

بسم الله الرحمن الرحيم



دانشگاه صنعتی امیرکبیر

(پلی‌تکنیک تهران)

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد

گرایش نرم افزار

داده کاوی ساختار وب با استفاده از اتماتای یادگیر توزیع شده و سلولی و  
کاربردهای آن

نگارش: سارا مطیعی

استاد راهنما: دکتر محمدرضا میبدی

# بسمه تعالی



دانشگاه صنعتی امیرکبیر

(پلی تکنیک تهران)

معاونت پژوهشی

## فرم اطلاعات پایان نامه

### کارشناسی ارشد و دکترا

تاریخ:

پیوست:

نام و نام خانوادگی:	مطیعی	دانشجوی آزاد	بورسیه	معادل
شماره دانشجوئی:	۸۴۱۳۱۰۴۵	دانشکده:	مهندسی کامپیوتر	رشته تحصیلی: مهندسی کامپیوتر
نام و نام خانوادگی استاد راهنما:	محمد رضا میبدی			
عنوان پایان نامه به فارسی:	داده کاوی ساختار وب با استفاده از اتماتای یادگیر توزیع شده و سلولی و کاربردهای آن			
عنوان پایان نامه به انگلیسی:	Web Structure Mining using Distributed and Cellular Learning Automata and its Applications			
نوع پژوهه:	دکترا	کارشناسی ارشد	■	نظری
تاریخ شروع:	مهر ۸۵	تاریخ خاتمه:	۸۶ اسفند	تعداد واحد:
سازمان تأمین کننده اعتبار:	مرکز تحقیقات مخابرات ایران			
واژه های کلیدی به فارسی:	داده کاوی ساختار وب، اتماتای یادگیر توزیع شده، اتماتای یادگیر سلولی، اتماتاهای یادگیر، پیمایشگر موضوعی، اجتماعات وب.			
واژه های کلیدی به انگلیسی:	Web Structure Mining , Distributed Learning Automata , Cellular Learning Automata , Focused Crawler , Web Community			
نظرها و پیشنهادها به منظور بهبود فعالیت های پژوهشی دانشگاه:				
استاد راهنما:				
دانشجو:	تشکیل گروه های تحقیقاتی			
امضاء استاد راهنما:				
تاریخ:				
نسخه ۱: معاونت پژوهشی				
نسخه ۲: کتابخانه و به انصمام دو جلد پایان نامه به منظور تسویه حساب با کتابخانه و مرکز استاد و مدارک علمی				

این پایان نامه توسط مرکز تحقیقات مخابرات ایران حمایت مالی شده است که از این طریق سپاسگزاری می‌گردد.

تقدیم به مادر عزیزم که از سپاسگزاری زحمات بی دریغش قاصرم.

و

تقدیم به خواهر عزیزم سیما، با آرزوی موفقیت های هر چه بیشتر.

در ابتدا لازم می‌دانم از استاد ارجمند جناب آقای دکتر محمدرضا میبدی به پاس راهنمایی‌های ارزنده ایشان در جهت پیشبرد و انجام هرچه بهتر این پایان‌نامه تشکر و قدردانی نمایم.

از آقای علی برادران هاشمی که در انجام این پایان‌نامه من را یاری دادند، سپاسگزارم.

## چکیده

در سال های اخیر، برای بهره برداری از حجم وسیع داده های وب روش های وب کاوی معرفی شده اند. وب کاوی، به کارگیری روش های داده کاوی برای کشف و استخراج خودکار اطلاعات از اسناد و سرویس های وب می باشد. یکی از انواع داده کاوی، داده کاوی ساختار وب است که با استفاده از پیوند ها اطلاعات جدیدی راجع به صفحات به دست می آورد، اما پیوند ها اطلاعات کافی راجع به ارتباط بین صفحات به دست نمی دهند. یک راهکار مناسب برای بهبود نتایج روش های داده کاوی ساختار وب، به کارگیری داده های استفاده از وب و نحوه پیمایش کاربران علاوه بر پیوند ها در این روش ها می باشد. در این پژوهه دو روش برای داده کاوی ساختار وب ارایه می شود که با استفاده از ترکیب پیوند ها و داده های استفاده از وب اطلاعات جدید راجع به صفحات و ارتباطشان به دست می آورند. روش اول مبتنی بر اتماتای یادگیر توزیع شده و روش دوم مبتنی بر اتماتای یادگیر سلولی است. هر دو روش پیشنهادی از دو مرحله کلی تشکیل شده است. در مرحله اول، با استفاده از اتماتای یادگیر (توزیع شده یا سلولی)، پیوند های بین صفحات و رفتار کاربران در مشاهده صفحات وب، ساختار ارتباطی صفحات وب به دست می آید. به آن معنی که صفحات مرتبط با یکدیگر و میزان ارتباط آنها تعیین می شود. در مرحله دوم، ساختار ارتباطی به دست آمده از مرحله قبل، در دو نوع از کاربردهای داده کاوی ساختار وب استفاده خواهد شد. کاربرد اول پیمایش موضوعی صفحات وب و کاربرد دوم تشخیص اجتماعات وب است. همچنین کارایی ساختار به دست آمده، پیمایشگر طراحی شده و اجتماعات وی که با استفاده از روش های پیشنهادی تشخیص داده می شوند، با روش های مشابه مقایسه شده و رفتار آنها در شرایط گوناگون مورد بررسی قرار می گیرد.

## فهرست مطالب

۱	۱	۱- مقدمه
۶		۱-۱- وب کاوی
۸		۲- داده کاوی ساختار وب
۱۰		۱-۲-۱- اصطلاحات داده کاوی ساختار وب
۱۱		۱-۲-۲-۱- مدل های بازنمایی ساختار وب
۱۱		۱-۲-۲-۲-۱- مدل های مبتنی بر گراف
۱۳		۱-۲-۲-۲-۱- مدل های مارکو
۱۴		۱-۳-۲-۲-۱- مدل رابطه ای احتمالی
۱۴		۱-۳-۲-۲-۱- معیارهای داده کاوی ساختار وب
۱۴		۱-۳-۲-۲-۱-۱- معیارهای صفحه منفرد
۱۵		۱-۳-۲-۲-۱-۲- معیارهای گروهی از صفحات
۱۵		۱-۳-۲-۲-۱-۳- معیارهای کل گراف وب
۱۶		۱-۴-۲-۱- الگوریتم های داده کاوی ساختار وب
۱۷		۱-۴-۲-۱- HITS
۱۸		۱-۴-۲-۱- Page Rank
۲۰		۱-۴-۲-۱-۳- الگوریتم جریان پیشینه
۲۰		۱-۴-۲-۱- Average Clicks
۲۱		۱-۵-۲-۱- کاربردهای داده کاوی ساختار وب
۲۳		۱-۳-۱- داده کاوی استفاده از وب
۲۳		۱-۳-۱-۱- منابع داده
۲۴		۱-۳-۱-۲- الگوریتم های داده کاوی استفاده از وب
۲۵		۱-۳-۱-۱- قواعد انجمنی
۲۵		۱-۳-۱-۲- الگوهای ترتیبی
۲۶		۱-۳-۱-۳- خوش بندی
۲۷		۱-۳-۳-۱- کاربردهای داده کاوی استفاده از وب
۲۸		۱-۴-۱- ترکیب داده کاوی ساختار وب و داده کاوی استفاده از وب
۲۸		۱-۴-۱-۱- الگوریتم PageRank آگاه از استفاده
۲۹		۱-۴-۱-۲- پیمایش در سایت وب
۲۹		۱-۴-۱-۳- جستجو در وب محدود
۳۰		۱-۴-۱-۴- HITS تغییر داده شده
۳۰		۱-۵-۱- انوماتیک یادگیر توزیع شده و سلولی

۳۰	۱-۵-۱- اتوماتای سلولی .....
۳۱	۲-۵-۱- اتوماتای یادگیر .....
۳۲	۳-۱-۲-۵-۱- الگوریتمهای یادگیری .....
۳۵	۲-۲-۵-۱- اتوماتای یادگیر با عمل های متغیر .....
۳۵	۳-۵-۱- اتوماتای یادگیر سلولی (CLA) .....
۳۷	۱-۳-۵-۱- اتوماتای یادگیر سلولی ناهمگام (ACLA) .....
۳۹	۴-۵-۱- اتوماتای یادگیر توزیع شده .....
۴۰	۶-۱- ساختار پایان نامه .....
۴۱	<b>۲ تشخیص ساختار ارتباطی صفحات وب .....</b>
۴۱	۱-۲- مقدمه .....
۴۳	۲-۲- تشخیص ساختار ارتباطی صفحات وب با استفاده از اتوماتای یادگیر توزیع شده .....
۴۹	۱-۲-۲- ارزیابی ساختار ارتباطی به دست آمده با استفاده از اتوماتای یادگیر توزیع شده .....
۴۹	۱-۱-۲-۲- مدل شبیه سازی .....
۵۱	۲-۱-۲-۲- معیار ارزیابی .....
۵۲	۳-۱-۲-۲- مقایسه با روش های موجود .....
۵۳	۴-۱-۲-۲- بررسی رفتار روش پیشنهادی .....
۵۶	۳-۲- تشخیص ساختار ارتباطی صفحات وب با استفاده از اتوماتای یادگیر سلولی ناهمگام .....
۶۳	۱-۳-۲- ارزیابی ساختار ارتباطی به دست آمده با استفاده از اتوماتای یادگیر سلولی .....
۶۳	۱-۱-۳-۲- مقایسه با روش های موجود .....
۶۴	۲-۱-۳-۲- بررسی رفتار روش پیشنهادی .....
۶۹	۴-۲- نتیجه گیری .....
۷۰	<b>۳ پیمایش موضوعی صفحات وب .....</b>
۷۰	۱-۳- مقدمه .....
۷۲	۲-۳- الگوریتم های پیمایش موضوعی .....
۷۲	۱-۲-۳- پیمایش موضوعی بدون استفاده از دانش پیش زمینه .....
۷۲	۱-۱-۲-۳- الگوریتم عرض اول .....
۷۳	۲-۱-۲-۳- BestFirst- الگوریتم .....
۷۴	۳-۱-۲-۳- الگوریتم SharkSearch .....
۷۵	۴-۱-۲-۳- الگوریتم PageRank .....
۷۶	۲-۲-۳- پیمایش موضوعی به کمک دانش پیش زمینه .....
۷۶	۱-۲-۲-۳- پیمایش مبتنی بر طبقه بندي موضوعی .....
۷۷	۲-۲-۲-۳- پیمایش با استفاده از گراف زمینه .....

۷۸ .....	- پیمایش مبتنی بر هسته شناسی ..... ۳-۲-۲-۳
۷۹ .....	- سایر الگوریتم ها ..... ۳-۲-۳
۷۹ .....	- پیمایش هوشمند ..... ۱-۳-۲-۳
۸۰ .....	- پیمایشگر یادگیر ..... ۲-۳-۲-۳
۸۱ .....	- الگوریتم پیشنهادی برای پیمایش موضوعی ..... ۳-۳
۸۵ .....	- ارزیابی پیمایشگر موضوعی وب ..... ۱-۳-۳
۹۱ .....	- نتیجه گیری ..... ۴-۳
۹۳ .....	<b>۴ تشخیص اجتماعات وب</b>
۹۳ .....	- مقدمه ..... ۱-۴
۹۴ .....	- الگوریتم های تشخیص اجتماعات وب ..... ۲-۴
۹۵ .....	- روش های مبتنی بر تحلیل پوندها ..... ۱-۲-۴
۹۵ .....	- HITS ..... ۱-۱-۲-۴
۹۶ .....	- RPA ..... ۲-۱-۲-۴
۹۷ .....	- روش های مبتنی بر تئوری گراف ..... ۲-۲-۴
۹۷ .....	- روش های مبتنی بر گراف کامل دویخشی ..... ۱-۲-۲-۴
۹۸ .....	- روش های مبتنی بر الگوریتم جریان پیشینه ..... ۲-۲-۲-۴
۹۹ .....	- الگوریتم پیشنهادی برای تشخیص اجتماعات وب ..... ۳-۴
۱۰۳ .....	- ارزیابی اجتماعات وب به دست آمده با استفاده از مدل شبیه سازی ..... ۱-۳-۴
۱۰۶ .....	- ارزیابی اجتماعات وب به دست آمده با استفاده از داده های واقعی وب ..... ۲-۳-۴
۱۰۶ .....	- مجموعه داده های واقعی وب ..... ۱-۲-۳-۴
۱۰۸ .....	- پردازش داده های واقعی وب ..... ۲-۲-۳-۴
۱۰۸ .....	- معیار ارزیابی ..... ۳-۲-۳-۴
۱۰۹ .....	- آزمایش ها ..... ۴-۲-۳-۴
۱۱۱ .....	- نتیجه گیری ..... ۴-۴
۱۱۲ .....	<b>۵ نتیجه گیری</b>
۱۱۵ .....	<b>۶ مراجع</b>
۱۲۱ .....	واژه نامه انگلیسی به فارسی
۱۲۶ .....	واژه نامه فارسی به انگلیسی
۱۳۱ .....	<b>ضمائم</b>
۱۳۱ .....	<b>الف - پیاده سازی</b>

---

۱۳۱ .....	- مقدمه .....	۱
۱۳۲ .....	- مولفه شبیه ساز وب .....	۲
۱۳۲ .....	- مولفه تعیین ساختار ارتباطی .....	۳
۱۳۳ .....	- ساختار مولفه .....	۴
۱۳۸ .....	- عملکرد مولفه .....	۵
۱۳۹ .....	- مولفه پیمایش موضوعی .....	۶
۱۳۹ .....	- ساختار مولفه .....	۷
۱۴۴ .....	- عملکرد مولفه .....	۸
۱۴۵ .....	- مولفه تشخیص اجتماعات وب .....	۹
۱۴۵ .....	- ساختار مولفه .....	۱۰
۱۴۸ .....	- عملکرد مولفه .....	۱۱
۱۵۰ .....	ب- کد برنامه .....	

## فهرست اشکال

..... شکل ۱-۱- طبقه بندی وب کاوی و کاوش ساختار وب ..... ۹
..... شکل ۲-۱- مدل های گراف تک گره ای ..... ۱۱
..... شکل ۳-۱- مدل های گراف چند گره ای ساده ..... ۱۱
..... شکل ۴-۱- مدل های گراف چند گرهای پیچیده ..... ۱۲
..... شکل ۵-۱- مدل گراف وب ..... ۱۲
..... شکل ۶-۱- الگوریتم جریان بیشینه ..... ۲۰
..... شکل ۷-۱- ارتباط بین اتماتاتی یادگیر و محیط ..... ۳۲
..... شکل ۸-۱- قانون ..... ۳۷
..... شکل ۹-۱- اتماتاتی یادگیر توزیع شده ..... ۳۹
..... شکل ۱۰-۱- شبکه کد الگوریتم به دست آوردن ساختار ارتباطی با استفاده از DLA ..... ۴۸
..... شکل ۱۰-۲- کورولیشن روش ارایه شده در [37] ..... ۵۲
..... شکل ۱۰-۳- کورولیشن روش پیشنهادی ..... ۵۳
..... شکل ۱۱-۲- مقایسه انواع کاربران در ساختار مبتنی بر DLA ..... ۵۴
..... شکل ۱۱-۳- حذف جریمه، پیوند و رابطه تراگذری ..... ۵۶
..... شکل ۱۱-۴- شبکه کد الگوریتم به دست آوردن ساختار ارتباطی با استفاده از ACLA ..... ۶۲
..... شکل ۱۱-۵- مقایسه کورولیشن ساختار مبتنی بر ACLA و DLA ..... ۶۴
..... شکل ۱۱-۶- مقایسه انواع کاربران در ساختار مبتنی بر ACLA ..... ۶۵
..... شکل ۱۱-۷- بررسی تاثیر مقدار R بر کورولیشن ساختار مبتنی بر ACLA ..... ۶۶
..... شکل ۱۱-۸- مقایسه اتماتاتی یادگیر سلولی همگام و ناهمگام ..... ۶۷
..... شکل ۱۱-۹- مقایسه استراتژی های فعال سازی سلول ..... ۶۸
..... شکل ۱۱-۱۰- الگوریتم عرض اول ..... ۷۳
..... شکل ۱۱-۱۱- الگوریتم BestFirst ..... ۷۴
..... شکل ۱۱-۱۲- الگوریتم SharkSearch ..... ۷۵
..... شکل ۱۱-۱۳- الگوریتم PageRank ..... ۷۶
..... شکل ۱۱-۱۴- شبکه کد الگوریتم پیمایش موضوعی مبتنی بر CLA / DLA ..... ۸۵
..... شکل ۱۱-۱۵- مقایسه نرخ حاصل برای چهار پیمایشگر ..... ۸۸
..... شکل ۱۱-۱۶- مقایسه فراخوان هدف برای چهار پیمایشگر ..... ۸۹
..... شکل ۱۱-۱۷- واپستگی به مجموعه آغازین ..... ۹۰
..... شکل ۱۱-۱۸- تاثیر به کارگیری امتیاز Hub در پیمایشگر مبتنی بر DLA ..... ۹۱
..... شکل ۱۱-۱۹- تاثیر به کارگیری امتیاز Hub در پیمایشگر مبتنی بر CLA ..... ۹۱

---

شکل ۴-۱- الگوریتم تشخیص اجتماعات وب با استفاده از گراف کامل دویخشی ..... ۹۷
شکل ۴-۲- مثالی از یک اجتماع وب (مجموعه گره های سمت چپ تصویر) ..... ۹۸
شکل ۴-۳- الگوریتم تشخیص اجتماعات وب با استفاده از الگوریتم جریان بیشینه ..... ۹۸
شکل ۴-۴- نحوه عملکرد الگوریتم جریان بیشینه ..... ۹۹
شکل ۴-۵- شبیه کد الگوریتم تشخیص اجتماعات وب مبتنی بر CLA / DLA ..... ۱۰۲
شکل ۴-۶- مقایسه روش پیشنهادی با الگوریتم HITS ..... ۱۰۴
شکل ۴-۷- مقایسه روش پیشنهادی با روش مبتنی بر گراف کامل ..... ۱۰۵
شکل ۴-۸- تاثیر مجموعه ریشه در تشخیص اجتماع و ب ..... ۱۰۷
شکل ۴-۹- مقایسه مدل شبیه سازی با داده های واقعی و ب ..... ۱۰۹
شکل ۴-۱۰- مقایسه با الگوریتم HITS و الگوریتم [42] با استفاده از داده های واقعی و ب ..... ۱۱۰
شکل ۱- مولفه های سیستم ..... ۱۳۱
شکل ۲- نمودار کلاس مولفه تعیین ساختار ارتباطی ..... ۱۳۳
شکل ۳- واسط کاربر برای تعیین ساختار ارتباطی ..... ۱۳۹
شکل ۴- نمودار دنباله برای مولفه تعیین ساختار ارتباطی ..... ۱۴۰
شکل ۵- نمودار کلاس مولفه پیمایش موضوعی ..... ۱۴۱
شکل ۶- واسط کاربر برای پیمایش موضوعی ..... ۱۴۴
شکل ۷- نمودار دنباله برای مولفه پیمایش موضوعی و ب ..... ۱۴۵
شکل ۸- نمودار کلاس مولفه تشخیص اجتماعات و ب ..... ۱۴۶
شکل ۹- واسط کاربر برای تشخیص اجتماعات و ب ..... ۱۴۸
شکل ۱۰- نمودار دنباله برای مولفه تشخیص اجتماعات و ب ..... ۱۴۹

## ۱ مقدمه

وب، محیطی وسیع، متنوع و پویا است که کاربران متعدد اسناد خود را در آن منتشر می کنند. وب طی یک فرآیند آشفته و غیر متمرکز رشد می کند و این روند منجر به تولید حجم وسیعی از مستندات متصل به یکدیگر گشته است که از هیچ گونه سازماندهی منطقی برخوردار نیستند. در حال حاضر بیش از ۱۱۵ بیلیون صفحه [34] در وب موجود است و این تعداد با نرخ ۷,۳ میلیون صفحه در روز افزایش می یابد. با توجه به حجم وسیع اطلاعات در وب، مدیریت آن با ابزارهای سنتی تقریباً غیر ممکن است و ابزارها و روش هایی نو برای مدیریت آن مورد نیاز است. به طور کلی کاربران وب در استفاده از آن با مشکلات زیر روبرو هستند:

- ۱- **یافتن اطلاعات مرتبط:** یافتن اطلاعات مورد نیاز در وب دشوار می باشد. روش های سنتی بازیابی اطلاعات که برای جستجوی اطلاعات در پایگاه داده ها به کار می روند، قابل استفاده در وب نمی

باشند و کاربران معمولاً از موتورهای جستجو که مهمترین و رایج‌ترین ابزار برای یافتن اطلاعات در وب می‌باشند، استفاده می‌کنند. این موتورها، یک پرس و جوی<sup>۱</sup> مبتنی بر کلمات کلیدی از کاربر دریافت کرده و در پاسخ لیستی از اسناد مرتبط با پرس و جوی وی را که بر اساس میزان ارتباط با این پرس و جو مرتب شده‌اند، به وی ارائه می‌کنند. اما موتورهای جستجو دارای دو مشکل اصلی هستند. اولاً دقت<sup>۲</sup> موتورهای جستجو پایین است، چراکه این موتورها در پاسخ به یک پرس و جوی کاربر صدها یا هزاران سند را بازیابی می‌کنند، در حالی که بسیاری از اسناد بازیابی شده توسط آنها با نیاز اطلاعاتی کاربر مرتبط نمی‌باشند. دوماً میزان فراخوان<sup>۳</sup> این موتورها کم می‌باشد، به آن معنی که قادر به بازیابی کلیه اسناد مرتبط با نیاز اطلاعاتی کاربر نیستند. چرا که حجم اسناد در وب بسیار زیاد است و موتورهای جستجو قادر به نگهداری اطلاعات کلیه اسناد وب، در پایگاه داده‌های خود نمی‌باشند.

**۲- ایجاد دانش جدید با استفاده از اطلاعات موجود در وب:** این مشکل در واقع بخشی از مشکل مطرح شده در قسمت قبل می‌باشد. در حال حاضر این سوال مطرح است که چگونه می‌توان داده‌های فراوان موجود در وب را به دانشی قابل استفاده تبدیل کرد، به طوری که یافتن اطلاعات مورد نیاز در آن به سادگی صورت بگیرد. همچنین چگونه می‌توان با استفاده از داده‌های وب به اطلاعات و دانشی جدید دست یافت.

**۳- خصوصی سازی<sup>۴</sup> اطلاعات:** از آن جا که کاربران متفاوت هر یک درباره نوع و نحوه بازنمایی اطلاعات سلیقه خاصی دارند، این مسئله باید مورد توجه تامین کنندگان اطلاعات در وب قرار بگیرد. برای این منظور با توجه به خواسته‌ها و تمایلات کاربران متفاوت، نحوه ارائه اطلاعات به آنها باید سفارشی گردد.

برای بهره برداری از حجم وسیع داده و کاهش مشکلات فوق الذکر، در سال‌های اخیر روش‌های وب کاوی<sup>۵</sup> معرفی شده‌اند. وب کاوی به کارگیری روش‌های داده کاوی<sup>۶</sup> برای کشف و استخراج خودکار اطلاعات از اسناد و سرویس‌های وب می‌باشد. البته روش‌های وب کاوی تنها ابزار موجود برای حل این مشکلات نیستند. بلکه روش‌های مختلفی از سایر زمینه‌های تحقیقاتی همچون پایگاه داده

<sup>1</sup> Query

<sup>2</sup> Precision

<sup>3</sup> Recall

<sup>4</sup> Personalization

<sup>5</sup> Web Mining

<sup>6</sup> Data Mining

ها، بازیابی اطلاعات، پردازش زبان طبیعی، ... قابل استفاده در این زمینه می باشند. همچنین وب کاوی با زمینه های مختلف تحقیقاتی علوم کامپیوتر همچون داده کاوی، پایگاه داده، بازیابی اطلاعات، هوش مصنوعی، یادگیری ماشین، پردازش زبان طبیعی، استخراج اطلاعات، انبار داده ها، طراحی واسط کاربر و ... در ارتباط تنگاتنگ است.

با توجه به گسترش روز افزون حجم اطلاعات در وب و ارتباط وب کاوی با تجارت الکترونیکی، وب کاوی به یک زمینه تحقیقاتی وسیع مبدل گشته است. تکنیک ها و روش های وب کاوی از کاربرد وسیعی در حوزه های مختلف همچون موتورهای جستجو، تجارت الکترونیکی، دولت الکترونیکی، آموزش الکترونیکی، آموزش از راه دور، سازمان های مجازی، مدیریت دانش، کتابخانه های دیجیتال، ... برخوردارند. اگرچه وب کاوی با چالش ها و محدودیت های متنوعی رو به رو است که از آن جمله می توان به داده های ناصحیح و نادقیق، حجم وسیع و رو به گسترش داده های وب، تغییر داده های وب در فواصل زمانی کوتاه، کاربران گوناگون با نیازهای مختلف، عدم وجود ابزارها و الگوریتم های مناسب اشاره کرد.

روش های وب کاوی بر اساس آن که چه نوع داده ای را مورد کاوش قرار می دهد، به سه دسته داده کاوی محتوای وب<sup>۱</sup>، داده کاوی ساختار وب<sup>۲</sup> و داده کاوی استفاده از وب<sup>۳</sup> تقسیم می شوند. داده کاوی محتوای وب، فرآیند استخراج اطلاعات مفید از محتوای مستندات وب است. داده کاوی ساختار وب به کشف اطلاعات جدید با استفاده از پیوندهای<sup>۴</sup> بین صفحات وب می پردازد. داده کاوی استفاده از وب نیز داده های مربوط به استفاده کاربران از وب را مورد کاوش قرار می دهد و الگوهای استفاده از وب را به منظور درک و برآوردن بهتر نیازهای کاربران استخراج می کند.

بخش عمدهی فعالیت ها و تحقیقات انجام شده در وب کاوی به محتوای صفحات وب می پردازند. اما در سال های اخیر داده کاوی ساختار وب و داده کاوی استفاده از وب نیز مورد توجه قرار گرفته اند. محور اصلی این پایان نامه داده کاوی ساختار وب می باشد. همان طور که گفته شد، الگوریتم های داده کاوی ساختار وب با استفاده از پیوندها اطلاعات جدیدی راجع به صفحات به دست می آورند. در

<sup>1</sup> Data Warehouse

<sup>2</sup> Web Content Mining

<sup>3</sup> Web Structure Mining

<sup>4</sup> Web Usage Mining

<sup>5</sup> Hyperlink

این نوع از وب کاوی، وب به صورت یک گراف مدلسازی می شود که در آن صفحات وب، گره های گراف و پیوندهای بین صفحات، یال های گراف هستند. الگوریتم های داده کاوی ساختار وب در کاربردهای متفاوتی همچون رتبه بندی<sup>۱</sup> صفحات وب، تشخیص اجتماعات وب<sup>۲</sup>، پیماش صفحات وب<sup>۳</sup>، تحلیل گراف وب، مدلسازی و شبیه سازی فرآیند تولید گراف وب به کار می روند. به عنوان مثال، الگوریتم HITS، صفحات مرتبط با یک موضوع و یا اجتماعات وب را تشخیص می دهد [47]. الگوریتم Page Rank، رتبه ای برای هر یک از صفحات محاسبه می کند که کیفیت صفحات را مستقل از یک موضوع خاص نشان می دهد [76]. الگوریتم جریان بیشینه<sup>۴</sup>، قسمت های متراکم گراف وب را تشخیص می دهد و آن را به عنوان اجتماع وب معرفی می کند [42]. الگوریتم Fish Search برای پیماش موضوعی صفحات وب<sup>۵</sup> به کار می رود و در حین پیماش صفحاتی را جمع آوری می کند که مرتبط به یک موضوع خاص هستند [25].

روش هایی که در داده کاوی ساختار وب به کار می رود، تنها از پیوندها استفاده می کنند، در حالی که پیوندها اطلاعات کافی راجع به ارتباط بین صفحات به دست نمی دهنند. چرا که بسیاری از پیوندهای بین صفحات، به دلایلی غیر از ارتباط صفحات، بین آنها ایجاد شده اند. از جمله این دلایل می توان به تسهیل پیماش توسط کاربران و تبلیغات تجاری و ... اشاره کرد. به همین جهت به روش هایی در داده کاوی ساختار وب نیاز است که تاثیر مشکلات فوق الذکر را کاهش دهند.

یکی از راهکار مناسب برای بهبود نتایج روش های داده کاوی ساختار وب، به کارگیری داده های استفاده از وب علاوه بر پیوندها در این روش ها می باشد. در واقع نحوه پیماش کاربران در وب و صفحاتی که مشاهده می کنند، می تواند اطلاعات مفیدی راجع به این صفحات به دست دهد. چرا که کاربران معمولاً از محتویات صفحه ای که می خواهند آن را در گام بعدی خود انتخاب کنند آگاهی نسبی دارند و بر اساس نیاز اطلاعاتی خود صفحه بعدی را انتخاب می کنند و حرکت کاربران در بین صفحات اتفاقی نیست. در واقع کاربر با استفاده از اطلاعات خود ارتباطی مجازی بین صفحات ایجاد کرده و آنها را مشاهده می کند و انتظار می رود که صفحات مشابه و مرتبط با یک موضوع با یکدیگر مورد استفاده

<sup>1</sup> Ranking

<sup>2</sup> Web Community

<sup>3</sup> Crawling

<sup>4</sup> Maximum Flow Algorithm

<sup>5</sup> Focused Crawling

قرار گیرند. به این ترتیب با تحلیل داده‌های استفاده از وب و بدون تلاش مضاعف کاربر یا افراد خبره، اطلاعات با ارزشی راجع به صفحات و ارتباطشان بدست می‌آید.

در این پروژه دو روش برای داده کاوی ساختار وب ارایه می‌شود که با استفاده از ترکیب پیوندها و داده‌های استفاده از وب اطلاعات جدید راجع به صفحات و ارتباطشان به دست می‌آورند. روش اول مبتنی بر اتوماتای یادگیر توزیع شده<sup>۱</sup> و روش دوم مبتنی بر اتوماتای یادگیر سلولی<sup>۲</sup> می‌باشد. هر دو روش پیشنهادی از دو مرحله کلی تشکیل شده است. در مرحله اول، با استفاده از اتوماتای یادگیر (توزیع شده یا سلولی)، پیوندهای بین صفحات و رفتار کاربران در مشاهده صفحات وب، ساختار ارتباطی صفحات وب به دست می‌آید. به آن معنی که صفحات مرتبط با یکدیگر و میزان ارتباط آنها تعیین می‌شود. در مرحله دوم، ساختار ارتباطی به دست آمده از مرحله قبل، در دو نوع از کاربردهای داده کاوی ساختار وب استفاده خواهد شد. کاربرد اول پیمایش موضوعی صفحات وب و کاربرد دوم تشخیص اجتماعات وب است که در ادامه هر یک به اختصار معرفی می‌شود.

پیمایش وب یکی از وظایف موتورهای جستجو می‌باشد. پیمایش وب برای جمع آوری صفحاتی به کار می‌رود که توسط موتور جستجو شاخص گذاری<sup>۳</sup> می‌شوند. موتور جستجو برای پاسخ دادن به پرس و جوی کاربران از این صفحات پیمایش شده و شاخص گذاری شده استفاده می‌کند. پیمایشگرهای وب کار خود را با یک مجموعه اولیه از صفحات آغاز می‌کنند و با استفاده از پیوندهای موجود در این صفحات، صفحات دیگر را پیمایش کرده و این روند را تا رسیدن تعداد صفحات پیمایش شده به حدی مشخص انجام می‌دهند. از آن جا که حجم صفحات وب بسیار بالا و همواره رو به افزایش است، موتورهای جستجو قادر به شاخص گذاری صفحات محدودی هستند. به همین جهت در سال‌های اخیر به جای پیمایش کل وب، الگوریتم‌های پیمایش موضوعی یا متمرکز مورد توجه محققین قرار گرفته است که در حین پیمایش به صورت انتخاب گر عمل نموده و صفحات مرتبط با یک موضوع خاص را جمع آوری می‌کنند.

کاربرد دومی که در این پروژه مورد بررسی قرار خواهد گرفت تشخیص اجتماع و ب است. اجتماع وب مجموعه‌ای از صفحات وب است که راجع به یک موضوع مشترک هستند و معمولاً توسط افراد یا سازمان‌های مختلف که علایق مشترک درباره یک موضوع خاص دارند، ایجاد می‌شوند. همچنین تعداد

<sup>1</sup> Distributed Learning Automata

<sup>2</sup> Cellular Learning Automata

<sup>3</sup> Indexing

اتصالات صفحات یک اجتماع وب با یکدیگر بیش از تعداد اتصالاتشان با سایر صفحات وب است. تشخیص اجتماعات وب از این جهت که می‌تواند به کاربران در بازیابی اطلاعات از وب کمک کند، اهمیت ویژه‌ای دارد. با تشخیص یک اجتماع وب درباره یک موضوع خاص، کاربران می‌توانند با استفاده از صفحات اجتماع، اطلاعات مفیدی درباره آن موضوع به دست آورند.

علت آنکه روش‌های پیشنهادی برای داده کاوی ساختار وب به دو مرحله کلی تقسیم شده آن است که با انجام مرحله اول و به دست آوردن ساختار ارتباطی صفحات (صفحات مرتبط و میزان ارتباط آنها با یکدیگر)، می‌توان از ساختار به دست آمده در کاربردهای گوناگون داده کاوی ساختار وب استفاده کرد. الگوریتم‌های داده کاوی ساختار وب تنها به پردازش گراف و ب می‌پردازنند، اما روش‌های پیشنهادی در این پروژه در واقع بر یک گراف اصلاح شده عمل می‌کنند. همچنین این رویکرد امکان استفاده مجدد از اطلاعات به دست آمده راجع به صفحات را در کاربردهای داده کاوی ساختار وب فراهم می‌کند. به عنوان مثال پس از آن که ساختار ارتباطی تعیین شد، می‌توان از آن به دفعات متعدد در تشخیص اجتماع وب استفاده کرد. به این ترتیب روش‌های ارایه شده نسبت به سایر روش‌هایی که محاسبات و پردازش‌های خود را در تشخیص اجتماع وب و یا هر کاربرد دیگری تکرار می‌کنند، به صرفه‌تر می‌باشند.

در ادامه این فصل به معرفی وب کاوی، داده کاوی ساختار وب، الگوریتم‌ها و کاربردهای آن پرداخته می‌شود. همچنین از آنجا که در روش‌های پیشنهادی ترکیب پیوندها و داده‌های استفاده از وب به کار گرفته شده است، به الگوریتم‌ها و کاربردهای داده کاوی استفاده از وب و الگوریتم‌هایی که در آنها از ترکیب این دو شاخه از وب کاوی استفاده شده است، اشاره می‌شود. از آنجا که در الگوریتم‌های پیشنهادی در این پایان نامه از اتوماتای یادگیر توزیع شده و اتوماتای یادگیر سلولی استفاده شده است، این مفاهیم نیز در انتهای فصل اول معرفی می‌شود.

## ۱-۱- وب کاوی

در [48] وب کاوی به صورت زیر تعریف شده است: وب کاوی به کارگیری تکنیک‌های داده کاوی برای کشف و استخراج خودکار اطلاعات از اسناد و سرویس‌های وب می‌باشد. وب کاوی شامل چهار مرحله اصلی می‌باشد:

۱. پیدا کردن منبع: این مرحله شامل بازیابی صفحات وب مورد نظر می‌باشد.

**۲. انتخاب اطلاعات و پیش پردازش:** در این مرحله به صورت خودکار اطلاعات خاصی از صفحات بازیابی شده انتخاب و پیش پردازش می شوند.

**۳. تعمیم<sup>۱</sup>:** در این مرحله به طور خودکار الگوهای عام در صفحات بازیابی شده کشف می شوند.

**۴. تحلیل:** در این مرحله الگوهای به دست آمده در مرحله قبل اعتبار سنجدی<sup>۲</sup> و تفسیر می شوند.

در مرحله اول داده ها از منابع موجود در وب مانند خبرنامه های الکترونیکی، گروه های خبری، اسناد HTML، پایگاه داده های متنی و ... بازیابی می شوند. مرحله انتخاب و پیش پردازش شامل هر گونه فرآیند تبدیل داده های بازیابی شده در مرحله قبل می باشد. این پیش پردازش می تواند کاهش کلمات به ریشه آنها<sup>۳</sup>، حذف کلمات زائد<sup>۴</sup>، پیدا کردن عبارات موجود در متن و تبدیل بازنمایی داده ها به قالب رابطه ای یا منطق مرتبه اول باشد. در مرحله سوم از تکنیک های داده کاوی و یادگیری ماشین برای تعمیم استفاده می شود. در مرحله آخر، الگوهای به دست آمده ارزیابی می گردند.

به این ترتیب وب کاوی، فرآیند کشف اطلاعات و دانش ناشناخته و مفید از داده های وب می باشد. این فرآیند به طور ضمنی شامل فرآیند کشف دانش در پایگاه داده ها (KDD<sup>۵</sup>) نیز می شود. در واقع وب کاوی گونه توسعه یافته KDD است که بر روی داده های وب عمل می کند.

روش های وب کاوی بر اساس آنکه چه نوع داده ای استفاده می کنند، به ۳ دسته تقسیم می شوند [87]:

**۱. کاوش محتوای وب:** کاوش محتوای وب فرآیند استخراج اطلاعات مفید از محتوای مستندات وب است. محتوای یک سند وب متناظر با مفاهیمی است که آن سند در صدد انتقال آن به کاربران است. این محتوا می تواند شامل متن، تصویر، ویدئو، صدا و یا رکوردهای ساخت یافته مانند لیست ها و جداول باشد. در این میان کاوش متن بیش از سایر زمینه ها مورد تحقیق قرار گرفته است. از جمله این تحقیقات می توان به تشخیص موضوع<sup>۶</sup>، استخراج قواعد انجمنی<sup>۷</sup>، خوشه بندی<sup>۸</sup> و طبقه بندی اسناد وب اشاره کرد. روش ها و تکنیک های موجود در این گروه، از تکنیک های بازیابی اطلاعات و پردازش زیان

<sup>1</sup> Generalization

<sup>2</sup> Validation

<sup>3</sup> Stemming

<sup>4</sup> Stop Words

<sup>5</sup> Knowledge Discovery in Data Base

<sup>6</sup> Topic Discovery

<sup>7</sup> Association Rule

<sup>8</sup> Clustering