

به نام خدا

متن کاوی به کمک
بادگیری ماشین

نویسنده‌ان: چاروی. آگراوال

مترجم:

دکتر مهدی اسماعیلی

متن کاوی به کمک یادگیری ماشین

مترجم: دکتر مهدی اسماعیلی

ناشر: انتشارات آتی نگر

ناشر همکار: وینا

صفحه‌آرایی و طراحی: همتا بیداریان

تیراز: ۵۰۰ نسخه

چاپ اول: ۱۳۹۸

قیمت: ۹۹۰,۰۰۰ ریال

شابک: ۷۸-۶۲۲-۶۱۰۲-۵۹-۹

ISBN: 978-622-6102-59-9

حق چاپ برای انتشارات آتی نگر حفظ است.

نشانی دفتر فروش: خیابان حمامه دبی، رو به روی کوچه رشتچی، پلاک ۱۴۴، واحد ۱

نامبر: ۶۶۵۶۵۳۳۷

تلفن: ۶۶۵۶۵۳۳۶-۸

www.ati-negar.com * info@ati-negar.com



سرشناسه: آگراوال، چارو سی، ۱۹۷۰ م. C. Aggarwal

متن کاوی به کمک یادگیری ماشین / نویسنده چارو سی آگراوال، مترجم: مهدی اسماعیلی

تهران: آتی نگر، وینا ۱۳۹۸

ص: مصور، جدول، نمودار. ۵۹۲

ISBN: 978-622-6102-59-9

و ضعیت فهرست نویسی: فیبا.

یادداشت: عنوان اصلی کتاب: Machine Learning for Text, 2018

یادداشت: کتابنامه.

یادداشت: تمايز.

موضوع: فرآگیری ماشینی -Machine learning - متن پردازی - Data mining - داده کاوی - Artificial intelligence - هوش مصنوعی

شناسه افزوده: اسماعیلی، مهدی، ۱۳۵۰، مترجم

شناسه افزوده: بیداریان، همتا، ۱۳۶۱، گرافیست

ردیبدی کنگره:

ردیبدی دیوبی:

شماره کتابشناسی ملی:

Q۳۲۵/۵

۰۰۶/۳۱

۵۸۱۵۴۴۳

فهرست مطالب

۹
۱۱
۱۱	۱-۱ مقدمه
۱۴	۱-۲ چه جزی درباره یادگیری از متن خاص است؟
۱۶	۱-۳ مدل های تحلیلی برای استناد متنی
۳۱	۴-۱ خلاصه
۳۱	۴-۲ کتاب شناختی
۳۲	۴-۳ تمرین
۳۵
۳۵	۵-۱ مقدمه
۳۶	۵-۲ استخراج متن و تبدیل آن به ته
۴۱	۵-۳ استخراج عبارات از توکن ها
۴۵	۵-۴ نمایش برداری و نرمال سازی
۴۷	۵-۵ محاسبه شباهت در متن
۵۰	۵-۶ خلاصه
۵۱	۵-۷ کتاب شناختی
۵۱	۵-۸ تمرین ها
۵۳
۵۳	۶-۱ مقدمه
۵۷	۶-۲ تجزیه مقدار منفرد
۶۵	۶-۳ تجزیه نامنفی ماتریس
۷۱	۶-۴ تحلیل معنایی نهفته احتمالاتی
۷۷	۶-۵ نگاهی اجمالی به تخصیص نهفته دیربکله
۸۲	۶-۶ تبدیل های غیرخطی و مهندسی ویژگی
۹۷	۶-۷ خلاصه

۹۸.....	۸-۳ کتاب‌شناختی
۹۹.....	۹-۳ تمرین‌ها

۱۰۳.....	۱-۴ مقدمه
۱۰۳.....	۲-۴ انتخاب و مهندسی ویژگی
۱۰۵.....	۳-۴ مدل‌سازی موضوعی و تجزیه ماتریس
۱۱۱.....	۴-۴ مدل‌های مولد آمیخته برای خوشه‌بندی
۱۱۶.....	۵-۱ "سوریتم k-means"
۱۲۲.....	۶-۴ الگوریتم خوشه‌بندی سلسله‌مراتبی
۱۲۵.....	۷-۴ خوشه‌بندی تا یعنی
۱۳۲.....	۸-۴ خوشه‌بندی مترازگار به خوشه‌بندی توالی‌ها
۱۳۴.....	۹-۴ تبدیل خوشه‌بندی یادگاری با ناظر
۱۴۱.....	۱۰-۴ ارزیابی خوشه‌بندی
۱۴۲.....	۱۱-۴ خلاصه
۱۴۸.....	۱۲-۴ کتاب‌شناختی
۱۴۹.....	۱۳-۴ تمرین‌ها
۱۵۰.....	

۱۵۳.....	۱-۵ مقدمه
۱۵۹.....	۲-۵ انتخاب و مهندسی ویژگی
۱۶۳.....	۳-۵ مدل بیز ساده
۱۷۹.....	۴-۵ رده‌بند نزدیک‌ترین همسایه‌ها
۱۸۸.....	۵-۵ درختان تصمیم و جنگل‌های تصادفی
۱۹۵.....	۶-۵ رده‌بندهای مبتنی بر قاعده
۲۰۱.....	۷-۵ خلاصه
۲۰۲.....	۸-۵ کتاب‌شناختی
۲۰۴.....	۹-۵ تمرین‌ها

۲۰۷.....	۱۶ مقدمه
----------	----------

۲۱۴	۶-۶ رگرسیون و رده‌بندی کمترین مربعات.....
۲۲۹	۶-۳ ماشین‌های بردار پشتیبان.....
۲۴۱	۶-۴ رگرسیون لجستیک.....
۲۴۸	۶-۵ تعمیم‌های غیرخطی از مدل‌های خطی
۲۶۰	۶-۶ خلاصه.....
۲۶۰	۶-۷ کتاب‌شناختی.....
۲۶۲	۶-۸ تمرین‌ها.....

۲۶۵	
۲۶۵	۱-۱ مقدمه.....
۲۶۶	۲-۷ موا به بایار د واریانس.....
۲۷۳	۳-۷ اثرات برنه بس و واریانس بر روی کارانی.....
۲۷۷	۴-۷ بهبود کارایی با متعدد: روش‌های تلفیقی
۲۸۱	۵-۷ ارزیابی رده‌بند.....
۲۹۵	۶-۷ خلاصه.....
۲۹۵	۷-۷ کتاب‌شناختی.....
۲۹۷	۸-۷ تمرین‌ها.....

۳۹۹	
۳۹۹	۱-۸ مقدمه.....
۴۰۲	۲-۸ ترفندهای تجزیه ماتریس مشترک
۳۱۶	۳-۸ ماشین‌های تجزیه.....
۳۲۱	۴-۸ تکنیک‌های مدل‌سازی احتمالاتی توأم
۳۲۳	۵-۸ تبدیل به تکنیک‌های گرافکاوی
۳۲۶	۶-۸ خلاصه.....
۳۲۶	۷-۸ کتاب‌شناختی.....
۳۲۸	۸-۸ تمرین‌ها.....

۳۲۹	
۳۲۹	۱-۹ مقدمه.....
۳۳۰	۲-۹ شاخص‌بندی و پردازش پرسش.....

۳۵۵	۳-۹ امتیازدهی با مدل‌های بازیابی اطلاعات
۳۶۳	۴-۹ خوش وب و کشف منع
۳۶۹	۵-۹ پردازش پرسش در موتورهای جستجو
۳۷۴	۶-۹ الگوریتم‌های رتبه‌بندی مبتنی بر لینک
۳۸۳	۷-۹ خلاصه
۳۸۴	۸-۹ کتاب‌شناختی
۳۸۶	۹-۹ تمرین‌ها
۳۸۹	
۳۸۹	۱-۱۰ نقدم
۳۹۲	۲-۱۰ مدل‌های آریان
۳۹۹	۳-۱۰ روش‌های ایرل
۴۰۰	۴-۱۰ مدل‌های تجزیه سازنده
۴۰۵	۵-۱۰ نمایش فواصل واژه‌ها با کم، گاف
۴۰۸	۶-۱۰ مدل‌های عصبی زبان
۴۳۵	۷-۱۰ شبکه‌های عصبی برگشتی
۴۵۳	۸-۱۰ خلاصه
۴۵۴	۹-۱۰ کتاب‌شناختی
۴۵۶	۱۰-۱۰ تمرین‌ها
۴۵۹	
۴۵۹	۱-۱۱ مقدمه
۴۶۲	۲-۱۱ روش‌های مبتنی بر واژه‌های موضوعی برای تشخیص استخراجی
۴۶۸	۳-۱۱ روش‌های نهفته برای تشخیص استخراجی
۴۷۲	۴-۱۱ یادگیری ماشین برای تشخیص استخراجی
۴۷۹	۵-۱۱ تشخیص چندسندی
۴۷۸	۶-۱۱ تشخیص چکیده‌ای
۴۸۱	۷-۱۱ خلاصه
۴۸۱	۸-۱۱ کتاب‌شناختی
۴۸۲	۹-۱۱ تمرین‌ها

۴۸۳	۱- مقدمه
۴۸۳	۲- شناسایی موجودیت نامدار
۴۸۹	۳- استخراج روابط
۵۰۶	۴- خلاصه
۵۱۷	۵- کتاب‌شناسی
۵۱۸	۶- تمرین‌ها
۵۲۰	
۵۲۳	۱- مقدمه
۵۲۳	۲- بهندی سیاست‌در سطح سند
۵۲۹	۳- ردیف اساسی در سطح جمله و عبارت
۵۳۳	۴- نظرکاوی مبتنی بر جمله، به عنوان استخراج اطلاعات
۵۳۶	۵- نظرات هرز
۵۴۱	۶- خلاصه‌سازی نظرات
۵۴۵	۷- خلاصه
۵۴۶	۸- کتاب‌شناسی
۵۴۷	۹- تمرین‌ها
۵۴۹	
۵۵۱	۱- مقدمه
۵۵۱	۲- تقطیع متن
۵۵۲	۳- کاوش جریان‌های متنی
۵۶۱	۴- تشخیص رویداد
۵۶۳	۵- خلاصه
۵۷۰	۶- کتاب‌شناسی
۵۷۱	۷- تمرین‌ها
۵۷۲	
۵۷۳	

ناظر روی تو صاحب نظری نیست که نیست
بوی گیسوی تو در هیچ سری نیست که نیست

پیشگفتار مترجم

در سال‌های اخیر به داشتن ایش متون در محیط‌های نظری و ب، رسانه‌های اجتماعی و دیگر پلتفرم‌ها، متن کاوی از اهمیت ویژه‌ای برخوردار است. در تحلیل متن، تکنیک‌هایی از حوزه‌های دیگر نظری بازیابی اطلاعات، یادگیری ماشین و پردازش ریاضی‌اند. این‌گویی دیده می‌شود، و برای هر یک از آن‌ها نیز کتاب‌های زیادی نوشته شده است. تمکن این کتاب بر روی اندکی از این‌ها یادگیری ماشین برای اسناد متنی است؛ هر چند در برخی از فصل‌ها رنگ روش‌های حوزه‌های دیگر پررنگتر مده است. به زعم مترجم، کتاب حاضر یکی از ارزش‌ترین کتاب‌هایی است که در این حوزه به رشته تحریر دنده است. به همین دلیل از میان کتاب‌های موجود در این حوزه، به انتخاب و ترجمه آن همت گماشتم.

مطلوب کتاب در چهارده فصل گردآوری و تهییه شده است که توان آن را در سه بخش گروه‌بندی کرد. بخش اول یعنی فصل‌های اول تا هشتم کتاب، شامل الگوریتم‌ها و مدل‌های پایه برای تحلیل متن است. تجزیه ماتریس، خوشبندی و رده‌بندی، موضوعات اصلی این فصل‌ها را تجییل می‌نمایند. طبیعی است روش‌های ارائه شده در این بخش، برای کار با متن، تطبیق داده شده‌اند. همان‌طور که قبل از این زمان به آن اشاره شد، تحلیل متن ارتباط نزدیکی با حوزه بازیابی اطلاعات دارد. در بخش دوم که تنها شامل صراحتاً مقاله است، به مروری اجمالی از روش‌های بازیابی اطلاعات از دیدگاه متن کاوی پرداخته شده است. فصل‌های دهم تا چهاردهم کتاب، بخش سوم کتاب را تشکیل می‌دهند. در این بخش موضوعات پیشرفته‌ای مانند یادگیری عمیق، استخراج اطلاعات، خلاصه‌سازی، نظرکاوی، تقطیع متن و تشخیص رویداد بررسی می‌شوند.

تمام تلاش خود را انجام داده‌ایم تا ترجمه کتاب به گونه‌ای انجام شود که خوانندگان محترم آن بتوانند مقاومتی آن را به راحتی درک کنند. بدون شک ممکن است برای برخی از واژه‌های انگلیسی بتوان معادل‌های بهتری یافت. آنچه مسلم است این است که شاید گاهی تبدیل واژه‌ها آنچنان که باید و شاید انجام نشده است؛ اما در ادای جملات و بیان موضوع تلاش فراوان شده است تا خوانندگان گرامی با متنی مبهم و غیج کننده رویه‌رو نشوند.

در اینجا لازم می‌دانم از همه اساتید و دانشجویان به خاطر راهنمایی‌هایی ارزشمندشان در حین آماده‌سازی این کتاب سپاسگزاری کنم. همچنین از مدیریت محترم انتشارات آتی نگر و دوست عزیزم جناب رامین مولاناپور نیز به خاطر آماده‌سازی، چاپ و پخش این کتاب تشکر می‌کنم. رهین محبت بی‌دریغ خانواده‌ام هستم که با فراهم‌سازی محیطی مناسب مرا یاری نمودند. اما با وجود همه سعی و تلاشی که در تمام مراحل آماده‌سازی این کتاب انجام گرفته است، یقین دارم که عاری از اشتباه نیست، چرا که تنها مکتوب بی‌نقص همان معجزه ـ جاوید قرآن کریم است. در آخر ضمن سپاسگزاری از همه کسانی که مرا یاری داده‌اند و با پذیرش مستولیت هرگونه کاستی احتمالی، امیدوارم که این اندک مفید افتند.

مهری اسماعیلی

۹۸ مردادماه