

به نام خدا



وزارت علوم، تحقیقات، و فناوری

پژوهشگاه علوم و فناوری اطلاعات ایران

(ایرانداک)

خلاصه گزارش برای پژوهشگران

شناسایی مسائل و ارائه راهکارهای زبانشناختی

سازمان‌دهی و بازیابی اطلاعات

در سامانه اطلاعات علمی ایران (گنج)

1397

مجری: زهرا محمودزاده

همکار: زهرا دهسرایبی

معرفی پژوهش

پردازش متن یکی از چالش‌هایی است که سیستم‌های بازیابی اطلاعات با آن روبرو هستند. ساختار زبانی متن شامل ویژگی‌های ساختارهای نحوی-معنایی و لغوی متن تأثیر بسزایی بر دقت بازیابی دارد. در سامانه گنج به عنوان مثال، مسائل خط فارسی از جمله عواملی است که سبب کاهش دقت بازیابی می‌شود. در این طرح هدف بر آن بوده است تا حد امکان مسائل بازیابی و سامان‌دهی اطلاعات به ویژه مسائلی که ریشه در پردازش متن و ساختار زبانی آن دارد، بر اساس ادبیات موجود شناسایی شده. پس از مطالعه و دسته‌بندی آنها، راه‌کارهای لازم بر اساس ادبیات موجود ارائه شود.

۱- مرحله سازمان‌دهی و نمایه‌سازی

در این مرحله نمایه‌ساز اطلاعات پایان‌نامه‌ها را وارد می‌کند. شناسایی مسائل نمایه‌سازی صورت‌های زبانی و ارائه راهکار برای حل یا بهبود آنها براساس ادبیات موجود از اهداف این مرحله است. سامانه اطلاعات علمی ایران (گنج) حاوی داده‌های متنی از نوع پایان‌نامه است. اطلاعات پایان‌نامه‌ها طی فرایند نمایه‌سازی (سازمان‌دهی) به صورت‌های زیر و با دو زبان فارسی و انگلیسی ذخیره می‌شود: چیکده (فارسی - انگلیسی)، عنوان (فارسی - انگلیسی)، کلیدواژه (فارسی - انگلیسی)، پدیدآوران شخص (پدیدآور، استاد راهنما، استاد مشاور)، پدیدآوران سازمان (دانشگاه، دانشکده)، اطلاعات تخصصی پایان‌نامه (رشته/گرایش، مکان تولید)، حوزه موضوعی (در دو سطح).

۲- مرحله بازیابی

بنابراین در فرایند بازیابی که بر روی داده‌های سامان‌دهی شده صورت می‌گیرد، می‌توان اطلاعات داده‌ها (یا پایان‌نامه‌ها) را براساس هفت صورت بالا بدست آورد. اما عوامل متعددی برخاسته از ساختار زبانی صورت‌های سامان‌دهی و ذخیره شده شامل ویژگی‌های نحوی، معنایی، صرفی،

لغوی و املائی بر دقت بازیابی تأثیر گذاشته، آن را کاهش می‌دهد. کاربرد برای بازیابی اطلاعات ذخیره شده به وسیله نمایه‌ساز با موانع، خطاها و کم و کاستی‌هایی روبرو است. شناسایی این موانع و ارائه راهکار بر اساس ادبیات موجود هدف دیگر این طرح است.

روش پژوهش

گام اول: شناسایی مسائل زبانشناختی سازمان‌دهی و بازیابی اطلاعات سامانه اطلاعات علمی ایران به روش مطالعه و مصاحبه و جمع‌آوری میدانی نمونه‌ها از پایگاه گنج. مشاهده گزارش‌های موجود در این زمینه، مطالعه ادبیات موجود در این زمینه، مصاحبه با نمایه‌سازان مرکز ایرانداک و بخش بازیابی اطلاعات و جمع‌آوری داده‌ها از پایگاه گنج.

گام دوم: ارائه راهکار برای مسائل بالا به روش کتابخانه‌ای: مطالعه ادبیات موجود

یافته‌ها

مسائل و راهکارهای ارائه شده در جدول‌های زیر به صورت خلاصه ارائه شده است. بدیهی است مسائل ساماندهی در بازیابی نیز ملاحظه می‌شود.

راهکار	مسائل ساماندهی
۱- بهره‌گیری از زبانشناسان و متخصصان زبان فارسی	ترجمه به فارسی یا انگلیسی اطلاعات کتابشناختی
۲- افزایش قابلیت نرم‌افزار نمایه‌سازی	استفاده از معادل‌های متعدد برای یک واژه لاتین
۳- تشکیل گروه تصمیم‌ساز متشکل از زبانشناس، متخصص حوزه‌های علمی و نمایه‌ساز	تعیین معادل فارسی کلیدواژه‌های لاتین
	تعیین اختصار فارسی برای یک واژه لاتین
	نوشتن صورت‌های مخفف لاتین به اشکال مختلف
نرم‌افزارهای خطایاب زبان فارسی موجود در بازار کمک‌چندانی به حل مسائل ناشی از خطاهای مختلف نگارشی و املائی و نیز درج فرمولها نمی‌کنند. هر چند می‌توانند مقداری از مسائل و خطاها را به نمایه‌ساز نشان دهند.	تعدد صورت‌های نوشتاری اسامی خاص افراد (پدیدآوران)
	انواع خطاهای املائی و نگارشی در هر دو سطح فارسی و لاتین
	مسائل مربوط به درج علائم ریاضی و فرمولهای شیمی در چکیده و عنوان
	مطالب نامربوط در چکیده

راهکار	مسائل بازیابی
<p>بررسی ادبیات موجود نشان می دهد که با کاربرد روش ها و تکنیک های پردازش زبان طبیعی در مدل های بازیابی می تواند دقت بازیابی را افزایش داد.</p> <p>بنابراین با مشخص کردن روش بازیابی اطلاعات در سامانه گنج و سپس افزودن تکنیک های NLP به مدل های موجود می توان دقت بازیابی در این سامانه را بهبود بخشید.</p>	<p>صورت های متفاوت نوشتاری</p> <p>۱- ناشی از کاربرد علامتهای مختلف نوشتاری: تشدید، مساد، همزه، تنوین، آوا، نقطه، خط تیره</p> <p>۲- صورت های متفاوت کسره اضافه</p> <p>۳- صورت های مختلف جمع</p> <p>۴- صورت های مختلف اختصارات</p>
	<p>مسائل معنایی-صرفی-نحوی</p> <p>۱- کلمات چندمعنا با املای یکسان (اما صورت آوایی یکسان یا متفاوت)</p> <p>۲- صورت های متفاوت واژگانی با معنای یکسان یا مرتبط</p> <p>۳- صورت های کوتاه / کامل اسامی یا سازمان ها</p>

نتیجه گیری و پیشنهادها

در این پژوهش ابتدا با بررسی میدانی پایگاه اطلاعات علمی ایران (گنج)، مسائل زبانشناختی بازیابی و ساماندهی اطلاعات در این پایگاه بدست آمد و سپس راهکارهایی ارائه شد. مسائل نمایه سازی شامل موارد زیر است: ترجمه به فارسی یا انگلیسی اطلاعات کتابشناختی، استفاده از معادل های متعدد برای یک واژه لاتین، یافتن معادل فارسی کلیدواژه های لاتین، تعیین اختصار فارسی برای یک واژه لاتین، نوشتن صورت های مخفف لاتین به اشکال مختلف، انواع خطاهای املایی و نگارشی در هر دو سطح فارسی و لاتین، مسائل مربوط به درج علائم ریاضی و فرمولهای شیمی در چکیده و عنوان، مطالب نامربوط در چکیده. برای حل مسائل زبانشناختی ساماندهی و نمایه سازی به نظر می رسد علاوه بر بهبود قابلیت نرم افزار نمایه سازی، لازم است دانش تخصصی در زمینه علوم مختلف و نیز زبان فارسی و زبانشناسی بکار گرفته شود.

مسائل بازیابی به دو دسته تقسیم می شوند الف) صورت های متفاوت نوشتاری ناشی از کاربرد علامتهای مختلف نوشتاری: تشدید، همزه، تنوین، آوا، نقطه، خط تیره، صورت های متفاوت کسره اضافه و صورت های مختلف جمع و صورت های مختلف اختصارات و ب) صورت های متفاوت ناشی از مسائل معنایی، صرفی و نحوی: کلمات چندمعنا با املای یکسان (اما صورت آوایی یکسان یا متفاوت)، صورت های متفاوت واژگانی با معنای یکسان یا مرتبط و صورت های کوتاه / کامل اسامی یا سازمان ها. بررسی ادبیات مربوط به بازیابی اطلاعات نشان می دهد که بکار گیری روشها و تکنیک های پردازش زبان طبیعی تأثیر معناداری در بهبود دقت و بازخوانی سیستم های بازیابی اطلاعات دارد.