

شماره اول- آذرماه ۱۴۰۰- سال اول
نشریه علمی دانشجویی انجمن علمی آمار دانشگاه اراک



در این شماره می خوانیم:

داده هارا از کجا پیدا کنیم؟

مقالات

فیشر

نرم افزار R

سرگرمی به زبان ۱۹۰

سخن سردبیر



دفتر را با نام او می‌گشایم که هر امر مبهمی بی‌یاد او
بی‌حاصل است.

شماره اول نشریه انيگما را در آذر ماه ۱۴۰۰ کلید
می‌زنیم این نشریه زیر نظر انجمن علمی آمار دانشگاه
اراک منتشر می‌شود. به همین بهانه چند سطrix را با
شما مخاطبان نشریه در رابطه با انتخاب نام این مجله
در میان می‌گذارم.

همه‌ی ما با گستره علم آمار در علوم مختلف آشنا
هستیم که بیان آن تکرار مکرات است. هر آماردان،
آمارخوان یا تحلیل‌گرداده امروزی نیازمند آشنایی با
علوم کامپیوتري به روز است چرا که در عصر
اطلاعات و ارتباطات دیجیتال این عنوان‌ها
را کسب می‌کند. به همین علت گذري
بر گذشته برنامه‌نويسی، الگوريتم‌نويسی
و هوش مصنوعی زديم و نام انيگما را
برگزيريم.

انيگما اولين دستگاه رمزگذاري و کدنگاري
بود که از هوش مصنوعی برای محافظت از
ارتباطات تجاری، دипلماتیک و نظامی به
طور گسترده در جنگ جهانی دوم توسط
ارتش آلمان نازی به کار گرفته می‌شد.

آلن تورینگ، پدر علم کامپیوت و هوش مصنوعی،
به کمک ماشین تورینگ، فرمولاسیون موثری برای
روش‌های الگوريتم‌نويسی و محاسبه تهیه کرد که موفق
به شکست کدنگاری‌های ماشین انيگما شد. به همین
جهت نام انيگما را برای این نشریه انتخاب کردیم تا
اهمیت علم آمار و نقش داده‌ها را در هوش مصنوعی
شفاف‌تر کنیم.

اميده است قدم کوچکی در جهت گسترش و معرفی
علم آمار برداریم.

یکی از عمدۀ مشکلاتی که نیازمندان تحلیل‌های آماری با آن سرو کار دارند یافتن یک مجموعه داده معتبر برای انجام تحقیقات اولیه و گزارش نویسی دانشگاهی است. در این بخش به معرفی وبسایت‌هایی می‌پردازیم که مجموعه داده‌های شناسه داری را رایگان در اختیار همگان قرار می‌دهند:

سایت مرکز ملی آمار ایران

در این سایت انواع داده‌های سرشماری‌های کشوری، مطالعات و تحقیقاتی که زیر نظر مرکز ملی آمار ایران انجام شده است در اختیار همگان قرار دارد و منبع اصلی دسترسی به داده‌های ملی است.

[data.world](http://www.data.world)

در Data world داده‌ها را کشف کنید و به اشتراک بگذارید. با مردم ارتباط برقراری کنید و برای حل مشکلات با یکدیگر همکاری کنید! این سایت داده‌های واقعی با قابلیت تجزیه و تحلیل و استفاده مجدد و قابل توسعه را در اختیار کاربران خود قرار می‌دهد. با یک کلیک موضوع خود را پیدا و داده دریافت و تحلیل کنید.

kaggle

اگر علاقه مند به حوزه پردازش و تحلیل داده هستید سری به مجموعه داده‌ای این وبسایت بزنید. این سایت مجموعه داده‌هایی را در اختیار کاربران می‌گذارد و مسائلی در حوزه تحلیل داده راجع به آن داده‌ها را مطرح می‌کند.

به علاوه با برگزاری مسابقات پردازش داده‌ها کمی چاشنی رقابت بین تحلیلگران هم اضافه کرده است.

هر تیم شرکت کننده در بازه زمانی معین باید برنامه یا الگوریتم بهینه مورد نظر را برای رسیدن به هدف تعیین شده ارایه کند. و گاه‌ها جایزه‌های چندین دلاری و استخدام تیم‌های برتر را به دنبال دارد.

البته این سایت برای پروکسی‌های ایران در دسترس نیست و نیازمند فیلترشکن می‌باشد.

machine learning uci repository

این سایت همانطور که از نامش پیداست یک مخزن دیتابست برای یادگیری ماشین می‌باشد که داده را رایگان در اختیار همگان با شناسه تحقیقاتی قرار می‌دهد.

این سایت جز برترین سایت‌های دسترسی به داده‌ها معرفی شده‌اند اما بسیاری سایت دیگر با اهداف مختلف برای دسترسی به داده وجود دارد کافیست هدف خود را گوگل کنید...

داده را از کجا پیدا کنیم؟



مقالات



در هر نشریه قصد داریم با ارایه مقالاتی کاربرد علم آمار را در علوم مختلف بیان کنیم. برای اولین نسخه کاربرد علم آمار در علوم اعصاب و علوم شناختی را در نظر گرفتیم: عنوان مقاله: پردازش غیر خطی در آنالیز آماری سیگنالهای EEG در بررسی خویشاوندی و علاقه مندی به انتخاب یک کالا

نویسنده‌گان: سعیده ریس دانا، سمانه صفری

در این مقاله یک طرح بازاریابی عصبی توسط پردازش سیگنالهای EEG انجام شده است که در آن میزان علاقه مندی افراد جامعه به خرید یک کالای تزیینی نسبتاً لوکس ارزیابی می‌شود.

استخراج احساسات و ترجیحات افراد از ناخودآگاه آنها که موضوع بازاریابی عصبی است می‌تواند به عنوان یک روش قابل اطمینان برای ارزیابی میزان علاقه مندی افراد به کالا و خدمات مورد بررسی قرار گیرد. نورو مارکتینگ یک علم میان رشته‌ای در حال ظهور در مرز بین علوم اعصاب، روانشناسی و بازاریابی است که به لزوم درک بهتر عملکرد مغز و بررسی رفتار مصرف کننده از دیدگاه مغز تاکید دارد. تحقیقات در حوزه بازاریابی عصبی با تجزیه و تحلیل سیگنالهای بیولوژیکی از جمله EMG, EOG, EEG و تصاویر FMRI و ردیابی چشم، آنالیز حالت چهره و ... انجام می‌گیرد.

در این مقاله روش‌های آماری مورد استفاده به ترتیب:

۱. ثبت دادگان

۲. ابزار پردازش و روش‌های محاسباتی

۳. آنالیز مولفه‌های مستقل با روش‌های محاسباتی و آماری برای جداسازی مولفه‌های مستقل، مجموعه‌ای از سیگنال‌ها، اندازه گیری، متغیرهای چندگانه است.

۴. شیوه نورو فازی: حل مسائل پیچیده بازشناسی الگوهای عصبی با کمک یادگیری ماشین و استنتاج مدل فازی در چارچوب شبکه عصبی چندگانه.

۵. آنالیزهای آماری: آزمون t، آنالیز واریانس anova

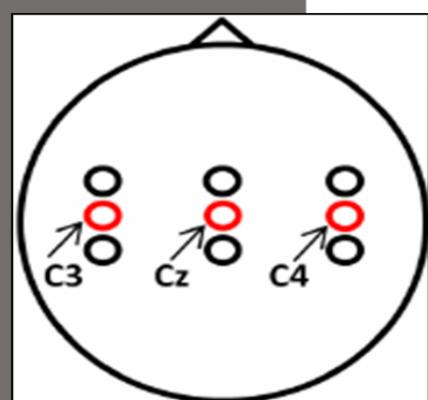
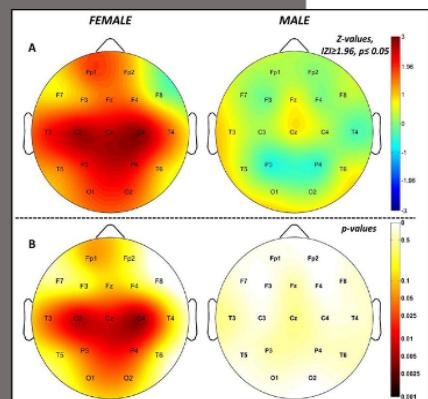
۶. پردازش تصویری داده‌ها و پردازش ICA

۷. استخراج ویژگی و طبقه بندی سیگنالها

تحلیل آماری مورد استفاده جهت اعتبارسنجی نتایج این تحقیق به کمک ابزارهای استاندارد آنالیز آماری، تجزیه و تحلیل شده‌اند. داده‌ها توسط پرسشنامه فراهم شده‌اند. از نرم افزار آماری SPSS هم برای آزمون‌های رگرسیونی و آنالیز داده‌ای ناخود آگاه استفاده شده است.

در پایان هم نتایج این مقاله با مقالات مشابه مقایسه شده و نتایج را اعلام نموده‌اند. برای اطلاع بیشتر لینک مقاله به آدرس زیر است:

<https://www.sid.ir/fa/journal/ViewPaper.aspx?ID=497594>



فیشر

فیشر فرزند کیتی هیث (وکیل) و جرج فیشر (صاحب یک شرکت فروش فوق العاده در خیابان کینگ لندن) بود. فیشر سه خواهر و یک برادر بزرگتر داشت. وقتی ۱۴ ساله بود مادرش از دنیا رفت و تنها ۱۸ ماه بعد پدرش بر اثر شکست قابل ملاحظه در معامله شغلش را از دست داد. رونالد در طول دوران تحصیل، استعداد درخشانی در ریاضیات از خود نشان داد. علی‌رغم از دست دادن مادرش در چهارده سالگی وی در همان سال، در رقابت ریاضیاتی که بین دانش‌آموزان تمام مدارس، در مدرسه «هارو» برگزار شد، موفق به کسب مدال شد. اگرچه رونالد فیشر بینایی ضعیفی داشت اما یک دانش‌آموز باهوش و پیشرو بود و در مسابقات ریاضی مدرسه هارو در ۱۶ سالگی مدال گرفت. به دلیل بینایی ضعیفیش برای آموزش خانگی ریاضی از کاغذ و قلم استفاده نمی‌کرد، طوری که توانایی او در تجسم موارد هندسی نیز گسترش یافت. او گسترش قابل توجهی در علم بیولوژی داد و سیر تکاملی خاصی در آن ایجاد کرد.

در سال ۱۹۰۹ او موفق به دریافت بورس کالج گانویل دانشگاه کمبریج لندن شد. در کمبریج فیشر به انجام مطالعاتی جدید درباره تئوری‌های مربوط به ژنتیک مندل دست زد و همچنین نظری بر بیومتری (کاربرد آمار در مبحث زیست‌شناسی) داشت و با تشکیل مجموعه‌ای از دست نوشته‌ها در مورد روش‌های آماری موجب شد که روشی بالقوه در تطبیق با طبیعت وارث مندل حاصل شود که در آن از واریانس‌های متوالی که رتبه تغییرات نهادین بواسطه اثر تغییرات چند ژن و محیط اطراف و تکامل تدریجی آن مشخصه بوجود آورده است استفاده می‌شود.

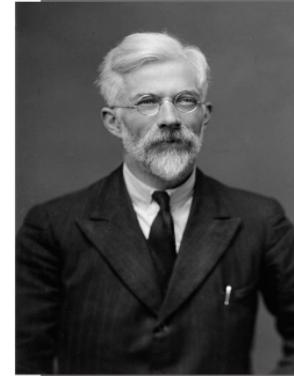
وی پس از دریافت جایزه ولستون به مطالعاتش بر روی تئوری خطاهای در دانشگاه کمبریج ادامه داد. در واقع این علاقه فیشر بود که آن‌ها را به سمت سرمایه‌گذاری بر روی مطالعات حوزه آمار سوق می‌داد. بعد از فارغ‌التحصیل شدن از کمبریج، او دیگر از هیچ‌گونه پشتیبانی مالی برخوردار نبود و برای چند ماه در شرکتی در کانادا مشغول به کار شد. سپس به لندن بازگشت و به عنوان یک متخصص علم آمار در یک شرکت بزرگ استخدام شد. علاقه‌اش به اصلاح نژاد و گیاهان و تجربیاتش در کمبریج، فیشر را علاقه‌مند کرد به اینکه برای خود مزرعه‌ای داشته باشد.

«کارل پیرسون» به او پیشنهاد کرد تا به عنوان کارشناس ارشد آمار در لابراتوار گالتون مشغول به کار شود. همچنین مسئولیت اجرایی مشابهی در ایستگاه آزمایش‌های کشاورزی «روتھامستد» به وی پیشنهاد شد که یکی از قدیمی‌ترین مؤسسات تحقیقات کشاورزی در انگلستان بود و در سال ۱۸۳۷ به منظور مطالعه بر روی اثرات تغذیه خاک و انواع مختلف خاک بر روی باروری گیاهان تأسیس شده بود. علاقه فیشر به کشاورزی سبب شد تا مسئولیت اجرایی پیشنهادی روتھامستد را بپذیرد. جایی که او به واسطه ارائه روش‌های آنالیز و تجزیه و تحلیل نتایج آزمایش‌ها، خدمات زیادی هم به علم آمار و هم به علم ژنتیک کرد. در آنجا بود که او برروی طراحی آزمایش‌هایی بوسیله معرفی مفهوم انتخاب تصادفی و آنالیز واریانس مطالعاتی انجام داد، روش‌هایی که هم‌اکنون نیز در تمام دنیا مورد استفاده قرار می‌گیرند. وی در سال ۱۹۲۱، مفهوم احتمال را معرفی کرد. احتمال یک پارامتر، متناسب است با امکان وجود داده‌ها و تابعی ارائه می‌دهد که معمولاً دارای یک واحد ارزش ماقریم است، که ماقریم احتمال نامیده می‌شود. فیشر در سال ۱۹۲۲، تعریف جدیدی از آمار ارائه کرد. هدف او این بود که ادعا کند، کاهش (کمبود) داده وجود دارد و سه مشکل اساسی را معرفی کرد:

۱. تنوع جمعیتی که داده‌ها از آنها استخراج می‌شوند

۲. تخمین

۳. توزیع



نرم افزار R

نرم افزار R یک نرم افزار محاسباتی آماری و یک زبان برنامه نویسی نیز می باشد که در زمانی کوتاهی مورد استقبال گسترده ای قرار گرفته است. با اطمینان می توان گفت که R هر روز از روز پیش کامل تر است، زیرا بیشتر کاربران آن با نوشتن برنامه هایی برای حل مساله های خود و به اشتراک گذاشتن آنها در پیشرفت نرم افزار نقشی اساسی دارند. امروزه ۱۸۲۷۶ پکیج مختلف در نرم افزار وجود دارد که قابل دانلود برای همگان است. این نرم افزار یک زبان برنامه نویسی و محیطی یک پارچه را در اختیار کاربران قرار می دهد تا به کمک آن بتوانند کارهای گوناگونی از جمله موارد زیر را انجام دهند:

- * وارد کردن، ویرایش، پیکره بندی و ذخیره های انواع داده ها در قالب های مناسب.

- * اجرای روش ها و تحلیل های رایج آماری با به کار بستن تابع های موجود در R

- * اجرای روش ها و تحلیل های جدید آماری با به کار گیری بسته های نرم افزاری موجود بر پایگاه های CRAN

- Github , BioConductor

- * اجرای محاسبات عددی مانند اعمال جبری روی بردارها، ماتریس ها و آرایه ها، حل معادله های خطی و ناخطی، یافتن ریشه های چند جمله ای ها، بهینه سازی، تقریب انتگرال، مشتق گیری و ...

- * رسم انواع نمودارهای آماری با انعطاف پذیری بالا

- * به کار بستن یک زبان برنامه نویسی ساده و کارا برای پیاده سازی الگوریتم ها و هر نوع تحلیل آماری و ریاضی دلخواه.

تاریخچه:

در اواسط قرن بیستم که انواع روش های اماری برای تحلیل داده ها توسط آماردانان و پژوهشگران دیگر رشته ها معرفی و در سطح گسترده ای از علوم و رشته های مهندسی فراگیر شده بودند، برای اجرای محاسبات آماری، بیشتر زبان های برنامه نویسی رایج آن دوران، به ویژه فورترن، به کار گرفته می شد. به همین خاطر انجام محاسبات مربوط به تحلیل داده ها به نوشتن برنامه های طولانی و مهارت بالای برنامه نویسی نیاز داشت.

در سال های ۱۹۷۵-۱۹۷۶ به نظر انجام ساده تر و سریع تر محاسبات آماری و عددی در مجموعه هی مشهور و با سابقه ای آزمایشگاه های بیل نسخه های آغازین زبان برنامه نویسی S توسط جان چمبرز و با همکاری ریک بکر و آلن ویک شکل گرفت. در سال ۱۹۸۸ تغییرات زیادی در زبان S ایجاد و ویژگی های جدیدی مانند به کار گیری تابع نویسی به آن افزوده شد و نسخه جدید S⁴ نامیده گرفت. سپس ویژگی های برنامه نویسی شیء گرا مانند به کار گیری کلاس ها و متد ها نیز به S افزوده و نسخه جدید S⁴ نامیده شد. در این میان آمار گیری به زبان S محدود به آزمایشگاه های بیل نشد و این زبان در اختیار دیگر افراد به ویژه پژوهشگران در دانشگاه ها و مرکز علمی قرار گرفت. سادگی برنامه نویسی با S، امکان اجرای روش های آماری رایج با دستور های ساده و کوتاه، رسم انواع نمودارهای گرافیکی و انتشار چند کتاب در مورد شیوه بکار گیری S توسط تیم توسعه دهنده آن زبان مورد توجه قرار گیرد و گسترش یابد. به دلیل استقبال از S، یک پیاده سازی تجاری از آن با عنوان S_PLUS یک شرکت تجاری منتشر و برای فروش به بازار عرضه شد. در سال ۱۹۹۳ راس ایه اکا و رابت جنتلمن از دانشگاه اوکلند نیوزیلند با افزودن ویژگی هایی از زبان برنامه نویسی Scheme به S، پروژه توسعه ای که پیاده سازی غیر تجاری و متن باز از S را آغاز کردند و آن را R نامیدند. این اسم، حرف اول اسم هردوی آنان بوده و حرف قبل از S در الفبای انگلیسی است. پس از آن متخصصان دیگری از سراسر جهان، از جمله جان چمبرز خالق S، به آن ها پیوستند و تیم هسته توسعه را R تشکیل دادند که تا به امروز مسؤولیت توسعه و به روز رسانی R را به عهده دارد.

برای دانلود آخرین نسخه نرم افزار R و R studio به لینک های زیر مراجعه کنید:

<https://cran.um.ac.ir/bin/windows/base/R-4.1.1-win.exe>

<https://download1.rstudio.org/desktop/windows/RStudio-121.09.0-351.exe>



چرا نرم افزار R؟

چرا R برای Data Science مهم است؟

چون می توانید کد خود را بدون هیچ کامپایلری اجرا کنید؛ یعنی R یک زبان تفسیر شده است. از این رو می توان بدون هیچ کامپایلری کد را اجرا کرد. R کد را تفسیر می کند و توسعه کد را آسان تر می کند. از طرفی دیگر بسیاری از محاسبات با بردارها انجام می شود؛ یعنی R یک زبان برداری است. بنابراین هر کسی می تواند توابع را به یک بردار اختصاص دهد، بدون اینکه حلقه ایجاد کند. از این رو، R قدرتمندتر و سریع تر از زبان های دیگر است.

چرا R برای تجارت (Business) خوب است؟

دلیل اصلی این است که R یک برنامه open-source است؛ بنابراین می تواند براساس نیاز کاربر اصلاح و توزیع شود. برای عملیات مصورسازی داده ها (visualization) عالی است و در مقایسه با سایر ابزارها قابلیت های بسیار بیشتری دارد. برای مشاغل داده محور، عدم وجود دانشمند داده (data scientist) یک نگرانی بزرگ است. شرکت ها از برنامه نویسی R به عنوان پلتفرم اصلی خود استفاده می کنند و برنامه نویسان آموزش دیده R را استخدام می کنند.

زبان R یک شغل پردرآمد

زبان R به طور گسترده ای در data science مورد استفاده قرار می گیرد که برخی از پردرآمدترین مشاغل روز جهان را در بردارد. اگر می خواهید وارد حوزه data science شوید و درآمد بالایی بست آورید، قطعاً باید R را یاد بگیرید.

منبع باز open-source

زبان R یک زبان منبع باز است که توسط جامعه ای از کاربران فعال توسعه داده می شود و به کار خود ادامه می دهد؛ از همین رو می توانید از R به صورت رایگان بهره مند شوید. شما می توانید توابع مختلف را در R اصلاح کرده و بسته های نرم افزاری خود را بازاریزید. از آنجا که R تحت مجوز عمومی (GNU) صادر شده است، هیچ محدودیتی در استفاده از آن وجود ندارد.

محبوبیت

زبان R به یکی از محبوب ترین زبان های برنامه نویسی در صنایع تبدیل شده است. به طور معمول، R بیشتر در دانشگاه استفاده می شده است اما با ظهور data scientist، نیاز به R در صنایع مشهود شد. به عنوان مثال، از در فیس بوک برای تجزیه و تحلیل شبکه های اجتماعی و یا در توییتر برای تحلیل مفهومی و همچنین مصورسازی داده ها (visualization) استفاده می شود.

داشتن کتابخانه قوی برای مصورسازی داده ها

زبان R شامل کتابخانه هایی مانند ggplot2 و plotly است که طرح های گرافیکی جذابی را به کاربران خود ارائه می دهند. R به دلیل داشتن قابلیت های مصورسازی فوق العاده نسبت به سایر زبان های برنامه نویسی بیشتر مورد توجه data science ها می باشد.

برخورداری R از جامعه پشتیبانی گسترده

برنامه نویسی R توسط جامعه وسیعی پشتیبانی می شود که R را توسعه می دهند و به روز می کنند. اگر هنگام کدنویسی در R با مشکلی روبرو شدید، می توانید از پشتیبانی و کمک افراد در مکان هایی مانند stack overflow بهره مند شوید. انجمنهای مختلفی در سراسر جهان وجود دارند که بوت کمپ و میتینگ های R را سازماندهی می کنند.

زبانی برای آمار و علوم داده

زبان R، زبان استاندارد علم آمار و علم داده است. برای آمار و توسط آمارشناسان ساخته شده است. این زبان حتی قبل از ایجاد کلمه data science نیز مورد استفاده بوده است. آمارشناسان و دانشمندان داده بیشتر از هر زبان برنامه نویسی با R آشنایی دارند چراکه R از طریق هزاران پکیج مختلف، انجام عملیات مختلف آماری را تسهیل می کند.

فراگیری زبان R در صنعت

امروزه R از پر کاربرد ترین زبان های برنامه نویسی در جهان است و تقریباً در هر صنعتی مورد استفاده قرار می گیرد؛ از مالی و بانکی گرفته تا دارو و تولید. R برای مدیریت سبد سرمایه یا داشتمان (portfolio management) و تجزیه و تحلیل ریسک در صنایع و بازارهای مالی، تجزیه و تحلیل کشف داروهای جدید، تجزیه و تحلیل ژنومیک در بیوانفورماتیک اجرای عملیات آماری مختلف برای بهینه سازی فرآیندهای صنعتی استفاده می شود.

سرگرمی به زبان R

برنامه نویسی با R

سوال ۱

کدام یک از رشته کدهای زیر فقط دو ردیف اول فایل csv را می خواند؟

- A) `csv('Dataframe.csv', header=TRUE, row.names=1, sep=',', nrows=2)`
- B) `csv2('Dataframe.csv', row.names=1, nrows=2)`
- C) `delim2('Dataframe.csv', header=T, row.names=1, sep=',', nrows=2)`
- D) `dataframe('Dataframe.csv', header=TRUE, row.names=1, sep=',', skip.last=2)`

سوال ۲

خروجی این کد چه خواهد بود؟ دلیل اجرای این کد را تفسیر کنید.

```
generic2 <- function(x) UseMethod("generic1")
generic2_a1 <- function(x) "a1"
generic2_a2 <- function(x) "a2"
generic2_b1 <- function(x) {
  class(x) <- "a1"
  NextMethod()
}

generic2(structure(list(), class = c("b1", "a2")))
```

در نسخه بعدی که منتشر خواهد شد پاسخ این سوالات قرار می‌گیرد و به علاقه‌مندانی که پاسخ سوالات مطرح شده را به شبکه اجتماعی انجمن علمی آمار ارسال کنند به قید قرعه جوایزی اهدا خواهد شد.

