

# تخمین دبی نفت تولیدی از چاه به وسیله روش های یادگیری ماشین با استفاده از داده های پمپ الکتریکی شناور (ESP)

محمدباقر صدیقی<sup>۱</sup>، مجید سیاوشی<sup>۲\*</sup>، روح الدین میری<sup>۳</sup>

۱- دانشکده مهندسی مکانیک، دانشگاه علم و صنعت ایران

۲- دانشیار، دانشکده مهندسی مکانیک، دانشگاه علم و صنعت ایران

۳- استادیار، دانشکده مهندسی شیمی، نفت و گاز، دانشگاه علم و صنعت ایران

\* نویسنده مسئول: [msiavashi@iust.ac.ir](mailto:msiavashi@iust.ac.ir)

## چکیده

تخمین دبی جریان در چاه های یک میدان نفتی، یک فرایند حیاتی و کاربردی است. با این حال جریان های استخراج شده از چاه های نفتی، چند فازی بوده و تخمین دقیق دبی آن ها، بسیار چالش برانگیز و پرهزینه است. دبی سنج های مجازی در مقایسه با دبی سنج های چند فازی و روش های چاه آزمایشی، از نظر اقتصادی گزینه بسیار مناسبی هستند که با بهره گیری از داده های موجود و استفاده از روش های هوش مصنوعی، قادر به پیش بینی دقیق دبی در آینده هستند؛ بنابراین، اخیراً به دبی سنج های مجازی داده محور توجه زیادی شده است. در این مقاله تخمین دبی تولیدی یک چاه با استفاده از سه روش یادگیری ماشین: ۱- k همسایه نزدیک تر (k-NN)؛ ۲- تقویت گرادیان (GBR)؛ و ۳- درخت تصمیم (DT)، با استفاده از داده های پمپ انجام شده است. به منظور انتخاب ویژگی های مناسب به عنوان ورودی روش ها، از تحلیل های آماری پیرسون و اسپیرمن استفاده شده است. مجموعه داده مورد بررسی مربوط به یکی از چاه های یک میدان نفتی در جنوب ایران است. مجموعه داده موجود دارای حجم کم و تنوع ناکافی است، اما با این وجود نتایج نشان می دهند که روش های پیشنهادی عملکرد مناسبی دارند. روش k-NN با دقت ۰.۹۴۹۴ نسبت به دو روش دیگر عملکرد بهتری در تخمین دبی نفت داشته است. برای بررسی عملکرد روش ها در برابر داده های دارای نویز، یک درصد انحراف معیار نویز به داده های ورودی اضافه شد. بررسی ها نشان داد که مدل k همسایه نزدیک تر با دقت ۰.۹۲۵۷ در مقایسه با دو روش دیگر عملکرد بهتری داشته و کمترین تأثیر را از نویزها گرفته است.

کلیدواژگان

دبی سنج مجازی داده محور؛ یادگیری ماشین؛ k همسایه نزدیک تر؛ تقویت گرادیان؛ درخت تصمیم.

## Estimation of oil flow production of well employing machine learning algorithms using electrical submersible pump (ESP) data

### Abstract

Estimating the flow rate in oil wells of a field is a vital and practical process. However, the flows extracted from oil wells are multiphase, and their accurate estimation is highly challenging and costly. Virtual flow meters, compared to multiphase flow meters and well-testing methods, are an economically viable option that can accurately predict future flow rates by leveraging existing data and artificial intelligence algorithms. Therefore, data-driven virtual flow meters have recently received significant attention. This paper estimates the production flow rate of a well using three machine learning algorithms: 1- k-nearest neighbors; 2- gradient boosting; and 3- decision tree, using pump data. Pearson and

Spearman statistical analyses were used to select appropriate features as the algorithm inputs. The dataset under investigation pertains to one of the wells of a southern oil field in Iran. The available dataset has a small volume and insufficient diversity, but despite this, the results show that the proposed algorithms perform well. The k-NN method, with an accuracy of 0.9494, performed better than the other two methods in estimating oil flow rate. To examine the performance of the algorithms against noisy data, one percent of standard deviation noise was added to the input data. The investigations showed that the k-NN model, with an accuracy of 0.9257, performed better than the other two methods and was least affected by the noise.

**Keywords**

*Data driven virtual flow meter, machine learning, k-nearest neighbor, gradient boosting, decision tree*

Accepted Paper

## ۱- مقدمه

آگاهی از نرخ تولید سیالات استخراج‌شده از یک چاه در صنعت نفت و گاز به‌منظور بهینه‌سازی تولید و مدیریت مخزن ضروری است [۱]. روش‌های چاه آزمایی و استفاده از دبی سنج چند فاز (MPFM<sup>۱</sup>) دو روش متداول برای اندازه‌گیری دبی در میدان‌های نفتی به‌حساب می‌آیند. با این حال این روش‌ها به دلیل نیاز به نصب دستگاه‌های فیزیکی از لحاظ عملیاتی پرهزینه هستند. پژوهشگران باهدف یافتن راهی برای کاهش هزینه‌های عملیاتی برداشت نفت و گاز، دبی سنج‌های مجازی (VFM<sup>۲</sup>) را توسعه دادند. دبی سنج مجازی در واقع یک بسته نرم‌افزاری است که بدون نیاز به تجهیزات فیزیکی و تنها با استفاده از داده‌های در دسترس چاه در بخش‌های مختلف (مانند فشار و دما)، دبی جریان را پیش‌بینی می‌نماید. هزینه نصب، راه‌اندازی، نگهداری و اجرای دبی سنج‌های مجازی در مقایسه با سایر روش‌ها بسیار کمتر است [۲]. استفاده از فناوری VFM می‌تواند به عملکرد بهتر مجموعه‌های تولید نفت و گاز کمک کند، زیرا این فناوری به کاربران این امکان را می‌دهد تا با دقت و سرعت بیشتری به داده‌های نرخ جریان دسترسی داشته باشند و عملیات تولید را برای بهبود بازدهی بهینه کنند. همچنین دبی سنج مجازی می‌تواند از طریق محاسبه‌ی برخط نرخ جریان سیالات به کاربران در تشخیص و پیش‌بینی مشکلات چاه کمک کرده، فرآیند نگهداری از چاه را ارتقا داده و به‌تبع آن، نیاز به تعمیرات و هزینه‌ها را کاهش دهد.

دبی سنج‌های مجازی از نظر رویکرد محاسبه دبی به دو دسته دبی سنج مجازی فیزیکی محور و دبی سنج مجازی داده محور تقسیم می‌شوند [۳]. یکی از روش‌های محاسبه دبی سیال به‌صورت مجازی، استفاده از مدل‌سازی فیزیکی بخش‌های مختلف یک مجموعه تولید است. مدل‌های ارائه‌شده برای جریان‌های چند فاز بسیار پیچیده هستند و حل آن‌ها بدون ساده‌سازی امکان‌پذیر نیست. از طرفی حل این معادلات نیازمند صرف هزینه محاسباتی و زمان اجرا بالاست [۲]. در روش فیزیکی محور، جریان سیال با استفاده از معادلات بقا تحلیل شده و دبی تخمین زده می‌شود. با توجه به اینکه معادلات حاکم بر جریان سیالات چند فاز پیچیده هستند و حل تحلیلی برای برخی از آن‌ها وجود ندارد، در نتیجه معمولاً برای حل این معادلات ساده‌سازی زیادی بر روی آن‌ها انجام می‌شود یا از روش‌های حل عددی استفاده می‌شود که همین موجب بروز خطا و افزایش هزینه محاسبات می‌شود. عملکرد این نوع دبی سنج‌ها وابستگی زیادی به مقدار فشار، حجم و دما دارد. تنظیم این نوع دبی سنج‌ها کار دشواری است زیرا در این روش کاربر باید دانش کافی در زمینه عملیاتی و نرم‌افزاری داشته باشد [۲].

دومین رویکرد در دبی سنج‌های مجازی استفاده از داده‌های میدانی برای تخمین دبی است. در این روش تلاش بر این است که با استفاده از داده‌های جمع‌آوری‌شده از میدان شامل داده‌های به‌دست‌آمده از حسگرهای درون‌چاهی و چاه آزمایی، یک مدل ریاضی به کمک روش‌های پیشرفته یادگیری ماشین جهت تخمین دبی ایجاد می‌شود [۲]. اگر مدل مبتنی بر داده به‌خوبی آموزش داده شود می‌تواند دبی فازها را به‌صورت آنی و دقیق پیش‌بینی کند. مزیت اصلی این روش نسبت به دبی سنج‌های مجازی فیزیکی محور این است که از مدل‌سازی‌های فیزیکی-دقیق سیستم‌ها استفاده نمی‌کنند زیرا این مدل‌سازی‌ها معمولاً غیرخطی بوده و حل دقیق آن‌ها بسیار دشوار است. روش‌های داده محور بر این اصل تکیه‌دارند که داده‌های تجربی و عملکردی، سیستم را به‌خوبی نمایندگی می‌کنند و سعی می‌کنند روابط فیزیکی سیستم را مستقیماً از داده‌ها یاد بگیرند.

<sup>۱</sup> Multi-Phase Flow Meter

<sup>۲</sup> Virtual Flow Meter

القطامی و همکاران [۴] یک سیستم دبی سنج مجازی مبتنی بر یادگیری گروهی ایجاد کردند. آن‌ها تلاش کردند با دست‌کاری داده‌های آموزشی، معماری شبکه عصبی و مسیر یادگیری یادگیرندگان، شبکه عصبی متنوعی را تولید کنند. هدف از این کار این بود که عملکرد دبی سنج‌های مجازی مبتنی بر داده در میدان‌هایی با داده‌های تولیدشده بسیار محدود بهبود یابد. دبی سنج مجازی توسعه‌یافته دقت بسیار مطلوبی دارد. این سیستم می‌تواند دبی فاز مایع و گاز را به ترتیب با میانگین خطا ۴.۷ و ۲.۴ درصد تخمین بزند. القطامی و همکاران [۵] یادگیری یک مجموعه ترکیبی را با ترکیب شبکه عصبی و رویکردهای درخت رگرسیون<sup>۱</sup> مورد بحث قرار دادند. ایشان رویکرد ترکیبی (NN-RTs)<sup>۲</sup> را با رویکردهای مجموعه همگن (RTs و NN) مقایسه کردند و نشان دادند روش ترکیبی عملکرد دقیق‌تری را نشان می‌دهد. احمدی و همکاران [۶] در یک مسئله جریان سنجی چند فاز، روش جدیدی برای پیش‌بینی نرخ نفت تولیدی چاه‌ها با استفاده از منطق فازی، شبکه‌های عصبی مصنوعی و روش رقابت امپراتوری (ICA)<sup>۳</sup> ارائه کردند. دماها و فشارهای خطوط به‌عنوان متغیر ورودی شبکه و نرخ جریان نفت به‌عنوان متغیر خروجی تنظیم شدند. یک مجموعه داده متشکل از ۱۶۰۰ داده‌ی مربوط به ۵۰ چاه دریکی از میدان‌های نفتی واقع در شمال خلیج فارس برای ساخت پایگاه داده استفاده شده است. نتایج به‌دست‌آمده کارآیی، قدرت و سازگاری مدل ICA-ANN<sup>۴</sup> را نشان دادند. بیک مخامدوف و جاشکه [۷] عملکرد روش تقویت‌گرادیان را به‌عنوان جایگزینی برای شبکه‌های عصبی در یک مجموعه دبی سنج مجازی بررسی کردند. آن‌ها نشان دادند که این سازوکار به‌عنوان پشتیبان دبی سنج چند فاز و یا به‌طور مستقل حتی برای مجموعه داده آموزشی کوچک عملکرد خوبی دارد. نتایج همچنین نشان داد که با ترکیب روش تقویت‌گرادیان با اندازه‌گیری‌های دبی آزمایش چاه در محدوده عملیاتی گسترده‌ای از چاه، می‌توان پیش‌بینی دقیق دبی را از مرحله تولید اولیه انجام داد.

گوئز و همکاران [۸]، هاتودت و همکاران [۹] و العجمی و همکاران [۱۰] از داده‌های شیر فشارشکن برای توسعه مدل داده محور استفاده کردند. مدل ایجادشده توسط گوئز و همکاران دبی هر چاه را به‌عنوان تابعی از داده‌های جمع‌آوری شده مانند دما و فشار شیر فشارشکن و ویژگی سیالات در میدان نفتی بیان می‌کند. نتایج نشان می‌دهد که برای یک دوره ۲۴ ساعته، خطاهای نسبی بین نرخ جریان پیش‌بینی شده و اندازه‌گیری شده به ترتیب زیر ۳.۵ و ۳ درصد برای کل نرخ جریان نفت و گاز بود. هاتودت و همکاران یک دبی سنج مجازی جدید با ترکیب روش فیزیک محور و روش داده محور ایجاد کردند. در این پژوهش از رویکرد ترکیبی برای مدل‌سازی شیر فشارشکن استفاده شده است. شیر فشارشکن با مجموعه‌ای ساده از معادلات پایه و یک شبکه عصبی برای تخمین ضریب جریان در چاه نشان داده می‌شود. آن‌ها نشان دادند که رویکرد ترکیبی عملکرد مناسبی دارد و ممکن است مزایایی نسبت به دبی سنج‌های مجازی فیزیک محور و داده محور داشته باشد. العجمی و همکاران علاوه بر فشار، دما، اندازه شیر فشارشکن و داده‌های برش آب از برخی پارامترهای اضافی برای ورودی‌ها (از جمله یک همبستگی تجربی برای جریان بحرانی شیر فشارشکن) استفاده کردند. در مقایسه با تخمین‌های دبی جریان به‌دست‌آمده با استفاده از همبستگی‌های تجربی شیر فشارشکن، شبکه عصبی (NN) عملکرد نسبتاً بهتری را نشان داد. لازم به ذکر است که مدل‌های شیر فشارشکن مورد استفاده در این مطالعه تجربی بودند.

<sup>1</sup> Decision Tree

<sup>2</sup> NN-RTE: Neural Network-Regression Trees

<sup>3</sup> Imperialist competitive algorithm

<sup>4</sup> ANN: Artificial Neural Networks

سندس و همکاران [۱۱] یک دبی سنج مجازی مبتنی بر داده با استفاده از معماری یادگیری چندوظیفه‌ای برای مطالعه ۵۵ حلقه چاه پیشنهاد کردند. نتایج با دو مدل پایه تک‌کاره مقایسه شده است. آن‌ها نشان دادند که الگوریتم یادگیری چند وظیفه‌ای استحکام را نسبت به روش‌های تک‌وظیفه‌ای، بدون به خطر انداختن عملکرد، بهبود می‌بخشد. این الگوریتم به‌طور متوسط برای دارایی‌هایی که در آن معماری تک‌وظیفه‌ای با مشکل مواجه هستند، خطا را به میزان ۲۵ تا ۵۰ درصد کاهش می‌دهد.

الجاسمی و همکاران [۱۲] یک مدل داده محور با استفاده از شبکه‌های عصبی برای پیش‌بینی دبی مایع و برش آب در یک مخزن کربناته دریایی با بیش از ۲۰ درصد برش آب را ارائه دادند. شبکه عصبی با استفاده از داده‌های تولید واقعی، داده‌های سطحی و ته-چاهی آموزش داده شد. داده‌هایی که برحسب زمان بودند به‌عنوان سری زمانی در نظر گرفته شدند تا به کاربران اجازه دهد سناریوهایی را با تغییر عملیات چاه ایجاد کنند. این رویکرد با تغییر متغیرهای کنترلی مانند فشار سر چاه و فرکانس پمپ، نتایج را شبیه‌سازی می‌کند. تغییر فشار سر چاه و فرکانس، به کاربران اجازه می‌دهد تولید را پیش‌بینی کرده و رویدادهای منفی پمپ را پیش‌بینی کنند. با وجود داده‌های محدود مخزن، نتایج نشان می‌دهد که شبکه‌های عصبی دقت قابل قبولی در پیش‌بینی دبی نفت و برش آب داشته است.

علاوه بر برنامه‌های کاربردی VFM توسط شبکه عصبی، شرکت بیکر هیوز<sup>۱</sup> نرم‌افزار نئورال فلو<sup>۲</sup> را توسعه داده است که مبتنی بر مدل شبکه عصبی است [۱۳]. این نرم‌افزار برای تخمین دبی در سیستم‌های دارای پمپ‌های شناور الکتریکی<sup>۳</sup> (ESP) با استفاده از رویکرد شبکه عصبی استفاده می‌شود. مشابه شبکه‌های عصبی‌ای که قبلاً مورد بحث قرار گرفت، این سازوکار فشار سر چاه، فشار ورودی و فشار تخلیه پمپ و همچنین سایر پارامترهای اندازه‌گیری شده مانند فرکانس و جریان در پمپ را به‌عنوان ورودی در نظر می‌گیرد و تخمین دبی را به‌عنوان خروجی شبکه تولید می‌کند.

در تحقیقات گذشته از پارامترهای مختلف درون چاهی مانند فشار و دمای سرچاهی و ته چاهی، فشار خط، کسر حجمی فازها و اطلاعات شیر فشارشکن مانند میزان بازشدگی شیر، فشار قبل و بعد شیر، دمای قبل و بعد شیر و ضریب تخلیه شیر ... برای تخمین نرخ دبی‌ها استفاده شده است و کمتر به تخمین دبی با استفاده از پمپ ESP به‌عنوان یک مدل داده محور پرداخته شده است. در حالی که تمامی این اطلاعات در چاه‌ها به‌صورت بر خط در دسترس نیست. از آنجاکه اکثر چاه‌های ایران وارد دوره کاهش برداشت شده و برای استخراج نفت نیاز به پمپ ESP وجود دارد، بنابراین اطلاعات این پمپ‌ها در اکثر چاه‌ها موجود است و تخمین دبی از این طریق می‌تواند بسیار ساده و مفید باشد.

لازم به ذکر است که در کارهایی که دبی سنجی با اطلاعات پمپ انجام شده (مانند پژوهش دنی و همکاران [۱۳] که موجب ایجاد یک نرم‌افزار تجاری شده است) نیز از مدل‌های شبکه عصبی استفاده شده است. این در حالی است که چنین روش‌هایی به حجم بالایی از داده با تنوع بالا نیاز دارد که فراهم کردن این مورد در بسیاری از چاه‌ها ممکن نیست. این مشکل یک چالش در استفاده از روش‌های شبکه عصبی در چاه‌هایی است که حجم بالا داده با تنوع زیاد را تأمین نمی‌کنند. این تحقیق باهدف تخمین دبی نفت در یکی از چاه‌های مربوط به یک میدان نفتی واقع در جنوب ایران انجام شده است. چاه مورد بررسی از حجم داده بالا و تنوع کافی برخوردار نیست بنابراین، مدل‌های شبکه عصبی در این چاه با خطا قابل توجهی

<sup>1</sup> Baker Hughes

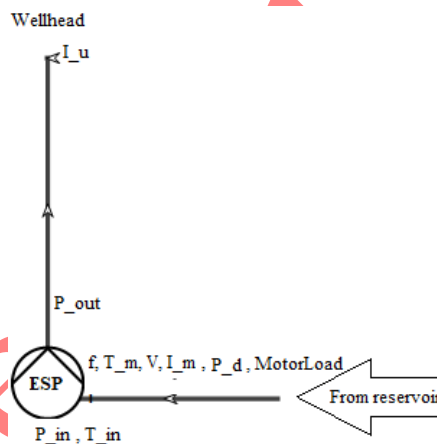
<sup>2</sup> NeuraFlow

<sup>3</sup> Electrical Submersible Pump

روبه‌رو هستند. به همین دلیل در این مطالعه، از روش‌های  $k$  همسایه نزدیک‌تر ( $k$ -NN)<sup>۱</sup>، درخت تصمیم (DT)<sup>۲</sup> و تقویت گرادیان (GBR)<sup>۳</sup> استفاده شده است. همچنین، داده‌های نویز دار که ممکن است به دلیل خرابی حسگرها در مجموعه داده پدیدار شوند، می‌توانند تأثیر معکوس قابل توجهی بر دقت و قابلیت اعتماد مدل‌ها داشته باشند. این نوع داده‌ها می‌توانند تأثیر منفی بر عملکرد دبی سنج مجازی داده محور داشته باشند؛ بنابراین، عملکرد دبی سنج مجازی داده محور در برابر داده‌های نویز دار نیز مورد بررسی قرار گرفته است. در بخش دوم مدل‌های مورداستفاده توضیح داده شده‌اند. بررسی مجموعه داده در بخش سه آورده شده است. نتایج در بخش ۴ آورده شده و در نهایت نتیجه‌گیری در بخش ۵ آمده است.

## ۲- تعریف مسئله

شماتیک چاه مورد مطالعه در شکل ۱ آمده است. در این تحقیق فقط پمپ ESP مدل‌سازی شده است. این منجر به پیچیدگی کمتر مدل می‌شود و نیاز به داده‌های ته چاه و سر چاه را از بین می‌برد. این کار برای دارایی‌هایی که اندازه‌گیری‌های ته چاه و یا بالای چاه وجود ندارد یا معیوب هستند، سودمند است. شکل ۱ پارامترهای اندازه‌گیری شده پمپ برای تخمین دبی را نشان می‌دهد. این پارامترها شامل جریان کنتور کنترل‌کننده ( $I_u$ )، فشار ورودی پمپ ( $P_{in}$ )، فشار خروجی پمپ ( $P_{out}$ )، فشار تخلیه ( $P_d$ )، دمای ورودی به پمپ ( $T_{in}$ )، دمای موتور ( $T_m$ )، جریان موتور ( $I_m$ )، ولتاژ ( $V$ )، فرکانس ( $f$ ) و بار موتور ( $MotorLoad$ ) هستند.



شکل ۱. شماتیک چاه مورد مطالعه

برای ساخت یک دبی سنج مجازی داده محور ابتدا باید داده‌های مختلف جمع‌آوری شوند سپس از بین این داده‌ها آن‌هایی که ارتباط قوی‌تری با تابع هدف<sup>۴</sup> دارند به‌عنوان ورودی انتخاب شده و آنگاه پیش‌پردازش روی آن‌ها انجام می‌گیرد. یعنی داده‌های پرت حذف و نویزها فیلتر می‌شوند. در گام بعدی داده‌های پیش‌پردازش شده وارد مرحله توسعه مدل می‌شود. در این قسمت مجموعه داده به دو قسمت آموزش و آزمون تقسیم می‌گردد؛ روش سعی می‌کند با استفاده از داده‌های آموزشی یک الگوریتم تکرارشونده را بدون دخالت انسان یاد بگیرد. در نهایت مدل با استفاده از داده‌های بخش آزمون اعتبار سنجی می‌شود. اطلاعات مربوط به مدل داده محور در بخش ۳ آمده است.

<sup>1</sup> k-Nearest Neighbor.

<sup>2</sup> Decision Tree.

<sup>3</sup> gradient boosting

<sup>4</sup> تابع هدف پارامتری است که قرار است مقدار آن تخمین زده شود

### ۳- مدل داده محور

پیشرفت‌های سریع و شگفت‌انگیز روش‌های یادگیری ماشین باعث شده این روش‌ها در صنایع مختلف مورد توجه پژوهشگران قرار گیرد. اصل اساسی یادگیری ماشینی این است که رایانه‌ها را قادر می‌سازد از داده‌ها یاد بگیرند و رفتار خود را بر این اساس تطبیق دهند. این شامل آموزش روش‌های حل با استفاده از مجموعه داده‌های بزرگ برای شناسایی الگوها، همبستگی‌ها و بینش‌هایی است که می‌توانند برای پیش‌بینی یا تصمیم‌گیری دقیق استفاده شوند. در پژوهش حاضر از سه روش DT، k-NN و روش GBR به‌عنوان مدل‌های داده محور برای پیش‌بینی نرخ جریان استفاده شده است. در ادامه مراحل ایجاد مدل داده محور و همچنین پیشینه هر یک از روش‌های انتخاب‌شده شرح داده می‌شود.

### ۳-۱- انتخاب ویژگی

در دبی سنج‌های داده محور پارامترهای ورودی که به‌عنوان ویژگی‌ها یا متغیرها شناخته می‌شوند، بلوک‌های سازنده مدل‌های یادگیری ماشین هستند. این پارامترها ویژگی‌های داده‌ها را نشان می‌دهند که روش حل برای پیش‌بینی یا طبقه‌بندی از آن‌ها استفاده می‌کند. انتخاب پارامترهای ورودی مناسب می‌تواند به بهبود عملکرد مدل کمک کند، درحالی‌که پارامترهای نادرست یا نامربوط می‌تواند منجر به نتایج غیر بهینه و برازش بیش‌ازحد شود. انتخاب پارامترهای ورودی مناسب برای روش یادگیری ماشین می‌تواند یک فرایند پیچیده و چالش‌برانگیز باشد، زیرا داده‌ها ممکن است حاوی متغیرهای متعددی باشند که برخی از آن‌ها ممکن است نامربوط، زائد یا نویز دار باشند. با انتخاب ویژگی‌های مرتبط‌تر، می‌توانیم از تأثیرپذیری مدل از نویز و اطلاعات نامربوط داده‌ها جلوگیری کنیم که می‌تواند به تعمیم بهتر داده‌های دیده نشده کمک کند. همچنین، با تعداد ویژگی‌های کمتر، مدل به زمان و منابع کمتری برای آموزش نیاز دارد که آن را برای کاربردهای دنیای واقعی کارآمدتر و مقیاس‌پذیرتر می‌کند.

علاوه بر این، انتخاب ویژگی مناسب به افزایش تفسیرپذیری مدل کمک می‌کند. با تمرکز بر مهم‌ترین ویژگی‌ها، می‌توان به بینش‌هایی در مورد عواملی دست‌یافت که بیشترین تأثیر را بر پیش‌بینی‌های مدل دارند و برای درک و توضیح رفتار مدل برای ذینفعان و کارشناسان حوزه بسیار مهم است. روش‌های مختلفی برای انتخاب ویژگی مناسب وجود دارد که هر کدام نقاط قوت و ضعف خاص خود را دارند. روش‌های فیلتر یک رویکرد محبوب برای انتخاب ویژگی هستند، زیرا این روش‌ها از معیارهای آماری برای رتبه‌بندی و انتخاب ویژگی‌ها بر اساس ارتباط بین آن‌ها با متغیر هدف استفاده می‌کنند معیارهای رایج مورد استفاده در روش‌های فیلتر عبارتند از انتخاب ویژگی مبتنی بر همبستگی<sup>۱</sup>، اطلاعات متقابل<sup>۲</sup> و آزمون مجذور کای<sup>۳</sup>. ویژگی‌ها بر اساس معیارهای از پیش تعریف‌شده مانند آستانه مشخصی از مقدار متریک انتخاب یا حذف می‌شوند. ماتریس همبستگی اسپیرمن و پیرسون روش‌های آماری هستند که برای اندازه‌گیری قدرت و جهت رابطه بین دو متغیر استفاده می‌شوند. در زمینه انتخاب ویژگی، این روش‌ها می‌توانند برای شناسایی مرتبط‌ترین ویژگی‌ها برای یک کار پیش‌بینی معین استفاده شوند.

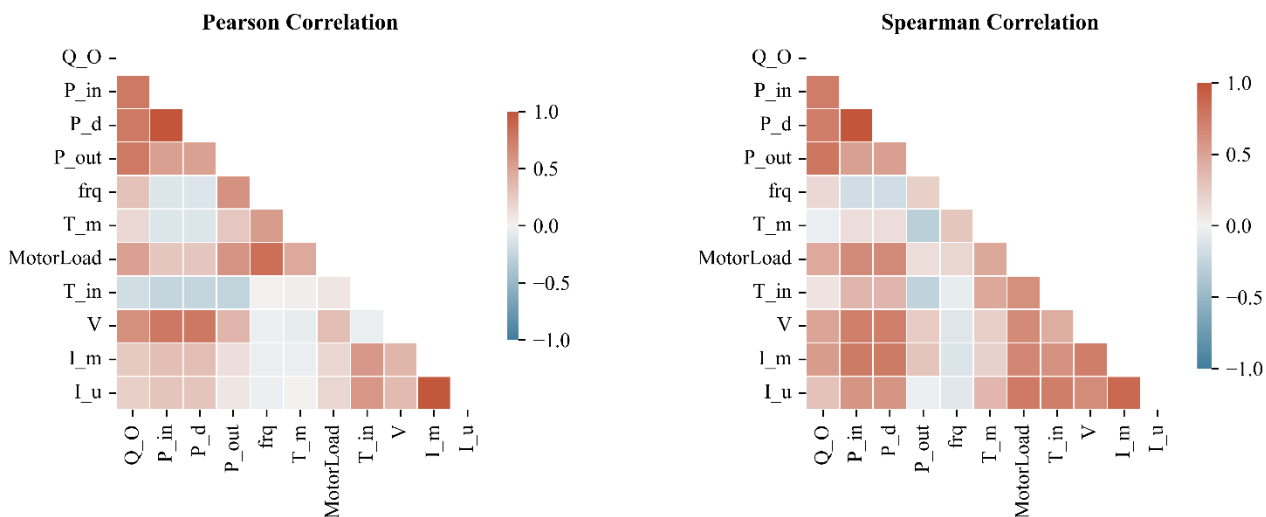
<sup>1</sup> Correlation-based

<sup>2</sup> mutual information

<sup>3</sup> Chi-Square

همبستگی پیرسون رابطه خطی بین دو متغیر پیوسته را اندازه‌گیری می‌کند، در حالی که همبستگی اسپیرمن رابطه یکنواخت بین دو متغیر را بدون توجه به خطی بودن یا نبودن آن اندازه‌گیری می‌کند. هر دو روش بینش‌های ارزشمندی را در مورد روابط بین ویژگی‌ها ارائه می‌دهند که می‌تواند برای شناسایی ویژگی‌های اضافی یا نامربوط برای پیش‌بینی استفاده شود.

هنگام استفاده از ماتریس همبستگی پیرسون و اسپیرمن برای انتخاب ویژگی، ویژگی‌هایی که با یکدیگر ضریب همبستگی بالا (نزدیک به ۱ یا -۱) داشته باشند ممکن است برای پیش‌بینی اضافی باشند. از بین این ویژگی‌ها تنها یک مورد انتخاب می‌شود چون اگر دو ویژگی با هم مرتبط باشند، مدل فقط به یکی نیاز دارد، زیرا دومی اطلاعات جدیدی برای بهبود عملکرد آموزش مدل اضافه نمی‌کند. از سوی دیگر، ویژگی‌هایی با ضرایب همبستگی پایین (نزدیک به ۰) ممکن است کمتر برای پیش‌بینی مرتبط باشند و همچنین می‌توانند برای حذف در نظر گرفته شوند. در شکل ۲ ارتباط بین ویژگی‌های مختلف با معیارهای اسپیرمن و پیرسون نشان داده شده است. با توجه به نتایج می‌توان نتیجه گرفت که فشار ورودی ( $P_{in}$ )، فشار خروجی ( $P_{out}$ )، بار موتور ( $MotorLoad$ )، جریان موتور ( $I_m$ ) و ولتاژ موتور ( $V$ ) به‌عنوان ویژگی‌های مناسب برای تخمین دبی در نظر گرفته می‌شوند. کامیلری و همکاران در سال ۲۰۱۰ از معادله بقای انرژی در پمپ برای تخمین دبی استفاده کردند [۱۴]. طبق این معادله توان جذب‌شده توسط پمپ برابر با توان تولیدشده توسط موتور است. توان پمپ تابعی از اختلاف فشار دبی و بازدهی بوده و از سوی دیگر توان موتور تابعی از ولتاژ موتور، جریان موتور، ضریب توان و بازدهی موتور است؛ بنابراین در یک پمپ ESP دبی تابعی از اختلاف فشار قبل و بعد از پمپ، جریان و ولتاژ موتور است. نتایج این مطالعه نشان می‌دهد که ویژگی‌های انتخاب‌شده با نرخ دبی ارتباط دارند.



شکل ۲. ارتباط بین نرخ جریان نفت و ویژگی‌ها با استفاده از معیارهای اسپیرمن (سمت راست) و پیرسون (سمت چپ)

### ۲-۳- روش‌های یادگیری ماشین

#### ۱-۲-۳- GBR

روش تقویت‌کننده گرادیان (GBR) یک شیوه یادگیری ماشین قدرتمند است که هم برای مسائل رگرسیون و هم برای طبقه‌بندی استفاده می‌شود. این روش به روش‌های یادگیری گروهی تعلق دارد که چندین مدل را برای بهبود عملکرد پیش‌بینی ترکیب می‌کنند. در مورد GBR، مدل‌ها معمولاً درخت تصمیم هستند.



ایده اصلی رگرسیون GBR، آموزش مکرر مدل‌های جدید است که باقی‌مانده یا خطاهای مدل‌های قبلی را پیش‌بینی می‌کنند. سپس این مدل‌های جدید به مجموعه اضافه می‌شوند و هرکدام بر روی خطاهایی تمرکز می‌کنند که توسط مدل‌های قبلی به خوبی پیش‌بینی نشده بودند. با انجام این کار به صورت مکرر، گروه به تدریج عملکرد پیش‌بینی خود را بهبود می‌بخشد [۱۵]. یکی از مزایای کلیدی GBR، توانایی آن در مدیریت روابط پیچیده و غیرخطی در داده‌ها است. علاوه بر قدرت پیش‌بینی آن، GBR معیاری از اهمیت ویژگی را نیز ارائه می‌دهد که می‌تواند برای درک اینکه کدام ویژگی‌ها در پیش‌بینی‌ها بیشترین تأثیر را دارند مفید باشد. این می‌تواند برای به دست آوردن بینش در مورد فرآیندهای اساسی که داده‌ها را هدایت می‌کند و برای تصمیم‌گیری آگاهانه بر اساس خروجی مدل ارزشمند باشد. یک مدل GBR با تعداد  $M$  درخت را می‌توان به صورت زیر تعریف کرد.

$$f_M(x_j) = \sum_m^M \gamma_m h_m(x_j) \quad (1)$$

در رابطه فوق  $h_m$  نشان‌دهنده یک یادگیرنده ضعیف است که به تنهایی عملکرد ضعیفی دارد و به عنوان یک عامل مقیاس برای اضافه کردن سهم یک درخت به مدل عمل می‌کند. روش GBR از تابع تلفات نزولی گرادینان برای به حداقل رساندن خطاها با به‌روزرسانی تخمین اولیه با تخمین جدید استفاده می‌کند؛ که نتیجه آن ایجاد یک مدل نهایی با ادغام تمام تخمین‌های اولیه با وزن‌های مناسب است.

### ۳-۲-۲- روش k-NN

روش  $k$  همسایه نزدیک‌تر ( $k$ -NN) اولین بار در سال ۱۹۷۶ توسط دودانی [۱۶] توسعه یافت در سال ۱۹۸۳ جوزویک [۱۷] نسخه پیشرفته‌تر این روش را معرفی کرد. در زمینه یادگیری ماشین،  $k$ -NN یک روش ناپارامتریک<sup>۱</sup> پرکاربرد برای طبقه‌بندی و رگرسیون است. این روش بر این اصل استوار است که نقاط داده مشابه در یک فضای چندبعدی به یکدیگر نزدیک هستند. روش  $k$ -NN با یافتن  $k$  همسایه نزدیک یک نقطه داده بر اساس فاصله (مانند فاصله اقلیدسی یا فاصله منهتن) پیش‌بینی یا طبقه‌بندی برای نقطه داده جدید را انجام می‌دهد. مقدار  $k$  تعیین می‌کند که چه تعداد همسایه در فرآیند پیش‌بینی در نظر گرفته شود. یکی از شاخص‌ترین مزایای این روش، عملکرد مناسب در مدل‌هایی است که حجم داده آن‌ها کم است. البته باید توجه داشت که این روش قادر نیست مقادیر بزرگ‌تر از بیشترین مقدار مشاهده‌شده و کوچک‌تر از کمترین مقدار مشاهده‌شده را تولید کند. به عبارت دیگر این روش تنها توانایی درون‌یابی بین داده‌ها را دارد و قادر به انجام برون‌یابی نیست [۱۸].

$k$ -NN با مقایسه یک نمونه آزمون داده‌شده  $(x, y)$  با یک مجموعه آموزشی  $D = [(x_i, y_i)]$  یاد می‌گیرد. روش  $k$ -NN ابتدا فاصله بین  $x$  و هر نمونه  $x_i$  در  $D$  را محاسبه می‌کند. فاصله در اینجا با نماد  $d$  بیان شده است. مقدار  $d_i$  از رابطه‌های زیر محاسبه می‌شود.

$$d_{Euc} = \sqrt{(x - x_i)^2} \quad (2)$$

<sup>۱</sup> منظور این است که در این روش رابطه پارامتری از پیش تعیین شده‌ای میان متغیرهای ورودی و خروجی برقرار نمی‌شود.

$$d_{Man} = |x - x_i| \quad (3)$$

که در آن  $d_{Ecu}$  فاصله اقلیدسی و  $d_{Man}$  فاصله منهن است. در مرحله بعد k-NN داده‌ها را بر اساس مقدار  $d$  از کوچک به بزرگ مرتب می‌کند؛ پس از رتبه‌بندی، نمونه مرتبط با  $d_i$  به‌عنوان نزدیک‌ترین همسایه  $i$  نامیده می‌شود و خروجی با  $y_i(x)$  نشان داده می‌شود. خروجی نهایی  $y_{pred}$  میانگین وزن دار خروجی k همسایه نزدیک آن است، همان‌طور که در معادله (۴) نشان داده شده است [۱۹].

$$y_{pred} = \frac{\sum_{i=1}^k w_i y_i(x)}{\sum_{i=1}^k w_i} \quad (4)$$

در رابطه بالا  $y_{pred}$  مقدار تخمین زده شده و  $y_{pred}$  وزن هر همسایه است که برابر با معکوس فاصله آن است.

### ۳-۲-۳- روش DT

درخت تصمیم (Decision Tree) یک روش یادگیری ماشین ساده و محبوب است که برای طبقه‌بندی داده‌ها و رگرسیون (پیش‌بینی مقادیر عددی) استفاده می‌شود. این روش از یک نمودار درخت مانند برای مدل‌سازی تصمیمات و پیامدهای مختلف استفاده می‌کند [۲۰].

در این روش، فضای ورودی به مناطق کوچک‌تر تقسیم می‌شود و در هر منطقه، یک مدل ساده برای پیش‌بینی ایجاد می‌شود. در رگرسیون با درخت تصمیم، هدف پیش‌بینی مقدار یک متغیر هدف با استفاده از قوانین تصمیم‌گیری ساده است. در واقع، درخت تصمیم یک ساختار درخت مانند دارد که هر گره داخلی آن یک تصمیم بر اساس یک ویژگی از داده‌ها را نشان می‌دهد و هر گره برگ نشان‌دهنده مقدار پیش‌بینی شده متغیر هدف است. قوانین تصمیم‌گیری از داده‌های آموزشی یاد گرفته می‌شوند و ساختار درخت بر اساس ویژگی‌هایی تعیین می‌شود که به بهترین شکل داده‌ها را در هر گره تقسیم می‌کند. [۲۱]. یکی از مزایای اصلی این روش، سادگی و تفسیرپذیری آن است. درخت‌های تصمیم را می‌توان به راحتی تجسم کرد و حتی توسط افراد غیرمتخصص قابل درک است. این درخت‌ها به‌عنوان ابزاری مفید برای تحلیل داده‌ها و فرآیندهای تصمیم‌گیری پیچیده مورد استفاده قرار می‌گیرند. همچنین، درخت‌های تصمیم قادر به مدیریت داده‌های عددی و مقوله‌ای هستند و نسبت به مقادیر پرت و گم شده نیز مقاومت دارند.

### ۳-۳- تنظیم پارامترها

تنظیم پارامترها در الگوریتم‌های یادگیری ماشین اهمیت بالایی دارد، زیرا بر کیفیت مدل‌سازی و دقت پیش‌بینی‌ها تأثیر می‌گذارد. در واقع تنظیم پارامتر مناسب برای اطمینان از اینکه مدل به‌طور مؤثر از داده‌های آموزشی یاد می‌گیرد و به‌خوبی به داده‌های دیده نشده تعمیم داده می‌شود؛ ضروری است. یکی از رویکردهای رایج برای تنظیم پارامتر، جستجوی شبکه‌ای است که در آن مجموعه‌ای از مقادیر از پیش تعریف شده برای هر پارامتر مشخص می‌شود و مدل برای هر ترکیبی از مقادیر، آموزش و ارزیابی می‌گردد. روش دیگر جستجوی تصادفی است که از فضای پارامتر به‌صورت تصادفی نمونه‌برداری می‌کند. هر دو روش مزایا و معایب خود را دارند. در این پژوهش، روش جستجوی شبکه‌ای برای تنظیم پارامترها انتخاب شده

است. برای الگوریتم  $k$ -NN، پارامترهای کلیدی شامل تعداد همسایه‌ها ( $k$ )، نحوه وزن دهی و معیار فاصله (مانند فاصله اقلیدسی) هستند. در درخت تصمیم، عمق درخت و تعداد نمونه‌ها در هر برگ از جمله پارامترهای مهم به شمار می‌روند. تقویت گرادیان نیز از پارامترهایی مانند نرخ یادگیری، تعداد درختان و عمق درختان استفاده می‌کند. این پارامترها هر کدام جنبه‌های مختلفی از الگوریتم را کنترل می‌کنند و برای دستیابی به بهترین عملکرد، باید با دقت تنظیم شوند. در جدول ۱، مقادیر پیشنهادی برای پارامترهای مذکور آورده شده است.

جدول ۱. مقدار پارامترهای مهم الگوریتم‌های یادگیری ماشین استفاده شده

مقدار پیشنهادی	پارامتر	الگوریتم
۵	تعداد همسایه‌ها ( $k$ )	$k$ -NN
فاصله اقلیدسی	نوع وزن دهی	DT
۷	نوع فاصله	
۵	عمق درخت	GBR
۰.۱	نمونه‌ها	
۵۰۰	نرخ یادگیری	
۲	تعداد درختان	
	عمق درختان	

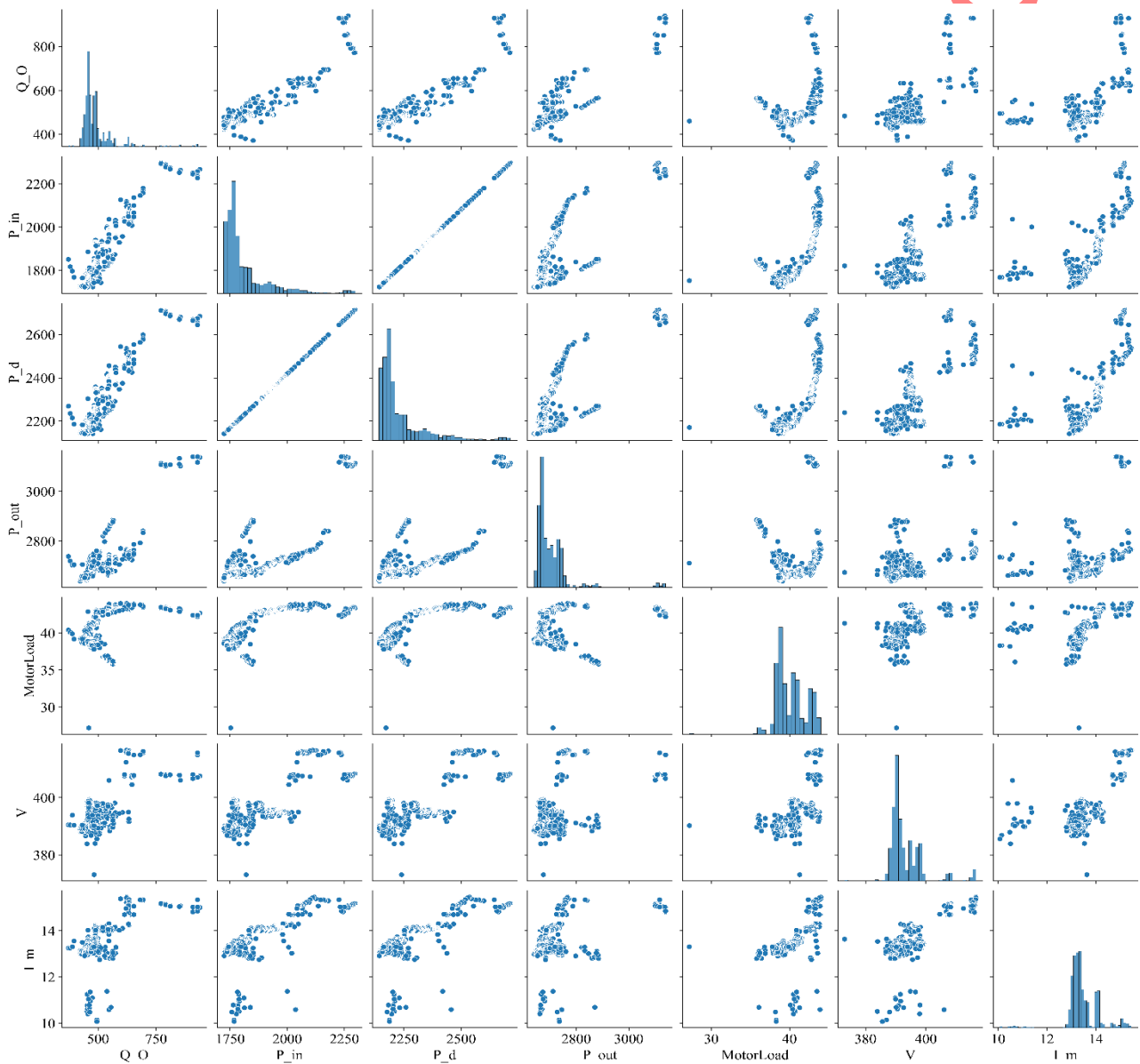
#### ۴- توضیحات داده‌ها

اطلاعات مورد استفاده در این مطالعه به صورت یک روز در میان یا دو روز در میان و از تاریخ ۱۰ اکتبر ۲۰۲۰ تا ۸ جولای ۲۰۲۳ جمع‌آوری شده است و تعداد کل داده‌ها ۹۸۰ عدد است. در این بازه زمانی تغییرات زیادی در چاه‌های میدان صورت نگرفته و همچنین عملیات تزریق آب و یا گاز و سایر فرآیندهایی که می‌تواند منجر به تغییر جدی در مکانیزم‌های تولیدی گردد، گزارش نشده است. در بخش پیش‌پردازش داده‌ها ابتدا داده‌های گمشده از مجموعه داده حذف گردید. سپس با بررسی بیشتر داده‌ها مشخص شد که برخی از داده‌ها در فاصله بسیار زیاد از میانگین قرار دارند. به عنوان مثال دبی نفت در بعضی از روزها صفر گزارش شده بود، از آنجاکه این بخش از داده‌ها فاصله قابل توجهی با سایر داده‌ها و میانگین داشتند از مجموعه داده حذف شدند. پس از حذف داده‌های پرت و گمشده حجم داده‌ها به ۸۵۰ مورد کاهش یافت. جدول ۲ اطلاعات آماری جامعی از مجموعه داده را نشان می‌دهد.

جدول ۲. اطلاعات آماری مجموعه داده پس از انجام پیش‌پردازش

پارامترها	فشار ورودی	فشار خروجی	بار موتور	جریان موتور	ولتاژ موتور	دبی نفت
نماد	$P_{in}$	$P_{out}$	$MotorLoad$	$I_m$	$V$	$Q_o$
کمترین مقدار	۱۷۲۲	۲۶۴۰	۲۷.۲	۱۰.۰۷	۳۷۳.۰۶	۳۷۲
بیشترین مقدار	۲۲۹۷	۳۱۴۰	۴۴	۱۵.۴۶	۴۱۶.۵۶	۹۳۹
میانگین	۱۸۱۳.۵۶	۲۷۰۴.۶۳	۴۰.۱۹	۱۳.۴۶	۳۹۳	۴۹۱.۶۷
انحراف معیار	۱۰۶.۸۳	۷۱.۶۷	۱.۸۷	۰.۶۵	۵.۷۷	۷۰.۹۹
واریانس	۱۱۳۹۹.۳۸	۵۱۳۶.۴۵	۳.۵۲	۰.۴۲	۳۳.۳۳	۵۰۴۰.۳۵

شکل ۳ رابطه بین ویژگی‌های انتخاب‌شده و تابع هدف و همچنین نمودار توزیع آن‌ها را نشان می‌دهد. در این شکل نمودارهای قطر اصلی توزیع فراوانی هر پارامتر را نشان می‌دهند. در دیگر ستون‌ها نیز پارامترها برحسب یکدیگر رسم شده‌اند. سطر اول مربوط به دبی نفت ( $Q_o$ ) اندازه‌گیری شده است. اولین نمودار در این سطر توزیع فراوانی دبی است همان‌طور که مشاهده می‌شود دبی نفت از توزیع نرمال برخوردار نیست و دبی در ناحیه ۴۰۰ تا ۶۰۰ بشکه در روز بیشترین فراوانی را دارند و حجم داده‌ها در ناحیه بالاتر از این نرخ کم است. همان‌طور که در شکل ۳ مشخص است دبی نفت با فشار تخلیه نیز ارتباط قوی دارد با این حال از آنجاکه بین فشار ورودی پمپ ( $P_{in}$ ) و فشار تخلیه ( $P_d$ ) خطی است بنابراین تنها یکی از این دو مورد به‌عنوان ورودی مدل داده محور انتخاب‌شده است. توزیع فراوانی دیگر ورودی‌ها نیز غیر نرمال هستند.



شکل ۳. رسم نمودارهای دوتایی بین ویژگی‌های انتخاب‌شده و نرخ جریان نفت

## ۵- ارائه نتایج و بحث و بررسی

در این بخش عملکرد هر کدام از روش‌های DT، GBR و k-NN در دو حالت داده‌های با نویز کم و داده‌های با نویز شدید مورد بررسی قرار می‌گیرد و با یکدیگر مقایسه می‌گردد. معیارهای آماری مختلفی برای بررسی عملکرد روش‌های یادگیری ماشین وجود دارد. این معیارها می‌توانند دقت مدل‌های ارائه‌شده را به خوبی نشان دهند. در اینجا از معیارهای جذر میانگین مربعات خطا، میانگین قدر مطلق خطاها، میانگین درصد قدر مطلق خطا و ضریب تعیین استفاده شده و روابط مربوط به آن‌ها، به ترتیب در زیر آمده است.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_{test,i} - y_{pred,i})^2} \quad (5)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_{test,i} - y_{pred,i}| \quad (6)$$

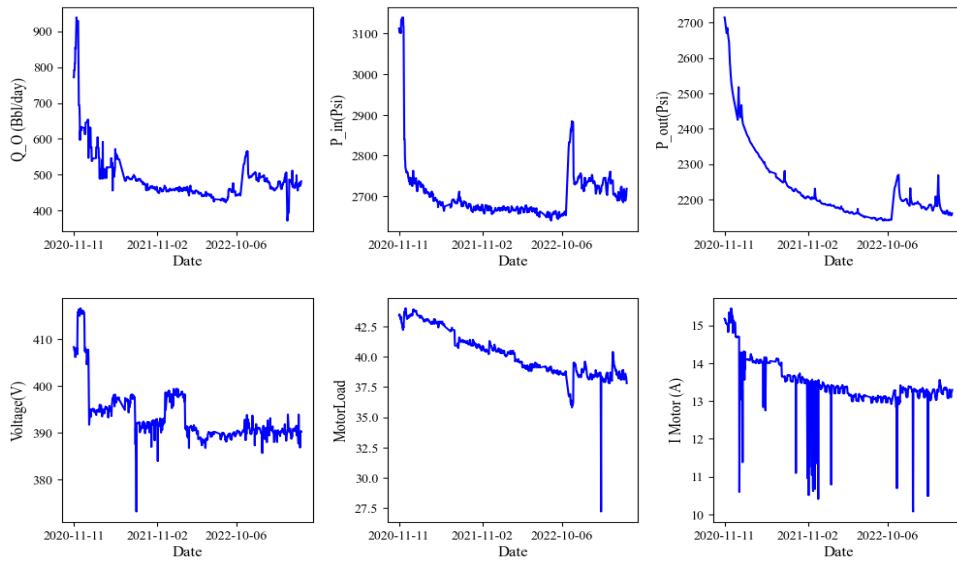
$$MAPE = \frac{100}{n} \sum_{i=1}^n \left| \frac{y_{test,i} - y_{pred,i}}{y_{test,i}} \right| \quad (7)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_{test,i} - y_{pred,i})^2}{\sum_{i=1}^n (y_{test,i} - \bar{y})^2} \quad (8)$$

در این روابط  $y_{test}$  مقادیر داده‌های آزمون،  $y_{pred}$  مقادیر پیش‌بینی‌شده،  $\bar{y}$  مقدار متوسط نمونه و  $n$  تعداد داده‌های آزمون است. در یادگیری ماشین معمولاً مجموعه داده به دو بخش داده‌های آموزش و آزمون تقسیم می‌شوند. مدل‌های ارائه‌شده روابط موجود بین ویژگی‌ها و تابع هدف را با استفاده از داده‌های آموزش یاد گرفته و سپس دقت آن‌ها با استفاده از داده‌های آزمون بررسی می‌گردد. در پژوهش حاضر نیز از ۷۵ درصد داده‌ها برای آموزش و ۲۵ درصد برای آزمون استفاده شده است.

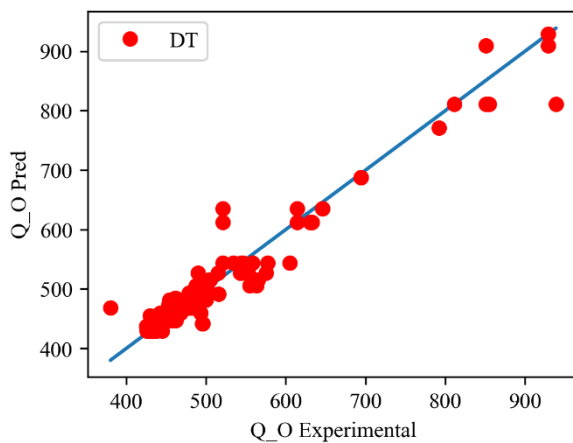
#### ۱-۵ - داده‌های با نویز کم

داده‌های واقعی مورد استفاده در مدل داده محور در شکل ۴ نشان داده شده است. شکل ۴ نشان می‌دهد ویژگی‌های ورودی دارای نویز شدیدی نیستند که این مورد می‌تواند به مدل‌های ارائه‌شده برای پیش‌بینی دقیق‌تر کمک کند. البته برخی از ویژگی‌ها مانند جریان و ولتاژ موتور دارای نویزند.

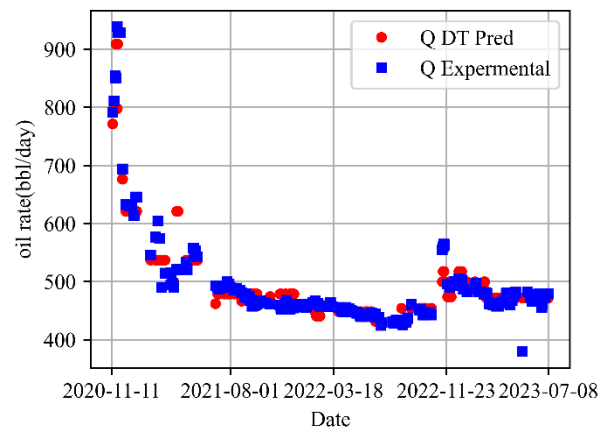


شکل ۴. رسم ورودی‌های مدل داده محور مجموعه داده کم نویز

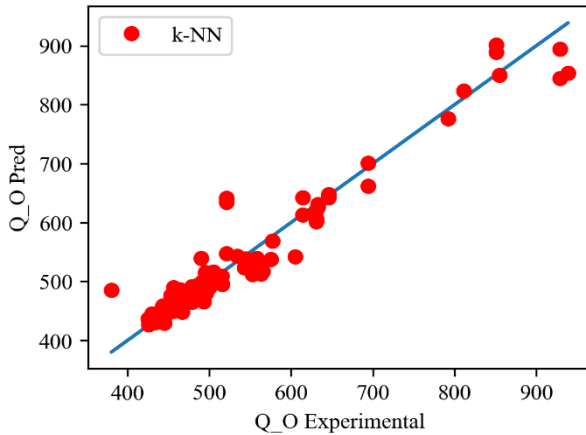
شکل ۵ (الف)، (ب) و (ج) سمت راست دبی تخمینی نفت توسط مدل‌های داده محور و دبی تجربی نفت را برحسب زمان نمایش می‌دهد. همان‌طور که مشاهده می‌شود مدل‌های دبی سنج داده محور در دبی‌های بیشتر از ۵۰۰ بشکه در روز نسبت به بازه کمتر از ۵۰۰ بشکه در روز خطای بالاتری دارند که دلیل آن حجم کمتر داده‌ها در این بازه نسبت به بازه کمتر از ۵۰۰ بشکه در روز است. از مقایسه شکل‌های ۵ (الف)، (ب) و (ج) سمت راست استنباط می‌شود که بین روش‌های ارائه‌شده با قاطعیت نمی‌توان درباره برتری یک روش نسبت به روش دیگر صحبت کرد. شکل‌های ۵ (الف)، (ب) و (ج) سمت چپ نشان می‌دهند که مدل‌های داده محور ارائه‌شده در محدوده کمتر از ۵۰۰ بشکه در روز عملکرد خوبی دارند و نمی‌توان گفت که کدام مدل عملکرد بهتری داشته است. با این حال در محدوده بالاتر از ۵۰۰ بشکه در روز به ترتیب مدل k-NN، GBR و DT بهترین عملکرد را دارند. برای بررسی دقیق‌تر عملکرد هر روش، دقت آن‌ها بر اساس معیارهای ارزیابی در جدول ۳ آمده است.



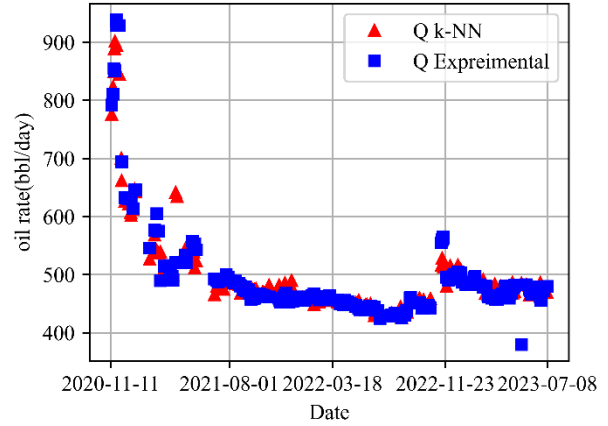
(الف) چپ



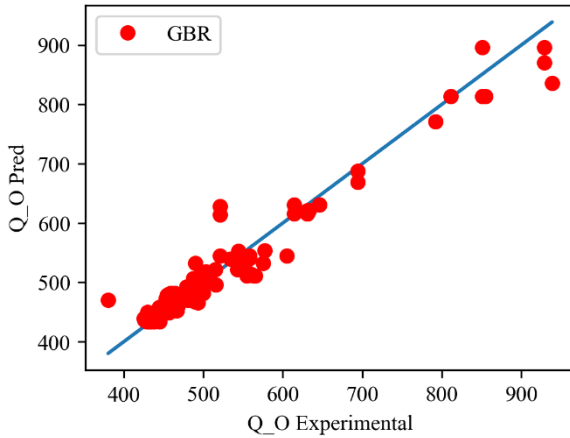
(الف) راست



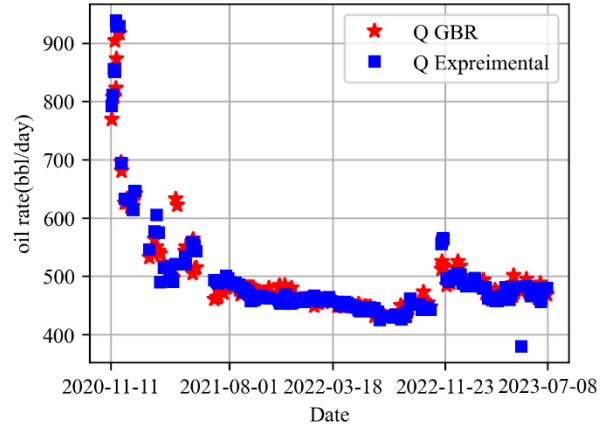
(ب) چپ



(ب) راست



(ج) چپ



(ج) راست

شکل ۵. دبی تخمینی نفت و دبی تجربی نفت نسبت به زمان (راست) و دبی تخمینی نفت در برابر دبی تجربی نفت (چپ) در

حالت ورودی با نویز کم توسط مدل‌های داده محور (الف) DT (ب) k-NN (ج) GBR

جدول ۳ نشان می‌دهد که مدل‌های داده محور ارائه شده عملکرد بسیار نزدیکی دارند با این حال مدل k-NN با اختلاف

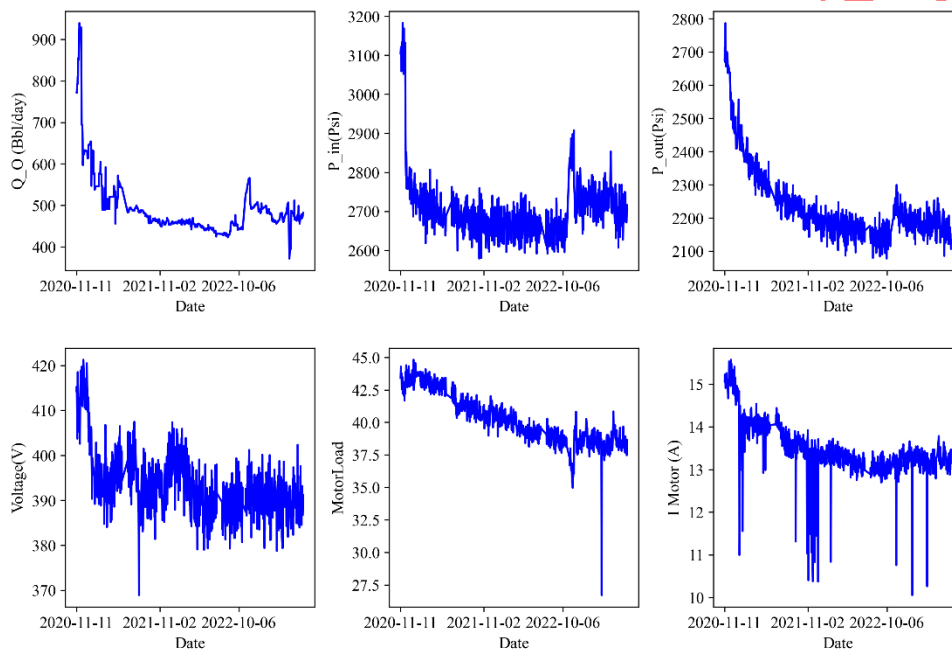
کمی نسبت به دو مدل دیگر عملکرد بهتری دارد. مدل DT نیز دقت پایین تری نسبت به دو مدل دیگر دارد.

جدول ۳. ارزیابی عملکرد مدل‌های ارائه شده در مجموعه داده با نویز کم برای داده‌ها آزمون

روش یادگیری ماشین	MAE	MAPE	RMSE	R <sup>2</sup>
DT	۱۲.۹۵۶۳	۲.۴۷۷۳	۲۱.۵۶۶۹	۰.۹۴۱۳
k-NN	۱۰.۹۵۶۵	۲.۰۵۹۲	۲۰.۳۱۴۰	۰.۹۴۷۹
GBR	۱۲.۱۸۳۴	۲.۳۱۹۱	۲۰.۲۸۱۵	۰.۹۴۷۷

## ۲-۵- داده‌های نويز دار

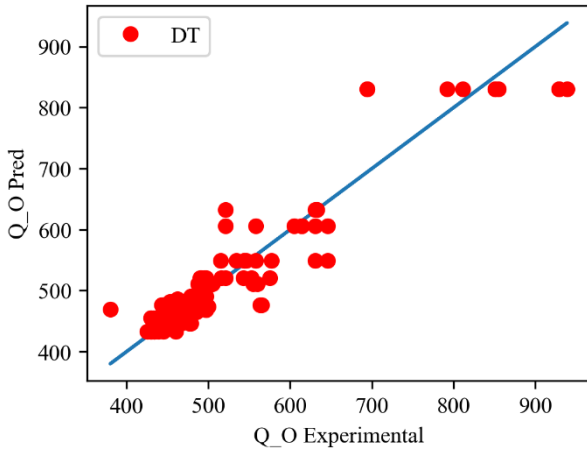
نويز موجود در مجموعه داده می‌تواند تأثیر مخربی بر عملکرد روش‌های یادگیری ماشین داشته باشد و منجر به سوگیری شده و دقت پیش‌بینی‌ها را کاهش دهد. همچنین می‌تواند منجر به برازش بیش‌ازحد شود. برازش بیش‌ازحد زمانی اتفاق می‌افتد که یک مدل از نويز موجود در داده‌ها به‌جای الگوهای زیربنایی یاد می‌گیرد و در نتیجه تعمیم ضعیفی به داده‌های جدید و نامرئی ایجاد می‌کند. این می‌تواند قابلیت اطمینان و اثربخشی مدل‌های یادگیری ماشین را به خطر بیندازد و منجر به نتیجه‌گیری و تصمیم‌گیری اشتباه شود. چندین منبع نويز در داده‌ها وجود دارند که از جمله آن‌ها می‌توان به خطاهای اندازه‌گیری، مقادیر از دست‌رفته، نقاط پرت و ویژگی‌های نامربوط اشاره کرد. این‌ها می‌توانند از عوامل مختلفی مانند نقص حسگرها، خطای انسانی یا تنوع ذاتی در فرآیند جمع‌آوری داده‌ها ناشی شوند. در این قسمت با افزودن مقدار یک درصد انحراف معیار نويز با توزیع گاوسی به پارامترهای اندازه‌گیری شده ورودی، عملکرد مدل‌ها در شرایط نويزی بررسی شده است. شکل ۶ داده‌های ورودی مورد مطالعه در این بخش از مقاله را نشان می‌دهد.



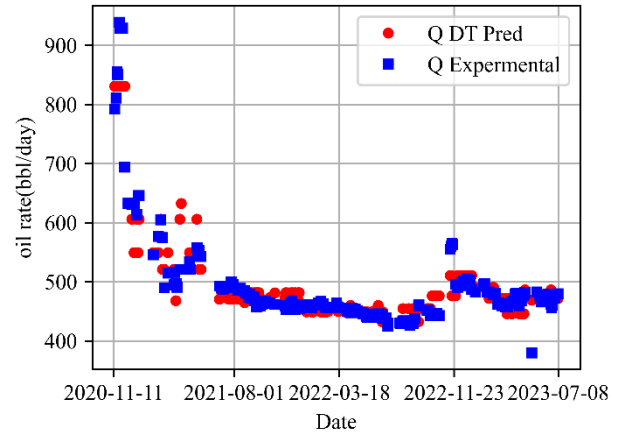
شکل ۶. داده‌های ورودی با یک درصد انحراف معیار نويز

در این حالت نیز مجموعه داده به دو قسمت بخش آموزش و بخش آزمون تقسیم شده است. شکل ۷ (الف)، (ب) و (ج) سمت راست دبی تخمینی نفت توسط مدل‌های داده محور و دبی تجربی نفت را برحسب زمان نمایش می‌دهد. همان‌طور که مشاهده می‌شود نويز موجود در داده‌ها تأثیر منفی بر عملکرد دبی سنج‌های مجازی داده محور داشته است. تأثیر نويزها در مناطقی که حجم داده کمتری دارند بیشتر است. شکل‌های ۷ (الف)، (ب) و (ج) سمت چپ نشان می‌دهند که مدل داده محور k-NN بهترین عملکرد را در برابر ورودی‌های دارای نويز دارد. دلیل این امر این است که این مدل از همسایه‌ها برای تخمین مقدار استفاده می‌کند. مدل DT نیز بیشتر از مدل‌های k-NN و GBR تحت تأثیر نويز موجود در ورودی‌ها بوده است. مدل GBR توانسته با ترکیب چندین مدل DT عملکرد این مدل را بهبود دهد.

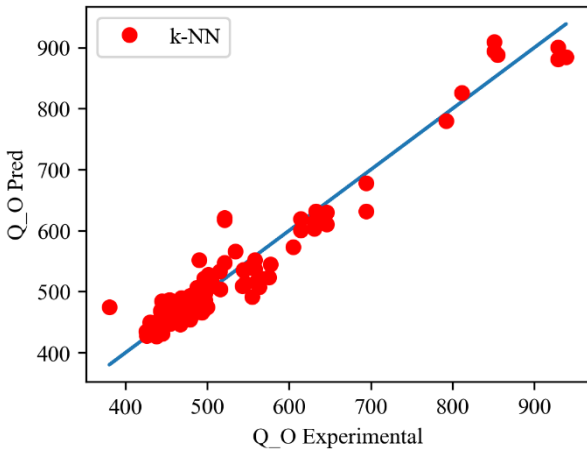




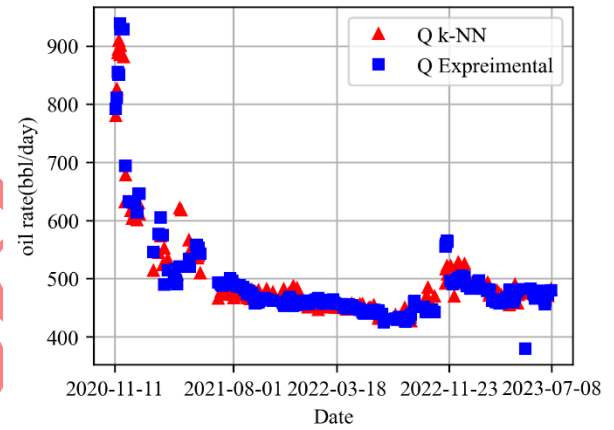
چپ (الف)



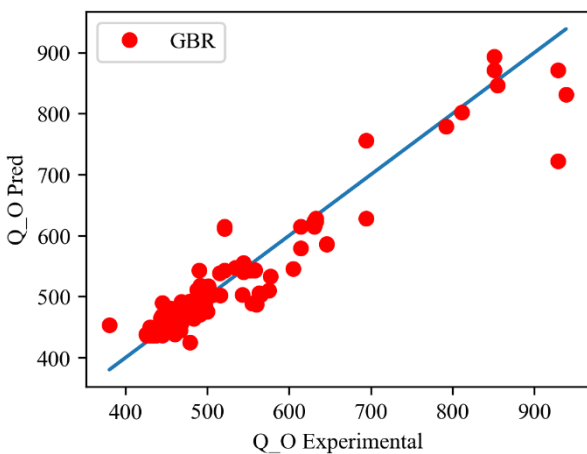
راست (الف)



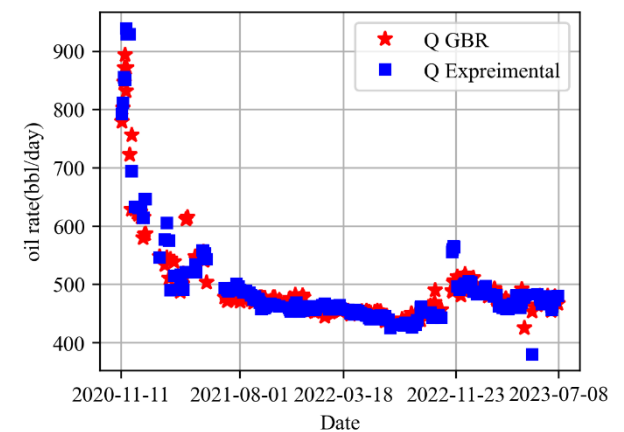
چپ (ب)



راست (ب)



چپ (ج)



راست (ج)

شکل ۷. دبی تخمینی نفت و دبی تجربی نفت نسبت به زمان (راست) و دبی تخمینی نفت در برابر دبی تجربی نفت (چپ) در

حالت ورودی با نویز کم توسط مدل‌های داده محور (الف) DT (ب) k-NN (ج) GBR

جدول ۴ عملکرد مدل‌های مختلف را در حالت با نویز نشان می‌دهد. جدول ۴ میزان دقت هر کدام از مدل‌ها را در برابر داده‌های نویز دار نشان می‌دهد. با توجه به نتایج ارائه‌شده در این جدول مشاهده می‌شود که مدل k-NN بهترین عملکرد را در برابر داده‌های نویز دار داشته است و نویز تأثیر کمی بر عملکرد این روش داشته است. مدل DT نیز بیشترین تأثیر را از نویز گرفته است.

جدول ۴. عملکرد روش‌های ارائه‌شده در مجموعه داده با یک درصد انحراف معیار نویز برای داده‌های آزمون

روش مورد بررسی	MAE	MAPE	RMSE	R <sup>2</sup>
DT	۱۶.۹۲۶۸	۳.۰۹۶۲	۳۲.۵۳۲۴	۰.۸۶۶۵
k-NN	۱۴.۹۱۹۷	۲.۸۳۰۵	۲۴.۲۵۷۹	۰.۹۲۵۷
GBR	۱۴.۶۶۱۱	۳.۴۲۵۳	۲۶.۸۰۸۸	۰.۸۹۶۵

بررسی نتایج به دست آمده برای ورودی‌های نویز دار و بدون نویز نشان داد که مدل DT یک مدل ساده است که در یک مجموعه داده با نویز کم، تنوع ناکافی و حجم محدود می‌تواند تخمین قابل قبولی داشته باشد با این حال این مدل در برابر نویز بسیار حساس است. الگوریتم GBR یک مدل پیشرفته‌تر از مدل DT است. این مدل با ایجاد چند درخت تصمیم مختلف سعی می‌کند که عملکرد مدل DT را بهبود دهد. نتایج به دست آمده در این مطالعه به خوبی نشان می‌دهد که مدل GBR در شرایط مجموعه داده با حجم محدود، تنوع ناکافی و کم نویز عملکرد نزدیک به عملکرد DT دارد. با این حال در مجموعه داده نویز دار عملکرد مدل GBR بهتر از مدل DT است. الگوریتم k-NN عملکرد بسیار خوبی در هر دو مجموعه داده کم نویز و نویز دار داشته است. این الگوریتم می‌تواند در شرایط مجموعه داده با حجم کم و تنوع ناکافی عملکرد خوبی داشته است. در این مطالعه سعی شده است که قابلیت‌های الگوریتم‌های DT، k-NN و GBR نشان داده شود. با وجود مزایایی که برای مدل دبی سنج مجازی داده محور ارائه شده یاد شده است، این مدل‌ها دارای محدودیت‌هایی نیز هستند. در این مدل تنها از اطلاعات پمپ استفاده شده است و تأثیرات دیگر پارامترهای درون چاهی در نظر گرفته نشده است. به همین دلیل نیاز است که دقت مدل در طول زمان به صورت مکرر امتحان شود. در این پژوهش دبی سنج مجازی داده محور برای تخمین دبی نیازمند اطلاع از شرایط درون چاه و مخزن نیست. با این حال در صورتی که شرایط چاه و مخزن تغییر کنند مدل داده محور نیازمند آموزش مجدد است.

## ۴- جمع بندی

در پژوهش حاضر از روش‌های یادگیری ماشین شامل روش تقویت گرادیان (GBR)، درخت تصمیم (DT) و k همسایه نزدیک‌تر (k-NN) برای تخمین نرخ دبی نفت خروجی از یک چاه نفتی واقع در جنوب ایران استفاده شد. انتخاب ویژگی‌های مناسب برای آموزش مدل‌ها با استفاده از معیارهای آماری اسپیرمن و پیرسون انجام شد. در این مطالعه از داده‌های ورودی و خروجی پمپ‌های شناور الکتریکی (ESP) شامل فشار ورودی ( $P_{in}$ )، فشار خروجی ( $P_{out}$ )، بار موتور (MotorLoad)، جریان موتور ( $I_m$ ) و ولتاژ موتور ( $V$ ) استفاده شد و از نتایج به دست آمده می‌توان به این نتیجه رسید که اطلاعات ورودی و خروجی پمپ‌ها می‌توانند به عنوان ورودی‌های دبی سنج‌های داده محور استفاده شوند. تعداد داده‌های موجود در مجموعه

داده ۹۸۰ عدد بود که پس از حذف داده‌های پرت و مقادیر گمشده تعداد داده‌ها به ۸۵۰ عدد کاهش یافت. توزیع فراوانی داده‌ها نشان می‌دهد که تعداد داده‌هایی که دبی نفت آن در محدوده بیش‌تر از ۵۰۰ بشکه در روز بوده کم است. همین امر چالشی در برابر استفاده از دبی سنج‌های مجازی داده محور ایجاد کرده است، با این حال نتایج نشان می‌دهد که مدل‌های ارائه‌شده عملکرد قابل قبولی داشتند. مدل k-NN بهترین عملکرد را نسبت به دو مدل دیگر نشان داد. در ادامه عملکرد مدل‌های ارائه‌شده در برابر داده‌های نویزی با اضافه کردن یک درصد انحراف معیار نویز با استفاده از تابع توزیع گاوسی به داده‌های ورودی مورد ارزیابی قرار گرفت. مشاهده شد که روش‌های یادگیری ماشین استفاده‌شده در اثر وجود نویز در اطلاعات ورودی، خصوصاً در محدوده‌های با داده ناکافی دچار خطا می‌شوند. روش k-NN بهترین عملکرد را در برابر داده‌های نویزی در مقایسه با دو روش دیگر به نمایش گذاشت که نشان از مدیریت مناسب این مدل در حضور داده‌های دارای نویز است. روش DT نیز در مقایسه با دو روش دیگر، بیشترین خطا را در حضور داده‌های دارای نویز نشان داد. نتایج این مطالعه نشان داد که انتخاب داده ورودی و مدل یادگیری ماشین مناسب می‌تواند عملکرد بسیار مناسبی را برای دبی سنج مجازی به ارمغان بیاورد. در ادامه پیشنهاد می‌شود که از داده‌های ته چاهی و سر چاهی در کنار داده‌های پمپ به‌عنوان دبی سنج مجازی داده محور استفاده شود.

### تقدیر و تشکر و پیوست‌ها

داده‌های این مطالعه توسط شرکت توسعه نفت و گاز پرشیا جمع‌آوری و ارائه‌شده است. لذا نویسندگان از این شرکت بابت همکاری ایشان، صمیمانه تشکر می‌نمایند.

### References

- [1] Dayev, Z. A. (2020). Application of artificial neural networks instead of the orifice plate discharge coefficient. *Flow Measurement and Instrumentation*, 71, 101674. .
- [2] Bismukhametov, T. & Jäschke, J. (2020). First principles and machine learning virtual flow metering: a literature review. *Journal of Petroleum Science and Engineering*, 184, 106487.
- [3] Mercante, R., & Netto, T. A. (2022). Virtual flow predictor using deep neural networks. *Journal of Petroleum Science and Engineering*, 213, 110338..
- [4] T. A. AL-Qutami, R. Ibrahim, I. Ismail, and M. A. Ishak, "Virtual multiphase flow metering using diverse neural network ensemble and adaptive simulated annealing," *Expert Syst. Appl.*, vol. 93, pp. 72–85, Mar. 2018, doi: 10.1016/j.eswa.2017.10.014..
- [5] Al-Qutami, T. A., Ibrahim, R., & Ismail, I. (2017, September). Hybrid neural network and regression tree ensemble pruned by simulated annealing for virtual flow metering application. In *2017 IEEE International Conference on Signal and Image Processing Applications (ICSIPA)* (pp. 304-309). IEEE.
- [6] Ahmadi, M. A., Ebadi, M., Shokrollahi, A., & Majidi, S. M. J. (2013). Evolving artificial neural network and imperialist competitive algorithm for prediction oil flow rate of the reservoir. *Applied Soft Computing*, 13(2), 1085-1098.
- [7] Bismukhametov, T., & Jäschke, J. (2019). Oil production monitoring using gradient boosting machine learning algorithm. *Ifac-Papersonline*, 52(1), 514-519.
- [8] Góes, M. R. R., Guedes, T. A., d'Avila, T. C., Vieira, B. F., Ribeiro, L. D., de Campos, M. C., & Secchi, A. R. (2021). Virtual flow metering of oil wells for a pre-salt field. *Journal of Petroleum Science and Engineering*, 203, 108586.

- [9] Hotvedt, M., Grimstad, B., Ljungquist, D., & Imsland, L. (2022). On gray-box modeling for virtual flow metering. *Control Engineering Practice*, 118, 104974.
- [10] AlAjmi, M. D., Alarifi, S. A., & Mahsoon, A. H. (2015, March). Improving multiphase choke performance prediction and well production test validation using artificial intelligence: a new milestone. In *SPE digital energy conference and exhibition* (p. D031S022R003). SPE.
- [11] A. T. Sandnes, B. Grimstad, and O. Kolbjørnsen, "Multi-task learning for virtual flow metering," *Knowl.-Based Syst.*, vol. 232, p. 107458, Nov. 2021, doi: 10.1016/j.knosys.2021.107458.
- [12] Al-Jasmi, A., Goel, H. K., Nasr, H., Querales, M., Rebeschini, J., Villamizar, M. A., ... & Saputelli, L. (2013, June). Short-term production prediction in real time using intelligent techniques. In *SPE Europec featured at EAGE Conference and Exhibition?* (pp. SPE-164813). SPE.
- [13] Denney, T., Wolfe, B., & Zhu, D. (2013, March). Benefit evaluation of keeping an integrated model during real-time ESP operations. In *SPE Digital Energy Conference and Exhibition* (pp. SPE-163704). SPE..
- [14] L. A. Camilleri, T. Banciu, G. Ditoiu, and S. A. Petrom, "First Installation of 5 ESPs Offshore Romania - A Case Study and Lessons Learned," presented at the *SPE Intelligent Energy Conference and Exhibition*, Mar. 2010, p. SPE-127593-MS. doi: 10.2118/127593-MS."
- [15] Rao, H., Shi, X., Rodrigue, A. K., Feng, J., Xia, Y., Elhoseny, M., ... & Gu, L. (2019). Feature selection based on artificial bee colony and gradient boosting decision tree. *Applied Soft Computing*, 74, 634-642..
- [16] Dudani, S. A. (1976). The distance-weighted k-nearest-neighbor rule. *IEEE Transactions on Systems, Man, and Cybernetics*, (4), 325-327.
- [17] Józwiak, A. (1983). A learning scheme for a fuzzy k-NN rule. *Pattern Recognition Letters*, 1(5-6), 287-289.
- [18] Song, Y., Liang, J., Lu, J., & Zhao, X. (2017). An efficient instance selection algorithm for k nearest neighbor regression. *Neurocomputing*, 251, 26-34. .
- [19] Devroye, L., Györfi, L., Krzyżak, A., & Lugosi, G. (1994). On the strong universal consistency of nearest neighbor regression function estimates. *The annals of Statistics*, 22(3), 1371-1385.
- [20] Xu, M., Watanachaturaporn, P., Varshney, P. K., & Arora, M. K. (2005). Decision tree regression for soft classification of remote sensing data. *Remote Sensing of Environment*, 97(3), 322-336.
- [21] L. Breiman, *Classification and regression trees*. Routledge, 2017.