

شناسایی خطاهای داده‌های خام بویه‌های موج‌نگار با استفاده از روش ضریب داده پرت محلی

کیومرث محمودی^۱، محمد جواد کتابداری^۲

ketabdar@aut.ac.ir

۱- کارشناس ارشد، دانشکده مهندسی دریا، دانشگاه صنعتی امیرکبیر

۲- دانشیار، دانشکده مهندسی دریا، دانشگاه صنعتی امیرکبیر

چکیده

استخراج مشخصات دریا معمولاً از طریق بویه‌های موج‌نگار انجام می‌شود. اما ثبت داده توسط موج‌نگارها معمولاً با خطاهایی همراه است. لذا قبل از استخراج هرگونه اطلاعاتی لازم است این خطاها را شناخت و آنها را حذف و یا تصحیح کرد. هدف از این تحقیق، شناسایی خطاهای موجود در برداشت داده‌های خام از بویه‌های موج‌نگار، با استفاده از روش ضریب داده پرت محلی (LOF) است. روشی قدرتمند جهت شناسایی ناهنجاری داده‌ها در یادگیری ماشین است، که در بسیاری از کاربردهای عملی از آن استفاده می‌شود. در این مقاله داده‌های روزانه بویه موج‌نگار سازمان بنادر و دریانوردی که در آب‌های خلیج فارس در منطقه عسلویه استان بوشهر به آب انداخته شده است در بازه ۹۲/۲/۱۷ تا ۹۲/۵/۱ مورد تحلیل خطا به این روش قرار گرفت. نتایج نشان می‌دهد LOF روشی کارآمد برای شناسایی خطاهای موجود در داده‌های موج‌نگار است.

واژگان کلیدی: موج‌نگار، داده موج، خلیج فارس، شناسایی خطا، روش ضریب داده پرت محلی (LOF).

تاریخ دریافت مقاله : ۹۳/۱۲/۱۸

تاریخ پذیرش مقاله : ۹۵/۰۱/۱۶

۱ - مقدمه

داده‌های حاصل از یک مشاهده وقتی مفید و کارآمد است که نماینده واقعی موضوع مورد مطالعه باشد. برای بدست آوردن نتایج مناسب و مطلوب از داده‌های جمع آوری شده، داده‌ها باید نماینده واقعی مشاهدات بوده، و به نحو مفیدی تهیه و تنظیم شده باشند تا قابل تجزیه و تحلیل بوده و بتوان به کمک آنها به نتیجه‌گیری صحیح رسید. اهمیت برداشت داده‌های مناسب و با کمترین خطا در مطالعات آزمایشگاهی از آنجا دارای اهمیت است که معمولاً از این داده‌ها، در صحت سنجی و کالیبراسیون مدل‌های عددی و ریاضی استفاده می‌شود [۱].

داده‌های خطا نتایج تجزیه و تحلیل داده‌ها را تحت تاثیر قرار داده و منجر به اختلال در نتیجه‌گیری از اطلاعات خواهند شد. داده‌های خطا می‌توانند بر اثر عوامل مختلفی به وجود آیند. دستگاه‌های فیزیکی برای تخمین اندازه‌گیری‌ها ممکن است با یک سری عوامل گذرا که موجب بد عمل کردن آنها می‌شود مواجه باشند. در نتیجه ممکن است خطایی در انتقال اطلاعات یا استنتاج نهایی بوجود آید. همچنین این داده‌ها ممکن است بر اثر تغییر در رفتار سیستم، خطاهای انسانی، خطاهای دستگاهی یا به طور ساده به خاطر انحرافهای طبیعی در جمعیت نمونه‌ها به وجود آیند. تعریف ریاضیاتی دقیقی برای تشخیص منبع داده خطا وجود ندارد. تشخیص اینکه یک مشاهده خطا است یا نه، صرفاً یک فرایند کیفی و ابتکاری است که به تجربه شخصی فرد و کاربردی که از داده‌ها داریم وابسته است.

برای اندازه‌گیری مشخصه‌های دریایی از دستگاه‌های متنوعی استفاده می‌شود. این دستگاه‌ها قادر به جمع‌آوری اطلاعات هواشناسی و اقیانوس شناسی مانند سرعت، جهت و مدت وزش باد، دما، فشار، رطوبت هوا، غلظت و شوری آب دریا و تاریخچه زمانی تغییرات تراز آب هستند [۲]. بویه‌های موج نگار و پلات فرم‌ها از وسایل متداول این اندازه‌گیری‌ها هستند. هر یک از این دستگاه‌ها می‌توانند اطلاعات مهم موج از جمله ارتفاع و پریود آن را از داده‌های خام استخراج کنند. ثبت داده توسط موج نگارها معمولاً با خطاهای مختلفی همراه است که قبل از استخراج هرگونه اطلاعاتی

لازم است این خطاها را شناخت و آنها را حذف و یا تصحیح نمود. هدف از این تحقیق، شناسایی داده‌های مشکوک به خطا موجود در داده‌های برداشت شده از بویه‌ها با استفاده از روش ضریب داده پرت محلی^۱ (LOF) است. در این مطالعه، داده‌ای به عنوان کاندیدای مشکوک به خطا در نظر گرفته شده است که به طور قابل توجهی متناقض با سایر اعضای نمونه ظاهر شده، و احتمالاً با یک فرآیند متفاوت تولید شده است.

۲- خطاهای احتمالی در رکورد موج

در این بخش، برخی از خطاهای احتمالی موجود در برداشت داده‌های خام از بویه‌های موج‌نگار معرفی شده است:

- وجود Trend در موج: گاهی تراز متوسط داده‌های زمانی ثبت شده توسط بویه‌ها شیبی به سمت بالا یا پایین دارد. عموماً این مسئله بواسطه اختلال در دستگاه‌ها و تجهیزات است. گاهی ممکن است ترند به دلیل تأثیر تغییرات جزر و مد بر روی سنسور فشار بویه بروز نماید.
- وجود Noise در موج: نویز به وجود پارازیت و ارتعاشات با فرکانس بالا در تاریخچه زمانی موج می‌گویند. نویز ممکن است به دلیل حرکت پروانه شناورهای عبوری، یا تشدیدهای موردی در سیستم موج نگار و مهار آن ایجاد شود.
- وجود MWL در موج: معمولاً در اثر عبور موج به طور طبیعی تراز سطح آب بالا^۲ یا پایین^۳ می‌رود. گاهی نیز این تغییرات تراز سطح آب در صورتی که شدید باشد به علت وجود اختلال Drift در دستگاه است [۲].

۳- مواد و روش‌ها

داده‌های مورد استفاده در این تحقیق، داده‌های گردآوری شده از یکی از بویه‌های موج‌نگار موجود در آب‌های خلیج فارس است که در عسلویه استان بوشهر به آب انداخته شده است. این بویه در طول جغرافیایی ۲۵,۵ و عرض جغرافیایی ۵۲,۵ واقع شده است. بویه مورد استفاده از نوع Wavescan است (شکل (۱)). این بویه یک ابزار چند کاره و یکی از پرکاربردترین بویه‌ها است در سرتاسر دنیا است.

³ Set down

¹ Local Outlier Factor (LOF)

² Set up

اندازه‌گیری جهت موج: این بویه مجهز به یک حسگر برای اندازه‌گیری‌های جهت موج مبتنی بر اصل اندازه‌گیری شیب است.

ارائه داده‌ها:

وضعیت پارامترهای سنسور انتخابی با استفاده از نرم‌افزار مخصوص به طور مداوم پایش شده بنابراین در صورت بروز هر گونه حادثه پیش‌بینی نشده‌ای یک پیام به آدرس از پیش تعریف شده ارسال می‌شود. گزارش‌های بویه به راحتی تولید و پیکربندی شده و برای خروجی به فرمت‌های مختلف آماده می‌شوند. این ویژگی امکان ارائه آئی داده‌ها روی اینترنت را فراهم می‌کند. همچنین نرم‌افزار مورد استفاده دارای روش‌هایی برای بررسی کیفیت داده‌ها است.

در این تحقیق داده‌های بویه از سایت سازمان بنادر و دریانوردی استخراج شده است [۳]. در جدول (۱)، داده‌های مورد بررسی به همراه فهرست مشخصه‌های اندازه‌گیری شده، آمده است. در این جدول، مقادیر مشخصه‌های ارتفاع موج شاخص، میزان رطوبت، فشار هوا، دمای آب و سرعت باد قرار دارد. هدف، شناسایی داده‌های مشکوک به خطا موجود در این اندازه‌گیری‌ها است.

۳-۱- روش ضریب داده پرت محلی

شناسایی داده‌های پرت (داده‌های مشکوک به خطا)، شاخه‌ای کاربردی و بسیار مهم از داده کاوی است که به شناسایی الگوهای متناقض با جمعیت نرمال نمونه، در یک جامعه آماری می‌پردازد. به علت اهمیت شناسایی داده‌های پرت در علوم مختلف، این موضوع به یکی از مهمترین موضوعات تحقیقاتی در آمار، یادگیری ماشین و داده‌کاوی تبدیل شده است. برخی از کاربردهای الگوریتم‌های شناسایی داده‌های پرت به صورت نمونه در زیر معرفی شده است:

- شناسایی تقلب: به عنوان مثال شناسایی تقلب در کارت‌های اعتباری و یا گوشی‌های تلفن همراه [۴].
- فرایندهای وام و وام دهی: شناسایی کلاهبرداری‌ها در امور وام و همچنین شناسایی مشتری‌هایی که پتانسیل ایجاد مشکل را دارند [۴].
- شناسایی نفوذ: به عنوان مثال شناسایی دسترسی‌های غیرمجاز به شبکه‌های کامپیوتری [۵].



شکل (۱) بویه موج نگار Wavescan عسلویه.

. از این دستگاه برای اندازه‌گیری و پایش داده‌های مختلف اقیانوس‌شناسی، هواشناسی و کیفیت آب استفاده می‌شود. در ادامه این بخش مشخصات این بویه به صورت مختصر معرفی شده است.

ویژگی‌ها:

- طراحی منحصر به فرد برای بهینه‌سازی اندازه‌گیری‌های جهت موج.
 - پردازش داده‌های اندازه‌گیری شده.
 - ارتباط دو طرفه برای انتقال و کنترل داده‌ها.
 - انتقال و ارائه آئی داده‌های اندازه‌گیری شده.
 - پیکربندی انعطاف‌پذیر سنسورها و جمع‌آوری داده‌ها.
 - طراحی شده برای حمل و نقل و استقرار امن و آسان.
 - حمل و نقل آسان و مونتاژ محلی.
 - ارائه لینک‌های صوتی آبی به رصدخانه اعماق اقیانوس.
 - طراحی ویژه به منظور به حداقل رساندن اثرات پهلوگیری و اغتشاش شناورها روی حرکت بویه.
- کاربردها:
- طراحی و عملیات دریایی.
 - مطالعات اقلیمی و هواشناسی.
 - مناسب برای عملیات در آبهای عمیق.
 - نظارت بر بندرگاه.
 - مهندسی سواحل.
 - مطالعات علمی.
 - مطالعات انرژی موج.
 - کنترل ترافیک دریایی.
 - مطالعات نظارت بر کیفیت آب.

جدول (۱) داده‌های گردآوری شده از بویه موج‌نگار عسلویه.

شماره	تاریخ برداشت	ارتفاع موج شاخص (m)	رطوبت (درصد)	فشار هوا (اتمسفر)	دمای آب (°C)	سرعت باد (m/sec)
۱	۱۳۹۲/۲/۱۷	۰,۳۸	۶۲,۲۴	۱۰۰۴,۲۳	۲۸,۱۳	۱,۶۳
۲	۱۳۹۲/۲/۱۸	۰,۲۳	۵۲,۹۳	۱۰۰۵,۴۷	۲۸,۶۳	۱,۸۹
۳	۱۳۹۲/۲/۱۹	۰,۲۲	۶۹,۵۹	۱۰۰۵,۴۹	۲۵,۷۰	۰,۸۴
۴	۱۳۹۲/۲/۲۰	۰,۲۷	۷۲,۰۹	۱۰۰۴,۹۱	۲۵,۸۸	۱,۰۹
۵	۱۳۹۲/۲/۲۱	۰,۲۲	۷۴,۰۲	۱۰۰۳,۱۴	۲۶,۱۶	۱,۰۵
۶	۱۳۹۲/۲/۲۲	۰,۱۹	۷۱,۶۷	۱۰۰۲,۷۳	۲۶,۹۲	۱,۰۶
۷	۱۳۹۲/۲/۲۳	۰,۱۶	۶۸,۶۸	۱۰۰۲,۷۰	۲۷,۶۵	۰,۶۵
۸	۱۳۹۲/۲/۲۴	۰,۱۷	۶۸,۰۵	۱۰۰۲,۸۰	۲۸,۴۲	۱,۰۴
۹	۱۳۹۲/۲/۳۱	۰,۱۸	۷۷,۴۲	۱۰۰۳,۷۸	۲۷,۷۴	۱,۱۵
۱۰	۱۳۹۲/۳/۱	۰,۵۵	۵۳,۴۴	۱۰۰۴,۰۵	۲۸,۵۹	۱,۰۵
۱۱	۱۳۹۲/۳/۲	۰,۳۷	۶۰,۹۰	۱۰۰۳,۱۹	۲۸,۱۷	۱,۲۲
۱۲	۱۳۹۲/۳/۳	۰,۳۷	۶۳,۶۰	۹۹۸,۸۲	۲۷,۱۹	۱,۱۵
۱۳	۱۳۹۲/۳/۴	۰,۵۶	۷۶,۵۹	۹۹۹,۹۷	۲۷,۷۵	۲,۰۲
۱۴	۱۳۹۲/۳/۵	۰,۷۸	۸۴,۴۳	۱۰۰۱,۲۹	۲۷,۸۶	۲,۵۱
۱۵	۱۳۹۲/۳/۶	۰,۶۵	۷۹,۷۰	۱۰۰۱,۵۳	۲۸,۶۸	۱,۸۱
۱۶	۱۳۹۲/۳/۷	۰,۵۸	۸۱,۷۷	۱۰۰۱,۷۹	۲۹,۰۷	۱,۶۶
۱۷	۱۳۹۲/۳/۸	۱,۳۱	۷۴,۹۴	۱۰۰۱,۰۸	۲۹,۴۱	۲,۸۱
۱۸	۱۳۹۲/۳/۹	۲,۱۳	۵۶,۱۸	۱۰۰۰,۰۰	۲۸,۸۵	۴,۸۸
۱۹	۱۳۹۲/۳/۱۰	۲,۱۵	۶۶,۹۲	۹۹۸,۷۳	۲۸,۱۴	۳,۱۷
۲۰	۱۳۹۲/۳/۱۱	۱,۱۱	۷۶,۹۲	۹۹۸,۰۸	۲۸,۵۱	۲,۶۷
۲۱	۱۳۹۲/۳/۱۲	۱,۱۵	۷۹,۱۶	۹۹۶,۶۹	۲۹,۵۲	۱,۹۷
۲۲	۱۳۹۲/۳/۱۳	۱,۶۴	۷۹,۰۵	۹۹۷,۰۴	۲۹,۰۳	۳,۷۲
۲۳	۱۳۹۲/۳/۱۴	۰,۶۲	۷۶,۴۰	۱۰۰۰,۱۴	۲۹,۷۲	۲,۲۵
۲۴	۱۳۹۲/۳/۱۵	۰,۶۷	۷۵,۵۳	۱۰۰۰,۸۴	۳۰,۰۸	۱,۷۴
۲۵	۱۳۹۲/۳/۱۶	۱,۰۷	۷۶,۴۶	۱۰۰۰,۳۷	۳۰,۶۴	۲,۴۹
۲۶	۱۳۹۲/۳/۱۷	۱,۶۶	۵۷,۳۵	۹۹۶,۰۸	۳۰,۳۴	۵,۲۸
۲۷	۱۳۹۲/۳/۱۸	۲,۰۴	۶۲,۲۸	۹۹۴,۱۸	۲۹,۲۶	۶,۴۷
۲۸	۱۳۹۲/۳/۱۹	۲,۰۰	۵۵,۹۷	۹۹۴,۳۵	۲۸,۵۵	۴,۹۱
۲۹	۱۳۹۲/۳/۲۰	۱,۵۰	۷۶,۹۹	۹۹۴,۲۷	۲۸,۷۸	۱,۴۹
۳۰	۱۳۹۲/۳/۲۱	۱,۵۹	۷۸,۴۹	۹۹۴,۱۴	۲۸,۸۱	۱,۸۴
۳۱	۱۳۹۲/۳/۲۲	۲,۰۱	۶۸,۹۳	۹۹۲,۷۳	۲۸,۸۴	۳,۵۷
۳۲	۱۳۹۲/۳/۲۳	۱,۶۵	۷۹,۶۵	۹۹۳,۷۱	۲۸,۷۹	۳,۱۴
۳۳	۱۳۹۲/۳/۲۴	۰,۹۲	۸۰,۷۱	۹۹۴,۳۱	۲۹,۴۶	۲,۴۱
۳۴	۱۳۹۲/۳/۲۵	۰,۶۹	۷۷,۵۶	۹۹۳,۸۸	۳۰,۱۰	۱,۸۴
۳۵	۱۳۹۲/۳/۲۶	۱,۷۳	۸۲,۰۵	۹۹۲,۳۳	۳۰,۳۹	۲,۹۲
۳۶	۱۳۹۲/۳/۲۷	۲,۰۷	۷۵,۴۰	۹۹۱,۴۲	۳۰,۶۳	۳,۲۳
۳۷	۱۳۹۲/۳/۲۸	۱,۸۹	۷۱,۹۹	۹۹۱,۳۷	۳۱,۰۰	۳,۲۸
۳۸	۱۳۹۲/۳/۲۹	۲,۱۵	۵۲,۹۶	۹۹۳,۳۹	۳۰,۸۳	۶,۲۰
۳۹	۱۳۹۲/۳/۳۰	۲,۱۷	۶۹,۱۲	۹۹۶,۷۷	۳۰,۴۱	۳,۳۶
۴۰	۱۳۹۲/۳/۳۱	۱,۸۹	۷۰,۶۹	۹۹۶,۷۰	۳۰,۶۵	۳,۱۴
۴۱	۱۳۹۲/۴/۱	۱,۵۳	۷۶,۴۳	۹۹۶,۲۶	۳۰,۸۱	۳,۶۷

۲,۱۰	۳۱,۳۴	۹۹۵,۹۸	۷۸,۶۲	۰,۷۵	۱۳۹۲/۴/۲	۴۲
۱,۸۲	۳۱,۸۵	۹۹۴,۴۵	۷۹,۶۰	۰,۷۹	۱۳۹۲/۴/۳	۴۳
۲,۵۷	۳۱,۷۷	۹۹۳,۳۹	۷۸,۷۶	۰,۹۵	۱۳۹۲/۴/۴	۴۴
۱,۷۹	۳۲,۲۰	۹۹۵,۱۰	۸۰,۶۵	۰,۸۷	۱۳۹۲/۴/۵	۴۵
۲,۱۷	۳۲,۲۶	۹۹۴,۷۷	۷۷,۲۷	۱,۱۱	۱۳۹۲/۴/۶	۴۶
۲,۸۸	۳۲,۲۴	۹۹۳,۳۹	۷۴,۹۴	۱,۶۰	۱۳۹۲/۴/۷	۴۷
۲,۹۲	۳۲,۳۴	۹۹۳,۸۷	۷۸,۴۹	۱,۲۰	۱۳۹۲/۴/۸	۴۸
۲,۱۳	۳۲,۶۴	۹۹۴,۷۱	۶۹,۴۷	۰,۹۰	۱۳۹۲/۴/۹	۴۹
۱,۶۱	۳۲,۶۳	۹۹۵,۰۸	۵۰,۰۸	۰,۴۴	۱۳۹۲/۴/۱۰	۵۰
۲,۶۲	۳۴,۷۸	۹۹۵,۰۸	۷۲,۶۰	۰,۷۴	۱۳۹۲/۴/۲۴	۵۱
۳,۰۴	۳۴,۷۱	۹۹۴,۲۱	۷۲,۵۳	۱,۰۷	۱۳۹۲/۴/۲۵	۵۲
۲,۳۰	۳۴,۵۶	۹۹۳,۸۸	۷۲,۴۸	۰,۸۵	۱۳۹۲/۴/۲۶	۵۳
۱,۹۳	۳۴,۷۱	۹۹۳,۳۸	۷۲,۸۰	۰,۸۳	۱۳۹۲/۴/۲۷	۵۴
۲,۳۹	۳۴,۹۲	۹۹۱,۵۴	۷۱,۲۹	۰,۶۵	۱۳۹۲/۴/۲۸	۵۵
۲,۷۱	۳۴,۹۶	۹۹۱,۱۳	۷۱,۰۰	۰,۷۹	۱۳۹۲/۴/۲۹	۵۶
۲,۴۵	۳۴,۹۶	۹۹۲,۰۴	۷۳,۸۶	۱,۵۸	۱۳۹۲/۴/۳۰	۵۷
۲,۰۹	۳۴,۹۶	۹۹۱,۸۱	۷۰,۶۳	۱,۱۶	۱۳۹۲/۴/۳۱	۵۸
۲,۳۴	۳۴,۹۶	۹۹۲,۹۸	۶۸,۲۶	۱,۵۹	۱۳۹۲/۵/۱	۵۹

روش ضریب داده پرت محلی، یکی از مهمترین روشهای داده کاوی است. این روش بسیار قدرتمند در شناسایی ناهنجاری‌های داده‌ها در یادگیری ماشین عمل می‌کند [۱۸]. در شماری از کاربردهای عملی مثل تشخیص تقلب در کارت‌های اعتباری [۱۹]، بازاریابی محصولات [۲۰]، امنیت شبکه [۲۱] و شناسایی خطاهای موجود در برداشت داده‌های آزمایشگاهی [۲۲] نیز کاربرد ویژه دارد. روش ضریب داده پرت محلی از مفهوم داده پرت محلی استفاده می‌کند. به این ترتیب که به هر داده درجه‌ای با توجه به چگالی همسایه‌های محلی آن نسبت می‌دهد. این درجه ضریب پرت بودن داده (LOF) نامیده می‌شود. ضریب LOF به میزان ایزوله بودن یک داده وقتی که با همسایه‌های محلی آن مقایسه می‌شود، بستگی دارد. در الگوریتم LOF تفاوت در چگالی بین یک داده و همسایه‌های آن، درجه پرت بودن را مشخص می‌کند. برای هر داده می‌توان مقدار ضریب LOF را محاسبه کرد که میزان پرت بودن آن را نشان می‌دهد. به طور شهودی مقادیر بالای ضریب LOF می‌تواند نماینده داده‌های مشکوک به خطا باشد؛ در حالی که مقدار پایین آن نشان دهنده نرمال بودن یک داده است، چون دلالت بر چگالی پایین همسایگی آن داده دارد. برای محاسبه مقادیر LOF تمامی اعضای یک مجموعه داده، احتیاج به یافتن k تا همسایه‌های نزدیک هر داده است.

- پایش فعالیت: به عنوان مثال شناسایی کلاهبرداری بوسیله گوشی‌های همراه با پایش نحوه فعالیت گوشی و یا خرید و فروش‌های مشکوک تجاری [۶].
- تشخیص عیب: به عنوان مثال پایش فعالیت‌ها برای شناسایی بروز مشکل در موتورها، ژنراتورها، خطوط لوله، پایش سازه‌های صنعتی مثلاً شناسایی وجود ترک در دیوارها یا تجهیزات فضایی در شاتل‌ها [۷، ۸، ۹].
- تحلیل تصاویر ماهواره‌ای: شناسایی پدیده‌های جدید و یا پدیده‌هایی که تا کنون دسته‌بندی نشده‌اند [۱۰، ۱۱].
- شناسایی موارد جدید در تصاویر: به عنوان مثال استفاده از آنها در سیستم‌های تجسسی [۱۲].
- شناسایی پدیده‌های متحرک: شناسایی پدیده‌هایی که نسبت به زمینه تصویر متحرک هستند، به عنوان مثال شناسایی ورود افراد به بانک [۱۳].
- پایش سری‌های زمانی: به عنوان مثال پایش ایمنی در امور مهمی مثل حفاری و یافتن الگوهای غیر عادی در سری‌های زمانی [۱۴].
- پایش وضعیت در امور پزشکی: مثلاً پایش میزان ضربان قلب بیمار [۱۵، ۱۶].
- شناسایی متن: به عنوان مثال شناسایی عناوین در اسناد خبری [۱۷].

برای هر داده p ، k -distance(p)، فاصله k امین همسایه نزدیک p از آن است. برای محاسبه k -distance(p) می توان از الگوریتم kNN استفاده کرد.

به این ترتیب که ابتدا فاصله k تا همسایه نزدیک p را پیدا کرد، سپس k امین فاصله را به عنوان k -distance(p) انتخاب کرد. داده هایی که در k تا فاصله همسایگی p واقع می شوند دارای دو ویژگی هستند:

(الف) برای حداقل k تا داده $o' \in D \setminus \{p\}$ $d(p, o') \leq d(p, o)$
 (ب) برای حداکثر $k-1$ داده $o' \in D \setminus \{p\}$ $d(p, o') < d(p, o)$
 بنابراین k -distance(p) تخمینی از چگالی همسایگی اطراف p را مشخص می کند.

گام دوم: پیدا کردن همسایگی k امین فاصله p ^۳
 هر داده ای که فاصله آن از p ، از k -distance(p) کمتر و یا مساوی باشد، در k امین فاصله همسایگی p قرار می گیرد.

این تعریف به صورت زیر ارائه می شود:

$N_{k\text{-distance}(p)}(p) = \{q \in D \setminus \{p\} \mid d(p, q) \leq k\text{-distance}(p)\}$
 گام سوم: محاسبه چگالی دسترسی داده p نسبت به داده o
^۴ برای هر داده o در k تا فاصله همسایگی p ، چگالی دسترسی

داده p نسبت به داده o به صورت زیر تعریف می شود:

$$\text{reachdist}_k(p, o) = \max\{k\text{-distance}(o), d(p, o)\}$$

شکل (۳) نشان دهنده مثالی از فاصله دسترسی به ازای $k=4$ است. همان طور که در این شکل نشان داده شده است، اگر داده p خارج از k -distance(o) واقع شده باشد (p_2 در شکل)، چگالی دسترسی فاصله بین آنها می شود یعنی $d(o, p_2)$. حال اگر فاصله آنها از k -distance(o) کمتر باشد (p_1 در شکل)، چگالی دسترسی برابر k -distance(o) است.

گام چهارم: محاسبه چگالی دسترسی محلی p ^۵
 چگالی دسترسی محلی p ، معکوس میانگین چگالی دسترسی k تا همسایه های نزدیک p است:

$$\text{lr}d_k(p) = \left[\frac{|N_k(p)|}{\sum_{o \in N_k(p)} \text{reach-dist}_k(p, o)} \right] \quad (1)$$

برای محاسبه چگالی داده p از چگالی دسترسی محلی آن استفاده می شود.

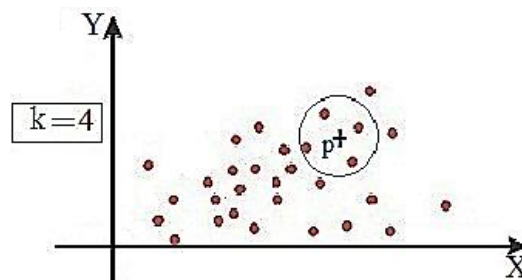
با استفاده از چگالی دسترسی محلی می توان ضریب LOF را محاسبه کرد.

گام پنجم: محاسبه ضریب داده پرت محلی (LOF)

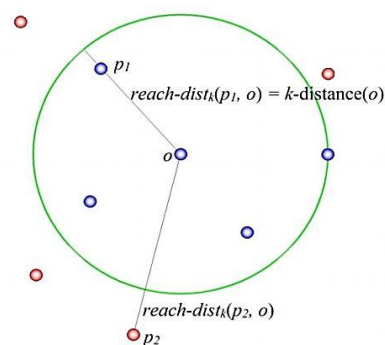
برای یافتن این همسایه ها از الگوریتم k تا نزدیکترین همسایه ها^۱ (kNN) استفاده می شود. قبل از تشریح الگوریتم، پارامترهای مورد استفاده معرفی می شوند. D مجموعه داده، o ، p و q چند داده نمونه که در D واقع شده اند و k یک عدد صحیح مثبت است که توسط کاربر انتخاب می شود. $d(p, q)$ تابع اندازه گیری فاصله بین داده های p و q است که می تواند هر تابع اندازه گیری فاصله باشد. برای محاسبه kNN داده p ، ابتدا فاصله آن با تمامی داده های موجود در مجموعه داده D بدست آید. پس از مرتب کردن فاصله های محاسبه شده به صورت صعودی، k تا داده ای که در فاصله کمتری نسبت به p واقع شده اند، در k تا فاصله همسایگی آن قرار می گیرند، که k تا همسایه محلی آن نامیده می شوند. به عنوان مثال، شکل (۲) نشان دهنده ۴ همسایگی داده p ، از مجموعه داده D است.

جهت محاسبه LOF داده p که در مجموعه داده D قرار دارد، انجام گام های زیر ضروری است:

گام اول: محاسبه k امین فاصله p ^۲



شکل (۲) kNN داده p ، به ازای $k=4$.



شکل (۳) مثالی از چگالی دسترسی به ازای $k=4$.

⁴ Reachability Distance of p w.r.t Object o

⁵ Local Reachability Density of p

¹ k -Nearest Neighbor (kNN)

² k -Distance of p

³ k -Distance Neighborhood of p

$$d_{st} = \sum_{j=1}^n |x_{sj} - y_{tj}| \quad (5)$$

۴- روش ماهالانوبیس^۴

$$d_{st}^2 = (x_s - y_t)C^{-1}(x_s - y_t)' \quad (6)$$

در این رابطه C کوواریانس ماتریس ورودی است.

۵- روش چبیشف^۵

$$d_{st} = \max_j \{|x_{sj} - y_{tj}|\} \quad (7)$$

۶- روش اقلیدسی استاندارد شده^۶

$$d_{st}^2 = (x_s - y_t)V^{-1}(x_s - y_t)' \quad (8)$$

V یک ماتریس قطری n در n است که ز امین المان قطری آن $S(j)^2$ است، که S برداری است که شامل وزنهای معکوس^۷ است.

۷- فاصله کسینوس^۸

$$d_{st} = \left(1 - \frac{x_s - y_t'}{\sqrt{(x_s x_s')(y_t y_t')}}\right) \quad (9)$$

۸- فاصله همبستگی^۹

$$d_{st} = \left(1 - \frac{(x_s - \bar{x}_s)(y_t - \bar{y}_t)'}{\sqrt{(x_s - \bar{x}_s)(x_s - \bar{x}_s)'(y_t - \bar{y}_t)(y_t - \bar{y}_t)'}}\right) \quad (10)$$

که

$$\bar{x}_s = \frac{1}{n} \sum_j x_{sj} \quad (11)$$

$$\bar{y}_t = \frac{1}{n} \sum_j y_{tj} \quad (12)$$

۹- فاصله همینگ^{۱۰}

$$d_{st} = \left(\frac{\#(x_{sj} \neq y_{tj})}{n}\right) \quad (13)$$

۱۰- فاصله جاکاردا^{۱۱}

$$d_{st} = \frac{\#[(x_{sj} \neq y_{tj}) \cap ((x_{sj} \neq 0) \cup (y_{tj} \neq 0))]}{\#[(x_{sj} \neq 0) \cup (y_{tj} \neq 0)]} \quad (14)$$

۱۱- فاصله اسپیرمن^{۱۲}

$$d_{st} = 1 - \frac{(r_s - \bar{r}_s)(r_t - \bar{r}_t)'}{\sqrt{(r_s - \bar{r}_s)(r_s - \bar{r}_s)'(r_t - \bar{r}_t)(r_t - \bar{r}_t)'}} \quad (15)$$

که

r_{sj} رنک x_{sj} می باشد که روی $x_{1j}, x_{2j}, \dots, x_{mj}$ گرفته

شده است. r_{tj} رنک y_{tj} است که روی $y_{1j}, y_{2j}, \dots, y_{mj}$ گرفته شده است.

r_s و r_t رنک بردارهای x_s و y_t است،

یعنی $r_s = (r_{s1}, r_{s2}, \dots, r_{sn})$ و $r_t = (r_{t1}, r_{t2}, \dots, r_{tn})$

$$\bar{r}_s = \frac{1}{n} \sum_j r_{sj} = \frac{n+1}{2} \quad (16)$$

از ضریب LOF برای تشخیص نرمال و یا عدم نرمال بودن یک داده استفاده می شود. LOF(p) میانگین نسبت های چگالی دسترسی محلی p و k تا همسایه های آن است:

$$LOF_k(p) = \frac{\frac{lad_k(o)}{\sum_{o \in N_k(p)} lad_k(p)}}{|N_k(p)|} \quad (2)$$

در ادامه برخی از توابع اندازه گیری فاصله که در این تحقیق در الگوریتم های kNN و LOF بکار رفته اند، به طور مختصر معرفی شده است.

۲-۳- توابع اندازه گیری فاصله

ماتریس $X_{mx \times n}$ داده شده است، می توان آن را به عنوان mx تا بردار سطری $(1 \times n)$ با عناصر $x_1, x_2, x_3, \dots, x_{mx}$ در نظر گرفت. همچنین ماتریس $Y_{my \times n}$ که می توان آن را به عنوان my تا بردار سطری $(1 \times n)$ با عناصر $y_1, y_2, y_3, \dots, y_{my}$ در نظر گرفت. توابع مختلف اندازه گیری فاصله بین عناصر بردار x_s و y_t به صورت زیر تعریف می شوند:

۱- روش مینکوسکی^۱

$$d_{st} = \sqrt[p]{\sum_{j=1}^n |x_{sj} - y_{tj}|^p} \quad (3)$$

p مرتبه فاصله مینکوسکی نامیده می شود که می تواند از ۱ تا بینهایت متغیر باشد. فاصله مینکوسکی برای حالت خاص $p=1$ فاصله بلاک شهری، برای $p=2$ فاصله اقلیدسی و برای حالت خاص $p=\infty$ ، فاصله چبیشف را نتیجه می دهد. n تعداد متغیرها، همچنین x_{sj} و y_{tj} به ترتیب مقادیر z امین متغیر در نقاط x_s و y_t است.

۲- روش اقلیدسی^۲

فاصله اقلیدسی حالت خاصی از فاصله مینکوسکی با $p=2$ است:

$$d_{st}^2 = (x_s - y_t)(x_s - y_t)' \quad (4)$$

۳- روش بلاک شهری^۳

فاصله بلاک شهری حالت خاص $p=\infty$ در رابطه مینکوسکی است:

⁸ Cosine Distance

⁹ Correlation Distance

¹⁰ Hamming Distance

¹¹ Jaccard Distance

¹² Spearman Distance

¹³ Rank

¹ Minkowski Metric

² Euclidean Distance

³ City Block Metric

⁴ Mahalanobis Distance

⁵ Chebychev Distance

⁶ Standardized Euclidean Distance

⁷ Inverse Weights

چگالی کمتر از میزان چگالی یکنواخت بقیه داده‌ها باشد، LOF آن داده از ۱ بیشتر می‌شود و هر چه تراکم داده‌ها حول آن داده کمتر باشد LOF عدد بزرگتری می‌شود و حتی می‌تواند بینهایت شود. بنابراین بسته به مسأله مورد نظر باید عددی بزرگتر از ۱ به عنوان آستانه تشخیص انتخاب شود. در این آزمایشات عدد ۱,۳ به عنوان آستانه انتخاب شده است. این بدان معنی است که داده‌هایی که چگالی اطراف آن‌ها بیشتر از ۳۰ درصد از چگالی یکنواخت داده‌ها کمتر است به عنوان داده‌های مشکوک به خطا در نظر گرفته می‌شوند. با توجه به رابطه LOF می‌توان دید که اگر مقدار آستانه از ۱,۳ بیشتر انتخاب شود برخی از داده‌ها از حوزه داده‌های خطا خارج شده و داده‌های نرمال به حساب خواهند آمد؛ همچنین اگر این آستانه از ۱,۳ کمتر انتخاب شود برخی از داده‌های نرمال وارد حوزه داده‌های خطا شده و به عنوان داده خطا شناخته خواهند شد.

بنابراین در این مطالعه با توجه به توزیع داده‌ها در صفحه دو بعدی و فاصله نسبی نقاط از هم و با توجه به دهک‌های بالا و پایین فواصل داده‌ها از میانگین، مقدار پارامتر آستانه ۱,۳ انتخاب شده است. نتایج آزمون LOF روی نمونه‌ها به ازای $k=7$ و $d=1,3$ ، و بر حسب توابع مختلف اندازه‌گیری فاصله در جداول ۲ تا ۶ ارائه شده است. برای انجام تمامی آزمون‌ها یک کد کامپیوتری در محیط نرم‌افزار MATLAB نوشته شد. همچنین جهت یکسان‌سازی شرایط آزمون از لحاظ زمان اجرا، تمامی آزمون‌ها بر روی یک سیستم با واحد پردازنده مرکزی ۲GHz، رم ۳GB و سیستم عامل ویندوز ۷ (۶۴ بیتی) انجام شده است. زمان اجرا، مدت زمان مورد نیاز جهت اجرای روش توسط کامپیوتر را مشخص می‌کند. هرچه زمان اجرای روش LOF با یک تابع اندازه‌گیری فاصله مشخص کمتر باشد، می‌توان در مدت زمان کمتری به نتایج دست یافت. از اینرو از هزینه‌های زمانی کاسته می‌شود.

همانطور که از نتایج موجود در جداول ۲ تا ۶ مشخص است، روش LOF به ازای توابع اقلیدسی، اقلیدسی استاندارد شده، بلاک شهری، مینکوسکی، چبیشف و ماهالانوبیس روی نمونه‌های مختلف، نتایج کاملاً مشابهی ارائه کرده است. و تنها اختلاف بین آنها، در زمان اجرای روش حین بکارگیری با توابع مختلف اندازه‌گیری فاصله است. بنابراین جهت شناسایی داده‌های مشکوک به خطا در داده‌های گردآوری

$$\bar{r}_t = \frac{1}{n} \sum_j r_{tj} = \frac{n+1}{2} \quad (17)$$

۴- نتایج

در این بخش، نتایج آزمون LOF به ازای توابع مختلف اندازه‌گیری فاصله روی نمونه مذکور ارائه شده است. جهت محاسبه ضریب LOF برای هر داده از پارامتر k استفاده می‌شود. مقدار k وابسته به مسأله مورد نظر و طبیعت داده‌هایی است که قرار است داده‌های مشکوک به خطا از بین آنها شناسایی شود. به طور کلی تمامی الگوریتم‌های شناسایی داده‌های پرت وابسته به یک یا چند پارامتر هستند که توسط کاربر و بر حسب نوع مسئله و شرایط فیزیکی داده‌ها انتخاب می‌شوند. تا کنون هیچ قانون مشخص و پذیرفته شده‌ای برای انتخاب این پارامترها ارائه نشده است و اگر ارائه شده باشد نیز نسبی است و در تمامی حالات معتبر نیست.

در الگوریتم نزدیک‌ترین همسایه نیز این حالت وجود دارد و انتخاب مقدار k از روی سعی و خطا و بسته به شرایط مسئله است. در این الگوریتم به طور کلی یک قانون برقرار است. اگر تعداد همسایه‌های انتخابی به تعداد قابل توجهی کم انتخاب شود آنگاه دقت روش بسیار پایین خواهد آمد در غیر این صورت اگر تعداد همسایه‌های انتخابی از یک حد مشخص بالاتر انتخاب شود آنگاه دقت روش تغییری نخواهد کرد و فقط حجم محاسبات و زمان اجرای برنامه کامپیوتری افزایش خواهد یافت. در این تحقیق، با توجه به اینکه هر داده برداشت شده مربوط به یک روز است، داده‌های برداشت شده طی یک هفته به عنوان داده‌های همسایه در نظر گرفته شده‌اند. بنابراین مقدار k برابر ۷ انتخاب شده است؛ این بدان معنی است که برای هر داده ۷ همسایه نزدیک آن به عنوان حوزه همسایگی آن داده تعیین شده است.

پس از محاسبه ضریب LOF برای هر داده، باید بر اساس ضریب محاسبه شده در مورد نرمال و یا نرمال نبودن آن داده تصمیم‌گیری شود. با توجه به فرمول LOF اگر همه داده‌ها به صورت مرتب و با فاصله دقیقاً یکسانی در صفحه قرار گیرند، آنگاه ضریب LOF همه داده‌ها (به جز داده‌های مرزی) ۱ خواهد بود. از طرفی اگر در اطراف یک داده فشردگی ایجاد شود، یعنی چگالی حول آن داده بیش از سایر داده‌ها شود، آنگاه ضریب LOF مربوط به آن داده از یک کمتر می‌شود و هرچه این چگالی بیشتر باشد این ضریب به صفر نزدیکتر می‌شود. همچنین اگر حول داده‌ای

مقادیر ۰,۳۸، ۰,۳۷ و ۰,۳۷ متر است و این مقادیر به عنوان ارتفاع موج شاخص طبیعی است. و تنها علت انتخاب این داده‌ها، دور افتادگی نسبی آنها از اعضای نمونه است. طبق نتایج جدول (۳)، هیچ داده‌ای به عنوان مشکوک به خطا شناسایی نشده است. با توجه به اینکه بیشتر توابع بر عدم وجود داده خطا در این نمونه تأکید دارند، بنابراین می‌توان گفت این نمونه فاقد خطا بوده و نتایج حاصل شده از بویه قابل اتکا است.

در جدول (۴)، داده‌های با اندیس ۲ و ۳ به عنوان کاندیدای مشکوک به خطا تشخیص داده شده‌اند. در شکل (۶)، این داده‌ها با یک علامت دایره به دور آنها مشخص شده‌اند. با توجه به فیزیک مسئله این داده‌ها نیز خطا نیستند، چون این اندیس‌ها به ترتیب مربوط به داده‌های با مقادیر ۱۰۰۵,۴۷ و ۱۰۰۵,۴۹ اتمسفر بوده و این مقادیر به عنوان فشار هوا طبیعی است. در شکل (۷) و (۸)، داده‌های مشکوک به خطای شناسایی شده در نمونه‌های مربوط به دمای آب و سرعت باد نشان داده شده است. در این نمونه‌ها نیز داده‌های علامت‌گذاری شده تفاوت چندانی با سایر اعضای نمونه نداشته و بر اساس فیزیک مسئله، این داده‌ها نمی‌توانند کاندیدای داده خطا باشند و صرفاً بر اساس تفاوت ناچیزی که با الگوی نرمال جامعه آماری دارند، توسط روش به عنوان کاندیدای خطا انتخاب شده‌اند. بنابراین می‌توان گفت در این مورد نیز داده‌های برداشت شده از بویه قابل اتکا بوده و می‌توان از آنها به منظور اهداف مختلف استفاده کرد. با توجه به نتایج موجود در جداول، در نمونه‌های مختلف داده‌هایی به عنوان کاندیدای مشکوک به خطا انتخاب شدند. و مشاهده شد که با توجه به فیزیک مسئله این داده‌ها نمی‌توانند خطا باشند. بنابراین می‌توان گفت داده‌هایی که به عنوان خطا توسط روش انتخاب می‌شوند، همواره به این معنا نیست که واقعاً خطایی حین گردآوری داده‌ها رخ داده و چه بسا ممکن است بر حسب شرایط خاص اندازه‌گیری به وجود آمده باشند و گویای اطلاعات بسیار مهمی در خصوص سیستم اندازه‌گیری بوده که تا کنون به آنها پرداخته نشده است و یا در حالت‌های بسیار نادر رخ خواهند داد. مثلاً وجود یک داده متناقض در نمونه میزان سرعت باد می‌تواند بیانگر وقوع طوفان‌های شدید در برخی از اوقات سال باشد.

شده از بویه‌ها، استفاده از هر یک از این توابع سبب تغییر نتایج نمی‌شود و می‌توان به صورت دلخواه از هر یک از توابع مذکور استفاده کرد. در اینجا اختلاف زمان اجرای توابع مختلف بسیار ناچیز بوده و قابل چشم‌پوشی است؛ اما هنگام کار روی نمونه‌های بسیار بزرگ، زمان اجرا حائز اهمیت است و می‌تواند سبب صرفه‌جویی در هزینه‌های زمانی و مالی شود. همچنین نتایج اجرای روش LOF بر حسب تابع مینکوسکی و با تغییر مقادیر p تغییر نکرده است.

اما همواره نمی‌توان گفت تغییر مقادیر p سبب تغییر نتایج نشده و می‌توان از یک مقدار p دلخواه استفاده کرد. مثلاً در نمونه سرعت باد اگر مقدار p بیشتر از ۳۶۵ انتخاب شود، مقادیر کلی LOF داده‌ها بیشتر شده و در نتیجه تعداد داده‌های بیشتری به عنوان کاندیدای مشکوک به خطا انتخاب خواهند شد. به عنوان مثال در شکل (۴)، مقادیر ضریب LOF داده‌های میزان سرعت باد به ازای $p=40$ و $p=500$ نشان داده شده است. همانطور که از این شکل مشخص است، مقادیر LOF داده‌ها به ازای $p=500$ نسبت به $p=40$ به طور قابل توجهی بیشتر است. اما در همه نمونه‌ها روش LOF بر حسب توابع کسینوس، همبستگی و اسپیرمن هیچ نتیجه‌ای ارائه نکرده و نتوانسته است با موفقیت اجرا شود.

بنابراین بکارگیری این توابع با روش LOF جهت شناسایی داده‌های مشکوک به خطا در نمونه‌های مذکور پیشنهاد نمی‌شود. همچنین بکارگیری روش LOF با توابع همینگ و جاکارد سبب شده تا تمامی داده‌ها به عنوان داده مشکوک به خطا تشخیص داده شوند. این امر به این دلیل است که ضریب LOF تمامی داده‌ها با استفاده از این توابع برابر ۷,۳۷۵۰ شده است. چون این مقدار از مقدار پارامتر آستانه بیشتر است، در نتیجه تمامی داده‌ها به عنوان داده مشکوک به خطا تشخیص داده شده‌اند. بنابراین می‌توان گفت بکارگیری این دو تابع جهت شناسایی داده‌های مشکوک به خطا در نمونه‌های مذکور نیز سبب نتایج اشتباه شده و استفاده از آنها با روش LOF پیشنهاد نمی‌شود.

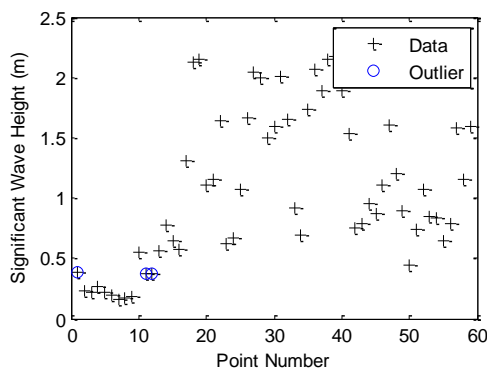
با توجه به نتایج جدول (۲) داده‌های با اندیس ۱، ۱۱ و ۱۲ به عنوان کاندیدای مشکوک به خطا تشخیص داده شده‌اند. در شکل (۵)، این داده‌ها با یک علامت دایره به دور آنها مشخص شده‌اند. با توجه به فیزیک مسئله این داده‌ها خطا نیستند، چون این اندیس‌ها به ترتیب مربوط به داده‌های با

جدول (۲) داده‌های مشکوک به خطا شناسایی شده توسط روش LOF، روی نمونه ارتفاع موج شاخص ($d=1,3$ و $k=7$).

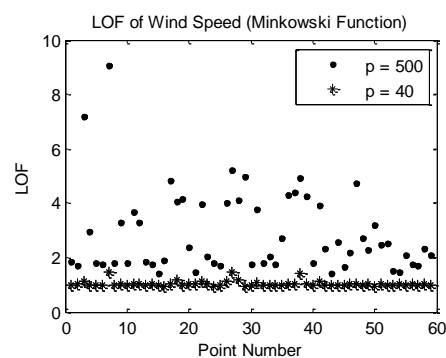
تابع اندازه‌گیری فاصله	شماره (اندیس) داده‌های مشکوک به خطا شناسایی شده	تعداد	زمان اجرا (ثانیه)
اقلیدسی	۱۲، ۱۱، ۱	۳	۰،۴۲۲۴۷۲
اقلیدسی استاندارد شده	۱۲، ۱۱، ۱	۳	۰،۵۴۴۲۲۲
بلاک شهری	۱۲، ۱۱، ۱	۳	۰،۳۲۳۶۴۲
مینکوسکی (p= ۵)	۱۲، ۱۱، ۱	۳	۰،۳۳۲۰۵۰
مینکوسکی (p= ۱۰)	۱۲، ۱۱، ۱	۳	۰،۳۳۲۰۰۶
مینکوسکی (p= ۲۰)	۱۲، ۱۱، ۱	۳	۰،۳۳۲۵۵۱
مینکوسکی (p= ۴۰)	۱۲، ۱۱، ۱	۳	۰،۳۳۶۰۰۱
چپیشف	۱۲، ۱۱، ۱	۳	۰،۳۲۳۹۷۴
ماهالانوبیس	۱۲، ۱۱، ۱	۳	۰،۶۷۷۹۲۵
کسینوس	بدون عملکرد	۰	۰،۳۹۴۳۵۸
همبستگی	بدون عملکرد	۰	بدون عملکرد
اسپیرمن	بدون عملکرد	۰	۰،۴۶۵۰۴۴
همینگ	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۸۶۵۸۴
جاکارد	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۸۸۲۲۱

جدول (۳) داده‌های مشکوک به خطا شناسایی شده توسط روش LOF، روی نمونه میزان رطوبت ($d=1,3$ و $k=7$).

تابع اندازه‌گیری فاصله	شماره (اندیس) داده‌های مشکوک به خطا شناسایی شده	تعداد	زمان اجرا (ثانیه)
اقلیدسی	---	۰	۰،۳۳۵۰۹۸
اقلیدسی استاندارد شده	---	۰	۰،۳۵۳۸۸۳
بلاک شهری	---	۰	۰،۳۱۳۲۳۱
مینکوسکی (p= ۵)	---	۰	۰،۳۲۴۹۳۱
مینکوسکی (p= ۱۰)	---	۰	۰،۳۳۳۶۰۸
مینکوسکی (p= ۲۰)	---	۰	۰،۳۲۷۸۹۴
مینکوسکی (p= ۴۰)	---	۰	۰،۳۲۵۳۶۴
چپیشف	---	۰	۰،۳۱۲۴۵۰
ماهالانوبیس	---	۰	۰،۳۶۲۵۸۰
کسینوس	بدون عملکرد	۰	۰،۳۹۱۹۵۵
همبستگی	بدون عملکرد	۰	بدون عملکرد
اسپیرمن	بدون عملکرد	۰	۰،۴۲۸۸۹۷
همینگ	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۹۳۴۶۴
جاکارد	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۸۷۰۶۹



شکل (۵) نمونه ارتفاع موج شاخص با داده‌های مشکوک به خطا.



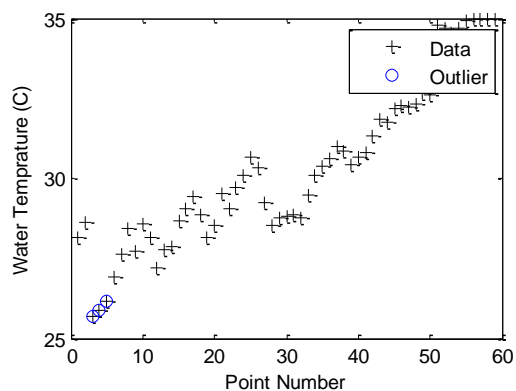
شکل (۴) مقادیر ضریب LOF نمونه سرعت باد با تابع مینکوسکی.

جدول (۴) داده‌های مشکوک به خطا شناسایی شده توسط روش LOF، روی نمونه میزان فشار هوا ($k=7$ و $d=1,3$).

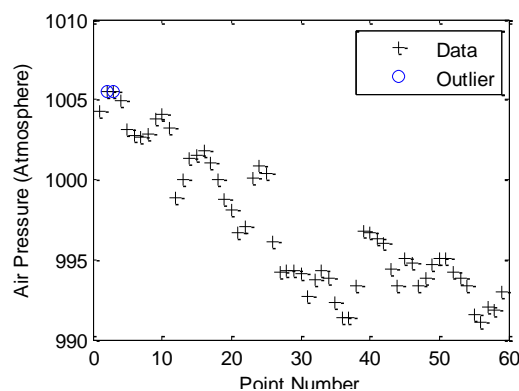
تابع اندازه‌گیری فاصله	شماره (اندیس) داده‌های مشکوک به خطا شناسایی شده	تعداد	زمان اجرا (ثانیه)
اقلیدسی	۳، ۲	۲	۰،۲۰۰۰۶۲
اقلیدسی استاندارد شده	۳، ۲	۲	۰،۳۷۱۹۷۴
بلاک شهری	۳، ۲	۲	۰،۳۳۵۰۷۷
مینکوسکی (۵ = p)	۳، ۲	۲	۰،۳۴۱۰۹۵
مینکوسکی (۱۰ = p)	۳، ۲	۲	۰،۳۳۹۳۷۷
مینکوسکی (۲۰ = p)	۳، ۲	۲	۰،۳۴۴۳۸۴
مینکوسکی (۴۰ = p)	۳، ۲	۲	۰،۳۴۳۰۴۴
چیشف	۳، ۲	۲	۰،۳۲۷۶۳۷
ماهالانوبیس	۳، ۲	۲	۰،۳۸۳۳۹۳
کسینوس	بدون عملکرد	۰	۰،۴۲۰۶۷۷
همبستگی	بدون عملکرد	۰	بدون عملکرد
اسپیرمن	بدون عملکرد	۰	۰،۴۳۲۸۶۵
همینگ	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۹۹۲۸۹
جاکارد	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۹۵۳۳۵

جدول (۵) داده‌های مشکوک به خطا شناسایی شده توسط روش LOF، روی نمونه دمای آب ($k=7$ و $d=1,3$).

تابع اندازه‌گیری فاصله	شماره (اندیس) داده‌های مشکوک به خطا شناسایی شده	تعداد	زمان اجرا (ثانیه)
اقلیدسی	۵، ۴، ۳	۳	۰،۳۴۱۲۱۲
اقلیدسی استاندارد شده	۵، ۴، ۳	۳	۰،۳۷۰۲۴۳
بلاک شهری	۵، ۴، ۳	۳	۰،۳۳۱۰۶۰
مینکوسکی (۵ = p)	۵، ۴، ۳	۳	۰،۳۴۵۷۱۸
مینکوسکی (۱۰ = p)	۵، ۴، ۳	۳	۰،۳۴۳۴۴۶
مینکوسکی (۲۰ = p)	۵، ۴، ۳	۳	۰،۳۴۲۲۵۷
مینکوسکی (۴۰ = p)	۵، ۴، ۳	۳	۰،۳۴۱۸۵۸
چیشف	۵، ۴، ۳	۳	۰،۳۳۱۲۷۴
ماهالانوبیس	۵، ۴، ۳	۳	۰،۳۸۲۴۸۷
کسینوس	بدون عملکرد	۰	۰،۳۹۸۴۳۵
همبستگی	بدون عملکرد	۰	بدون عملکرد
اسپیرمن	بدون عملکرد	۰	۰،۴۳۹۳۳۵
همینگ	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۹۵۰۳۲
جاکارد	۱، ۲، ۳، ۴، ۵، ...، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹	۵۹	۰،۳۹۶۷۲۶



شکل (۷) نمونه میزان دمای آب با داده‌های مشکوک به خطا.



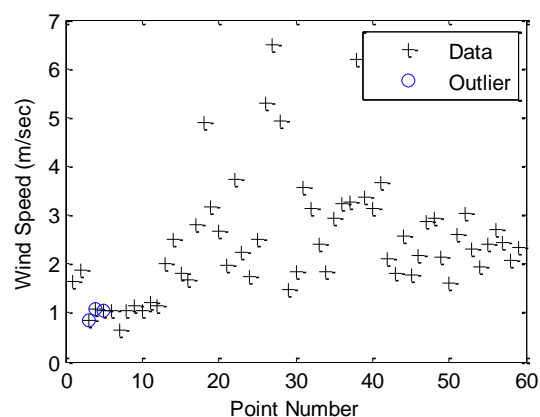
شکل (۶) نمونه میزان فشار هوا با داده‌های مشکوک به خطا.

جدول (۶) داده‌های مشکوک به خطا شناسایی شده توسط روش LOF، روی نمونه میزان سرعت باد ($k=7$ و $d=1.3$).

تابع اندازه‌گیری فاصله	شماره (اندیس) داده‌های مشکوک به خطا شناسایی شده	تعداد	زمان اجرا (ثانیه)
اقلیدسی	۳۸، ۲۷، ۷	۳	۰،۳۳۱۲۲۷
اقلیدسی استاندارد شده	۳۸، ۲۷، ۷	۳	۰،۳۷۱۵۶۷
بلاک شهری	۳۸، ۲۷، ۷	۳	۰،۳۲۵۳۶۰
مینکوسکی ($p=5$)	۳۸، ۲۷، ۷	۳	۰،۳۳۸۲۳
مینکوسکی ($p=10$)	۳۸، ۲۷، ۷	۳	۰،۳۴۰۵۹۵
مینکوسکی ($p=20$)	۳۸، ۲۷، ۷	۳	۰،۳۴۱۸۲۳
مینکوسکی ($p=40$)	۳۸، ۲۷، ۷	۳	۰،۳۵۴۶۳۰
چبیشف	۳۸، ۲۷، ۷	۳	۰،۳۲۷۹۶۳
ماهالانوبیس	۳۸، ۲۷، ۷	۳	۰،۳۷۶۵۸۸
کسینوس	بدون عملکرد	۰	۰،۴۰۱۲۴۱
همبستگی	بدون عملکرد	۰	بدون عملکرد
اسپیرمن	بدون عملکرد	۰	۰،۴۳۲۵۰۱
همینگ	۵۹، ۵۸، ۵۷، ۵۶، ۵۵، ...، ۵، ۴، ۳، ۲، ۱	۵۹	۰،۳۹۸۴۰۲
جاکارد	۵۹، ۵۸، ۵۷، ۵۶، ۵۵، ...، ۵، ۴، ۳، ۲، ۱	۵۹	۰،۳۹۶۰۰۱

فیزیک مسئله کمک قابل توجهی کند. در این مقاله عوامل بروز خطا در داده‌های برداشت شده از بویه‌های موج‌نگار بررسی شد و از روش ضریب داده پرت محلی که یکی از قدرتمندترین روش‌های یادگیری ماشین جهت تشخیص ناهنجاری در داده‌ها است، جهت شناسایی داده‌های مشکوک به خطا در یک نمونه واقعی از داده‌های برداشت شده از بویه‌های موج‌نگار استفاده شد. مهمترین نتایج این تحقیق به شرح زیر است:

- بررسی‌ها حاکی از عملکرد مناسب روش ضریب داده پرت محلی جهت شناسایی داده‌های مشکوک به خطا در نمونه دیتای موج مورد بررسی است.
- در نمونه‌های مورد بررسی روش ضریب داده پرت محلی به ازای توابع اقلیدسی، اقلیدسی استاندارد شده، بلاک شهری، مینکوسکی، چبیشف و ماهالانوبیس نتایج کاملاً مشابهی ارائه کردند. بنابراین می‌توان به دلخواه از هر یک از این توابع استفاده کرد. اما توابع کسینوس، همبستگی، اسپیرمن، همینگ و جاکارد نتوانسته‌اند با موفقیت اجرا شوند. بنابراین بکارگیری این توابع با روش ضریب داده پرت محلی پیشنهاد نمی‌شود.
- داده‌های علامت‌گذاری شده توسط روش همواره بیانگر وجود خطا در داده‌ها نیست و چه بسا ممکن است به دلیل شرایط خاص مسئله ایجاد شده باشند. بنابراین پس از شناسایی این داده‌ها، باید در خصوص علت



شکل (۸) نمونه میزان سرعت باد با داده‌های مشکوک به خطا.

از اینرو توجه مسئولین و مهندسان می‌تواند به آن جلب شود تا در صورت وقوع دوباره آن، تدابیر لازم را بکار گیرند. بنابراین می‌توان گفت که از تمامی این روش‌ها فقط به منظور تسهیل در شناسایی داده‌های مشکوک به خطا استفاده می‌شود و تشخیص نهایی خطا و یا صحیح بودن داده بر عهده فرد متخصص است. اوست که با توجه به ماهیت مسئله و شرایط داده‌ها، تصمیم‌گیری نهایی در خصوص خطا و یا نرمال بودن آن را انجام می‌دهد.

۵- نتیجه‌گیری

حین انجام آزمایش‌های مختلف، عوامل متعددی سبب بروز خطا در اندازه‌گیری‌ها می‌شوند. شناسایی این خطاها و در صورت امکان برطرف کردن آنها می‌تواند به بیان واقعیت

- Biomedical Image Analysis, IEEE Computer Society, Washington, DC, USA, Vol. 3, 2001.
- [13] Pokrajac, D., Lazarevic, A., and Latecki, L. J. "Incremental Local Outlier Detection for Data Streams", In Proceedings of IEEE Symposium on Computational Intelligence and Data Mining, pp.1-11, 2007.
- [14] Yankov, D., Keogh, E., and Rebbapragada, U. "Disk Aware Discord Discovery: Finding Unusual Time Series in Terabyte Sized Datasets", Knowledge and Information Systems, Vol.17, No.2, pp.241-262, 2008.
- [15] Solberg, H. E., and Lahti, A. "Detection of Outliers in Reference Distributions: Performance of Horn's algorithm", Clinical Chemistry, Vol.51, No.12, pp.2326-2332, 2005.
- [16] Suzuki, E., Watanabe, T., Yokoi, H., and Takabayashi, K. "Detecting Interesting Exceptions From Medical Test Data With Visual Summarization", In Proceedings of the 3rd IEEE International Conference on Data Mining, pp.315-322, 2003.
- [17] Chen, D., Odobez, J. M., and Boulard, H. "Text Detection and Recognition in Images and Video Frames", Pattern Recognition, Vol.37, No.3, pp.595-608, March 2004.
- [18] Breuning, M., Kriegel, H-P., Ng, R., and Sander, J., "LOF: Identifying Density Based Local Outliers", In Proc. of ACM SIGMOD International Conference on Management of Data (SIGMOD'00), Dallas, Texas, pp.93-104, 2000.
- [19] Chen, M. C., Wang, R. and Chen, A. P., "An Empirical Study for the Detection of Corporate Financial Anomaly Using Outlier Mining Techniques", In ICCIT'07: Proc. of the International Conference on Convergence Information Technology, pp.612-617, 2007.
- [20] Xi, J., "Outlier Detection Algorithms in Data Mining", Intelligent Information Technology Applications", Vol.1, pp.94-97, 2008.
- [۲۱] محمودی، کیومرث، کتابداری، محمد جواد، و سایبانی، مصباح. "شناسایی نفوذ به شبکه‌های کامپیوتری سیستم‌های نظامی، به روش تشخیص ناهنجاری"، اولین همایش ملی فناوری‌های نوین دریایی، نوشهر، دانشگاه امام خمینی (ره)، ۱۲ تا ۱۳ شهریور ماه ۱۳۹۲.
- [۲۲] محمودی، کیومرث، سایبانی، مصباح، و مرادی، عباس. "شناسایی خطاهای موجود در برداشت داده‌های آزمایشگاهی، با استفاده از الگوریتم‌های تشخیص داده خطا"، پنزدهمین کنفرانس دینامیک شاره‌ها FD2013، بندرعباس، دانشگاه هرمزگان، ۲۷-۲۹ آذر ۱۳۹۲.
- وقوع آنها تحقیق کرده و در صورتی که خطا باشند نسبت به حذف و یا تصحیح آنها اقدام کرد.
- ### ۶- مراجع
- [۱] واقفی، محمد، قدسیان، مسعود، و ادیب، آرش. "نگرشی بر خطاهای موجود در برداشت داده‌های آزمایشگاهی"، نهمین کنفرانس هیدرولیک ایران، دانشگاه تربیت مدرس، آبان ماه ۱۳۸۹.
- [۲] کتابداری، محمد جواد، و شهریبجاری، علی نظری. "مدل کامپیوتری جامع جهت آنالیز دیتاهای خام ناشی از بویه‌های موج نگار"، هشتمین همایش صنایع دریایی، بوشهر، ۹ لغایت ۱۱ آبانماه ۱۳۸۵.
- [3] <http://marinedata.pmo.ir>
- [4] Bolton, R., and Hand, D. "Unsupervised Profiling Methods for Fraud Detection", In Credit Scoring and Credit Control VII, 1999.
- [5] Eskin, E., "Outlier Detection Over Noisy Data using Learned Probability Distributions", In Proceedings of the 7th International Conference on Machine Learning, Morgan Kaufmann Publishers Inc., pp.255-262, 2000.
- [6] Aggarwal, C., "On Abnormality Detection in Spuriously Populated Data Streams", In Proceedings of 5th SIAM Data Mining, pp.80-91, 2005.
- [7] Manson, G. "Identifying Damage Sensitive, Environment Insensitive Features for Damage Detection", In Proceedings of the IES Conference, Swansea, UK, pp. 1-5, 2002.
- [8] Hollier, G., Austin, J. "Novelty Detection for Strain-Gauge Degradation using Maximally Correlated Components", In Proceedings of the European Symposium on Artificial Neural Networks, pp.262-539, 2002.
- [۹] محمودی، کیومرث، و سایبانی، مصباح. "سلامت‌سنجی سیستم‌های دریایی به روش تشخیص ناهنجاری در یادگیری ماشین"، اولین همایش ملی فناوری‌های نوین دریایی (MMT2013)، مازندران، دانشگاه علوم دریایی امام خمینی (ره)، ۱۶ الی ۱۷ تیر ۱۳۹۲.
- [10] Augusteijn, M., Folkert, B. "Neural Network Classification and Novelty Detection", International Journal on Remote Sensing, Vol. 23, No. 14, pp. 2891-2902, 2002.
- [11] Theiler, J., and Cai, D. M. "Re-sampling Approach for Outlier Detection in Multispectral Images", In Proceedings of SPIE, Vol.5093, pp.230-240, 2003.
- [12] Spence, C., Parra, L., and Sajda, P. "Detection, Synthesis and Compression in Mammographic Image Analysis with a Hierarchical Image Probability Model", In Proceedings of the IEEE Workshop on Mathematical Methods in