

# مقدمه‌ای بر بازیابی اطلاعات

تألیف:

کریستوفر دی. مینیگ  
پرباکر رَگون  
هنریک شوتزه

ترجمه:

دکتر هدیه ساجدی (عضو هیأت علمی دانشگاه تهران)  
مهندس زهرا سادات تقی  
فرناز سادات تقی

نیاز دانش

## مقدمه مترجمین

بازیابی اطلاعات به فن آوری و دانش پیچیده جستجو و استخراج اطلاعات، داده‌ها و فرآداده‌ها در انواع گوناگون منابع اطلاعاتی اطلاق می‌شود. با افزایش روزافزون حجم اطلاعات ذخیره شده در منابع مختلف، فرآیند بازیابی و استخراج اطلاعات اهمیت ویژه‌ای پیدا کرده است. برخلاف پایگاه داده‌ها، اطلاعات ذخیره شده در منابع اطلاعاتی بزرگ مانند وب و زیرمجموعه‌های آن مانند شبکه‌های اجتماعی، از ساختار مشخصی پیروی نمی‌کنند و عموماً دارای معانی تعریف شده و مشخصی نیستند. هدف بازیابی اطلاعات در چنین شرایطی، کمک به کاربر برای یافتن اطلاعات مورد نظر در انبوهی از اطلاعات بدون ساختار است.

کتاب مقدمه‌ای بر بازیابی اطلاعات اولین کتاب با درک منسجمی از بازیابی اطلاعات سنتی و وب است که شامل جستجوی وب و حوزه‌های مرتبط با دسته‌بندی و خوشبندی متن است. یکی از دلایل ترجمه این کتاب، انسجام خوب آن در این حیطه بود که توسط اساتید متخصص دانشگاهی نوشته شده است و به خوبی می‌تواند نیاز خوانندگان را مرتყع کند. از طرفی تلاش ما بر این بود که کتاب درسی منسجمی به زبان فارسی ارائه شود که به تحولات نسبتاً زیاد حیطه بازیابی اطلاعات پرداخته باشد.

در بخش پیشگفتار، مؤلفان کتاب، ویژگی‌های کتاب را به خوبی مطرح کرده‌اند و نیازی به تکرار آنها نیست. شایان توجه است که این کتاب در مجموع، با قلمی روان طیف قابل قبولی از مطالب را در زمینه بازیابی اطلاعات مطرح کرده است. به طوریکه طبق مصوبات شورای عالی برنامه‌ریزی وزارت علوم، تحقیقات و فن آوری تقریباً کلیه سرفصل‌های درس کارشناسی و کارشناسی ارشد بازیابی اطلاعات پیشرفته را پوشش می‌دهد.

در ترجمه این کتاب بر امانت‌داری، صحت و روانی ترجمه تاکید شده است. از ابداع واژه‌های جدید و مترادف واژه‌های موجود پرهیز و سعی شده از منابع موجود استفاده شود. مجموعه واژگان (فارسی به انگلیسی و انگلیسی به فارسی) در انتهای کتاب آورده شده است.

در اینجا لازم می‌دانیم تشکر ویژه‌ای از اعضای خانواده داشته باشیم. بدیهی است انجام اینکار بدون حمایت و همراهی آنها عملی نبود.

به یقین ترجمه حاضر خالی از اشکال و خطای نیست. بنابراین از تمامی عزیزانی که این کتاب را مطالعه می‌کنند، تقاضا داریم که ما را از نظرات و انتقادات خود از طریق ناشر مطلع نمایند. در انتها امیدواریم که مطالب این کتاب بتواند پاسخگوی نیاز اساتید، دانشجویان و افراد علاقه‌مند به حیطه بازیابی اطلاعات بوده و موجبات تعالی دانش آنها را فراهم آورد.

مترجمین

زمستان ۱۳۹۳

## پیشگفتار

مطالعات دهه ۱۹۹۰ نشان دادند که اکثر افراد ترجیح می‌دادند که اطلاعات را از یکدیگر دریافت کنند تا اینکه از سیستم‌های بازیابی اطلاعات بگیرند. البته، در آن دوره زمانی، اکثر افراد از دفاتر آژانس‌های مسافرتی نیز برای رزرو سفر خود استفاده می‌کردند. هر چند، در طول دهه گذشته، بهینه‌سازی‌ها در راستای اثربخشی بازیابی اطلاعات، موتورهای جستجوی وب را ارتقاء داده، طوری که اکثر افراد بیشتر اوقات از آن راضی بوده و جستجوی وب به یک استاندارد تبدیل شده و اغلب مرجع یافتن اطلاعات گشته است. برای مثال، Pew Internet Survey (2004) (Fallows 2004) بیان کرد که «۹۲٪ کاربران اینترنت بیان می‌کنند که اینترنت جای خوبی برای گرفتن اطلاعات روزانه است». حوزه بازیابی اطلاعات برای بسیاری از افراد جالب توجه است و از اطلاعات سیستم دانشگاهی تا مبانی اساسی دستیابی به اطلاعات را تحت تأثیر قرار داده است. این کتاب پشتیبانی علمی در این زمینه را در سطحی که برای دانشجویان و مقاطع بالاتر قابل استفاده باشد، فراهم می‌آورد.

بازیابی اطلاعات هم‌زمان با وب آغاز نشد: در پاسخ به چالش‌های متعدد در فراهم آوردن دسترسی به اطلاعات، حوزه بازیابی اطلاعات برای ارائه روش‌های جستجو برای صورت‌های مختلف محتوا مطرح شد. این حوزه با انتشارات علمی و رکوردهای کتابخانه‌ای آغاز گشت، اما خیلی زود به صورت‌های دیگر محتوا، به ویژه آن صورت‌هایی که مورد استفاده متخصصان است مانند، روزنامه‌نگاران، موکلان، و پژوهان است بسط پیدا کرد. اکثر تحقیقات علمی در بازیابی اطلاعات روی این محتواها رخ داده است و اکثر کاربردهای ادامه‌دار از بازیابی اطلاعات، به فراهم آوردن دسترسی به اطلاعات بدون ساختار در شرکت‌ها و دامنه‌های دولتی مختلف سروکار داشته و این کاربردها و مباحث بخش اساسی این کتاب را تشکیل می‌دهد.

با این حال، در سال‌های اخیر، محرك اصلی بازیابی اطلاعات، شبکه جهانی وب بوده است که در مقیاس‌دها میلیون تولیدکننده محتوا عمل کرده است. اگر اطلاعات بازیابی و تفسیر و تحلیل نمی‌شد طوری‌که کاربر بتواند به سرعت اطلاعاتی مرتبط به نیازش را بیابد این انفجار اطلاعات منتشر شده مطرح نمی‌شد. در اوخر دهه ۱۹۹۰، افراد زیادی دریافتند که ادامه شاخص‌گذاری کل وب، به دلیل رشد نمایی وب از نظر اندازه، به سرعت غیرممکن خواهد شد. اما بخش اعظم نوآوری‌های علمی و مهندسی، کاهش سریع قیمت سخت افزار کامپیوتر و افزایش تجارت تحت جستجوی وب، همگی به

قدرت امروزی اکثر موتورهای جستجو افزودند. این موتورها قادر هستند، نتایج با کیفیت بالا را در زمان پاسخ کمتر از یک ثانیه برای هزاران میلیون جستجو در روز، روی میلیاردها صفحه وب فراهم کنند.

### سازماندهی کتاب و توسعه دوره

این کتاب نتیجه یک سری دوره‌هایی است که در دانشگاه استنفورد و دانشگاه اشتوتگارت در بازه زمانی یک نیم ترم، یک ترم کامل و دو نیم ترم تدریس کردایم. این دوره‌ها بیشتر برای دانشجویان در اوایل دوره کارشناسی ارشد علوم کامپیوتر برگزار شده است. اما دانشجویان وکالت، پژوهشکی، آمار، زبان‌شناسی و گرایشات مختلف مهندسی را نیز شامل شده است. از این‌رو، اصل کلیدی طراحی این کتاب، پوشش مطالبی بود که اعتقاد داریم در دوره یک ترم تحصیلی در بازیابی اطلاعات مهم است. اصل دیگر ایجاد هر فصل براساس مطالبی است که می‌توان آن را در ارائه درس ۷۵ تا ۹۰ دقیقه‌ای پوشش داد.

هشت فصل اول کتاب، به مبانی بازیابی اطلاعات و به ویژه قلب موتورهای جستجو تخصیص داده شده، که ما این مطالب را اساس هر دوره‌ای در بازیابی اطلاعات می‌دانیم. فصل ۱، شاخص‌های وارونه را معرفی می‌کند و نشان می‌دهد که چگونه پرس و جو های بولی ساده می‌تواند با استفاده از چنین شاخص‌هایی پردازش شود. فصل ۲، این مقدمه را با ارائه جزئیات روشی که در آن استناد از شاخص‌گذاری پیش‌پردازش شده و با بحث راجع به بکارگیری شاخص وارونه به شیوه‌های مختلف برای افزایش سرعت و عملیاتی بودن، بسط می‌دهد. فصل ۳، ساختارهای جستجو برای لغتنامه‌ها و چگونگی پردازش پرس و جوهایی که خطاهای املایی داشته و دیگر تطبیق‌های غیردقیق به مجموعه واژگان در مجموعه استنادی که در حال جستجو است را مورد بحث قرار می‌دهد. فصل ۴، شماری از الگوریتم‌ها را برای ساخت شاخص وارونه از مجموعه متن با توجه ویژه به الگوریتم‌های توزیع شده و مقیاس‌پذیر توصیف می‌کند، که می‌توانند در مجموعه‌های بسیار بزرگ به کار گرفته شوند. فصل ۵ روش‌های فشرده‌سازی لغتنامه و شاخص‌های وارونه را پوشش می‌دهد. این روش‌ها، برای دستیابی به زمان‌های پاسخ کمتر از یک ثانیه برای پرس و جوهایی که در موتورهای جستجوی بزرگ، ضروری هستند. شاخص‌ها و پرس و جوهایی که در فصول ۱ تا ۵ مورد توجه قرار گرفتند، تنها با بازیابی بولی سروکار دارند که در آن یک سند یا با یک پرس و جو تطبیق دارد یا تطبیق ندارد. تمایل به اندازه‌گیری میزانی که یک سند با یک پرس و جو تطبیق دارد، یا نمره یک سند برای یک پرس و جو، موجب طرح وزن‌دهی عبارت و محاسبات نمره‌ها در فصول ۶ و ۷ می‌شود و به ایده‌ی لیستی از استناد منتهی می‌شود که برای پرس و جو رتبه بندی شده‌اند. فصل ۸، بر ارزیابی سیستم بازیابی اطلاعات مبتنی بر ربط سندی که آن را بازیابی می‌کند، تمرکز کرده و به ما اجازه می‌دهد تا عملکردهای نسبی سیستم‌های مختلف را براساس مجموعه‌های سند و پرس و جو که برای محک و ارزیابی سیستم‌های بازیابی اطلاعات ارائه شده‌اند مقایسه کنیم.

فصل ۹ تا ۲۱ بر مبنای هشت فصل اول، انواع موضوعات پیشرفت‌تر را پوشش می‌دهند. فصل ۹ روشن‌هایی را بررسی می‌کند که توسط آن بازیابی می‌تواند از طریق روش‌هایی مانند، بازخورد ربط و توسعه پرس‌وجو، که برای افزایش درستنمایی بازیابی اسناد مرتبط تلاش می‌کنند، بهبود یابد. فصل ۱۰، بازیابی اطلاعات از استنادی را در نظر می‌گیرد که با زبان‌های نشانه‌گذاری مانند XML و HTML ساختار یافته‌اند. بازیابی ساخت‌یافته را به وسیله کاهش آن به روش‌های نمره‌گذاری فضای بردار که در فصل ۶ مطرح شده ارائه می‌دهیم. فصول ۱۱ و ۱۲، به نظریه احتمالاتی برای محاسبه نمره‌های سند در پرس و جو استناد می‌کند. فصل ۱۱، بازیابی اطلاعات احتمالاتی سنتی را توسعه می‌دهد که چارچوبی را برای محاسبه احتمال ربط یک سند با معلوم بودن مجموعه عبارات پرس و جو فراهم می‌آورد. از این‌رو، این احتمال می‌تواند به عنوان نمره‌ای در رتبه بندی استفاده شود. فصل ۱۲، راه حل دیگری را شرح می‌دهد که در آن برای هر سندی در مجموعه، یک مدل زبانی می‌سازیم که از آن می‌توان احتمالی که مدل زبانی یک پرس و جوی مشخص را ایجاد کند، تخمین بزنیم. این احتمال کمیت دیگری است که با آن می‌توان اسناد را رتبه بندی کرد.

فصل ۱۳ تا ۱۷ صورت‌های مختلف یادگیری ماشین و روش‌های عددی در بازیابی اطلاعات را بیان می‌کند. فصول ۱۳ تا ۱۵ مسئله دسته بندی اسناد را با دانستن مجموعه اسناد و دسته‌هایی که به آن تعلق دارند، بیان می‌کند. فصل ۱۳، دسته بندی آماری را به عنوان یکی از تکنولوژی‌های کلیدی برای موتور جستجو مطرح کرده و Naïve Bayes را که روش ساده و کارآمد دسته بندی متن است، معرفی می‌کند و روش استانداردی را برای ارزیابی دسته بندی متن طرح ریزی می‌کند. فصل ۱۴، مدل فضای بردار را از فصل ۶ اتخاذ کرده و دو روش دسته بندی Rocchio و k-nzدیکترین همسایه (kNN) را معرفی می‌کند. همچنین، مصالحه بایاس سواریانس را به عنوان ویژگی مهمی از مسائل یادگیری مطرح می‌کند که ضابطه‌ای را برای انتخاب یک روش مناسب برای مسئله دسته بندی متن فراهم می‌آورد. فصل ۱۵، ماشین‌های بردار پشتیبان را معرفی می‌کند که بسیاری از محققان آن را در حال حاضر به عنوان کارآمدترین روش دسته بندی معرفی می‌کنند. همچنین، روابطی را بین مسئله دسته‌بندی و موضوعات به ظاهر مجزا، مانند استنتاج توابع نمره‌گذاری از مجموعه مثال‌های آموزشی، ارائه می‌دهیم.

فصل ۱۶ تا ۱۸ مسئله ایجاد خوشه‌ها را از اسناد مرتبط در یک مجموعه در نظر می‌گیرند. در فصل ۱۶، ابتدا مروری بر شماری از کاربردهای مهم خوشه‌بندی در بازیابی اطلاعات خواهیم داشت. سپس دو الگوریتم خوشه‌بندی را معرفی می‌کنیم: الگوریتم K-means که به طور گسترده و مؤثر برای خوشه‌بندی اسناد به کار می‌رود؛ و الگوریتم بیشینه‌سازی امیدریاضی<sup>۱</sup> که از نظر ریاضی پر هزینه‌تر بوده اما منعطف‌تر می‌باشد. فصل ۱۷، خوشه‌بندی سلسله مراتبی ساخت‌یافته را در بسیاری از کاربردهای بازیابی اطلاعات مطرح کرده و شماری از الگوریتم‌های خوشه‌بندی که سلسله مراتب خوشه‌ها را تولید می‌کنند، معرفی می‌کند. این فصل، همچنین مسئله دشوار محاسبه خودکار برچسب‌ها

<sup>1</sup> Expectation Maximization

را برای خوشها مطرح می‌کند. فصل ۱۸، به طرح روش‌های جبرخطی که در خوشبندی می‌تواند مورد استفاده قرار گیرند پرداخته و همچنین کاربردهای روش‌های جبری در بازیابی اطلاعات را مطرح می‌کند که در روش شاخص‌گذاری معنایی نهان بکار گرفته شده‌اند.

فصل ۱۹ تا ۲۱، مسئله جستجوی وب را در نظر می‌گیرند. در فصل ۱۹، بطور خلاصه به چالش‌های پایه در جستجوی وب به همراه مجموعه‌ای از روش‌ها در بازیابی اطلاعات وب پرداخته می‌شود. سپس، فصل ۲۰، معماری و الزامات پیمایشگر وب را توصیف می‌کند. در نهایت فصل ۲۱، قدرت تحلیل پیوند<sup>۱</sup> را در جستجوی وب در نظر می‌گیرد که در این روال، چندین روش را از جبر خطی گرفته تا نظریه احتمالات پیشرفت‌به کار می‌گیرد.

این کتاب، تمامی عناوین مرتبط با بازیابی اطلاعات را بطور جامع پوشش نمی‌دهد. ما شماری از موضوعات را کنار گذاشتیم زیرا خارج از محدوده مطالعی است که تمایل به پوشش آنها در مقدمه‌ای بر بازیابی اطلاعات داشتیم. با این حال، برای افرادی که علاقه‌مند به این موضوعات هستند، چندین کتاب معرفی می‌کیم.

**بازیابی اطلاعات بین زبانی:** Grossman and Frieder (2004)، فصل ۴؛ و Oard and Dorr (1996).  
**بازیابی اطلاعات تصویر و چندسانه‌ای:** Grossman and Frieder (2004)، فصل ۴؛  
Baeza-Yates and Ribeiro-Neto (1999)، فصل ۶؛  
Baeza-Yates and Ribeiro-Neto (1999)، فصل ۱۲؛ Lew (1999) del Bimbo؛  
Baeza-Yates and Ribeiro-Neto (2001)، فصل ۱۱؛ و Smeulders et al. (2000).

**بازیابی صوت:** Coden et al. (2002).

**بازیابی موسیقی:** Downie (2006) و http://www.ismir.net/.

**واسطه‌ای کاربری برای بازیابی اطلاعات:** Baeza-Yates and Ribeiro-Neto (1999)، فصل ۱۰.  
**بازیابی اطلاعات موازی و نظیر به نظیر:** Grossman and Frieder (2004)، فصل ۷؛  
Baeza-Yates and Ribeiro-Neto (1999)، فصل ۹؛ و Aberer (2001).

**کتابخانه‌های دیجیتال:** Baeza-Yates and Ribeiro-Neto (1999)، فصل ۱۵؛ و Lesk (2004).

**چشم‌انداز علم اطلاعات:** Meadowet et al. (1997) Korfhage (1997) و Ingwersen and Järvelin (2005).

**روش‌های مبتنی بر منطق، برای بازیابی اطلاعات:** van Rijsbergen (1989).

**روش‌های پردازش زبان طبیعی:** Jurafsky and Martin (1999) Manning and Schütze (2008) و Lewis and Jones (1996).

## پیش‌نیازها

دوره‌های مقدماتی در ساختمان داده‌ها و الگوریتم‌ها، جبر خطی و نظریه احتمالات به عنوان پیش‌نیاز برای تمامی ۲۱ فصل کفایت می‌کند.

<sup>1</sup> Link Analysis

فصل ۱ تا ۵، دوره مبانی الگوریتم و ساختمندانه را به عنوان پیش نیاز در نظر می‌گیرد. به علاوه، فصول ۶ و ۷ به دانش جبرخطی پایه شامل بردارها و ضرب نقطه‌ای نیاز دارند. هیچ پیش نیاز دیگری تا فصل ۱۱ لازم نیست. در فصل ۱۱ دوره مبانی در نظریه احتمالات مورد نیاز می‌باشد؛ بخش ۱-۱۱ مرور کوتاهی دارد بر مفاهیم مورد نیاز فصول ۱۱ تا ۱۳. فصل ۱۵ فرض می‌کند که خواننده با مفهوم بهینه‌سازی غیرخطی آشنا است، اگر چه این فصل ممکن است بدون دانش دقیق در مورد الگوریتم‌های بهینه‌سازی غیرخطی نیز مطالعه شود. فصل ۱۸ نیازمند دوره مقدماتی در جبرخطی شامل آشنایی با مفاهیم رتبه ماتریس و بردار مشخصه است که مرور مختصری بر آن در بخش ۱-۱۸ ارائه شده است. دانش مرتبط با مقادیر ویژه و بردارهای ویژه نیز در فصل ۲۱ مورد نیاز است. سطح دشواری تمرینات با آسان (★)، متوسط (★★)، و یا دشوار (★★★) مشخص می‌گردد.

## تقدیر و تشکر

از انتشارات دانشگاه کمبریج تشکر می‌کنیم که به ما اجازه دادند پیش‌نویسی از کتاب را به صورت آنلاین ایجاد کنیم و بازخورده از کتاب در حین نگارش آن به دست آوریم. همچنین از Lauren Cowles تشکر می‌کنیم. او ویراستار برجسته‌ای است و چندین دوره نظراتی در مورد هر فصل ازنظر شیوه نگارش، سازماندهی و پوشش مطالب، ارائه کرد. وی در رسیدن ما به اهداف نگارش این کتاب نقش ویراستاری بسزایی داشته است.

ما مدیون افراد بسیاری هستیم که نظرات، پیشنهادات و اصلاحاتی بر مبنای نسخه پیش‌نویس این کتاب ارائه کرده‌اند.

افراد بسیاری بازخورد دقیق در مورد هر یک از فصول یا بنا به درخواست و یا براساس خواست خود ایجاد کرده‌اند و برای این امر قدردان آنها هستیم.

و در آخر از داوران که بر حسب کیفیت و کمیت نظراتی را فراهم آورده‌اند به خاطر تأثیر قابل توجه آنها بر محتوا و ساختار کتاب تشکر می‌کنیم. ما قدردانی خود را از Pavel Berkhin، Andrew Trotman، Byron Dom، Jamie Callan، Stefan Büttcher، Stefan Suel، Torsten Suel، Byron Dom، Ray Mooney و Ray در این سه فصل و به ویژه توصیف پیچیدگی زمانی در تمامی الگوریتم‌های دسته‌بندی متن قدردانی می‌کنیم.

مؤلفان از دانشگاه استنفورد و اشتوتگارت برای فراهم آوردن محیط دانشگاهی برای به بحث گذاشتن ایده‌ها و فرصت تدریس دوره‌هایی که از آن این کتاب نتیجه شده و در آن محتواش اصلاح شده است، تشکر می‌کنند. کریستوفر دی. منینگ از خانواده‌اش برای ساعات بسیاری که وی صرف کار بر روی این کتاب کرده است، تشکر کرده و امیدوار است که سال آینده وقت آزاد بیشتری را آخر هفته‌ها با آنها سپری کند. پریاکر رگون از خانواده‌اش برای پشتیبانی صبورانه‌شان در طول نگارش این کتاب قدردانی کرده و به! Yahoo برای فراهم آوردن محیط پژوهشی که در آن این کتاب کار شده

است، مدييون است. هنريک شوتزه از والدين، خانواده و دوستانش برای حمایتشان در طول نگارش اين كتاب تشکر می کند.

### اطلاعات تماس و وب

این کتاب، وب سایتی با آدرس <http://informationretrieval.org> دارد. برروی این وب سایت مجموعه‌ای از اسلایدها که برای هر فصل ایجاد شده قرار داده شده است و می‌تواند برای تدریس درس بازیابی اطلاعات بکار رود. ما از بازخوردها، اصلاحات و پیشنهادات استقبال می‌کنیم که می‌تواند از طریق [informationretrieval \(at\) yahoogroups \(dot\) com](mailto:informationretrieval(at)yahoogroups(dot)com) به تمامی مؤلفان فرستاده شود.

# فهرست مطالب

۲۱	فصل ۱ بازیابی بولی
۲۳	۱-۱ مثالی از مسئله بازیابی اطلاعات
۲۷	۲-۱ اولین برداشت در ساخت شاخص وارونه
۳۰	۳-۱ پردازش پرس‌وجوهای بولی
۳۵	۴-۱ مدل بسط یافته بولی در مقابل بازیابی رتبه‌بندی شده
۳۸	۵-۱ مراجع و مطالعات آتی
۴۱	فصل ۲ مجموعه واژگان عبارات و لیست پست‌ها
۴۲	۱-۲ شرح و توصیف سند و کدگشایی دنباله‌ی کاراکتر
۴۲	۱-۱-۲ دستیابی به دنباله کاراکتر در یک سند
۴۳	۲-۱-۲ انتخاب واحد سند
۴۴	۲-۲ تعیین مجموعه واژگان عبارت
۴۴	۱-۲-۲ نشانه‌گذاری
۴۹	۲-۲-۲ حذف عبارات متعارف: کلمات توقف
۵۰	۳-۲-۲ نرمال‌سازی (دسته‌کردن همارزی عبارات)
۵۵	۴-۲-۲ ریشه‌گیری و مدخل‌گیری
۵۹	۳-۲ اشتراک سریعتر لیست پست‌ها از طریق پرش اشاره‌گرها
۶۲	۴-۲ پست‌های موقعیتی و پرس‌وجوهای اصطلاح
۶۲	۱-۴-۲ شاخص‌های دو کلمه‌ای
۶۴	۲-۴-۲ شاخص‌های موقعیتی
۶۶	۳-۴-۲ طرح‌های ترکیب
۶۹	۵-۲ مراجع و مطالعات آتی

### فصل ۳ لغت‌نامه‌ها و بازیابی مقاوم

۷۳.	۱-۳ ساختار جستجو برای لغت‌نامه
۷۴.	۲-۳ پرس‌وجوهای جایگزین
۷۶.	۱-۲-۳ پرس‌وجوهای جایگزین کلی
۷۷.	۲-۲-۳ شاخص‌های k-گرمی برای پرس‌وجوهای جایگزینی
۷۹.	۳-۳ تصحیح املایی
۸۱.	۱-۳-۳ پیاده‌سازی تصحیح املایی
۸۱.	۲-۳-۳ صورت‌های تصحیح املایی
۸۲.	۳-۳-۳ فاصله ویرایشی
۸۴.	۴-۳-۳ شاخص‌های k-گرمی برای تصحیح املایی
۸۶.	۵-۳-۳ تصحیح املایی حساس به متن
۸۸.	۴-۳ تصحیح آوایی
۸۹.	۵-۳ مراجع و مطالعات آتی

### فصل ۴ ساخت شاخص

۹۱.	۱-۴ مبانی ساخت افزاری
۹۲.	۲-۴ شاخص‌گذاری بلوکی مبتنی بر مرتب سازی
۹۳.	۳-۴ شاخص‌گذاری درون حافظه‌ای تک گذره
۹۷.	۴-۴ شاخص‌گذاری توزیع شده
۹۹.	۵-۴ شاخص‌گذاری پویا
۱۰۳.	۶-۴ انواع دیگر شاخص‌ها
۱۰۵.	۷-۴ مراجع و مطالعات آتی

### فصل ۵ فشرده‌سازی شاخص

۱۱۱.	۱-۵ ویژگی‌های آماری عبارات در بازیابی اطلاعات
۱۱۲.	۱-۱-۵ قانون Heaps: تخمین تعداد عبارات
۱۱۴.	۲-۱-۵ قانون Zipf: مدل‌سازی توزیع عبارات
۱۱۵.	۲-۵ فشرده‌سازی لغت‌نامه
۱۱۷.	۱-۲-۵ لغتنامه به صورت یک رشته

۱۱۹.	۲-۲-۵ ذخیره‌سازی بلوکی
۱۲۲.	۳-۵ فشرده‌سازی فایل پست‌ها
۱۲۲.	۱-۳-۵ کدهای بایت متغیر
۱۲۴.	۲-۳-۵ کدهای ۲
۱۳۳.	۴-۵ مراجع و مطالعات آتی

## فصل ۶ نمره‌گذاری، وزن‌دهی عبارات و مدل فضای بردار

۱۳۵.	۱-۶ شاخص‌های ناحیه‌ای و پارامتری
۱۳۶.	۱-۱-۶ نمره‌گذاری وزن‌دار ناحیه
۱۳۸.	۲-۱-۶ یادگیری وزن‌ها
۱۳۹.	۳-۱-۶ وزن بهینه‌ی $g$
۱۴۱.	۲-۶ وزن‌دهی و فراوانی عبارت
۱۴۳.	۱-۲-۶ فراوانی وارونه سند
۱۴۳.	۲-۲-۶ وزن‌دهی tf-idf
۱۴۵.	۳-۶ مدل فضای بردار برای نمره‌گذاری
۱۴۶.	۱-۳-۶ ضرب نقطه‌ای
۱۴۶.	۲-۳-۶ پرس‌وجو به عنوان بردار
۱۴۹.	۳-۳-۶ محاسبه‌ی نمره‌های بردار
۱۵۱.	۴-۶ توابع گوناگون tf-idf
۱۵۳.	۱-۴-۶ مقیاس‌گذاری زیرخطی tf
۱۵۳.	۲-۴-۶ نرمال‌سازی بیشینه tf
۱۵۴.	۳-۴-۶ طرح‌های وزن‌دهی سند و پرس‌وجو
۱۵۵.	۴-۴-۶ نرمال‌سازی محوری طول سند
۱۵۹.	۵-۶ مراجع و مطالعات آتی

## فصل ۷ محاسبه‌ی نمره‌ها در یک سیستم کامل جستجو

۱۶۱.	۱-۷ رتبه‌بندی و نمره‌گذاری کارآمد
۱۶۱.	۱-۱-۷ بازیابی غیردقیق $K$ سند برتر
۱۶۳.	۲-۱-۷ حذف شاخص
۱۶۴.	۳-۱-۷ لیست‌های قهرمان

۱۶۴.	۴-۱-۷ مرتب‌سازی و نمره‌های کیفیت ایستا
۱۶۶.	۵-۱-۷ مرتب‌سازی برخورد
۱۶۷.	۶-۱-۷ هرس خوش‌های
۱۶۹.	۲-۷ مؤلفه‌های یک سیستم بازیابی اطلاعات
۱۷۰.	۱-۲-۷ شاخص‌های لایه‌ای
۱۷۰.	۲-۲-۷ تقریب عبارت - پرس‌وجو
۱۷۱.	۳-۲-۷ طراحی توابع نمره‌گذاری و تجزیه
۱۷۳.	۴-۲-۷ کنار هم قرار دادن همه بخش‌ها
۱۷۴.	۳-۷ نمره‌گذاری فضای بردار و تعامل عملگر پرس‌وجو
۱۷۶.	۴-۷ مراجع و مطالعات آتی

## الفصل ۸ ارزیابی در بازیابی اطلاعات

۱۷۷.	۱-۸ ارزیابی سیستم بازیابی اطلاعات
۱۷۸.	۲-۸ مجموعه‌های آزمایشی استاندارد
۱۷۹.	۳-۸ ارزیابی مجموعه‌های بازیابی رتبه بندی نشده
۱۸۱.	۴-۸ ارزیابی نتایج بازیابی رتبه بندی شده
۱۸۴.	۵-۸ ارزیابی ربط
۱۹۳.	۱-۵-۸ انتقادات و توجيهات مفهوم ربط
۱۹۵.	۶-۸ یک چشم انداز وسیع‌تر: کیفیت سیستم و مطلوبیت کاربر
۱۹۵.	۱-۶-۸ مسائل سیستمی
۱۹۶.	۲-۶-۸ مطلوبیت کاربر
۱۹۷.	۳-۶-۸ تصحیح سیستم ساخته شده
۱۹۸.	۷-۸ بخش‌هایی از نتایج
۲۰۰.	۸-۸ مراجع و مطالعات آتی

## الفصل ۹ بازخورد ربط و گسترش پرس‌وجو

۲۰۴.	۱-۹ بازخورد ربط و شبه بازخورد ربط
۲۰۷.	۱-۱-۹ الگوریتم Rocchio برای بازخورد ربط
۲۰۹.	۲-۱-۹ بازخورد ربط احتمالاتی
۲۰۹.	۳-۱-۹ چه زمانی بازخورد ربط کار می‌کند؟

۲۱۱.	بازخورد ربط روی وب	۴-۱-۹
۲۱۲.	ارزیابی استراتژی‌های بازخورد ربط	۵-۱-۹
۲۱۳.	شبیه بازخورد ربط	۶-۱-۹
۲۱۴.	بازخورد ربط غیر مستقیم	۷-۱-۹
۲۱۵.	خلاصه	۸-۱-۹
۲۱۶.	روش‌های سراسری برای فرموله کردن پرس‌وجو	۲-۹
۲۱۷.	ابزارهای مجموعه واژگان برای فرموله کردن مجدد پرس‌وجو	۱-۲-۹
۲۱۸.	گسترش پرس‌وجو	۲-۲-۹
۲۱۹.	تولید خودکار فرهنگ لغات جامع	۳-۲-۹
	مراجع و مطالعات آتی	۳-۹

## فصل ۱۰ بازیابی XML

۲۲۱.	مفاهیم پایه‌ای XML	۱-۱۰
۲۲۵.	چالش‌ها در بازیابی XML	۲-۱۰
۲۲۸.	یک مدل فضای بردار برای بازیابی XML	۳-۱۰
۲۳۳.	ارزیابی بازیابی XML	۴-۱۰
۲۳۸.	بازیابی XML متن-محور در مقابل داده-محور	۵-۱۰
۲۴۲.	مراجع و مطالعات آتی	۶-۱۰

## فصل ۱۱ بازیابی اطلاعات احتمالاتی

۲۴۷.	مروری بر نظریه احتمالات پایه	۱-۱۱
۲۴۸.	اصل رتبه‌بندی احتمالاتی	۲-۱۱
۲۴۹.	مورد اتفاف ۰/۱	۱-۲-۱۱
۲۴۹.	اصل رتبه‌بندی احتمالاتی با هزینه‌های بازیابی	۲-۲-۱۱
۲۵۰.	مدل استقلال دودویی	۳-۱۱
۲۵۰.	به دست آوردن تابع رتبه‌بندی برای عبارات پرس‌وجو	۱-۳-۱۱
۲۵۲.	برآوردهای نظری احتمالاتی	۲-۳-۱۱
۲۵۴.	برآورد عملی احتمالاتی	۳-۳-۱۱
۲۵۵.	روش‌های احتمالاتی برای بازخورد ربط	۴-۳-۱۱
۲۵۶.	یک ارزیابی و ارائه تعدادی نسخه	۴-۱۱

۲۵۹	ارزیابی مدل‌های احتمالاتی	۱-۴-۱۱
۲۶۰	ساختار درختی برای وابستگی‌های بین عبارات	۲-۴-۱۱
۲۶۰	Okapi BM25: یک مدل غیردودویی	۳-۴-۱۱
۲۶۳	روش‌های شبکه بیزی برای بازیابی اطلاعات	۴-۴-۱۱
۲۶۴	مراجع و مطالعات آتی	۵-۱۱

## فصل ۱۲ مدل‌های زبانی برای بازیابی اطلاعات

۲۶۵	۱-۱۲ مدل‌های زبانی
۲۶۵	۱-۱-۱۲ آناماتای متناهی و مدل‌های زبانی
۲۶۸	۲-۱-۱۲ انواع مدل‌های زبانی
۲۶۹	۳-۱-۱۲ توزیع چندجمله‌ای روی کلمات
۲۷۰	۲-۱۲ مدل درستنمایی پرس‌وجو
۲۷۰	۱-۲-۱۲ استفاده از مدل‌های زبانی درستنمایی پرس‌وجو در بازیابی اطلاعات
۲۷۲	۲-۲-۱۲ برآورد احتمال تولید پرس‌وجو
۲۷۴	۳-۲-۱۲ آزمایشات Croft و Ponte
۲۷۷	۳-۱۲ مدل‌سازی زبانی در مقابل روش‌های دیگر در بازیابی اطلاعات
۲۷۸	۴-۱۲ نظریه‌های مدل‌سازی زبانی توسعه یافته
۲۸۰	۵-۱۲ مراجع و مطالعات آتی

## فصل ۱۳ دسته‌بندی متن و Naïve Bayes

۲۸۳	۱-۱۳ مسئله دسته‌بندی متن
۲۸۶	۲-۱۳ دسته‌بندی متن
۲۸۸	Naïve Bayes
۲۹۳	۱-۲-۱۳ رابطه با مدل زبانی یک-گرمی چندجمله‌ای
۲۹۴	۳-۱۳ مدل برنولی
۲۹۶	۴-۱۳ ویژگی‌های Naïve Bayes
۳۰۱	۱-۴-۱۳ یک نوع از مدل چندجمله‌ای
۳۰۲	۵-۱۳ انتخاب ویژگی
۳۰۳	۱-۵-۱۳ اطلاعات متقابل
۳۰۵	۲-۵-۱۳ انتخاب ویژگی $\chi^2$
۳۰۸	۳-۵-۱۳ انتخاب ویژگی مبنی بر فراوانی

۳۰۹	۴-۵-۱۳ انتخاب ویژگی برای دسته‌بندهای چندگانه
۳۰۹	۵-۵-۱۳ مقایسه روش‌های انتخاب ویژگی
۳۱۰	۶-۱۳ ارزیابی دسته‌بندی متن
۳۱۸	۷-۱۳ مراجع و مطالعات آتی

۳۲۱	<b>فصل ۱۴</b> دسته‌بندی فضای بردار
۳۲۳	۱-۱۴ نمایش سند و معیارهای وابستگی در فضاهای بردار
۳۲۴	۲-۱۴ دسته‌بندی Rocchio
۳۲۹	۳-۱۴ k-نزدیکترین همسایه
۳۳۲	۱-۳-۱۴ پیچیدگی زمانی و بهینگی k-نزدیکترین همسایه
۳۳۴	۴-۱۴ دسته‌بندهای خطی در برایر دسته‌بندهای غیرخطی
۳۳۹	۵-۱۴ دسته‌بندی با بیشتر از دو دسته
۳۴۱	۶-۱۴ مصالحه بایاس-واریانس
۳۵۰	۷-۱۴ مراجع و مطالعات آتی

۳۵۳	<b>فصل ۱۵</b> ماشین‌های بردار پشتیبان و یادگیری ماشین روی اسناد
۳۵۴	۱-۱۵ ماشین‌های بردار پشتیبان: حالت قابل جداسازی به صورت خطی
۳۶۱	۲-۱۵ توسعه مدل ماشین بردار پشتیبان
۳۶۱	۱-۲-۱۵ دسته‌بندی حاشیه‌ای نرم
۳۶۴	۲-۲-۱۵ دسته‌بندهای ماشین بردار پشتیبان چند دسته‌ای
۳۶۵	۳-۲-۱۵ دسته‌بندهای ماشین بردار پشتیبان غیرخطی
۳۶۸	۴-۲-۱۵ نتایج آزمایشگاهی
۳۶۹	۳-۱۵ مسائل در دسته‌بندی اسناد متنی
۳۶۹	۱-۳-۱۵ انتخاب نوع دسته‌بند برای استفاده
۳۷۱	۲-۳-۱۵ بهبود کارایی دسته‌بند
۳۷۵	۴-۱۵ روش‌های یادگیری ماشین در بازیابی اطلاعات موردنی
۳۷۶	۱-۴-۱۵ مثال ساده‌ای از نمره‌گذاری یادگرفته شده توسط ماشین
۳۷۸	۲-۴-۱۵ رتبه بندی نتایج یادگیری ماشین
۳۸۰	۵-۱۵ مراجع و مطالعات آتی

## فصل ۱۶ خوشبندی مسطح

۳۸۵	۱-۱۶ خوشبندی در بازیابی اطلاعات
۳۸۷	۲-۱۶ بیان مسئله
۳۹۰	۳-۱۶ کاردینالیتی - تعداد خوشها
۳۹۲	۴-۱۶ ارزیابی خوشبندی
۳۹۶	۴-۱۶ الگوریتم K-means
۴۰۱	۱-۴-۱۶ کاردینالیتی خوش در K-means
۴۰۴	۵-۱۶ خوشبندی مبتنی بر مدل
۴۱۱	۶-۱۶ مراجع و مطالعات آتی

## فصل ۱۷ خوشبندی سلسله مرتبی

۴۱۵	۱-۱۷ خوشبندی سلسله مرتبی تجمعی
۴۱۶	۲-۱۷ خوشبندی Complete-Link و Single-Link
۴۲۰	۳-۱۷ پیچیدگی زمانی
۴۲۳	۳-۱۷ خوشبندی تجمعی Group-Average
۴۲۶	۴-۱۷ خوشبندی Centroid
۴۲۸	۵-۱۷ بهینگی خوشبندی سلسله مرتبی تجمعی
۴۳۰	۶-۱۷ خوشبندی تقسیم‌کننده
۴۳۴	۷-۱۷ برچسب‌زنی خوش
۴۳۶	۸-۱۷ نکات پیاده‌سازی
۴۳۹	۹-۱۷ مراجع و مطالعات آتی

## فصل ۱۸ تجزیه‌های ماتریس و شاخص‌گذاری معنایی نهان

۴۴۱	۱-۱۸ مروری بر جبر خطی
۴۴۲	۱-۱-۱۸ تجزیه‌های ماتریس
۴۴۴	۲-۱۸ ماتریس‌های عبارت-سنده و تجزیه‌های مقدار منفرد
۴۴۶	۳-۱۸ تقریب‌های رتبه پایین
۴۴۹	۴-۱۸ شاخص‌گذاری معنایی نهان
۴۵۱	۵-۱۸ مراجع و مطالعات آتی
۴۵۷	

۴۵۹

**فصل ۱۹ مبانی جستجوی وب**

۴۵۹	۱-۱۹ پس زمینه و تاریخچه
۴۶۰	۲-۱۹ مشخصات وب
۴۶۱	۱-۲-۱۹ گراف وب
۴۶۲	۲-۲-۱۹ هرزنامه
۴۶۳	۳-۱۹ تبلیغات به عنوان مدل اقتصادی
۴۶۴	۴-۱۹ تجربه کاربر جستجو
۴۶۵	۱-۴-۱۹ نیازهای مربوط به پرس‌وجوی کاربر
۴۶۶	۵-۱۹ اندازه و برآورد شاخص
۴۶۷	۶-۱۹ استاد دونسخه‌ای و پوشاندن
۴۶۸	۷-۱۹ مراجع و مطالعات آتی

۴۸۳

**فصل ۲۰ پیمایش وب و شاخص‌ها**

۴۸۳	۱-۲۰ مرور کلی
۴۸۴	۱-۱-۲۰ ویژگی‌هایی که یک پیمایشگر موظف است فراهم آورده
۴۸۴	۲-۱-۲۰ ویژگی‌هایی که یک پیمایشگر باید فراهم بیاورد
۴۸۵	۲-۲۰ پیمایش کردن
۴۸۵	۱-۲-۲۰ معماری پیمایشگر
۴۹۰	۲-۲-۲۰ تحلیل DNS
۴۹۱	۳-۲-۲۰ URL آتی
۴۹۵	۳-۲۰ توزیع شاخص‌ها
۴۹۶	۴-۲۰ سرویس‌دهنده‌های اتصال
۴۹۹	۵-۲۰ مراجع و مطالعات آتی

۵۰۱

**فصل ۲۱ تحلیل پیوند**

۵۰۲	۱-۲۱ وب به عنوان یک گراف
۵۰۲	۱-۱-۲۱ لنگر متن و گراف وب
۵۰۴	۲-۲۱ رتبه صفحه
۵۰۶	۱-۲-۲۱ زنجیره‌های مارکوف

۵۰۸-	محاسبه رتبه صفحه	۲-۲-۲۱
۵۱۱-	رتبه صفحه موضوع خاص	۳-۲-۲۱
۵۱۶-	قطبیت و نفوذ	۳-۲۱
۵۱۹-	انتخاب زیرمجموعه‌ای از وب	۱-۳-۲۱
۵۲۲-	مراجع و مطالعات آتی	۴-۲۱
۵۲۵	مراجع	
۵۵۳	لیست اختصارات	
۵۵۵	واژه‌نامه فارسی به انگلیسی	
۵۷۰	واژه‌نامه انگلیسی به فارسی	