

۱۴۹۹۸۲۶



تحلیل داده با پایتون

www.ketab.ir

تألیف:

وحید قربانی

۱۴۰۳

سر شناه	قربانی ، وحید، ۱۳۶۱
عنوان و پدیدآور	تحلیل داده با پایتون / تألیف: وحید قربانی.
مشخصات نشر	تهران: نشر نور علم، ۱۴۰۲.
مشخصات ظاهري	۳۴۸ ص: مصور (بخشی رنگی)، جدول، نمودار (رنگی).
شابک	۹۷۸-۶۰۰-۶۲۹-۹
موضوع	داده کاوی -
Data mining	
Data structures (Computer science)	ساختار داده ها -
پایتون (زبان برنامه نویسی کامپیوتر) -	
Python (Computer program language)	
ردہ بندی کنگره	QA ۷۶۱۸
ردہ بندی دیوبی	۰۰۶۱۳۱۲

نشر نورعلم و پخش قلم سینا: تهران انقلاب خ ۱۲ فروردین پلاک ۲۸۶ تلفن ۰۹۱۲۳۴۶۲۸۱۱
۰۹۱۲۳۳۴۲۲۹-۰۹۱۲۲۰۷۹۸۴۹

وب سایت <https://www.modiranketab.ir> پیج اینستاگرام - @modiranketab

تحلیل داده با پایتون تألیف: وحید قربانی

ناشر: نور علم
شابک: ۹۷۸-۶۰۰-۶۲۹-۹
چاپ و صحافی: سورنا
نوبت چاپ: اول ۱۴۰۳
شمارگان: ۲۰۰ جلد
قیمت: ۲۹۵۰۰۰ تومان

از طریق تماس با ۰۹۱۲۳۳۴۲۲۹ کتاب ها به تمام نقاط ایران ارسال می شود.

پیشگفتار

تحلیل داده در دنیای امروز از اهمیت بسیاری برخوردار است. تحلیل داده‌ها به افراد و سازمان‌ها کمک می‌کند تا تصمیم‌های بهتری بگیرند. در داده‌کاوی و به طور کلی در علم داده از یک متداول‌وزیری به نام CRISP استفاده می‌شود. CRISP یک روش استاندارد برای تحلیل داده است که در فرآیند استخراج اطلاعات از داده‌ها و مدل‌سازی استفاده می‌شود. این روش از شش مرحله کلیدی تشکیل شده است که به ترتیب شامل شناخت کسب‌وکار، فهم داده‌ها، آماده‌سازی داده‌ها، مدل‌سازی، ارزیابی و نهایتاً استقرار می‌شود.

به طور کلی، روش CRISP به سازمان‌ها در افزایش بهره‌وری از داده‌ها، افزایش قابلیت پیش‌بینی، بهبود فرایند تصمیم‌گیری، کاهش ریسک و هزینه‌ها و افزایش رقابت‌پذیری کمک می‌کند.

با تحلیل داده‌ها و استفاده از روش CRISP می‌توان روندها و الگوهای موجود در داده‌ها را شناسایی کرده و روابط‌های آینده را پیش‌بینی کرد. این امر به سازمان‌ها کمک می‌کند تا استراتژی‌ها و تصمیمات خود را بر اساس اطلاعات قابل اعتماد و مبتنی بر داده و پیش‌بینی‌های دقیق انجام داده و تصمیم‌های بهتری درباره مسائل مختلف مانند بازاریابی، مدیریت مشتریان، بهبود عملکرد و... بگیرند.

با تحلیل داده‌ها و شناسایی الگوها، سازمان‌ها می‌توانند از تصمیماتی که ممکن است باعث خطاها و ضررها مالی شوند، جلوگیری نموده و از این طریق ریسک‌هایی مرتبط با تصمیمات نادرست را کاهش داده و در طولانی‌مدت بهبود عملکرد سازمان را تضمین نمایند. تحلیل داده به سازمان‌ها کمک می‌کند تا بازار رقابتی را بهتر بشناسد و استراتژی‌هایی را برای افزایش بازدهی، بهبود رقابت‌پذیری و جایگاه خود در بازار انتخاب کنند. هم چنین به سازمان‌ها این امکان را می‌دهد تا از داده‌های بزرگ و پیچیده بهره‌برداری کنند، روابط مخفی را شناسایی و اطلاعات ارزشمندی را استخراج نمایند که همین امر موجب اتخاذ تصمیمات استراتژیک و اجرایی بهتر سیاست‌ها خواهد شد.

زبان برنامه‌نویسی پایتون به عنوان یکی از زبان‌های برنامه‌نویسی محبوب و قدرتمند در زمینه علوم داده شناخته می‌شود. به دلایل زیادی از جمله سادگی، فهم و خوانایی بالا، دارا بودن بستر وسیع، داشتن کتابخانه‌ها و ابزارهای قدرتمندی که به طور خاص برای تحلیل داده‌ها و علوم داده

طراحی شده و توسعه یافته‌اند، جامعیت و قابل اجرا بودن بر روی اکثر سیستم عامل‌ها، وجود منابع آموزشی غنی و پشتیبانی از برنامه‌نویسان و توسعه‌دهندگان و سازگاری و هماهنگی با سایر ابزارهای مورد استفاده در حوزه علوم داده آن را به یک ابزار بسیار مناسب برای تحلیل داده تبدیل نموده که در حال حاضر توسط افراد و سازمان‌های زیادی مورد استفاده قرار می‌گیرد.

با مطالعه این کتاب خواهید آموخت که چگونه با به کارگیری زبان برنامه‌نویسی قدرتمند پایتون و با استفاده از چهار کتابخانه پرکاربرد آن یعنی Seaborn و Matplotlib، Pandas و NumPy به شناخت و فهمی از داده‌ها برسید، آن‌ها را پاکسازی، تبدیل، بصری‌سازی، و تحلیل نمایید.

در طول این کتاب، با مفاهیم اصلی تحلیل داده آشنا خواهید شد و با استفاده از متدها و پارامترهای موجود در این چهار کتابخانه، می‌توانید داده‌های خود را تفسیر کنید و نتایج دقیق و قابل اعتمادی را به دست آورید. هم‌چنین، نمونه‌های عملی و تمرین‌هایی ارائه شده است تا بتوانید مهارت‌های خود را تقویت نموده و در طول مسیر یادگیری، تجربه‌های جدیدی کسب نمایید.

کتاب پیش‌رو شامل مطالبی درباره متداول‌ترین CRISP در علم داده، کتابخانه NumPy، کتابخانه Matplotlib و کتابخانه Pandas است.

در مقدمه گام‌های متند CRISP به تفصیل آمدند. در فصل اول کتاب، به کتابخانه NumPy پرداخته شده است که امکان انجام عملیات جبری بر روی آرایه‌ها، توابع ریاضی و عملیات ضرب ماتریس‌ها را فراهم می‌کند.

در فصل دوم، کتابخانه Pandas معرفی می‌گردد که برای شناخت داده‌ها، آماده سازی داده‌ها و عملیات مختلف بر روی داده‌ها از این کتابخانه استفاده می‌شود. شما با ساختارهای داده مانند سریز و دیتاframes در Pandas آشنا می‌شوید و یاد می‌گیرید چگونه داده‌ها را از منابع مختلف مانند SQL و فایل‌های اکسل بخوانید و با آن‌ها کار کنید.

در فصل سوم، کتابخانه Matplotlib معرفی می‌شود که برای ترسیم و نمایش داده‌ها در قالب نمودارها از آن استفاده می‌شود. شما با انواع نمودارها مانند histplot، scatterplot و countplot آشنا می‌شوید و یاد می‌گیرید چگونه داده‌ها را در قالب پلات نشان دهید و شکل و رنگ پلات‌ها را تنظیم نمایید. در فصل چهارم، کتابخانه Seaborn معرفی می‌شود که برای ترسیم داده‌ها بر

پایه کتابخانه Matplotlib است. این کتابخانه امکان ترسیم انواع نمودارها مانند `rugplot` و `heatmaps` و `displot` را فراهم می‌کند. امیدوارم این کتاب به عنوان یک راهنمای جامع در مبحث تحلیل داده با پایتون برای خوانندگان، مفید واقع شده و به افزایش دانش و مهارت آنها در این حوزه کمک کند.

وحید قربانی

-  VahidGhorbani@hotmail.com
-  <http://www.linkedin.com/in/VahidGhorbani>
-  <https://www.youtube.com/@DataLand01>

فهرست

۱۳	مقدمه
۱۳	متدولوژی CRISP در علم داده
۱۳	گام‌ها در متا CRISP
۲۱	فصل ۱: کتابخانه NumPy
۳۰	توزيع نرمال
۳۸	ایندکس گذاری ماتریس‌ها
۴۱	عملیات جبری روی آرایه‌ها
۴۲	توابع ریاضی و آرایه‌ها
۴۴	ضرب ماتریس‌ها
۴۶	ضرب دو آرایه دو بعدی
۴۸	ضرب خارجی بردارها
۵۳	فصل ۲: کتابخانه Pandas
۵۳	شناخت داده‌ها و آماده‌سازی داده‌ها در Pandas
۵۳	درباره کتابخانه Pandas
۵۴	انواع داده ساختار در کتابخانه Pandas
۵۴	دستور وارد شدن به کتابخانه Pandas
۵۵	داده ساختار سریز
۵۶	تفاوت و شباهت لیست و سریز
۵۷	تفاوت سریزها با آرایه‌های numpy
۵۸	تغییر ایندکس در Series
۵۸	مرتب کردن داده‌ها

۶۶	اضافه کردن عنصر به سریز
۶۷	ضرب و تقسیم مقادیر سریز
۷۰	داده ساختار دیتافریم
۷۲	اتصال دیتافریم‌ها (متد concat)
۷۵	اتصال دیتافریم‌ها با چند ستون
۸۲	انواع Merge
۸۷	گام‌های شناخت داده
۹۱	خواندن داده از منابع داده‌ای مختلف
۹۷	دسترسی به یک ستون در دیتافریم
۱۰۰	خواندن داده‌ها از SQL
۱۰۱	اتصال به SQL با استفاده از کتابخانه SQLAlchemy
۱۰۲	کار با فایل‌های اکسل با استفاده از کتابخانه openpyxl
۱۰۳	خواندن داده‌ها از فایل اکسل با متدهای ExcelFile
۱۰۴	ایجاد فایل اکسل با متدهای ExcelWriter
۱۰۵	ذخیره داده‌ها در فایل اکسل: متدهای to_excel
۱۰۵	مشاهده سطرهای ابتدایی: متدهای head
۱۰۷	مشاهده سطرهای انتهایی: متدهای tail
۱۱۲	تنظیم ایندکس روی ستون دلخواه با متدهای set_index
۱۱۳	بازگردانی تنظیمات ایندکس به حالت اولیه با متدهای reset_index
۱۱۵	پارامتر inplace
۱۱۵	خواندن داده از صفحه‌های وب
۱۱۹	خواندن داده‌ها از فایل جی‌اسون با متدهای read_json

۱۲۰	متدهای <code>to_json</code>
۱۲۰	دسترسی به سطرها در دیتا فریم (متدهای <code>loc</code>)
۱۲۱	دسترسی به سطر و ستون در دیتا فریم (متدهای <code>iloc</code>)
۱۲۴	متدهای <code>sort_index</code>
۱۲۷	داده های مفقود
۱۳۰	محاسبات آماری (متدهای <code>mode</code> , <code>mean</code> , <code>min</code> , <code>max</code> و غیره)
۱۳۰	متدهای <code>mean</code>
۱۳۷	شمارش تعداد مقادیر معلوم (متدهای <code>count</code>)
۱۳۷	محاسبه فراوانی: متدهای <code>value_counts</code>
۱۴۳	تغییر نوع داده ها با متدهای <code>astype</code>
۱۴۴	مرتب کردن داده ها با متدهای <code>sort_values</code>
۱۴۵	فیلتر کردن داده ها
۱۴۵	شرط ها
۱۵۱	متدهای <code>between</code>
۱۵۵	فیلتر کردن داده ها در دیتا فریم (متدهای <code>isin</code>)
۱۵۷	متدهای <code>where</code>
۱۵۸	متدهای <code>query</code>
۱۶۱	شمارش تعداد مقادیر <code>Nan</code> (متدهای <code>isna()</code> , <code>sum</code>)
۱۶۳	حذف سطر و ستون (متدهای <code>drop</code>)
۱۶۶	انواع تحلیل داده
۱۶۷	دسته بندی داده ها (متدهای <code>groupby</code>)
۱۷۷	فیلتر کردن و دسته بندی کردن داده ها

۱۷۹	تعیین متاد برای دسته‌بندی داده‌ها (متاد <code>agg</code>)
۱۸۱	مدل‌سازی
۱۸۶	تغییردادن نام ستون‌ها با متاد <code>rename</code>
۱۸۷	اضافه کردن ستون با مقدار یکسان به دیتا فریم
۱۸۷	اضافه کردن یک سطر به دیتا فریم با متاد <code>_append</code>
۱۸۹	اضافه کردن ستون در محل لخواه به دیتا فریم با متاد <code>insert</code>
۱۹۳	فیلتر کردن سطرهای دارای مقادیر مفقود
۱۹۵	فیلتر کردن سطرهای دارای مقدار معلوم
۱۹۶	شناسایی سطرهای تکراری در دیتا فریم متاد <code>duplicated</code>
۱۹۷	حذف مقادیر تکراری در یک ستون با متاد <code>duplicated</code>
۱۹۷	حذف سطرهای تکراری با متاد <code>drop_duplicates</code>
۱۹۹	شناسایی مقادیر یکتای یک ستون (متاد <code>unique</code>)
۱۹۹	شناسایی تعداد مقادیر یکتای یک ستون (متاد <code>nunique</code>)
۲۰۰	حذف و جایگزینی داده‌های مفقود
۲۰۳	حذف همه سطرهای دارای مقادیر مفقود
۲۱۲	جایگزین کردن مقادیر با متاد <code>map</code>
۲۱۳	متاد <code>interpolate</code>
۲۱۵	درون‌یابی به روش <code>polynomial</code>
۲۲۰	روش درون‌یابی <code>spline</code>
۲۲۰	تابع
۲۲۸	جایگزین کردن مقادیر با متاد <code>replace</code>
۲۳۱	متاد <code>nlargest</code>

۲۲۲ nsallest متدها
۲۲۳ sample نمونه گیری دادهها با متدها
۲۲۵ Matplotlib فصل ۳: کتابخانه
۲۲۶ ترسیم و نمایش نمودار
۲۲۶ histplot
۲۴۱ محدود کردن پلاتها
۲۴۴ countplot
۲۴۵ scatterplot
۲۴۸ subplots
۲۵۰ pairplot
۲۵۱ regplot
۲۵۴ plot
۲۸۸ تنظیم شکل و رنگ پلات
۳۰۲ Seaborn فصل ۴: کتابخانه
۳۱۳ rugplot
۳۱۵ displot
۳۱۷ histplot
۳۲۳ countplot
۳۲۷ barplot
۳۳۰ jointplot
۳۳۳ pairplot
۳۴۰ heatmap

۳۴۶	clustermapper
۳۵۱	پیوست
۳۵۳	منابع و مأخذ

www.ketab.ir

مقدمه

متدولوژی CRISP^۱ در علم داده

می خواهیم با متodi به نام CRISP آشنا شویم. ابتدا علم داده را تعریف می کنیم.

علم داده^۲ چیست؟

در حوزه های مختلف علم و یا در کسب و کار سوال ها و چالش هایی به وجود می آید که خبرگان و صاحبان آن حوزه و کسب و کار نمی توانند به آن پاسخ دهند. علم داده با استفاده از داده های موجود در همان حوزه علمی و کسب و کار به این سوال ها پاسخ می دهد.

برای اینکه پروژه های علم داده نظم پیدا کند از یک متدولوژی به نام CRISP استفاده می کنیم.

برای پروژه های داده کاوی نیز از همین روش استفاده می نماییم.

گام ها در متod CRISP

گام اول: شناخت کسب و کار ^۳ در گام اول لازم است فهم و درک درستی نسبت به آن حوزه و کسب و کار پیدا کنیم. از آن جایی که کسب و کارها بیعیدگی های زیادی دارند معمولاً شناخت درست و دقیق آن ها مستلزم صرف زمان زیادی است: بوابی حل این مشکل می توانیم از خبرگان و صاحبان آن حوزه و کسب و کار کمک بگیریم. این افراد به عنوان مشاور^۴ در کنار دانشمند داده^۵ حضور دارند و به او کمک می کنند تا بتواند مسائل را درست متوجه شود، مسائل درستی را مطرح کند و برای رسیدن به راه حل آن مسائل مسیر درستی نیز انتخاب نماید.

گام دوم: شناخت داده ها^۶: در این گام تلاش می کنیم تا فهم و درک درستی نسبت به داده ها پیدا کنیم. یعنی اگر یک مجموعه داده^۷ در اختیار داشته باشیم، لازم است همه ستون های آن را بشناسیم. برای مثال در مجموعه داده های کشتی تایتانیک ستون Survived یعنی نجات یافتن یا غرق شدن مسافر و ستون Embarked یعنی مسافر در کدام ایستگاه سوار کشته شده است.

¹ The CRoss Industry Standard Process for Data Mining (CRISP-DM)

² Data science

³ Business Understanding

⁴ Business Experts

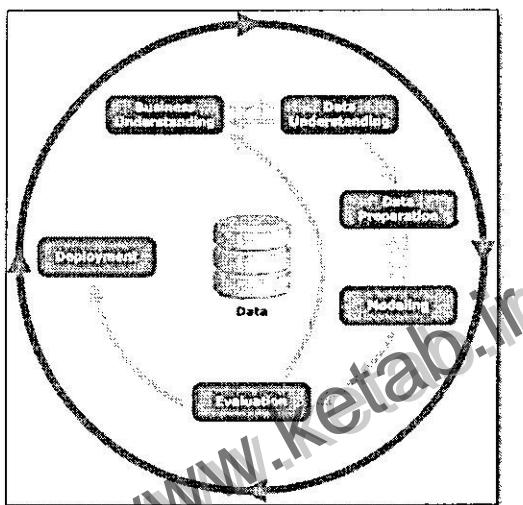
⁵ Data Scientist

⁶ Data Understanding

⁷ Dataset

از طرفی باید نسبت به هر یک از ستون‌ها در مجموعه داده‌ها شناخت داشته باشیم که چه تعداد داده مفقود^۱، داده دور افتاده^۲ و Noise Data وجود دارد.

به عبارتی در مرحله شناخت داده‌ها باید یک دید آماری نسبت به کل داده‌ها پیدا کنیم، نوع داده‌ای آن‌ها را بدانیم، داده‌هایی که نیاز به تغییر و تبدیل دارند و داده‌های Valid را شناسایی کنیم.



تصویر ۱-۱. گام‌ها در متداول‌ترین CRISP

برای اینکه بتوانیم به این شناخت از مجموعه داده‌ها برسیم، بهترین روش، استفاده از کتابخانه‌های پایتون مانند Seaborn، Pandas و Matplotlib است. (کتابخانه Seaborn توسط کتابخانه Matplotlib توسعه پیدا کرده است).

یکی از مواردی که در شناخت داده‌ها به ما کمک می‌کند بصری‌سازی داده‌ها^۳ است. اگر دیتا را بصری کنیم می‌توانیم در کمترین زمان ممکن به چندین مورد در شناخت داده برسیم از جمله نرمال بودن داده‌ها یعنی آیا توزیع نرمال در این داده‌ها وجود دارد؟ با ترسیم یک نمودار به نام

¹ Missing data

² Outlier data

³ Data Visualization

در کتابخانه **Histogram Plot** از **Matplotlib** می‌توانیم به این هدف برسیم. پس برای بصری‌سازی داده‌ها از کتابخانه **Matplotlib** استفاده می‌کنیم.

برای اینکه نوع داده‌ها را بشناسیم، گزارش‌های آماری در مورد داده‌ها داشته باشیم، داده‌های دور افتاده، داده‌های مفقود، **Noise data** و داده‌هایی که **Valid** نیستند را پیدا کنیم از کتابخانه **Pandas** استفاده می‌کنیم.

ممکن است گاهی یک رفت و برگشت بین شناخت داده‌ها و شناخت کسب و کار انجام شود. برای مثال ممکن است مساله‌ای مطرح کنیم ولی داده‌ای برای بررسی آن مساله وجود نداشته باشد. در چنین شرایطی باید مساله را تغییر دهیم یا با اقلام اطلاعاتی دیگری سعی کنیم مساله را حل نماییم. پس ممکن است بارها و بارها رفت و برگشت بین داده‌ها و کسب و کار اتفاق بیافتد.

گام سوم: آماده‌سازی داده^۱: در این مرحله داده‌ها را تغییر می‌دهیم. وقتی در مرحله شناخت داده‌ها، برای مثال متوجه می‌شویم چه تعداد داده گمشده داریم، در مرحله آماده‌سازی داده‌ها تصمیم می‌گیریم که به جای داده‌هایی که مقدارشان وجود ندارد چه مقداری قرار دهیم، این داده‌ها را چگونه جایگزین و آماده کنیم. همچنین اگر بخواهیم یک ستون را به دو ستون یا بر عکس دو ستون را به یک ستون تبدیل کنیم یا نوع داده‌ای یک ستون را تغییر دهیم مثلاً مقادیر رشته‌ای را به عددی تبدیل کنیم همه این موارد را در مرحله آماده‌سازی داده‌ها و با استفاده از کتابخانه **Pandas** انجام می‌دهیم.

حدود ۶۰ تا ۷۰ درصد وقت یک دانشمند داده صرف دو مرحله شناخت داده و آماده‌سازی داده می‌شود. پس این دو مرحله خیلی مهم هستند.

گام چهارم: مدل‌سازی^۲: در این مرحله از مجموعه الگوریتم‌های یادگیری ماشین^۳ مانند **Sklearn** استفاده می‌کنیم. برای این کار الگوریتم‌های زیادی وجود دارند که بینه شده و وارد مرحله یادگیری عمیق شده‌اند.

¹ Data Preparation

² Modeling

³ Machine Learning

در این مرحله نیز بین آماده‌سازی داده‌ها و مدل‌سازی دادها رفت و برگشت وجود دارد. برای مثال وقتی که داریم الگوریتم را پیاده‌سازی می‌کنیم متوجه می‌شویم یک ویژگی کم داریم، باید برگردیم به مرحله آماده‌سازی داده، داده‌ها را آماده کنیم. یا وقتی مدل در حال ساخته شدن است با پیغام خطای مواجه می‌شویم و متوجه می‌شویم که نوع داده‌ای یک ستون را باید تغییر دهیم، بعضی الگوریتم‌ها مانند الگوریتم خوشبندی^۱ مقداری Text نمی‌پذیرند. هنر یک دانشمند داده انتخاب الگوریتم مناسب برای مدل‌سازی داده و انتخاب پارامترهای مناسب است. مرحله مدل‌سازی محبوب‌ترین مرحله است و حدود ۱۰ تا ۱۵ درصد وقت دانشمند داده را می‌گیرد.

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1,0,3,"Braund, Mr. Owen Harris",male,jfhsdjfgdshfgds,1,0,A/5 21171,7.25,,S											
2,1,1,"Cumings, Mrs. John Bradley (Florence Briggs Thayer)",female,38,1,0,PC 17599,71.2833,C85,C											
3,1,3,"Heikkinen, Miss. Laina",female,26,0,0,STON/O2. 3101282,7.925,,S											
4,1,1,"Futrelle, Mrs. Jacques Heath (Lily May Peel)",female,35,1,0,113803.53,1,C123,S											
5,0,3,"Allen, Mr. William Henry",male,35,0,0,373450,8.05,,S											
6,0,3,"Moran, Mr. James",,,0,0,330877,8.4583,,Q											
7,0,1,"McCarthy, Mr. Timothy J",male,54,0,0,17463,51.3625,E46,S											
8,0,3,"Palsson, Master. Gosta Leonard",male,2,3,1,349909,21.075,,S											
9,1,3,"Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)",female,27,0,2,347742,11.1333,S											
10,1,2,"Nasser, Mrs. Nicholas (Adele Achem)",female,14,1,0,237736,30.0708,,C											
11,1,3,"Sandstrom, Miss. Marguerite Rut",female,4,1,1,PP 9549,16.7,G6,S											
12,1,1,"Bonnell, Miss. Elizabeth",female,58,0,0,113783,26.55,C103,S											
13,0,3,"Saunderscock, Mr. William Henry",male,20,0,0,A/5. 2151,8.05,,S											
14,0,3,"Andersson, Mr. Anders Johan",male,39,1,5,347082,31.275,,S											
15,0,3,"Vestrom, Miss. Hulda Amanda Adolfsina",female,14,0,0,350406,7.8542,,S											
16,1,2,"Hewlett, Mrs. (Mary D Kingcome)",female,55,0,0,248706,16,,S											
17,0,3,"Rice, Master. Eugene",male,2,4,1,382652,29.125,,Q											
18,1,2,"Williams, Mr. Charles Eugene",male,,0,0,244373,13,,S											
19,0,3,"Vander Planke, Mrs. Julius (Emelia Maria Vandemoortele)",female,31,1,0,345763,18,,S											
20,1,3,"Masselmani, Mrs. Fatima",female,,0,0,2649,7.225,,C											
21,0,2,"Fynney, Mr. Joseph J",male,35,0,0,239865,26,,S											
22,1,2,"Beesley, Mr. Lawrence",male,34,0,0,248698,13,D56,S											
23,1,3,"McGowan, Miss. Anna ""Annie""",female,15,0,0,330923,8.0292,,Q											
24,1,1,"Sloper, Mr. William Thompson",male,28,0,0,113788,35.5,A6,S											
25,0,3,"Palsson, Miss. Torborg Danira",female,8,3,1,349909,21.075,,S											
26,1,3,"Asplund, Mrs. Carl Oscar (Selma Augusta Emilia Johansson)",female,38,1,5,347077,31.3875,S											

¹ Clustering

گام پنجم: درستی مدل^۱: بعد از اینکه مدل ساخته شد آن را مورد تحلیل و بررسی قرار می‌دهیم و مقدار صحت یا درستی مدل را به دست می‌آوریم. معمولاً در مجموعه داده‌هایی که انتخاب می‌کنیم اگر داده‌ها درست باشند، الگوریتمی که برای مدل‌سازی انتخاب کردیم درست باشد و پارامترهای مناسبی انتخاب کرده باشیم، میزان درستی مدل بالای ۸۰ درصد می‌شود. در دنیای واقعی احتمال وقوع هر پدیده‌ای ۵۰ درصد است. پس اگر درستی مدل را ۵۰ درصد به دست آوریم در واقع هیچ کاری انجام نداده‌ایم. از طرفی درستی مدل ۱۰۰ درصد نخواهد شد مگر اینکه تعداد جامعه آماری و نمونه‌ها کم باشد و ورودی‌ها شباهت خیلی زیادی به جامعه آماری داشته باشند. در چنین شرایطی می‌توانیم نزدیک به ۱۰۰ درصد یا دقیقاً ۱۰۰ درصد درستی مدل داشته باشیم. در مجموعه داده‌هایی که تعداد موردها در مدل‌سازی زیاد باشد به دست آوردن درستی مدل کار پیچیده‌ای می‌شود. گاهی ممکن است آن قدر میزان درستی مدل پایین باشد که مجبور شویم به مرحله اول یعنی شناخت کسب و کار برگردیم و مراحل را دوباره تکرار کنیم. به همین دلیل است که مدیران علاقه‌ای به پژوهش‌های حوزه داده‌کاوی، یادگیری ماشین و علم داده ندارند. اما اگر پژوهه جواب دهد آنگاه وارد مرحله بعد می‌شویم.

نرم افزارهایی که تولید می‌کنیم و استفاده می‌کنیم را سامانه‌های توصیه گر یا DSS می‌نامیم. اعضای تیم در متدهای CRISP: یک مثلث را تشکیل می‌دهند که اخلاع آن متخصص کسب و کار^۲، متخصص فنی^۳ و Statistics Experts هستند. وقتی یک مجموعه داده داریم و می‌خواهیم روی آن کار تحلیلی انجام دهیم، در درجه اول باید یک شناختی روی آن پیدا کنیم. پس به یک Data Dictionary نیاز داریم. یعنی ضروری است یک سری اطلاعات راجع به مجموعه داده‌ها و هر یک از ستون‌ها کسب کنیم. Data Dictionary را می‌توانیم از وب سایتهای مانند GitHub، Kaggle و غیره نیز جستجو کنیم.

¹ Evaluation

² Business Experts

³ Technical Experts

در وب سایت Kaggle گفته شده است که هر مجموعه داده برای چه کاری مناسب است. در این وب سایت کدهای پایتون را نیز می‌توانیم ببینیم.

Dataset Information

Additional Information

The data is related with direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls. Often, more than one contact to the same client was required, in order to access if the product (bank term deposit) would be ('yes') or not ('no') subscribed.

There are four datasets:

- 1) bank-additional-full.csv with all examples (41188) and 20 inputs, ordered by date (from May 2008 to November 2010), very close to the data analyzed in [Moro et al., 2014]
- 2) bank-additional.csv with 10% of the examples (4119), randomly selected from 1), and 20 inputs.
- 3) bank-full.csv with all examples and 17 inputs, ordered by date (older version of this dataset with less inputs).
- 4) bank.csv with 10% of the examples and 17 inputs, randomly selected from 3) (older version of this dataset with less inputs). The smallest datasets are provided to test more computationally demanding machine learning algorithms (e.g., SVM).

The classification goal is to predict if the client will subscribe (yes/no) a term deposit (variable y).

SHOW LESS ^

Has Missing Values?

No

تصویر ۲-۱

از وب سایت UCI Machine Learning Repository اطلاعاتی در مورد مجموعه داده‌های بازاریابی بانکی^۱ به دست می‌آوریم.

در این وب سایت به جای عبارت Dataset Information از عبارت Data Dictionary استفاده می‌شود. با مطالعه توضیحات مربوط به این مجموعه داده متوجه می‌شویم که این داده‌ها مربوط به کمپین بازاریابی مستقیم است. از طریق شماره تماس یا به صورت حضوری افراد را مقنused کنند که از آنها وام بگیرند. نتایج کمپین‌های گذشته را نیز به ما نشان می‌دهند. در مورد این کمپین برای مثال گفته شده ۴۱۰۰۰ نفر از افراد جامعه نمونه شده‌اند. این ۴۱ هزار داده، اطلاعات افراد با ویژگی‌های مختلف است. همین‌طور نشان می‌دهد داده‌هایی که ما در اختیار داریم ۱۰ درصد از کل داده‌ها است. این مجموعه داده ۱۷ ورودی و یک خروجی دارد.

از این منبع داده‌ای با هدف Classification می‌توانیم استفاده کنیم. برای هر کدام از ستون‌ها توضیحاتی نوشته شده است که به بعضی از مهم‌ترین آن‌ها اشاره می‌نماییم:

¹ Bank Marketing

سن افراد (Age)

شغل (Job): مدیر (admin)، کارآفرین (entrepreneur)، خانهدار (housemaid)، بازنشسته (retired)، کارگر (blue-collar)، تکنسین (technician)، بیکار (unemployed)، برای خودش کار می‌کند (self-employed)

توضیحات مربوط به وضعیت تأهل (marital): جدا شده (divorced)، متاهل (married)، مجرد (single)

توضیحات مربوط به تحصیلات (education): تحصیلات مقدماتی (high school)، بی سواد (illiterate)، تحصیلات دانشگاهی (professional)، تحصیلات حرفه‌ای (university)

Variables Table				
Variable Name	Role	Type	Demographic	Description
age	Feature	Integer	Age	
job	Feature	Categorical	Occupation	type of job [categorical]: 'admin'; 'Blue-collar'; 'entrepreneur'; 'housemaid'; 'manually employed'; 'services'; 'student'; 'technician'; 'unemployed'; 'unknown'
marital	Feature	Categorical	Marital Status	Marital status [ategorical]: 'divorced'; 'married'; 'single'; 'unknown'; note: 'divorced' > 'married'
education	Feature	Categorical	Education Level	[categorical]: 'basic.4y'; 'basic.6y'; 'high.school'; 'illiterate'; 'professional.course'; 'university'
default	Feature	Binary		has credit in default?
balance	Feature	Integer		average yearly balance
housing	Feature	Binary		has housing loan?
loan	Feature	Binary		has personal loan?
contact	Feature	Categorical		contact communication type [categorical]: 'cellular'; 'telephone'
day_of_week	Feature	Date		last contact day of the week

تصویر ۱-۳

Default: یعنی برای ما مسجّل شده که این شخص وام می‌گیرد. اعتبار سنجی کردیم و مطمئنیم که این شخص می‌تواند وام بگیرد. این ستون از نوع Binary یعنی true و false یا صفر و یک است؛

Balance: میانگین درآمد سالانه بر مبنای euros

Housing: مشخص می‌کند این شخص برای خانه وام گرفته است یا خیر؛
نوع برقراری ارتباط با مشتری (Contact): تلفن (telephone)، cellular، (cellular)، Duration: یعنی آخرین تماسی که با مشتری داشتیم چند ثانیه صحبت کردیم؛
Pdays: یعنی از آخرین تماس با این شخص چند روز گذشته است؛
Previous: چند بار کلا با این مشتری تماس گرفتیم؛
. (nonexistent) Putcome: پاسخ دادن به تماس (success)، پاسخ نداده است (failure).