



A COMPARATIVE STUDY ON THE RELATIONSHIP BETWEEN STATISTICS AND DATA SCIENCE

AMIR HOSSEIN GHATARI^{1*}, ELHAM TABRIZI², AND EHSAN BAHRAMI SAMANI³

¹Department of Statistics, Faculty of Mathematics and Computer Science, Amirkabir University of Technology, Tehran, IRAN

a.h.ghatari@aut.ac.ir

²Department of Mathematics, Faculty of Mathematical Science and Computer, Kharazmi University, Tehran, Iran

elham.tabrizi@khu.ac.ir

³Department of Statistics, Faculty of Mathematical Science, Shahid Beheshti University, Tehran, Iran

ehsan_bahrani_samani@yahoo.com

Abstract. The study of the relation between data science as an emergent phenomenon in 21st century and statistics, an original and applicable science, is one the most challenged topics by different researchers. Data science is dependent to some theoretical basics such as linear models, clustering and etc, which are from statistics. Some of the scientists believe that data science is begotten from statistics and consider it as statistics' offspring. Some others consider data science as a branch of statistics using modern computational techniques. Also, there is an uncommon idea that considers data science as a separated scientific field from statistics. In this paper, we recall theoretical basics of data science and bring up comments and interpretations of some well-known scientists about the fact that data science is not separate from statistics. Finally, we consider Statistics⁺ as another name of data science. Additionally, this study aims to explore the practical implications of this relationship by examining how advancements in data science continue to reshape and augment statistical methodologies. Through a comprehensive review of contemporary applications, we aim to illustrate the dynamic interplay between data science and statistics, illuminating their evolving synergy.

2010 Mathematics Subject Classification. 97K70, 97D10.

Keywords. Statistics, Statistics⁺, Data Science, Machine Learning.

Date: Received 20-1-2024 Revised 29-2-2024 Accepted 28-5-2024 Available Online 30-7-2024

©Ferdowsi University of Mashhad.



مطالعه‌ای تطبیقی بر ارتباط بین آمار و علم داده‌ها

امیرحسین قطاری^{1*}، الهام تبریزی²، و احسان بهرامی سامانی³

¹ گروه آمار، دانشکده ریاضی و علوم کامپیوتر، دانشگاه صنعتی امیرکبیر، تهران، ایران.

a.h.ghatari@aut.ac.ir

² گروه ریاضی، دانشکده علوم ریاضی و کامپیوتر، دانشگاه خوارزمی، تهران، ایران.

elham.tabrizi@khu.ac.ir

³ گروه آمار، دانشکده علوم ریاضی، دانشگاه شهید بهشتی، تهران، ایران.

ehsan_bahrami_samani@yahoo.com

چکیده. ارتباط میان علم داده‌ها به‌عنوان پدیده‌ای نوظهور در قرن ۲۱ با دانش ریشه‌دار و کاربردی آمار همواره مورد بحث و تبادل نظر میان پژوهشگران مختلف است. علم داده‌ها در میدان عمل، متکی به مبانی نظری مختلفی از جمله احتمالات، مدل‌های خطی و غیرخطی، رده‌بندی، روش‌های بازنمونه‌گیری، خوشه‌بندی و... است که جملگی از مبانی نظری علم آمار هستند. برخی از دانشمندان بر این باورند که علم داده‌ها برخاسته از آمار بوده و مولود این دانش سنتی است. برخی علم داده‌ها را شاخه‌ای از آمار می‌دانند که تکنیک‌های محاسباتی نرم افزاری نوین را به‌کار بسته است. همچنین، یک نظر نامرسوم نیز این است که علم داده‌ها فارغ از آمار بوده و بر آمار استوار نشده است. در این نوشتگان، ضمن بررسی پایه‌های نظری علم داده‌ها، بررسی شباهت مبانی نظری آن با آمار و نقد و بررسی نظرات برخی از صاحب‌نظران عصر حاضر این گزاره را که «علم داده‌ها عملاً از آمار جدا نیست و همان آمار است با ابزارهای نوین» به‌عنوان خروجی این پژوهش مطرح کرده و لقب آمار⁺ را برای علم داده‌ها در نظر می‌گیریم.

2010 Mathematics Subject Classification. 97K70, 97D10.

واژگان کلیدی. آمار، آمار⁺، علم داده‌ها، یادگیری ماشین.

تاریخ: دریافت ۱۴۰۲/۱۰/۳۰ بازنگری ۱۴۰۲/۱۲/۱۰ پذیرش ۱۴۰۳/۳/۸ انتشار برخط ۱۴۰۳/۵/۹

نحوه ارجاع به این مقاله: ا.ح. قطاری، ا. تبریزی، ا. بهرامی سامانی، مطالعه‌ای تطبیقی بر ارتباط بین آمار و علم داده‌ها، به

سوی علوم ریاضی، ۴ (۱۴۰۳)، شماره ۱، ۱۸-۳۰.

© دانشگاه فردوسی مشهد.

۱. مقدمه

میان علم آمار و علم داده‌ها^۱ ارتباط تنگاتنگی وجود دارد و بسیاری از روش‌های تحلیل داده از حوزه‌ی علم آمار اخذ شده است؛ تا جایی که ارتباط علم داده‌ها و علم آمار مورد بحث برخی پژوهشگران در قرن ۲۱ شده است. این‌که آیا علم داده‌ها شاخه‌ای از علم آمار است؟ یا این‌که بستگی به آمار ندارد و تنها از آمار استفاده می‌کند؟ همچنین، آیا علم داده‌ها بیشتر به علوم کامپیوتر نزدیک است تا آمار؟ از جمله سوالاتی هستند که در ضمن این مساله مطرح شده‌اند. همواره نظرات مختلفی در مورد پاسخ به سوالات فوق مطرح شده است. بهترین آغاز برای بررسی مساله و پاسخ به سوالات فوق و البته سوالاتی از این قبیل، اشاره به تعاریف علم داده‌ها و آمار است که در نشریه انجمن آمار سلطنتی انگلستان و در مقاله دیگل^۲ [۱۲] به کار رفته‌اند:

- آمار: این علم، دانش مطالعه‌ی مجموعه داده‌ها، تحلیل، تفسیر، نمایش و سازماندهی آن‌هاست.
- علم داده‌ها: استخراج اطلاعات از داده‌ها است. تکنیک‌ها و نظریه‌های مطرح شده در زمینه‌های مختلف علوم ریاضی، آمار و همچنین فناوری اطلاعات^۳ را برای رسیدن به این هدف به کار می‌بندد.

به وضوح از تعریف این دو علم درمی‌یابیم که تقریباً هدف هر دو یک مقصد است با این تفاوت که آمار با ظاهری که امروزه مشاهده می‌کنیم در آغاز قرن بیستم و با مشارکت دانشمندانی چون فرانسیس گالتون^۴، رونالد فیشر^۵، کارل پیرسون^۶ و.... شکلی نظام‌مند به خود گرفته است و در طی سالیان متمادی رشدی نظری و کاربردی داشته است؛ درحالی‌که علم داده‌ها با اهدافی مشابه آمار، در قرن ۲۱ پای به عرصه علوم نهاده است. به بیان دیگر، علم داده‌ها به طور مستقیم با فن‌آوری‌های نوین و پیشرفته ارتباط داشته و به استخراج مفاهیم از داده‌ها و تولید محصولات داده‌محور می‌پردازد.

برخی از دانشمندان، علم داده‌ها را زیر مجموعه‌ی آمار خوانده‌اند. از جمله پژوهش‌های هم‌راستا با این نظر، می‌توان از گودمن^۷ و همکاران [۱۷] و سینگپرووالا^۸ [۲۸] نام برد. همچنین کسلا^۹ و همکاران [۷] معتقدند که آمار ریشه‌ی پیدایش دانشی است که به‌عنوان علم ساختاردهی، طبقه‌بندی و تحلیل داده شناخته می‌شود. در

¹Data Science

²Diggle

³Information Technology

⁴Francis Galton

⁵Ronald Fisher

⁶Karl Pearson

⁷Goodman

⁸Singpurwalla

⁹Casella

یک گزارش اینترنتی برای نشریه‌ی پرایس اکونومیکس^{۱۰}، بهاردواج^{۱۱} [۲] به بررسی تفاوت‌های علم داده‌ها و آمار از دیدگاه‌های مختلف پرداخته است.

از جمله موارد دیگری که از علل پیدایش علم داده‌ها محسوب می‌شود؛ نیاز مبرم به گسترش و تعمیم تکنیک‌های محاسباتی آماری در تحلیل داده‌ها در دنیای امروز است. چمبرز^{۱۲} [۸] از قُدمای در تحقیق و بررسی درباره‌ی نیاز به این مساله بود. وی از جمله پژوهشگرانی است که در قرن ۲۱م در زمره‌ی موافقین این گزاره است که علم داده‌ها و مفاهیم مورد استفاده در آن ریشه‌ی آماری دارند. در پژوهشی به‌روزتر در مورد ارتباط میان علم داده‌ها و آمار، وایز^{۱۳} و ایکشتات^{۱۴} [۳۰]، به بررسی تأثیرات آمار در پیدایش علم داده‌ها پرداخته‌اند. آن‌ها معتقدند که آمار نقشی حیاتی در تشکیل علم داده‌ها و فراهم آوردن مفاهیم اولیه‌ی آن دارد. دیگر پژوهش توسط کارمیشل^{۱۵} و مارون^{۱۶} [۶] انجام شده است. آن‌ها این سوال را که علم داده‌ها ریشه‌ی آماری دارد یا خیر را در ابعاد مختلف از جمله نظرات سایر دانشمندان و حتی عملکرد در مسائل کاربردی، مورد بررسی قرار دادند. پژوهشگران متعددی درباره‌ی این‌که علم داده‌ها در مسیر روزرسانی و تکامل آمار به‌وجود آمده است؛ تحقیق و بررسی کرده‌اند. برای اطلاعات بیشتر در این باره به [۲۳، ۱۰، ۱۸، ۲۰، ۲۶] مراجعه شود.

در این مقاله به واکاوی جزئی‌تر اهداف علم داده‌ها، ابزار مورد استفاده در این علم و مقایسه‌ی آن با مفاهیم و ابزارهای آماری پرداخته و نظرات برخی از دانشمندان حال حاضر را درباره‌ی ارتباط این دو گرایش از علم مورد نقد و بررسی قرار خواهیم داد. در نهایت می‌خواهیم به این سوال پاسخ دهیم که آیا علم داده‌ها چیزی جدای از علم آمار است و یا تنها یک پدیده‌ی پیشرفته در ضمن آمار و متناسب با نیاز عصر حاضر است؟ هدف ثانویه‌ی این مقاله ارائه‌ی یک معرفی عمومی از آمار و علم داده‌ها به همراه برخی از مشاهیر و فعالان در این زمینه‌ها به نحوی است که فعالان علمی پژوهشی خارج از این دو حوزه را نیز همراه خود کرده و اطلاعات و منابع سودمندی درباره‌ی آمار و علم داده‌ها در اختیار این طیف از خوانندگان مقاله قرار دهد.

۲. علم داده‌ها، ظهور، اهداف و ابزارها

۱.۲. ظهور علم داده‌ها. عبارت علم داده‌ها بیش از دو دهه است که موجودیت دارد. با استناد به گزارش منتشر شده توسط جیل پرس^{۱۷} [۲۵]، ویلیام کلیولند^{۱۸} اولین کسی است که اصطلاح علم داده‌ها را در سال

¹⁰priceonomics

¹¹Bhardwaj

¹²Chambers

¹³Weihs

¹⁴Ickstadt

¹⁵Carmichael

¹⁶Marron

¹⁷Press

¹⁸William Cleveland

۲۰۰۱ مطرح کرده است. وی در مقاله [۹] «علم داده‌ها: برنامه‌ای برای گسترش جنبه‌های فنی در رشته آمار» پیشنهاد کرد که علم داده‌ها به عنوان یک رشته مستقل شناخته شود. کلیوند این رشته‌ی جدید را مرتبط با علوم کامپیوتر و داده‌کاوی می‌داند. وی بر این باور بود که منافع استفاده از یک تحلیلگر داده محدود است. چون مهندسين کامپیوتر شناخت کمی از روش‌های کار با داده دارند و دانش محاسباتی متخصصین آمار هم محدود است؛ بنابراین تلفیق این دو گروه می‌تواند منجر به نوآوری‌های زیادی شود. گروه‌های علم داده‌ها باید استادانی داشته باشند که بتوانند دانش داده‌ها را با دانش محاسبات تلفیق کنند.

۲.۲. اهداف علم داده‌ها. پیشرفت تکنولوژی در روش‌های ذخیره‌سازی داده‌ها در عصر حاضر، منجر به پیدایش پایگاه‌های داده‌های بزرگ^{۱۹} شده است. از طرف دیگر، پردازش و تحلیل داده‌های حجیم^{۲۰} نظیر داده‌های هواشناسی، زیستی یا مالی به دلیل حجم و معمولا بعد بالا، امر ساده‌ای نیست.

تحلیل داده‌های حجیم یک نیاز حیاتی در دنیای تجارت و همچنین علوم اجتماعی است [۲۴]. در واقع از نیازهای اساسی که می‌توان آن را یکی از علل پیدایش علم داده‌ها عنوان کرد؛ همین مسأله‌ی تحلیل داده‌های حجیم است. براساس نظر کلیوند، آمار سنگ بنای نظری برای علم داده‌ها را فراهم می‌کند و تکنولوژی نیز به کمک آمار آمده و تحلیل‌های دقیق‌تر و متناسب با نیاز روز را فراهم می‌آورد. براساس همین نظر نیز می‌توان گفت که علم داده‌ها در واقع آمار است که بازوان اجرایی آن تکنیک‌های مدرن در دنیای مهندسی کامپیوتر است.

۳.۲. ابزارهای علم داده‌ها. در این زیربخش به مروری بر ابزارها و تکنیک‌های مرسوم در علم داده‌ها پرداخته می‌شود. از یک سو، مشخص می‌شود که علم داده‌ها خاستگاه نظری خود را از چه منبعی تهیه کرده و از طرف دیگر نشان می‌دهد که علم داده‌ها عملاً تلفیق چند رشته‌ی علمی نیست و در واقع آمار است که با استفاده از ابزارهای محاسباتی پیشرفته در عصر حاضر دست به ارائه‌ی نتایج می‌زند. این ابزارها عبارتند از:

- رگرسیون خطی^{۲۱}، مبحثی که تیم تحقیقاتی گالتون نخستین بار در اواخر قرن ۱۹ میلادی بدان پرداخته‌اند.

- مدل‌های خطی تعمیم یافته^{۲۲}، این مقوله در ذیل مبحث کلاس‌بندی داده‌ها بوده که نخستین بار تحلیل این نوع از داده‌ها توسط کارل پیرسون در آغاز قرن ۲۰ میلادی مطرح گشت (اگرستی^{۲۳} [۱] فصل ۱۶، بخش اول).

¹⁹Big Data

²⁰Huge Data

²¹Linear Regression

²²Generalized Linear Models

²³Agresti

- درخت تصمیم^{۲۴}، این مبحث نیز یک شاخه‌ی نظری در علم آمار است که در مسائل رده‌بندی و رگرسیونی می‌توان از آن استفاده کرد. برای اطلاعات بیشتر در مورد درخت تصمیم به بریمن^{۲۵} و همکاران [۴] و ون وینترفلت^{۲۶} و ادواردز^{۲۷} [۲۹] مراجعه شود.
- در نوعی پیشرفته‌تر از درخت تصمیم، جنگل تصادفی^{۲۸}، کیسه‌بندی^{۲۹} و روش تقویتی^{۳۰}، سه روش مبتنی بر یادگیری گروهی^{۳۱} هستند که برای ساختن مدل‌های پیشگویی‌کننده از درخت‌های تصمیم استفاده می‌کنند. در اکثر روش‌های سنتی آمار فرض می‌شود که تعداد متغیرهای توضیحی کمتر از تعداد نمونه‌های مستقل است. بنابراین، ظهور مجموعه داده‌های با بعد بالا^{۳۲} که در آنها تعداد متغیرهای توضیحی بیشتر از تعداد نمونه‌های مستقل است، به‌طور قابل توجهی جریان اصلی بسیاری از تحقیقات آماری را به خود اختصاص داده است. جنگل‌های تصادفی یکی از جدیدترین روش‌های یادگیری با ناظر^{۳۳} است که محققین، مطالعات زیادی را پیرامون این روش برای بهبود عملکرد پیش‌بینی متغیر پاسخ، هنگام مواجهه با داده‌های با بعد بالا، انجام داده‌اند. مجموعه اطلاعات در مورد این نوع از یادگیری آماری در هستی^{۳۴} و همکاران [۱۹] در دسترس است.
- ماشین بردار پشتیبان^{۳۵}، این نیز یک مبحث در کلاس‌بندی داده‌هاست که عملاً مباحثی نظری در آمار دارد. مجموعه اطلاعات در مورد این نوع از کلاس‌بندی، منشا و معرّف آن در هستی^{۳۶} و همکاران [۱۹] در دسترس است.

²⁴Decision Tree

²⁵Breiman

²⁶Von Winterfeldt

²⁷Edwards

²⁸Random Forest

²⁹Bagging

³⁰Boosting

³¹Ensemble Learning

³²High Dimensional Data

³³Supervised Learning

³⁴Hastie

³⁵Support Vector Machine

³⁶Hastie

- خوشه‌بندی^{۳۷}، نخستین بار ارائه شده توسط درپور^{۳۸} و کروبر^{۳۹} [۱۴]، بحثی مجزا با گونه‌های بسیار متنوع در علم آمار است. تفاوت آن با کلاس‌بندی آمار این است که در کلاس‌بندی برچسب طبقات داده‌ها مشخص است و در خوشه‌بندی با اتکا به خصوصیات موجود در داده‌ها اقدام به دسته‌بندی می‌کنیم. خوشه‌بندی عضوی از مجموعه‌ی یادگیری بدون ناظر^{۴۰} است که شامل تمام مطالعاتی است که در آن‌ها فرض قبلی در مورد برچسب طبقات داده‌ها برقرار نیست. برای اطلاعات کامل در مورد انواع خوشه‌بندی و یادگیری بدون ناظر به کاسامبرا^{۴۱} [۲۱] مراجعه شود.
- کاهش بُعد^{۴۲}، مبحثی در زمینه‌ی تلخیص داده‌هاست. برخی داده‌های بُعد بالا، ایجاد نقایصی در خواص ریاضیاتی مدل‌ها را در بر خواهند داشت. کاهش بعد یک مبحث مشترک در آمار و جبرخطی است که کاربرد در تلخیص داده‌ها دارد.
- یادگیری ماشین^{۴۳}، اساساً ابزار محاسباتی کامپیوتر محور در علم داده‌ها محسوب می‌شود که با ارائه‌ی الگوریتم‌های محاسباتی دقیق و پرسرعت موجب تسهیل دریافت خروجی می‌شود. در واقع، یادگیری ماشین کاربردی از هوش مصنوعی است که به سیستم، توانایی یادگیری خودکار و بهبود تجربه را می‌دهد. جهت اطلاعات بیشتر در مورد تکنیک‌های یادگیری ماشین به بیشاپ^{۴۴} [۳] مراجعه شود. آنچه در بالا مطرح شد ابزارهای مورد استفاده در علم داده‌ها بود. در ادامه از دیدگاه آماری به این ابزارها نگاه خواهیم کرد.

۳. آمار و نقش مبانی نظری آن در علم داده‌ها

براساس موارد ذکر شده در بخش ۲، علم داده‌ها از رگرسیون، کلاس‌بندی آماری، خوشه‌بندی، کاهش بعد و یادگیری ماشین برای تحلیل، پیش‌بینی و ارائه خروجی در مورد داده‌ها استفاده می‌کند. همانطور که در زیربخش ۳.۲ ذکر شد، تمامی موارد فوق (غیر از الگوریتم‌های رایانه‌ای در یادگیری ماشین) اساساً مباحثی ریشه‌ای در علم آمار هستند. در واقع این چنین استنتاج می‌شود که در علم داده‌ها از یادگیری ماشین استفاده می‌کنند تا مباحث نظری آماری را پیاده‌سازی کنند. اگر تاکنون به پاسخ نرسیده‌اید، اکنون سوال دیگری مطرح می‌کنیم که

³⁷Clustering

³⁸Driver

³⁹Kroeber

⁴⁰Unsupervised Learning

⁴¹Kassambara

⁴²Dimension Reduction

⁴³Machine Learning

⁴⁴Bishop

در یافتن پاسخ کمک خواهد کرد.

فرض کنید مسافرانی سوار هواپیما شوند، آیا مسافران را هواپیما به مقصد رسانده است یا خلبان؟ بهتر مطرح کنیم: آیا هواپیما ابزاری در دست خلبان برای جابجایی مسافران بوده یا خود آن‌ها را به مقصد رسانده است؟ مسافران همان داده‌ها، خلبان همان نظریات آماری و علم داده‌ها همان هواپیمایی است که آمار، داده‌ها را به کمک آن به مقصد که همان تحلیل و نتایج است می‌رساند.

هوش مصنوعی به‌عنوان یکی از ابزارهای پیشرفته‌ی مورد استفاده‌ی علم داده‌ها شناخته می‌شود. گودفلوو^{۴۵} و همکاران [۱۶] در مورد نقش هوش مصنوعی در علم داده‌ها اینگونه تعبیر می‌کنند که، تکنیک‌های مرتبط با تحلیل و محاسبات در هوش مصنوعی معمولاً برپایه‌ی یادگیری عمیق^{۴۶} طراحی شده‌اند و شدیداً تشنه‌ی داده هستند.

نقش هوش مصنوعی در علم داده‌ها را می‌توان تشبیه به یک رایانه کرد که وظایف یک کارشناس آمار را انجام می‌دهد. در واقع مباحث نظری آمار به وسیله علم داده‌ها روی داده‌های حجیم تولید شده در عصر حاضر پیاده‌سازی می‌شوند. می‌توان گفت هوش مصنوعی نکته‌ای شبیه به خلبان اتوماتیک در هواپیماست. یعنی از ماهیت موضوع که فاعل اصلی آمار است، چیزی کاسته نمی‌شود. در ادامه، به بررسی بیشتر پاسخ‌های ارائه شده در مورد رابطه‌ی بین آمار و علم داده‌ها می‌پردازیم. در نهایت می‌خواهیم به پاسخی تجمیعی برای سوالی که از ابتدا مطرح کرده‌ایم، برسیم.

۴. آیا علم داده‌ها چیزی بیش از آمار است یا تنها نسخه‌ی پیشرفته‌ی آمار می‌باشد؟

در این بخش، به بررسی و نقد نظرات برخی از افراد صاحب نظر در دنیای آمار و علم داده‌ها درباره‌ی ارتباط این دو رشته می‌پردازیم:

- چمبرز [۸] در ۱۹۹۳، در انتقاد به واکنش آماردان‌ها به پیشرفت فناوری اینگونه بیان داشته که: آمار در پاسخ به فناوری‌های جدید تغییرات قابل توجهی نداشته است. این رشته بر تئوری تاکید دارد و دروس آمار مقدماتی بیشتر تمرکز بر آزمون‌های آماری دارد تا محاسبات آماری. نقد و بررسی: با توجه به این‌که این نظرات در دهه ۱۹۹۰ بوده است؛ نقدی قابل‌پذیرش و نوعی هشدار در جهت همسو شدن فعالان رشته‌ی آمار با پیشرفت تکنولوژی است. بصورت غیرمستقیم و البته آنچه در حقیقت رخ داد نیز این بود که در شرایط احساس نیاز به تحلیل داده‌های حجیم، علمی به اسم علم داده‌ها در بستر نظری دانش آمار ظهور کرد.

⁴⁵Goodfellow

⁴⁶Deep Learning

- کیلیوند [۹] در مورد ظهور علم داده‌ها اینگونه بیان می‌کند که: علم داده‌ها را باید بر اساس میزان توانمندسازی یک تحلیلگر در تفسیر داده‌ها قضاوت کرد. ابزارهایی که توسط تحلیلگر داده‌ها استفاده می‌شود، اثربخشی مستقیم دارند. نظریه‌هایی که به عنوان مبنایی برای توسعه ابزارها عمل می‌کنند، فواید غیر مستقیم دارند.
- نقد و بررسی: نظر کیلیوند [۹] تبیین‌کننده‌ی این نکته است که علم داده‌ها در بستر نیاز به گسترش روش‌های تحلیل داده‌ها به وجود آمده است. وی ابزارهای تحلیل را مستقیماً موثر بر علم داده‌ها دانسته و مبنای نظری را بصورت غیرمستقیم برای این علم سودمند دیده است. این نظر منافاتی با این‌که علم داده‌ها در واقع گسترش روش‌های آماری و مبنای نظری آن در جهت حل مسائل جدید است، ندارد.
- لوکیدز^{۴۷} [۲۲] معتقد است که تنها تفاوت میان آمار و علم داده‌ها این است که علم داده‌ها از رویکردهای کلی‌نگرتری استفاده می‌کند.
- نقد و بررسی: نظر لوکیدز [۲۲] همسو با این عقیده است که علم داده‌ها از دیدگاه مبنای تئوری از آمار جدا نیست. در واقع وی تنها روش‌های حل مساله میان این دو رشته را متفاوت می‌داند.
- سیلور^{۴۸} آماردان و ژورنالیست مشهور آمریکایی درباره ارتباط علم داده‌ها با آمار اینگونه بیان می‌کند که علم داده‌ها تنها می‌تواند نام دیگر آمار باشد نه یک رشته‌ی علمی جداگانه [۲۷].
- نقد و بررسی: نظر سیلور [۲۷] بر این اساس است که علم داده‌ها از همان مبنای نظری مورد استفاده‌ی آمار برای تحلیل نتایج استفاده می‌کند. او معتقد است که آنچه از آمار نیاز داریم چیزی جدا از علم داده‌ها نیست.
- برومن^{۴۹} [۵] در مقاله‌ای تحت عنوان علم داده‌ها همان آمار است اینگونه بیان کرده است که وقتی در فیزیک از ریاضی استفاده می‌کنیم فیزیکدانان نمی‌توانند بگویند در حال انجام علم اعداد هستند، آن‌ها دارند از ریاضی استفاده می‌کنند. بنابراین، هرگاه ما تجزیه و تحلیل انجام دهیم، در حال بهره بردن آمار هستیم. شما می‌توانید آن را علم داده‌ها یا انفورماتیک یا تحلیل یا هر چیز دیگری بنامید، اما هنوز هم آمار هست و این نباید باعث شود اصطلاح «آمار» را کنار بگذارید.
- نقد و بررسی: دیدگاه برومن [۵] هم راستا با نظر سیلور [۲۷] است. در واقع هر دو بر این باورند که علم داده‌ها چیزی جدای از آمار نیست و تحلیل داده‌ای که علم داده‌ها بدان می‌پردازد در واقع استفاده از آمار است.

⁴⁷Loukides⁴⁸Silver⁴⁹Broman

● گلمان^{۵۰} [۱۵] بر این باور است که آمار کمترین نقش را دارد و آمار را یک بخش غیر اساسی علم داده‌ها برشمارده است.

نقد و بررسی: براساس موارد قید شده در بخش ۲، نقش به سزای مباحث نظری آمار در ساختن پیکره‌ی ابزارها و تکنیک‌های مورد استفاده برای تحلیل داده در علم داده‌ها غیرقابل انکار است. در واقع علم داده‌ها بدون رگرسیون، سری زمانی، کلاس‌بندی و خوشه‌بندی چیزی بیشتر از یک ماشین حساب پرسرعت نخواهد بود و این تکنیک‌های مذکور تماماً از سرفصل‌های شناخته‌ی شده‌ی علم آمار هستند. نظر نگارنده بر آن است که گلمان [۱۵] از دیدگاه توان محاسباتی و نقطه‌نظر فناوری‌های رایانه‌ای به مساله پرداخته‌اند. در کل نقطه‌ی قوت علم داده‌ها استفاده از الگوریتم‌های پیچیده و پیشرفته‌ی کامپیوتری است که همین الگوریتم‌ها نیز در واقع کدنویسی از تکنیک‌های آماری روی کاغذ هستند. بدون مباحث نظری آمار روی کاغذ این الگوریتم‌ها وجود خارجی ندارند.

● دهار^{۵۱} [۱۱] دانشمند هندی، معتقد است که آمار تاکید بر داده‌های کمی و توصیف و تفسیر آن‌ها دارد و علم داده‌ها در تعامل با داده‌های کمی و کیفی سعی و تاکید بر تحلیل و پیش‌بینی این داده‌ها دارد. نقد و بررسی: نظرات آقای دهار عملاً تعابیری نادرست از دانش آمار مطرح می‌کند. براساس هستی و همکاران [۱۹]، پیش‌گویی و پیش‌بینی براساس مدل‌بندی داده‌ها به کمک روش‌های آماری چه برای داده‌های کمی و چه کیفی امکان‌پذیر است. مساله‌ی رگرسیون عملاً مدل‌بندی و پیش‌بینی داده‌هایی از متغیر هدف است که براساس ویژگی‌های موجود و تعیین مدل مناسب این هدف در دسترس قرار می‌گیرد. همچنین ارائه‌ی مدلی که داده‌ها را به کلاس‌های گوناگون طبقه‌بندی کند و یک ورودی جدید را به یک طبقه‌ی خاص نسبت دهد عملاً یک رویه برای تحلیل داده‌های کیفی است که در مساله‌ی کلاس‌بندی آماری انجام می‌شود. تکنیک‌های یادگیری ماشین و الگوریتم‌های موجود در آن، این امکان را به روش‌های مذکور می‌دهد که نتایج و تحلیل‌ها را با سرعت بهتر انجام دهند. همچنین امکان تحلیل داده‌های حجیم با استفاده از زبان برنامه‌نویسی پایتون و به کار بستن روش‌های آمار مذکور وجود دارد. شایان ذکر است که در دنیای امروز به کار بستن یک رویکرد علمی آماری در یک ماشین یادگیری و گرفتن خروجی نهایی را از وظایف علم داده‌ها می‌دانند. حال آنکه براساس آنچه گفته شد تنها یک ابزار محاسباتی قوی در خدمت روش‌ها و مفاهیم نظری آمار داریم.

● دونوهو^{۵۲} [۱۳] در یک کارگاه تخصصی در متنی تحت عنوان ۵۰ سال از علم داده‌ها، اذعان داشته است که علم داده‌ها به واسطه‌ی بهره‌جستن از تکنیک‌های محاسباتی پیشرفته و حجم داده‌های مورد

⁵⁰Gelman

⁵¹Dhar

⁵²Donoho

تحلیل آن که از آمار حجیم‌تر است؛ از آمار متمایز نمی‌شود. او بر این باور است که علم داده‌ها تنها یک زمینه‌ی کاربردی است که در بستر آمار سنتی رشد کرده است.

نقد و بررسی: براساس نظر دونوهو [۱۳]، علم داده‌ها پدیده‌ای است که براساس دانش آمار طرح ریزی شده و علم داده‌ها را یک زمینه‌ی کاربردی می‌داند که از دانش آمار نشات می‌گیرد. این تعبیر دونوهو چیزی دور از این جمله نیست که علم داده‌ها فرزند آمار است و عملاً پدیده‌ای است که آمار را آماده‌ی پاسخ‌دهی به نیازهای متناسب با زمانه یعنی تحلیل داده‌های حجیم و بهره‌جستن از الگوریتم‌های رایانه‌ای کرده است.

در واقع خالی از ارزش نیست اگر بگوییم، علم داده‌ها تنها نسخه‌ای جدید از آمار است با ابزارهای پیشرفته‌ی محاسباتی که متناسب با نیازهای روز پژوهشگران در علوم کاربردی عمل می‌کند. این نسخه‌ی جدید از آمار را می‌توان آمار⁺ یا همان Statistics⁺ در نظر گرفت.

۵. نتیجه‌گیری

براساس موارد گفته شده، علم داده‌ها و آمار هر دو یک مقصد دارند. پیش‌بینی و تحلیل و تفسیر داده‌ها هدف اصلی هر دو است. داده‌های حجیم و نیاز مبرم به تکنیک‌های محاسباتی پیشرفته برای تحلیل و مدل‌بندی آن‌ها از جمله اصلی‌ترین انگیزه‌های ظهور گرایشی علمی به نام علم داده‌ها است. علم داده‌ها از همان مفاهیم و مبانی نظری موجود در آمار بهره می‌برد. بزرگترین وجه تمایز آمار و علم داده‌ها ابزارهای نوین نرم افزاری در تحلیل و مدل‌بندی داده‌ها است. علم داده‌ها با اتکا به مبانی آمار و ابزار نرم افزاری مدرن به تحلیل داده‌هایی می‌پردازد که در قرن ۲۱ به واسطه پیشرفت صنعت و تکنولوژی تولید شده‌اند. اگر اینگونه مطرح کنیم که علم داده‌ها همان آمار است که مبانی نظری آن را در نرم افزارهای پیشرفته پیاده‌سازی شده است؛ در واقع می‌توان نتیجه گرفت که علم داده‌ها در واقع نسخه‌ی پیشرفته‌ی علم آمار یا همان آمار⁺ است.

۶. تشکر و قدردانی

نویسندگان این مقاله، از داوران محترم بابت ارتقاء سطح کیفی مقاله با نظرات سازنده و مفیدشان و همچنین از دبیر رابط و سردبیر محترم نشریه بابت مد نظر قرار دادن مقاله کمال قدردانی و تشکر خود را ابراز می‌دارند.

مراجع

- [1] A. Agresti, *Categorical data analysis*, (Vol. 792), John Wiley & Sons, 2012.
- [2] A. Bhardwaj, *What's the Difference Between Data Science and Statistics?* Priceonomics. [https://priceonomics.com/whats-the-difference-between-data-science-and/\(2015\)](https://priceonomics.com/whats-the-difference-between-data-science-and/(2015)).
- [3] C.M. Bishop, *Pattern recognition and machine learning*, Springer, New York, 2006.
- [4] L. Breiman, J. Friedman, R. Olshen and C. Stone, *Classification and regression trees*, Chapman and Hall/CRC, 1984.

- [5] K. Broman, *Data science is statistics*, The Stupidest Thing: Statistics, Genetics, Programming, Academics, 2013.
- [6] I. Carmichael and J.S. Marron, *Data science vs. statistics: two cultures?* Jpn. J. Stat. Data Sci, **1**, (2018) 117–138.
- [7] G. Casella, S. Fienberg and I. Olkin, *Springer Texts in Statistics, Design* (Vol. 102) 2006.
- [8] J.M. Chambers, *Greater or lesser statistics: a choice for future research*, Statist. Comput., **3** (1993), no. 4, 182–184.
- [9] W.S. Cleveland, *Data science: an action plan for expanding the technical areas of the field of statistics*, International statistical review, **69** (2001) no. 1, 21–26.
- [10] R.D. De Veaux, M. Agarwal, M. Averett, B.S. Baumer, A. Bray, T.C. Bressoud, L. Bryant, L.Z. Cheng, A. Francis, R. Gould and A.Y. Kim, *Curriculum guidelines for undergraduate programs in data science*. Annu. Rev. Stat. Appl., **4** (2017) 15–30.
- [11] V. Dhar, *Data science and prediction*, *Communications of the ACM*, **56** (2013) no. 12, 64–73.
- [12] P.J. Diggle, *Statistics: a data science for the 21st century*, J. R. Stat. Soc., A: Stat. Soc. **178** (2015) no. 4, 793–813.
- [13] D. Donoho, *50 Years of Data Science*, J. Comput. Graph. Stat. **26** (2017) no. 4, 745–766.
- [14] H.E. Driver and A.L. Kroeber, *Quantitative expression of cultural relationships* Berkeley, University of California Press, 1932.
- [15] A. Gelman, *Statistics is the least important part of data science*, Statistical Modeling, Causal Inference, and Data Science blog. 2013.
- [16] I. Goodfellow, Y. Bengio and A. Courville, *Deep learning* (Vol. 1) USA: MIT Press, 2017.
- [17] A. Goodman, C. Kamath and V. Kumar, *Data analysis in the 21st century*, Statistical Analysis and Data Mining: The ASA Data Science Journal, **1** (2008) no. 1, 1–3.
- [18] J. Hardin, R. Hoerl, N.J. Horton, D. Nolan, B. Baumer, O. Hall-Holt, P. Murrell, R. Peng, P. Roback, D. Temple Lang and M. Ward, *Data science in statistics curricula: Preparing students to “think with data”*, The American Statistician **69**(2015) no. 4, 343–353.
- [19] T. Hastie, R. Tibshirani and J. Friedman, *The elements of statistical learning (2nd Edition)*, New York: Springer, 2009.
- [20] S.C. Hicks and R.A. Irizarry, *A guide to teaching data science*, The American Statistician, **72**(2018) no.4, 382–391.
- [21] A. Kassambara, *Practical guide to cluster analysis in R: Unsupervised machine learning*, France: Statistical Tools For High-Throughput Data Analysis, 2017.
- [22] M. Loukides, *What is data science?* O’Reilly Media, Inc. 2011.
- [23] D. Nolan and D. Temple Lang, *Computing in the statistics curricula*, The American Statistician, **64**(2010) no. 2, 97–107.

- [24] P. Pham, *The Impacts of big data that you may not have heard of*, Forbes. [https://www.forbes.com/sites/peterpham/2015/08/28/the-impacts-of-big-data-that-you-may-not-have-heard-of/\(2015\)](https://www.forbes.com/sites/peterpham/2015/08/28/the-impacts-of-big-data-that-you-may-not-have-heard-of/(2015)).
- [25] G. Press, *A very short history of big data*, Forbes Tech Magazine. [https://www.forbes.com/sites/gilpress/2013/05/28/a-very-short-history-of-data-science/\(2013\)](https://www.forbes.com/sites/gilpress/2013/05/28/a-very-short-history-of-data-science/(2013)).
- [26] N. Reid, *Statistical science in the world of big data*, Stat. Probab. Lett. **136**, (2018) 42–45.
- [27] N. Silver, *What I need from statisticians*, Stats & Data Science Views. [https://www.statisticsviews.com/article/nate-silver-what-i-need-from-statisticians/\(2013\)](https://www.statisticsviews.com/article/nate-silver-what-i-need-from-statisticians/(2013)).
- [28] D. Singpurwalla, *A handbook of statistics: An overview of statistical methods*, 2013.
- [29] D. Von Winterfeldt, and W. Edwards, *Decision analysis and behavioral research*, Cambridge University Press, Cambridge, 1986.
- [30] C. Weihs and K. Ickstadt, *Data science: the impact of statistics*, Int. J. Data Sci. Anal. **6** (2018) 189–194.