوی دیگری که در متون فینگلیش مشاهده می‌شود، استفاده از برخی کلمات انگلیسی است. این کلمات عمدتاً در دسته‌ی کلمات فنی قرار می‌گیرند و رایج‌ترین حوزه‌ی آن‌ها، حوزه‌ی کلمات تخصصی کامپیوتر است. برخی از این کلمات در جدول ۳ گردآوری شده‌اند. (لازم بذکر است که این جدول شامل تمامی این کلمات نمی‌باشد.)

جدول ۳ - برخی از واژه‌های انگلیسی رایج در متون بینگلیش

معادل	واژه‌ی انگلیسی
کپی	copy
سی‌دی	cd
سایت	site
صفحه	page
آی‌پی	IP
تایپ	type
چک	check
آزاد، خالی	free
دایال‌آپ،	dialup
کامپیوتر	comp, computer
عنوان	title
سی‌پی‌یو، پردازنده	CPU
حافظه	memory
پروژه	project
ایمیل	mail, email
الصاق کردن	paste
ابزار	tool
سلام	hi
خداحافظ	bye
خوب	good
باشه، بسیارخوب،	ok
آدرس	address
تلویزیون	tv

۲.۵ استفاده از حروف و کاراکترهای ویژه در کلمات

از این الگو در برخی موارد برای کوچکتر کردن طول کلمات و خلاصه‌سازی استفاده می‌شود:

جدول ۴

مثال	معادل	حرف – کاراکتر انگلیسی
zood -> zood	o	0
1doone -> yedoone	ye, yek	1
beznin -> betooni, bedoonin 2a -> doa	too, doo, do	2
3tar -> setar	se	3
4kerim -> chakerim	chahar, char, cha, for	4
w8-> weit (صبرکن)	eit	8
sa@ -> saat	at	@
mer30 -> merci	si, ci	30
raft -> rafti	ti	T
merC -> merci	si, ci	C

۳ مبدل پینگلیش

در بررسی روش‌های متداول تبدیل خودکار متون پینگلیش به فارسی نکته‌ای که جلب توجه می‌کند، عدم توانایی این روش‌ها در پوشش دادن همه‌ی الگوهای ذکر شده در بخش قبل است. برخی از روش‌ها نیز کلمات پینگلیش را در صورتی بدرستی تبدیل می‌کنند که در نوشتن آن کلمه، قوانین خاصی رعایت شده باشد (بطور مثال حرف «ع» بصورت «'» نوشته شده باشد، یا حرف «ا» بصورت «aa» در کلمه‌ی e'temaad).

بنابراین در این پروژه روش جدیدی برای تبدیل متون پینگلیش ایجاد و پیاده‌سازی شد. ایده‌ی اصلی این روش، یادگیری و کشف الگوهای تبدیل توسط خود نرم‌افزار است. این روش ۲ مرحله دارد:

۱. کشف و یادگیری الگوهای تبدیل، از طریق تعدادی کلمات نمونه: در این مرحله، تعدادی کلمه‌ی پینگلیش و معادل آنها بعنوان ورودی به نرم‌افزار داده می‌شود. نرم‌افزار طبق مراحل، معادل‌های حروف انگلیسی را یاد گرفته و آنها را ذخیره می‌کند.

۲. استفاده از الگوهای یاد گرفته شده در مرحله‌ی قبل: در این مرحله، کلمات پینگلیش بعنوان ورودی به نرم‌افزار داده می‌شود و نرم‌افزار با استفاده از الگوهای کشف شده در مرحله‌ی قبل، معادل‌های محتمل برای آن کلمات را ارائه می‌کند.

یکی از مزیت‌های این روش در این است که می‌توان الگوهای جدید را نیز در هر لحظه به نرم‌افزار آموخت و نیازی به بازنویسی نرم‌افزار نیست.

در ادامه، توضیح هر یک از مراحل فوق می‌آید.

۳.۱ کشف و یادگیری الگوهای تبدیل

در این مرحله، تعدادی کلمه‌ی پینگلیش و نحوه‌ی نگاشت آنها به فارسی به نرم‌افزار داده می‌شود. بعنوان مثال ورودی‌های زیر را در نظر بگیرید:

n --> ن
o --> و

a --> ا
m --> م
a --> ا
d --> د
e --> _

g	-->	گ
i	-->	ی

c	-->	چ
h	-->	
e	-->	ِ
s	-->	ش
h	-->	
m	-->	م
e	-->	ه

با در دست داشتن نگاشت‌هایی مانند مثال‌های فوق، نرم افزار سعی می‌کند الگوهای نگاشت را کشف کند. بدین ترتیب که برای تک‌تک حرف در کلمه‌ی پینگلیش، نگاشت‌های زیر را یافته و ذخیره می‌کند:

➤ نگاشت ۲-۲ (۲ حرف قبلی + حرف موردنظر + ۲ حرف بعدی) - «حرف فارسی معادل

مثال: حرف e در cheshme

کسره -> chesh

➤ نگاشت ۱-۲ (۱ حرف قبلی + حرف موردنظر + ۲ حرف بعدی) - «حرف فارسی معادل)

➤ نگاشت ۱-۱ (۱ حرف قبلی + حرف موردنظر + ۱ حرف بعدی) - «حرف فارسی معادل)

➤ نگاشت ۰-۱ (حرف موردنظر + ۱ حرف بعدی) - «حرف فارسی معادل)

مثال: حرف c در cheshme

چ -> ch

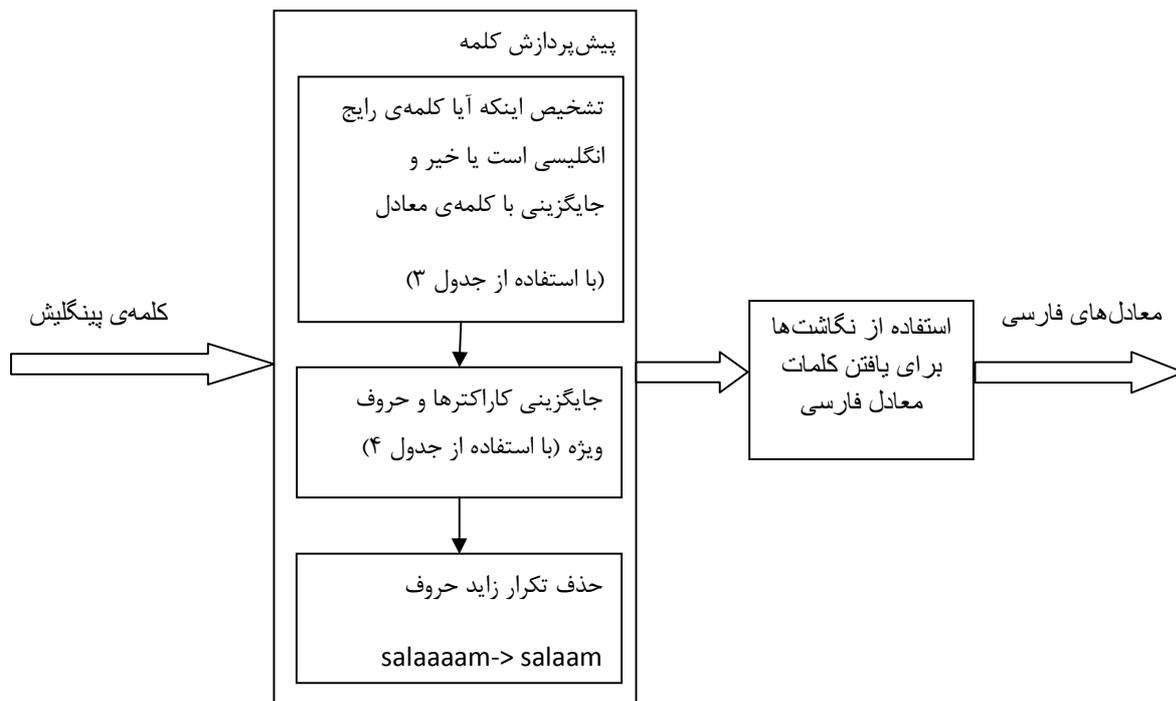
➤ نگاشت ۰-۰ (حرف مورد نظر) - «حرف فارسی معادل)

۳.۲ استفاده از الگوهای یادگرفته شده

در این مرحله، نرم‌افزار از نگاشت‌های ذخیره شده در مرحله‌ی قبل استفاده می‌کند تا معادل فارسی کلمات جدید را تولید کند. برای این منظور، الگوریتم زیر را اجرا می‌کند:

۱. برای تمامی حروف کلمه‌ی فینگلیش، نگاشت‌های ۲-۲، ۱-۲، ۱-۱، ۰-۱ و ۰-۰ را بترتیب جست‌جو کن. اگر نگاشت مشابهی یافت شد، حرف فارسی معادل آن را استخراج کن.
- اگر نگاشت مشابهی یافت شد، جست‌جو را خاتمه بده و به مرحله‌ی بعد برو.
۲. با استفاده از حروف یافت شده در مرحله‌ی قبل، کلمه (یا کلمات) فارسی معادل را تشکیل بده

۳.۳ الگوریتم کلی تبدیل کلمات پینگلیش



۳.۴ دقت الگوریتم

دقت الگوریتم ارائه شده در این مستند، ارتباط مستقیمی با مجموعه داده‌های اولیه‌ی آن دارد. هرچه تعداد کلمات پینگلیش و تنوع نگارشی آن‌ها در مجموعه داده‌های اولیه بالاتر باشد، دقت الگوریتم در تبدیل کلمات جدید بیشتر خواهد بود.

بهمین علت، ابزار جداگانه‌ای برای تولید مجموعه داده‌های اولیه توسعه یافته‌است. این ابزار یک متن پینگلیش را بعنوان ورودی دریافت کرده و به‌ازای تک‌تک کلمات موجود در آن، همه‌ی معادل‌های فارسی آن کلمه را تولید می‌کند و به کاربر نمایش می‌دهد. (ایجاد تمامی معادل‌های فارسی یک کلمه‌ی پینگلیش، با استفاده از جدول ۲ صورت می‌گیرد. بدیهی است که تعداد این معادل‌ها در برخی حالات به چندصد عدد می‌رسد). کاربر از میان کلمات نمایش داده‌شده، کلمات صحیح را تایید می‌کند و نرم‌افزار بصورت خودکار نگاشت‌های آن کلمه را یافته و ذخیره می‌کند. توسط این ابزار براحتی می‌توان یک مجموعه داده‌ی غنی از کلمات و نگاشت‌های آن‌ها ایجاد کرده و دقت نرم‌افزار را افزایش داد.

البته رسیدن به دقت کامل نیز ممکن نیست، بعنوان مثال می‌توان از کلماتی نام برد که بیش از ۲ معادل درست فارسی دارند و تعیین معادل صحیح آن‌ها، جز با پردازش معنایی متن امکان پذیر نیست. *madar* و *dar* از جمله‌ی این کلمات هستند که اولی ممکن است «مادر» یا «مدار» بوده و دومی نیز ممکن است «دار» یا «در» باشد.