

# ارزیابی یادگیری ساختواژی زبان فارسی با استفاده از یک الگوریتم تخمین بیزی

<b>مزدک انوشه</b> استادیار دانشکده ادبیات و علوم انسانی دانشگاه تهران mazdakanushe@ut.ac.ir	<b>مصطفی صالحی</b> استادیار دانشکده علوم و فنون نوین دانشگاه تهران mostafa_salehi@ut.ac.ir	<b>الهام سادات سبط</b> کارشناس ارشد زبان‌شناسی رایانشی دانشگاه تهران e.sebt@ut.ac.ir
---	--	--

## چکیده

پردازش ساختواژی با استفاده از الگوریتم‌های داده‌محور، رویکردی است که طیف وسیعی از پژوهش‌ها را به خود اختصاص داده است و می‌تواند در پردازش متون تخصصی، کهن و غیررسمی زبان فارسی مانند شبکه‌های اجتماعی، کاربرد داشته باشد. علی‌رغم گستردگی این حوزه، پژوهش‌ها در زمینه ساختواژه زبان فارسی بیشتر به روش‌های قاعده‌مند اختصاص دارد و کارایی این الگوریتم‌ها تاکنون بر روی داده‌های زبان، ارزیابی نشده است. در این پژوهش یادگیری ساختواژی بی‌نظارت و نیمه‌نظارتی با استفاده از تخمین بیزی بررسی شده، کارایی آن بر روی داده‌هایی از مجموعه داده‌های استاندارد وبلاگ‌های ایران، پیکره همشهری و پیکره متنی زبان فارسی ارزیابی شد. به دلیل دردسترس نبودن داده استاندارد آموزش، با شناسایی و اصلاح پارامترهای تأثیرگذار مانند نویز و طول واژه‌ها در داده و کنترل نقش نیم‌فاصله در فرآیند آموزش و تقطیع، نتایج در هر دو معیار صحت و فراخوانی بهبود یافته و معیار F به نتایجی که پژوهش‌ها برای زبان انگلیسی ثبت کرده‌اند، نزدیک شده است.

**کلیدواژه‌ها:** پردازش ساختواژی، یادگیری ماشین، تخمین بیزی، الگوریتم داده محور

## ۱. مقدمه

پردازش ساختواژی به عنوان پردازش پایه در طیف وسیعی از ابزارهای پردازش زبان طبیعی از جمله سیستم‌های ترجمه ماشینی، پرسش و پاسخ، بازیابی اطلاعات و غلطیاب املائی کاربرد دارد. اهمیت بهره‌گیری از این پردازش در کوچک کردن حجم واژگان، پوشش دادن ساخت‌های گوناگون واژه‌ها و امکان پردازش کلمه‌های خارج از واژگان است.

پردازشی که پردازشگر ساختواژی بر روی واژه انجام می‌دهد، می‌تواند بسته به نوع زبان و کاربرد آن در سایر ابزارهای پردازش زبان طبیعی متفاوت باشد و غالباً به‌منظور تحقق یکی از اهداف زیر به کار می‌رود (ویرپیوا<sup>۱</sup> و همکاران، ۲۰۱۱):

- تقطیع ساختواژی: هدف تقطیع ظاهر واژه و استخراج تک‌واژه‌ها به‌صورت فهرستی از زیررشته‌های سازنده واژه یا صورت ظاهری تک‌واژه‌ها است. هرچند در این رویکرد تک‌واژگونی و فرآیندهای واژه-واجی<sup>۲</sup> نادیده گرفته می‌شود، به دلیل کاربرد در حوزه‌های گوناگون پردازش زبان طبیعی از جمله ترجمه

<sup>1</sup> Virpioja

<sup>2</sup> Morphophonemic

ماشینی، بازناسی گفتار و بازیابی اطلاعات و نیز توسعه الگوریتم‌های بی‌نظارت<sup>۳</sup> و نیمه‌نظارتی<sup>۴</sup> دارای اهمیت است.

• ریشه‌یابی و خوشه‌بندی ظاهری واژه‌ها: شامل ریشه‌یابی و دسته‌بندی واژه‌ها برحسب ریشه یا ستاک است و خروجی الگوریتم برای هر کلمه یک برچسب یا ریشه است.

• تجزیه و تحلیل ساختواژی: خروجی این پردازشگرها فهرستی از برچسب‌های تعیین‌کننده ریشه و وندها است، به طوری که ساختار واژه را تحلیل و تک‌واژگونه‌ها را شناسایی می‌کند.

در حوزه پردازش ساختواژی دو رویکرد عمده استفاده از روش‌های قاعده‌مند و الگوریتم‌های داده‌محور<sup>۵</sup> است. روش‌های قاعده‌مند نیازمند استخراج تک‌واژها و روابط ساختواژی زبان است. اما در رویکردهای یادگیری ماشین<sup>۶</sup> مبتنی بر الگوریتم‌های داده‌محور، الگوریتم با استفاده از واژگانی که برای آموزش در اختیار دارد ساختواژه زبان را می‌آموزد و تک‌واژها را در واژه‌های جدید شناسایی می‌کند. علی‌رغم رشد این حوزه پژوهشی در دهه‌های اخیر و توانایی الگوریتم‌ها در شناسایی و تقطیع ساخت‌های پیچیده زبان، وندهای تصریفی و اشتقاقی، واژه‌بست‌ها، ساخت‌های مرکب و نیز کلمات خارج از واژگان، تاکنون بر روی داده‌های زبان فارسی ارزیابی نشده است. با توجه به تنوع واژه‌ها در حوزه‌های مختلف علوم و محدودیت پردازشگرهای قاعده‌مند در پردازش واژه‌های خارج از واژگان، همچنین گسترش پژوهش‌ها در حوزه پردازش زبان طبیعی به‌منظور تحلیل متون محاوره‌ای و غیررسمی پیام‌رسان‌ها و شبکه‌های اجتماعی، توجه به الگوریتم‌های داده‌محور در پردازش ساختواژی می‌تواند رویکردی بهینه‌تر و کم‌هزینه‌تر باشد.

در این پژوهش، تقطیع ساختواژی داده متنی زبان فارسی با استفاده از مدل تخمین بیزی بررسی شده، چالش‌های داده متنی و دستیابی به بهترین نتیجه با توجه به عدم دسترسی به داده آموزش استاندارد و کاستی‌هایی در داده‌های موجود، از جمله حجم و نویز داده و گوناگونی روش‌های نگارشی و کاربرد نویسه غیرزبانی نیم‌فاصله، مورد ارزیابی قرار گرفت. برای ارزیابی دقیق‌تر نتایج به‌دست آمده از الگوریتم‌های بی‌نظارت و نیمه‌نظارتی، داده حاشیه‌نویسی‌شده<sup>۷</sup> و استاندارد ساختواژی زبان فارسی طراحی و تهیه شد.

## ۲. پیشینه موضوع

در پردازش ساختواژی، دو رویکرد غالب استفاده از روش‌های قاعده‌مند و روش‌های یادگیری ماشین است. در روش قاعده‌مند، قواعد ساختواژی و واژگان زبان به یک ماشین حالت<sup>۸</sup> یا مبدل حالت محدود<sup>۹</sup> داده می‌شود و سامانه می‌تواند با دریافت یک واژه در ورودی، تجزیه و تحلیل ساختواژی کامل آن را در خروجی نمایش دهد (جورافسکی و مارتین، ۲۰۰۷).

<sup>3</sup> Unsupervised

<sup>4</sup> Semi-supervised

<sup>5</sup> Data-driven algorithm

<sup>6</sup> Machine learning

<sup>7</sup> Annotated data

<sup>8</sup> Finite-state automaton (FSA)

<sup>9</sup> Finite -state transducer (FST)

در حوزه یادگیری ماشین نخستین پژوهش‌ها به رویکردهای بانظارت<sup>۱۰</sup> بر می‌گردد که به حجم بالای داده برچسب‌خورده نیاز دارد و بسیار هزینه‌بر است. بنابراین گرایش‌ها به سمت حوزه پژوهشی بی‌نظارت سوق پیدا کرد که برای زبان‌هایی با منابع زبانی اندک کارایی بالایی دارد. در پژوهش‌های اخیر روش‌های نیمه‌نظارتی به دلیل بایاس شدن برای نوع خاصی از زبان و کاربرد، رشد چشمگیری را نشان می‌دهند. این روش‌ها را می‌توان در چهار دسته کلی طبقه‌بندی کرد:

۱- الگوریتم‌های خانواده مورفسور<sup>۱۱</sup>: خانواده مورفسور، گروهی از پردازشگرهای ساختوازی برپایه روش‌های یادگیری ماشین آماری است که از سال ۲۰۰۲ تاکنون تکامل و بهبود یافته است. یادگیری در این پردازشگرها به عنوان یک روش یادگیری پارامتری با سه مؤلفه توصیف می‌شود: مدل، تابع هزینه، الگوریتم یادگیری و تقطیع (کوریمو<sup>۱۲</sup> و همکاران، ۲۰۱۰a؛ ویرپیوا و اسمیت، ۲۰۱۵).

در بیشتر پردازشگرهای این خانواده از مدل‌های احتمالاتی زایشی<sup>۱۳</sup> استفاده شده است که می‌تواند احتمال توأم<sup>۱۴</sup>  $P(A = a, W = w | \theta)$  را برای هر کلمه  $w$  و تقطیع ساختوازی آن، که به صورت لیستی از برچسب‌های معرف تک‌واژه‌ها است  $a = (m_1, m_2, \dots, m_n)$ ، به شرط پارامترهای مدل  $\theta$  محاسبه کند.

$$P(A = a, W = w | \theta) = \prod_{i=1}^{|a|} P(m_i | \theta) \quad (1)$$

در این پردازشگرها تابع هزینه وابسته به پارامترها، مدل و مجموعه داده است. در نسخه‌های گوناگون عمدتاً برپایه تخمین بیزی<sup>۱۵</sup> یا تخمین بیشینه احتمال پسین<sup>۱۶</sup> (MAP) تعریف شده است. در نسخه‌های نخستین از تخمین بیشینه شباهت<sup>۱۷</sup> و تخمین کوتاه‌ترین طول توصیفی<sup>۱۸</sup> نیز استفاده شده است.

۲- دستور تطبیق‌دهنده<sup>۱۹</sup>: دستور تطبیق‌دهنده یک چارچوب مدل‌سازی است که از مدل بیزی غیرپارامتری<sup>۲۰</sup> استفاده و تلاش می‌کند ساخت‌های درختی را بر روی رشته نویسه‌های داده ورودی جستجو کرده، قواعد ساختوازی با پیچیدگی‌های گوناگون را بیاموزد (سرتز و گلدواتر<sup>۲۱</sup> ۲۰۱۳).

۳- میدان تصادفی شرطی: روشی برپایه مدل گرافیکی احتمالاتی غیرمستقیم تمایزی<sup>۲۲</sup> است که در آن مدل به‌منظور یادگیری ارتباط بین خروجی الگوریتمی بی‌نظارت مانند مورفسور پایه‌ای با بافت زیررشته‌ای پیرامونش آموزش می‌بیند. در این رویکرد ایده کلی تلاش و تمرکز بر مدل کردن مرزهای تک‌واژه به جای شناسایی رشته نویسه‌ها است که در چارچوب مسئله‌ای از دنباله برچسب‌ها مطرح می‌شود (روکولینن<sup>۲۳</sup> و همکاران، ۲۰۱۴).

<sup>10</sup> Supervised

<sup>11</sup> Morfessor

<sup>12</sup> Kurimo

<sup>13</sup> Generative probabilistic models

<sup>14</sup> Joint probability

<sup>15</sup> Bayesian estimation

<sup>16</sup> Maximum a posteriori

<sup>17</sup> Maximum likelihood

<sup>18</sup> Minimum description length

<sup>19</sup> Adaptor grammar

<sup>20</sup> Non-parametric Bayesian model

<sup>21</sup> Sirts & Goldwater

<sup>22</sup> Discriminative undirected probabilistic graphical model

<sup>23</sup> Ruokolainen

۴- گوناگونی پسین حروف<sup>۲۴</sup> (LSV) و انتروپی پسین حروف<sup>۲۵</sup> (LSE): مبنای این روش استفاده از خواص توزیع حروف در داخل واژه است. در این مدل گوناگونی پسین به حروفی اشاره دارد که پس از هر حرف در واژه قرار می‌گیرد و گوناگونی پیشین، حروفی است که پیش از آن قرار می‌گیرد. بنابراین نقطه‌ای که دارای پتانسیل تقطیع بالایی است را می‌توان با توجه به تغییرات قابل ملاحظه‌ای که در شماره پسین‌های حروف، در داخل واژه رخ می‌دهد، شناسایی کرد. LSE نیز روشی برپایه استفاده از انتروپی حروف پسین و پیشین به جای شمارش آنها است. مرزهای تکواژها نوعاً دارای LSE بالاتری نسبت به سایر نقاط در واژه هستند (بردگ<sup>۲۶</sup>، ۲۰۰۷).

در حوزه پردازش ساختوازی زبان فارسی اغلب پژوهش‌ها بر روش‌های قاعده‌مند متمرکز بوده و از داده واژگان زایا (اسلامی و همکاران، ۱۳۸۳) استفاده کرده‌اند. این پردازشگرها قادر به پردازش کلمات خارج از واژگان از جمله کلمات تخصصی و ساخت‌های گونه محاوره نیستند (مگردومیان، ۲۰۰۰؛ شمس‌فرد و جعفری، ۲۰۱۰؛ مواجی و همکاران، ۱۳۹۰؛ سرابی و همکاران، ۲۰۱۳). به‌طور محدود مواجی و همکاران (۱۳۹۱) پیاده‌سازی یک الگوریتم یادگیری ساختوازی بی‌نظارت از خانواده مورفسور را بر روی داده‌های زبان فارسی گزارش کرده‌اند، اما اشاره‌ای به داده آموزش و آزمون مورد استفاده و نحوه ارزیابی و دقت به‌دست آمده نداشته‌اند.

### ۳. تقطیع ساختوازی داده‌های زبان فارسی با استفاده از تخمین بی‌زی

با توجه به کاربرد گسترده تخمین بی‌زی در حوزه پردازش ساختوازی داده‌محور و دقت‌های بالایی که در پژوهش‌های گوناگون ثبت کرده است (روکولینن و همکاران، ۲۰۱۶؛ کوریمو و همکاران، ۲۰۱۰b)، یادگیری ساختوازی الگوریتم مورفسور پایه‌ای گسترش‌یافته (ویرپیوا و اسمیت، ۲۰۱۵؛ ویرپیوا و همکاران، ۲۰۱۳) بر روی داده‌های زبان فارسی ارزیابی شد.

این الگوریتم مطابق رابطه (۲) از تابع هزینه MAP استفاده و تلاش می‌کند به مدلی بهینه برسد که تابع هزینه را کمینه کند.

$$L(\theta, D_W) = -\log P(\theta) - \log P(D_W|\theta) \quad (2)$$

در تخمین MAP فرآیند یادگیری مستلزم یک گمانه یا پیش‌فرض در مورد پارامترهای مدل است که احتمال پیشین<sup>۲۷</sup>  $P(\theta)$  خوانده می‌شود. تخمین پارامترهای مدل یا احتمال پسین<sup>۲۸</sup>، در چارچوب آمار بی‌زی، محصول احتمال پیشین و تخمین شباهت داده به شرط داشتن پارامترهای مدل  $P(D_W|\theta)$  است.

<sup>24</sup> Letter successor variety

<sup>25</sup> Letter successor entropy

<sup>26</sup> Bordag

<sup>27</sup> Prior probability

<sup>28</sup> Posterior probability

مدل  $\theta$  فهرستی از تک‌واژه‌های زبان است که در داده  $D_W$  شناسایی و استخراج شده است. پارامترهایی که در این مدل تعریف و ارزیابی می‌شود، شامل خواص ظاهری هر تک‌واژه و خصوصیات کلی واژگان یا مدل است که به صورت زیر تعریف می‌شود:

- تعداد گونه‌های تک‌واژه‌ها یا به عبارتی اندازه واژگان  $\mu$
- تعداد قطعات مشاهده شده در داده  $\gamma$
- رشته‌های حروف در تک‌واژه‌ها  $(\sigma_1, \sigma_2, \dots, \sigma_\mu)$
- تعداد کل تک‌واژه‌ها و رخداد آن‌ها  $(\tau_1, \tau_2, \dots, \tau_\mu)$

برای توزیع احتمال تک‌واژه‌ها از پیش فرضی استفاده می‌شود که حاصل اشتراک دو متغیر مستقل بر مبنای شکل ظاهری<sup>۲۹</sup> تک‌واژه، رابطه (۳) و کاربرد<sup>۳۰</sup> آن در پیکره، رابطه (۴) است. به این ترتیب احتمال بالاتر به تک‌واژه‌هایی با طول کمتر اختصاص می‌یابد.

$$P(\sigma_i) = P(|\sigma_i|) \prod_{j=1}^{|\sigma_i|} P(C = \sigma_{ij}) \quad (۳)$$

$$P(\tau_1, \dots, \tau_\mu | \mu, \gamma) = 1 / \binom{\gamma-1}{\mu-1} \quad (۴)$$

در آموزش، استخراج تک‌واژه‌ها و تخمین هزینه مدل، از یک الگوریتم حریصانه بازگشتی<sup>۳۱</sup> استفاده می‌شود که در هر جستجو پس از دستیابی به هر نقطه بهینه محلی، پارامترهای مدل را به روزرسانی می‌کند. می‌توان با وزن‌دهی نقشی را که تابع پیشین در هزینه نهایی مدل ایفا می‌کند، کنترل کرد.

$$L(\theta, D_W) = -\alpha \log P(\theta) - \log P(D_W | \theta) \quad (۵)$$

در یادگیری نیمه نظارتی علاوه بر داده برچسب‌نخورده به حجم اندکی داده برچسب‌خورده نیاز است و لازم است در تخمین هزینه مدل، علاوه بر لگاریتم شباهت داده برچسب‌نخورده  $D_W$ ، لگاریتم شباهت داده برچسب‌خورده  $D_{W-An}$  نیز افزوده شود.

$$L(\theta, D_W, D_{W-An}) = -\log P(\theta) - \log P(D_W | \theta) - \log P(D_{W-An} | \theta) \quad (۶)$$

### ۳-۱. تهیه داده حاشیه‌نویسی شده ساختواژی

از چالش‌های اصلی پردازش ساختواژی زبان فارسی با استفاده از الگوریتم‌های داده‌محور، تهیه داده مناسب است. به منظور آموزش نیمه نظارتی و آزمون الگوریتم‌ها نخست می‌بایست داده استاندارد حاشیه‌نویسی شده برای تقطیع ساختواژی زبان فارسی فراهم شود که تا حد امکان چشم‌انداز قابل قبولی از واژه‌ها و تجزیه و تحلیل ساختواژی هر کدام ارائه دهد و ساخت‌ها و اشکال واژه‌ای<sup>۳۲</sup> زبان اعم از تصریفی و اشتقاقی و ترکیبی را به خوبی پوشش دهد. به این منظور برای تهیه داده استاندارد ساختواژی زبان فارسی از داده‌های مجموعه استاندارد وبلاگ‌های ایران (آل‌احمد و همکاران، ۲۰۱۶) و پیکره متنی زبان فارسی (بی‌جن‌خان و همکاران، ۲۰۱۱) استفاده شد. به دلیل تنوع نویسه‌های مورد استفاده در متون الکترونیکی فارسی و نیز گوناگونی قواعد

<sup>29</sup> Form

<sup>30</sup> Usage

<sup>31</sup> Recursive, greedy search

<sup>32</sup> Word forms

و روش‌های املائی در نگارش ساخت‌های مرکب و وندهای اشتقاقی و تصریفی، لازم است پیش از تقطیع واحدها، متون بهنجارسازی<sup>۳۳</sup> شود. با توجه به اینکه، بهنجارساز<sup>۳۴</sup>‌های موجود مانند هضم (۱۳۹۳) و ویراستیار (۱۳۹۳) هرکدام تعداد محدودی از نویسه‌های یونیکد، ساخت‌های اشتقاقی و ترکیبی و نیز وندهای تصریفی و واژه‌بست‌های احتمالی متصل به آن‌ها را بهنجار می‌کنند، در این پژوهش برای بهنجارسازی متون، بهنجارساز هضم بهبود داده شد و با اضافه کردن موارد زیر کاستی‌های مورد نظر رفع گردید:

- صد و ده نویسه یونیکد غیراستاندارد که در گذشته برای تایپ فارسی استفاده می‌شد و وبلاگ‌نویسان استفاده می‌کنند و معادل استاندارد آن‌ها به بهنجارساز افزوده شد. مثلا حرف «ی» در زبان فارسی دارای یک نویسه استاندارد یونیکد «\u06cc» و نویسه‌های غیراستانداردی مثل «\ufbaf» «\u06d2» «\ufbff» و غیره، است.

- بهنجارسازی پی‌بست‌های «ی» نکره، کسره اضافه، ربطی و کمکی و ضمیری، مانند: «خانه اش، خانه ست، خانه ی، خانه ای» و وندهای تصریفی متصل به اسم و صفت «ها - تر - ترین - جات» و قطعاتی<sup>۳۵</sup> که شامل وندهای تصریفی و پی‌بست‌های متصل به آن‌ها باشد، مانند: «های: کتابخانه های - هایمان: کتابخانه هایمان و...».

- بهنجارسازی فعل‌های چندجزئی که بیانگر زمان و نمودهای گوناگون فعل هستند مانند ماضی استمراری.

- بهنجارسازی بیش از هفتاد وند و ساخت ترکیبی با توجه به پژوهش‌های سراجی (۲۰۱۳) و اضافه کردن برخی وندهای دیگر از جمله «وش، سالگی، گانگی و...» به همراه وندهای تصریفی و پی‌بست‌های احتمالی متصل به هرکدام مانند: «ده سالگیش - تلاش گری».

- اصلاح درج اشتباه نویسه نیم‌فاصله در ابتدا یا انتهای واژه یا بیش از یک نویسه در فاصله تک‌واژها. اصلاح واژه‌های دارای تکرار حروف مانند «سلاااام» که به عنوان تاکید در صفحات وبلاگی و شبکه‌های اجتماعی رایج است.

بعد از بهنجارسازی و حذف علائم نگارشی و نویسه‌ها و سایر حروف غیرفارسی و نیز جملات و پست‌های دارای ساختارهای محاوره‌ای، فهرست واژه‌ها از متون پیکره متنی زبان فارسی و داده‌های وبلاگی ایران به همراه فراوانی هر واژه استخراج شد. همچنین بخشی از ساخت‌های غیرفارسی از زبان‌هایی مانند ترکی، کردی و عربی که دارای نویسه‌های مشترک با فارسی هستند، به صورت نیمه‌خودکار حذف شد. با توجه به حجم بالای غلط‌های املائی و تایپی، در واژگان فراوانی کمتر از ۵ که بیشترین فراوانی غلط‌های املائی را در خود داشت حذف شده، در نهایت از فهرست واژه‌های یکتا با حجم ۴۸۲۰۳۰ قطعه، ۵۰۰۰ واژه به صورت تصادفی استخراج و به صورت نیمه‌خودکار حاشیه‌نویسی شد. جهت حفظ تنوع و تناسب واژه‌های داده ساختواژی با واژه‌های موجود در پیکره، نخست واژه‌ها برحسب فراوانی دسته‌بندی شدند و سپس از هر دسته با توجه به اینکه چند درصد از کل پیکره را شامل می‌شود، تعدادی واژه به صورت تصادفی، با در نظر گرفتن تناسب مذکور، استخراج گردید. به طوریکه داده ساخت‌های پرتکرار، کم‌تکرار و با تکرار متوسط را متناسب

<sup>33</sup> Normalize

<sup>34</sup> Normalizer

<sup>35</sup> Tokens

با فراوانیشان در پیکره پوشش دهد. جدول ۱، درصد حجمی هر دسته از فراوانی‌ها نسبت به حجم کل قطعات یکتا در پیکره را نشان می‌دهد.

جدول ۱. درصد حجمی هر دسته از فراوانی‌ها در مجموع دو پیکره

فراوانی قطعات یکتا	۵ تا ۲۰	۲۰ تا ۵۰	۵۰ تا ۱۰۰	۱۰۰ تا ۱۰۰۰	بیش از ۱۰۰۰
درصد حجمی هر دسته در دو پیکره	۶۱/۶٪	۳۱/۵٪	۳/۱۶٪	۲۱٪	۶۴/۴٪

### ۲-۳. تهیه داده آموزش ساختواژی بی نظارت

داده آموزش بی نظارت فهرستی از اشکال و ساخت‌های واژه‌ای گوناگون و صحیح زبان است که حجمی بین چندصد هزار تا بیش از یک میلیون واژه دارد (روکولینن و همکاران، ۲۰۱۶؛ کوریمو و همکاران، ۲۰۱۰b). با توجه به اینکه چنین فهرستی برای زبان فارسی موجود نیست، داده آموزش از پیکره‌های موجود زبان تهیه شد.

- پیکره متنی زبان فارسی با حجمی حدود ۱۰۰ میلیون قطعه
  - مجموعه داده استاندارد وبلاگ‌های ایران با حجمی حدود ۴۰۰ میلیون قطعه
  - پیکره همشهری (آل احمد و همکاران، ۲۰۰۹) با حجمی حدود ۱۵۰ میلیون قطعه
- پس از پیش‌پردازش اولیه شامل جداسازی ساخت‌های محاوره، حذف نویسه‌های غیرفارسی، اعداد و علائم نگارشی، بهنجارسازی کدگذاری متون، اصلاح نویسه نیم‌فاصله و حذف برخی ساخت‌های پرتکرار زبان‌های دیگر مانند کردی و عربی و غیره در داده، که در بخش پیش توضیح داده شد، واژگان هر پیکره استخراج شد. داده‌های به‌دست آمده از پیکره‌ها، نسبت به داده استاندارد آموزشی که در الگوریتم‌های یادگیری ساختواژی استفاده می‌شود، کاستی‌هایی دارد که می‌توان در موارد زیر دسته‌بندی کرد:

- حجم داده صحیح: تعداد اشکال و ساخت‌های واژه‌ای گوناگون و صحیح زبان. این حجم نسبت به حجم داده آموزش مورد استفاده در الگوریتم‌ها کمتر است.
- نویز داده: غلط‌های املائی و ساخت‌های نادرست زبان در داده یا واژگان. در داده‌های مورد بررسی نویز داده بیشتر ناشی از سرهم‌نویسی گروه‌های اسمی و حرف‌اضافه‌ای، و نیز ساختارهایی از دیگر زبان‌ها است که نویسه‌های مشابه با زبان فارسی دارند مانند ترکی، عربی و غیره.
- تنوع شیوه‌های املائی: گوناگونی روش‌های املائی و شیوه‌های نگارش واژه‌های قرضی و ساخت‌های تصریفی، اشتقاقی و ترکیبی به عنوان مثال: «کتابخانه، کتاب‌خانه؛ خوبترین، خوب‌ترین».
- نویسه نیم‌فاصله: نویسه پرکاربرد غیرزبانی و صرفاً تایپی در متون الکترونیکی زبان فارسی که در الگوریتمی مانند مورفسور به عنوان یک نویسه زبانی ارزیابی می‌شود. این نویسه هم در واژه‌های غیربسیط فارسی و هم در واژه‌های بسیط قرضی به کار می‌رود مانند: «جهان‌دیده» و «اس‌ام‌اس». در واژه‌های غیربسیط نمایانگر مرز تک‌واژه است، اما در ساخت‌ها قرضی بسیط، جزئی از تک‌واژه به حساب می‌آید.

#### ۴. ارزیابی

کارایی مدل‌های ساخته‌شده با سه معیار صحت<sup>۳۶</sup>، فراخوانی<sup>۳۷</sup> و معیار<sup>۳۸</sup> F، برمبنای BPR<sup>۳۹</sup> یا تشخیص درست یا نادرست مرز تک‌واژ در واژه، ارزیابی شد. این روش، از رایج‌ترین روش‌های ارزیابی الگوریتم‌های پردازش ساختواژی داده‌محور است و به‌صورت زیر محاسبه می‌شود (کوریمو، ۲۰۰۷):

- صحت (P): نسبت تعداد مرزهای تک‌واژی صحیحی که توسط مدل تشخیص داده شده، به کل مرزهای تک‌واژی یافت شده

- فراخوانی (R): نسبت تعداد مرزهای تک‌واژی صحیحی که توسط مدل تشخیص داده شده به کل مرزهای درست تک‌واژی در داده

- معیار (F1) F:  $2 * \frac{P * R}{P + R}$

در ارزیابی مدل‌های ساخته‌شده از داده ثابت آزمون استفاده شده که مجموعه‌ای ثابت از ۳۰۰۰ واژه است و به‌طور تصادفی از داده حاشیه‌نویسی شده استاندارد انتخاب شده است.

#### ۴-۱. ارزیابی یادگیری ساختواژی بی‌نظارت

به‌منظور یافتن پارامترهایی که در بهبود یادگیری الگوریتم با توجه به داده‌های موجود موثر باشد، نقش حجم، نویز، طول واژه‌ها و نویسه نیم‌فاصله در داده، مورد ارزیابی قرار گرفت.

#### ۴-۱-۱. حجم و نویز داده

جدول ۲، نتایج یادگیری بی‌نظارت الگوریتم را برای هر سه مجموعه داده با بررسی نویز و حجم در هر مجموعه، نشان می‌دهد.

برای تخمین نویز داده از هر مجموعه با فراوانی بیشتر از یک، ۱۰۰۰ واژه تصادفی به‌صورت ۸۰ درصد با فراوانی ۱۵ و کمتر و ۲۰ درصد با فراوانی بیش از ۱۵ انتخاب شد و با شمارش تعداد اشتباهات املائی و تایپی و ساختارهای غیرفارسی برای هر مجموعه، درصد نویز تخمین زده شد. هدف از این تخمین صرفاً بررسی نسبی داده‌ها است و معیاری برای سنجش قطعی نویز داده نیست. با توجه به اینکه واژه‌ها و قطعاتی با فراوانی ۱ بیشترین درصد نویز داده را به خود اختصاص داده‌اند، در آموزش الگوریتم از واژه‌ها یا قطعات یکتایی استفاده شد که در واژگان فراوانی بیش از ۱ داشتند.

<sup>36</sup> Precision

<sup>37</sup> Recall

<sup>38</sup> F-score

<sup>39</sup> Boundary precision and recall



جدول ۲. یادگیری ساختوازی بی نظارت بر روی سه مجموعه داده.

صحت:  $P$ ، فراخوانی:  $R$ ، معیار  $F1$

F1	R	P	واژه کمتر از ۵ نویسه	نویز داده	قطعات یکتا با فراوانی بیشتر از ۱	حجم قطعات یکتا	
۰/۶۵	۰/۵۷	۰/۷۵	٪۲۰/۷	٪ ۲۸	۴۰۴۹۰۷	۸۱۱۹۱۱	پیکره متنی
۰/۶۳	۰/۵۴	۰/۷۶	٪۲۳/۳	٪ ۳۳	۷۲۴۳۸۵	۱۸۲۸۹۶۲	داده‌های وبلاگی
۰/۶۴	۰/۵۲	۰/۸۳	٪۲۸/۵	٪۲۲	۲۵۳۷۷۱	۵۸۰۸۷۷	پیکره همشهری

هرچند حجم داده همشهری به طور قابل ملاحظه‌ای از دو داده دیگر کمتر است در معیار  $F$  نتایجی مشابه دو داده دیگر و در معیار صحت بیش از ۷ درصد بالاتر ثبت کرده است. جدول ۲، درصد واژه‌ها با طول ۵ نویسه و کمتر را برای هر سه مجموعه داده، نشان می‌دهد این حجم در داده همشهری بیش از ۵ درصد از داده وبلاگی و نزدیک ۸ درصد از داده پیکره متنی بیشتر است. طول ۵ نویسه با توجه به متوسط طول تک‌واژه‌های زبان در داده حاشیه‌نویسی شده استاندارد انتخاب شده است. الگوریتم در فرآیند مدل‌سازی نخست تمام واژه‌ها را به عنوان تک‌واژه به مدل وارد کرده، سپس تلاش می‌کند با جستجوی محلی بر روی امکانات تقطیع واژه‌ها، مدل کم‌هزینه‌تر را انتخاب کند. به دلیل هزینه کم‌تری که به مدل تحمیل می‌شود، تمایل الگوریتم به تقطیع واژه‌های با طول کمتر نسبت به واژه‌های با طول بیشتر، کمتر است و این موضوع درصد تقطیع‌های نادرست در مدل را کاهش می‌دهد.

#### ۴-۱-۲. نویسه نیم‌فاصله

به منظور بررسی کارایی نیم‌فاصله، در آزمون بعد، این نویسه از داده حذف و در مرحله پیش‌پردازش با نویسه فاصله جایگزین شد. جدول ۳، حجم هر مجموعه داده بدون نیم‌فاصله و نتایج حاصل از یادگیری الگوریتم را نشان می‌دهد. با توجه به اینکه داده همشهری فاقد نویسه نیم‌فاصله است و وندها و ساخت‌های ترکیبی یا متصل به ستاک نوشته شده، مانند «کتابها» یا با نویسه فاصله از ستاک جدا شده‌اند، مانند «کتاب‌ها»، در این مرحله داده همشهری و نتایج آن تغییری با مرحله قبل ندارد. در مورد داده‌های دو پیکره دیگر، بهبود در نتایج معیار  $F$  مشاهده می‌شود که بیشتر ناشی از بهبود در فراخوانی است.

نکته قابل توجه در نتایج به دست آمده در این مرحله این است که حذف نیم‌فاصله هم سبب کاهش بخشی از نویز داده، مربوط به استفاده از نویسه نیم‌فاصله در بین اجزای یک گروه نحوی مانند «آسمان وزمین» می‌شود و هم سبب می‌شود برخی از وندهای پرکاربرد اشتقاقی و تصریفی از ستاک‌ها جدا شده، با فراوانی بالا در مدل اولیه قرار گیرد. به این ترتیب احتمال شناسایی و قرارگیری آن‌ها را در مدل نهایی افزایش می‌یابد.

نقش دیگر نویسه نیم‌فاصله در مرحله تقطیع واژه است. در صورتی که بین تک‌واژه‌ها نویسه نیم‌فاصله درج شده باشد و الگوریتم تک‌واژه را به درستی تشخیص دهد، این نویسه به عنوان بخشی از تک‌واژه قبلی یا بعدی یا به صورت یک تک‌واژه مستقل شناسایی می‌شود. این مسئله مغایرتی در خروجی الگوریتم نسبت به داده استاندارد آزمون ایجاد می‌کند. با توجه به نبود این نویسه در داده آموزش لازم است در مرحله تقطیع نیم‌فاصله به عنوان یک نویسه غیرزبانی به الگوریتم معرفی و از درج اشتباه آن در خروجی جلوگیری شود.

معرفی الگوی نیم‌فاصله، نتایج را در معیار F و صحت به‌طور متوسط نزدیک به ۷ درصد بهبود داده، که در جدول ۳، مشاهده می‌شود.

جدول ۳. ارزیابی یادگیری الگوریتم با داده فاقد نویسه نیم‌فاصله و معرفی الگوی نیم‌فاصله

صحت: P، فراخوانی: R، معیار F: F1

حذف نیم‌فاصله و معرفی الگوی نیم‌فاصله			حذف نیم‌فاصله از داده آموزش				
F1	R	P	F1	R	P	حجم داده	
۰/۷۵	۰/۶۶	۰/۸۶	۰/۶۹	۰/۶۳	۰/۷۷	۲۴۵۸۴۶	پیکره متنی
۰/۷۰	۰/۵۸	۰/۸۸	۰/۶۶	۰/۵۷	۰/۷۷	۵۴۵۸۹۸	داده‌های وبلاگی
۰/۷۴	۰/۶۷	۰/۸۳	—————				پیکره همشهری

#### ۴-۱-۳. وزن دهی به تابع احتمال پیشین

جدول ۴، بهبود نتایج را با وزن دهی به پارامتر  $\alpha$  نشان می‌دهد.

جدول ۴. نتایج حاصل از یادگیری الگوریتم با وزن بهینه  $\alpha$  و معرفی الگوی نیم‌فاصله

معیار F	فراخوانی	صحت	وزن	
۰/۷۸	۰/۷۱	۰/۸۷	۰/۷۳	پیکره همشهری
۰/۷۴	۰/۶۷	۰/۸۲	۰/۷۹	داده وبلاگی بدون نیم‌فاصله
۰/۷۹	۰/۷۲	۰/۸۶	۰/۷۴	پیکره متنی بدون نیم‌فاصله

به‌منظور تعیین وزن بهینه از داده توسعه‌یافته<sup>۴۰</sup> استفاده شد که شامل ۷۵۰ واژه است که به‌صورت تصادفی و بدون هم‌پوشانی با داده آموزش و آزمون، از داده حاشیه‌نویسی شده انتخاب شده است. پس از یادگیری با وزن بهینه، نتایج با استفاده از داده آزمون ارزیابی شد که به‌طور متوسط ۶/۳ درصد بهبود در فراخوانی و ۴ درصد بهبود در معیار F مشاهده می‌شود.

#### ۴-۲. ارزیابی یادگیری ساختوازی نیمه‌نظارتی

در یادگیری نیمه‌نظارتی الگوریتم با ۱۰۰۰ واژه که به‌صورت تصادفی و بدون هم‌پوشانی با داده آزمون، از داده حاشیه‌نویسی شده استخراج شد، آموزش دید. جدول ۵، بهبود نتایج یادگیری نیمه‌نظارتی را نسبت به یادگیری بی‌نظارت، نشان می‌دهد. در یادگیری نیمه‌نظارتی نیز معرفی الگو و حذف نیم‌فاصله با بهبود صحت و فراخوانی نتایج را در معیار F به‌طور چشمگیری بهبود داده است. با توجه به اینکه حجم ۱۰۰۰ واژه از داده حاشیه‌نویسی شده تقریباً تمام وندهای تصریفی، واژه‌بست‌ها و وندهای اشتقاقی پرکاربرد زبان را پوشش می‌دهد، سبب بهبود مدل ساخته‌شده و افزایش قابل توجه فراخوانی می‌شود. بهبود متوسط ۱۰ درصدی در فراخوانی و ۲/۷ درصدی در معیار F برای یادگیری نیمه‌نظارتی با معرفی الگو و حذف نیم‌فاصله نسبت به یادگیری بی‌نظارت ثبت شده است.

<sup>40</sup> Held out set

جدول ۵. ارزیابی یادگیری نیمه نظارتی. صحت:  $P$ ، فراخوانی:  $R$ ، معیار  $F1$

یادگیری نیمه نظارتی			یادگیری با حذف نیم فاصله از داده و معرفی الگوی نیم فاصله		
F1	R	P	F1	R	P
۰/۷۱	۰/۷۳	۰/۷۲	۰/۸۲	۰/۸۰	۰/۸۵
۰/۶۵	۰/۷۵	۰/۶۹	۰/۸۱	۰/۸۰	۰/۸۲
۰/۷۳	۰/۷۹	۰/۷۶	۰/۸۳	۰/۸۱	۰/۸۵

## ۵. جمع بندی و نتیجه گیری

در این پژوهش به ارزیابی یادگیری ساختارهای زبان فارسی با استفاده از تخمین بیزی پرداخته و داده استاندارد ساختارهای جهت ارزیابی مدل‌ها و آموزش نیمه نظارتی برای زبان فارسی طراحی و تهیه شد. با توجه به در دسترس نبودن داده استاندارد آموزش، داده برچسب نخورده از سه پیکره بزرگ زبان فارسی انتخاب و یادگیری الگوریتم با توجه به حجم، نویز، طول واژه و نقش نویسه نیم فاصله در داده مورد ارزیابی قرار گرفت. حذف نیم فاصله از داده و معرفی الگوی نیم فاصله به الگوریتم تقطیع، برای جلوگیری از درج اشتباه نیم فاصله در تک واژه‌ها، نتایج را در هر سه معیار صحت، فراخوانی و معیار  $F$  به طور قابل ملاحظه‌ای بهبود می‌دهد و متوسط نتایج ثبت شده در معیار  $F$  برای سه مجموعه داده با اختلاف ۳ درصد در یادگیری بی نظارت و ۲ درصد در یادگیری نیمه نظارتی به نتایجی که پژوهش‌های دیگر برای زبان انگلیسی ثبت کرده‌اند (روکولین و همکاران، ۲۰۱۶) نزدیک می‌شود. از دلایل بهبود یادگیری الگوریتم، کوتاه شدن طول واژه‌های داده، معرفی تک واژه‌های بیشتر به الگوریتم، کاهش تنوع الگوهای نگارشی و کاهش نویز در مواردی است که نیم فاصله در ساخت گروه‌های نحوی مانند «آسمان وزمین» وارد شده است. با توجه به نتایج به دست آمده، ارزیابی یادگیری الگوریتم با داده نوشتار غیر رسمی مانند متون شبکه‌های اجتماعی، می‌تواند زمینه‌ای برای پژوهش بیشتر باشد.

## منابع

- اسلامی، محرم،، شریفی آتشگاه، مسعود،، علیزاده، صدیقه،، زندی، طاهره. (۱۳۸۳). «واژگان زبانی فارسی». اولین کارگاه آموزش زبان فارسی و رایانه، دانشگاه تهران.
- مواجی، وحید،، اسلامی، محرم،، وزیرنژاد، بهرام. (۱۳۹۰). «پارس مورف: تحلیلگر ساختارهای زبان فارسی». دوفصلنامه پردازش علائم و داده‌ها (۱۰).
- مواجی، وحید،، اسلامی، محرم،، وزیرنژاد، بهرام. (۱۳۹۱). «پارس مورف: تحلیلگر ساختارهای زبان فارسی». پایان نامه کارشناسی ارشد، گروه زبان‌شناسی رایانشی، دانشگاه صنعتی شریف
- ویراستیار. (۱۳۹۳). [نرم افزار]. دبیرخانه شورای عالی اطلاع رسانی. برگرفته از: <http://virastyar.ir>.
- هضم (۱۳۹۳). [نرم افزار]. برگرفته از: <http://www.sobhe.ir/hazm>.
- AleAhmad, A., Amiri, H., Darrud, E., Rahgozar, M., & Oroumchian, F. (2009). "Hamshahri: A standard Persian text collection". *Knowledge-Based Systems*, 22 (5), 382-387.

- AleAhmad, A., Zahedi, M., Rahgozar, M., & Moshiri, B. (2016). "irBlogs: A standard collection for studying Persian bloggers". *Computers in Human Behavior*, 57, 195-207.
- Bijankhan, M., Sheykhzadegan, J., Bahrani, M., & Ghayoomi, M. (2011). "Lessons from building Persian written corpus: Peykare". *Language Resources and Evaluation*, 45(2), 143-164.
- Bordag, S. (2007, September). "Unsupervised and knowledge-free morpheme segmentation and analysis. In *Workshop of the Cross-Language Evaluation Forum for European Languages*", 881-891. Berlin, Springer.
- Jurafsky, D., & Martin, H. (2007). *Speech and language processing An introduction to natural language processing, computational linguistics, and speech recognition*. New Jersey: Prentice-Hall.
- Kurimo, M., Creutz, M., & Varjokallio, M. (2007, September). Morpho challenge evaluation using a linguistic gold standard. "*Workshop of the Cross-Language Evaluation Forum for European Languages*" 864-872. Berlin, Heidelberg: Springer.
- Kurimo, M., Virpioja, S., & Turunen, V. T. (2010a). "Proceedings of the Morpho challenge 2010 workshop". *Morpho Challenge Workshop; 2010; Espoo*. Aalto University School of Science and Technology.
- Kurimo, M., Virpioja, S., Turunen, V., & Lagus, K. (2010b). "Morpho Challenge competition 2005-2010: evaluations and results". *Proceedings of the 11th Meeting of the ACL Special Interest Group on Computational Morphology and Phonology*, 87-95. Association for Computational Linguistics.
- Megerdooian, K. (2000). "Persian computational morphology: A unification-based approach". *Computing Research Laboratory, New Mexico State University*.
- Ruokolainen, T., Kohonen, O., Sirts, K., Grönroos, S. A., Kurimo, M., & Virpioja, S. (2016). "A comparative study minimally supervised morphological segmentation". *Computational Linguistics*, 42(1), 91-120
- Ruokolainen, T., Kohonen, O., Virpioja, S., & Kurimo, M. (2014). "Painless semi-supervised morphological segmentation using conditional random fields". *Proceedings of the 14<sup>th</sup> Conference of European Chapter of the Association for Computational Linguistics*, 84-89.
- Sarabi, Z., Mahyar, H., & Farhoodi, M. (2013, October). "ParsiPardaz: Persian language processing toolkit". *Computer and Knowledge Engineering (ICCKE)*, 73-79. IEEE.
- Seraji, M. (2013). "Preper: A pre-processor for Persian, doctoral dissertation", *Department of Linguistics and Philology, Uppsala University*.
- Shamsfard, M., Jafari, H. S., & Ilbeygi, M. (2010, May). "STeP-1: A set of fundamental tools for Persian text processing". *LREC*.
- Sirts, K., & Goldwater, S. (2013). "Minimally-supervised morphological segmentation using adaptor grammars". *Transactions of the Association for Computational Linguistics*, 1, 255-266.
- Virpioja, S., & Smit, P. (2015). *Morfessor documentation*. Retrieved from the Web May, 2017. <https://morfessor.readthedocs.io/en/latest>.
- Virpioja, S., Smit, P., Grönroos, S., & Kurimo, M. (2013). *Morfessor 2.0: Python implementation and extensions for Morfessor Baseline*. Retrieved from the Web May, 2017. <http://www.cis.hut.fi/projects/morpho>.
- Virpioja, S., Turunen, V. T., Spiegler, S., Kohonen, O., & Kurimo, M. (2011). "Empirical comparison of evaluation methods for unsupervised learning of morphology". *TAL*, 52(2), 45-90.