

لا اله الا الله محمد رسول الله

سمینار

عنوان

مروری جامع بر متن کاوی و تکنیک‌های آن

فهرست مطالب

چکیده	۵
فصل اول: مقدمه	
۱-۱- مقدمه	۷
فصل دوم: مروری بر متن کاوی	
۱-۲- مقدمه	۱۰
۲-۲- سیستم‌های استخراج اطلاعات	۱۱
۳-۲- نیاز به متن کاوی	۱۲
۱-۳-۲- آماده سازی متن	۱۴
۲-۳-۲- پردازش متن	۱۵
۳-۳-۲- تحلیل متن	۱۵
۴-۲- چگونه متن کاوی را انجام دهیم	۱۵
۵-۲- چارچوب متن کاوی	۱۸
۶-۲- یافتن روابط	۱۹
۷-۲- کاربردهای متن کاوی	۱۹
۱-۷-۲- شناسایی spam	۱۹
۲-۷-۲- نظارت	۲۰
۳-۷-۲- شناسایی نامهای مستعار	۲۰
۴-۷-۲- خلاصه سازی	۲۰
۵-۷-۲- روابط میان مفاهیم	۲۰
۶-۷-۲- یافتن و تحلیل ترندها	۲۱
۷-۷-۲- گروه بندی و طبقه بندی داده	۲۱
۸-۷-۲- برچسب زدن نحوی	۲۱
۸-۲- کاربرد متن کاوی در کتابخانه‌ها	۲۲
۹-۲- جمع بندی	۲۳

فصل سوم: مروری بر تکنیک های متن کاوی

۲۵	۱-۳-۱- مقدمه
۲۶	۲-۳-۲- تکنیک های متن کاوی
۲۶	۱-۲-۳-۱- خلاصه سازی متن
۲۶	۲-۲-۳-۲- طبقه بندی
۲۷	۱-۲-۲-۳-۱- طبقه بندی کننده ی ساده ی بیزی
۲۷	۲-۲-۲-۳-۲- طبقه بندی کننده ی K نزدیکترین همسایه
۲۸	۳-۲-۳-۳- خوشه بندی
۲۸	۱-۳-۲-۳-۱- خوشه بندی سلسله مراتبی
۲۸	۱-۱-۳-۲-۳-۱- روش خوشه بندی سلسله مراتبی پایین به بالا
۲۹	۲-۱-۳-۲-۳-۲- روش خوشه بندی سلسله مراتبی بالا به پایین
۲۹	۲-۳-۲-۳-۲- خوشه بندی افزاینده
۲۹	۱-۲-۳-۲-۳-۲- الگوریتم k means
۳۰	۴-۲-۳-۴- استخراج اطلاعات
۳۰	۵-۲-۳-۵- بصری سازی
۳۱	۳-۳-۳- جمع بندی

فصل چهارم: نتیجه گیری و پیشنهادات

۳۳	۱-۴-۱- نتیجه گیری و پیشنهادات
۳۴	مراجع

فهرست اشکال

شکل (۱-۲) نحوه عملکرد متن کاوی ۱۸

شکل (۲-۲) فرآیند متن کاوی ۱۸

چکیده

رشد فزاینده پایگاه داده‌ها در زمینه‌های مختلف از فعالیت انسان باعث شده است که نیاز به ابزارهای قدرتمند جدید، برای تغییر دادن داده به دانش مفید، افزایش یابد. جهت برآوردن این نیاز، محققان به کاوش در زمینه‌های مختلف برای یافتن روش‌ها و ایده‌های مناسب پرداختند. متن کاوی یکی از زمینه‌های است که به دنبال استخراج اطلاعات مفید، از داده‌های متنی بدون ساختار، به وسیله شناسایی و اکتشاف الگوها می‌باشد. ایده اصلی متن کاوی، یافتن قطعات کوچک اطلاعات از حجم زیاد داده‌های متنی، بدون نیاز به خواندن تمام آن است. در این سمینار با توجه به اهمیت این روش مختصراً به متن کاوی، زمینه‌های مرتبط با آن و برخی روش‌های رایج طبقه‌بندی و خوشه‌بندی پرداخته شده است. اگرچه بیان همه روش‌ها و کاربردها ممکن نیست، اما این سمینار می‌تواند دید کلی از متن کاوی را در ذهن خواننده ایجاد کرده و در صورت علاقه برای مطالعه بیشتر، فرد را به منابع مناسب هدایت کند.

واژه های کلیدی: متن‌کاوی، اطلاعات، دانش، طبقه بندی، خوه بندی.

فصل اول

مقدمه

داده‌ها، نخستین شکل اطلاعات هستند که به منظور ایجاد دانش، مدیریت و کاویده می‌شوند. داده‌ها دارای چندین مشخصه هستند: حجم، سرعت بر حسب زمان، تنوع، صحت، دوام پذیری، اعتبار، ارزش و مدت اعتبار. "حجم"، به مقدار زیاد داده‌ها بر می‌گردد. "سرعت بر حسب زمان"، نرخ تولد داده‌ها در هر واحد زمانی را نشان می‌دهد. "تنوع"، بر شکل‌های مختلف داده، مانند متن (گزارش سلامت)، عدد (آمار و ارقام بازار سهام)، تصاویر (عکسبرداری ماهواره ای)، صوت (تماس‌های تلفنی)، ویدئو و هر فرم و شکل دیگری که بتوان تصور کرد، دلالت دارد. "صحت"، با انحرافات، اختلالات و مویز در داده‌ها سروکار دارد. "دوام پذیری"، به معنای بررسی ارتباط یک متغیر در آرایه‌ی وسیعی از متغیرهای مربوط به داده‌های چند بعدی است و ارتباطات میان متغیرها است. "اعتبار"، این پرسش را درباره‌ی داده‌ها مطرح می‌کند که آیا آن داده، برای استفاده و کاربرد در نظر گرفته شده، قابل اعتماد و دقیق است؟. "ارزش" حاکی از اهمیت کلیدی داده‌ها است. برخی از داده‌ها میتوانند بسیار مهم باشند؛ در حالی که بعضی دیگر از ارزش کم‌تری برخوردارند. آخرین مورد، یعنی "مدت اعتبار" در این باره است که داده‌ها چه مدت اعمار دارند و می‌بایست ذخیره شوند [1].

هدف از داده کاوی، کشف ضمنی الگوها و روند ناشناخته قبلی از پایگاه داده‌ها است. داده کاوی شامل تکنیک‌های بسیاری چون طبقه‌بندی، خوشه‌بندی، شبکه‌های عصبی و درخت‌های تصمیم است. اگر داده‌ها به اندازه تمامی اب سطح زمین باشند، آنگاه داده‌های متنی همانند اقیانوس بیش‌ترین بخش آن را تشکیل میدهند [2]. متن ممکن است در اندازه زیاد و فرم‌های متفاوتی هم چون زبان‌های مختلف، با استفاده از نمادهای مختلف و قالب‌های متفاوت موجود باشد. از این رو، این پرسش ایجاد می‌شود که چگونه اطلاعات را می‌توان از این متن خارج کرد، در این جاست که متن کاوی به ایفای نقش می‌پردازد [1]. متن کاوی، کاربردی از داده کاوی است. تفاوت اصلی این دو، آن است که در متن کاوی، الگوها از متنی با زبان طبیعی استخراج می‌شوند، این در حالی است که داده کاوی بر روی پایگاه داده‌های ساخت یافته عمل می‌کند. بنابراین، داده کاوی، بازیابی اطلاعات، پردازش زبان طبیعی و استخراج اطلاعات از زمینه‌های مرتبط با متن کاوی هستند.

داده کاوی: روشی بسیار کارا برای کشف اطلاعات از داده‌های ساخته یافته است. متن کاوی مشابه داده کاوی است، اما ابزارهای داده کاوی طراحی شده اند تا داده‌های ساخت یافته از پایگاه داده را به کار ببرند. می‌توان گفت، متن کاوی یک راه حل بهتر برای شرکت‌ها است.

بازیابی اطلاعات: معمولا در بازیابی اطلاعات با توجه به نیاز مطرح شده از سوی کاربر، مرتبط ترین متون و مستندات و یا در واقع کیسه ی کلمه از میان دیگر مستندات یک مجموعه بیرون کشیده می شود. بازیابی اطلاعات یافتن دانش نیست بلکه تنها آن مستنداتی را که مرتبط تر به نیاز اطلاعاتی جستجوگر تشخیص داده به او تحویل می دهد. این روش در واقع هیچ دانش و حتی هیچ اطلاعاتی را به ارمغان نمی آورد.

پردازش زبان طبیعی: هدف کلی آن رسیدن به یک درک بهتر از زبان طبیعی توسط کامپیوترهاست. تکنیک های مستحکم و ساده ای را برای پردازش سریع متن به کار می برد. همچنین از تکنیک های آنالیز زبان شناسی نیز برای پردازش متن استفاده می کنی. نقش پردازش زبان طبیعی در متن کاوی فراهم کردن یک سیستم در مرحله استخراج اطلاعات با داده های زبانی است.

استخراج اطلاعات: هدف استخراج اطلاعات خاص از سندهای متنی است و می تواند به عنوان یک فاز پیش پردازش در متن کاوی به کار رود. استخراج اطلاعات عبارتند از نگاشت متن های زبان طبیعی به یک نمایش ساخت یافته و از پیش تعریف شده یا قالب هایی که وقتی پر می شوند، منتخبی اط اطلاعات کلیدی از متن اصلی را نشان می دهند. این سیستم های استخراج اطلاعات به شدت بر داده های تولید شده توسط سیستم های پردازش زبان طبیعی تکیه دارند.