

## ارزیابی عملکرد روش‌های مختلف یادگیری ماشین در شناسایی نوع حمل و نقل با استفاده از داده‌های خط سیر

مر ترضی طیبی<sup>۱</sup>، پرهام پهلوانی<sup>۲\*</sup>

۱- دانشجوی دکتری سیستم‌های اطلاعات مکانی - دانشکده مهندسی نقشه‌برداری و اطلاعات مکانی - دانشکده فنی - دانشگاه تهران  
۲- دانشیار دانشکده مهندسی نقشه‌برداری و اطلاعات مکانی - دانشکده فنی - دانشگاه تهران

تاریخ دریافت مقاله: ۱۴۰۱/۰۲/۲۶ تاریخ پذیرش مقاله: ۱۴۰۱/۱۱/۰۱

### چکیده

با ظهور گسترده گوشی‌های هوشمند که به سامانه موقعیت‌یاب جهانی (GPS) مجهز هستند، حجم انبوهی از داده‌های مکانی خط سیر کاربران ایجاد شده است. مطالعه بر روی این داده‌ها در راستای تسهیل مدیریت شهری و ارائه مناسب خدمات به کاربران به‌عنوان یک زمینه تحقیقاتی گسترده مطرح شده و در حال رشد است. در این تحقیق به شناسایی نوع حمل‌ونقل خطوط سیر کاربران بر مبنای داده‌های خام GPS آن‌ها پرداخته شده است. این داده‌ها غالباً دارای خطاهایی هستند که در این تحقیق با اعمال یک فرآیند پیش‌پردازش چندمرحله‌ای سعی شده است مقدار خطا به حداقل برسد. سپس به‌منظور شناسایی نوع حمل‌ونقل شامل پیاده‌روی، دوچرخه، قطار، اتوبوس و رانندگی ویژگی‌های متعددی استخراج می‌شود. در ادامه به‌منظور ساختن مدل پیش‌بینی‌کننده از چهار روش طبقه‌بندی درخت تصمیم، شبکه عصبی پرسپترون چندلایه، بیز ساده و ماشین بردار پشتیبان استفاده می‌شود. در جهت بهبود عملکرد روش‌های پیاده‌سازی، از درصد حضور نقاط هر خط سیر در فاصله یک انحراف معیار از میانگین کل سرعت نوع‌های حمل‌ونقل به‌عنوان یک ویژگی جدید استفاده شده است. پیاده‌سازی چهار مدل یادشده به ازای پارامترهای تنظیم‌کننده مختلف انجام شده و پس از یک جستجوی جامع شبکه‌ای پارامترهای مختلف موجود در این روش‌ها در بهینه‌ترین مقدار تنظیم می‌شوند. در ادامه از دو شاخص کاپا و دقت کلی برای ارزیابی روش‌های مختلف استفاده می‌شود. نتایج حاصل از این تحقیق نشان داد که شبکه عصبی پرسپترون چندلایه با دقت کلی ۰/۸۸ توانست بهترین نتایج را نسبت به سایر مدل‌ها از خود نشان دهد.

کلیدواژه‌ها: داده‌های خط سیر، تعیین نوع حمل‌ونقل، طبقه‌بندی، یادگیری ماشین.

\* نویسنده مکاتبه‌کننده: خیابان کارگر شمالی، دانشکده مهندسی نقشه‌برداری و اطلاعات مکانی، دانشکده فنی، دانشگاه تهران.

تلفن: ۰۲۱۶۱۱۱۴۵۲۴

Email: [pahlavani@ut.ac.ir](mailto:pahlavani@ut.ac.ir)

## ۱- مقدمه

امروزه گسترش محیط‌های شهری و افزایش روزافزون جمعیت در کلان‌شهرها چالش‌های زیادی مانند ترافیک، آلودگی هوا، تأمین تسهیلات شهری و غیره با خود به همراه داشته است. از سوی دیگر منابعی نظیر گوشی‌های تلفن همراه، شبکه‌های اجتماعی، سامانه موقعیت‌یابی جهانی<sup>۱</sup> (GPS) خودروها و غیره موجب تولید داده‌های مکانی انبوه و ارزشمند شده است. دولت‌مردان در راستای حل چالش‌های یادشده و اتخاذ تصمیم‌های حوزه مدیریت شهری نیاز به شناخت محیط‌های شهری و درک رفتار حرکتی شهروندان دارند که می‌تواند از منابع مذکور حاصل شود. محاسبات شهری<sup>۲</sup> اصطلاحی است که به این حوزه از مطالعات مکانی شهرها اختصاص داده شده است [۱].

پیشرفت در تکنیک‌های محاسبه حرکت و دریافت موقعیت، داده‌های خط سیر مکانی حجیمی را تولید کرده است که حرکت اشیاء متحرک گوناگونی مانند مردم، وسایل نقلیه و حیوانات را ارائه می‌کند. تکنیک‌های زیادی برای پردازش، مدیریت و کاویدن داده‌های خط سیر در دهه گذشته پیشنهاد شده است که توسعه‌دهنده دامنه وسیعی از کاربردها است. هر خط سیر مکانی یک دنباله تولیدشده توسط یک شی متحرک در فضای جغرافیایی است که معمولاً به وسیله یک سری نقاط مرتب بر اساس زمان ارائه می‌شود. مثلاً  $p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_n$  که هر نقطه شامل یک مجموعه مختصات مکانی و یک لحظه زمانی است مانند  $p = (x, y, t)$  [۲].

در چند سال اخیر، تحقیقات متعددی در حوزه داده‌های خط سیر مکانی انجام شده است. این تحقیقات از یک دیدگاه کلی به سه‌شاخه اصلی تقسیم می‌شوند:

۱- مدل‌های داده: تعریف و توسعه نوع داده‌های خط سیر مانند نقاط و نواحی مشترک [۳ و ۴].

۲- مدیریت داده: ذخیره‌سازی مؤثر داده‌های حرکت با نمایه‌گذاری و تکنیک‌های پرس و جوی ویژه [۵، ۶] که منجر به توسعه چندین سیستم مدیریت پایگاه داده مانند [۷] *SECONDO*، [۸] *HERMES* و [۹] *DOMINO* شده است.

۳- داده‌کاوی: طراحی و پیاده‌سازی الگوریتم‌های داده‌کاوی نظیر خوشه‌بندی، طبقه‌بندی، کشف خطا و الگو کاوی بر روی داده‌های خط سیر مانند ژئولایف<sup>۳</sup>، *GeoPKDD* و *MoveMine* [۲].

داده‌کاوی خطوط سیر منجر به حصول نتایجی می‌شود که کاربردهای فراوانی نظیر شناسایی کاربران مشابه، تشخیص ناهنجاری‌های ترافیکی، شناسایی نقاط و مناطق جذاب و غیره دارند. تعیین نوع حمل‌ونقل<sup>۴</sup> یکی از کاربردهایی است که می‌توان با بهره‌گیری از تکنیک‌های داده‌کاوی به آن دست‌یافت. نوع حمل‌ونقل یکی از مهم‌ترین ویژگی‌ها در رفتارهای حرکتی کاربران است که در تحلیل تقاضای سفر، مدیریت ترافیک و برنامه‌ریزی حمل‌ونقل تأثیرگذار است [۱۰]. با استفاده از توزیع نوع سفر، آژانس‌های حمل‌ونقل می‌توانند استراتژی‌های مناسبی را برای کاهش زمان سفر کاربران، کاهش ترافیک و کاهش آلودگی هوا ایجاد کنند. به‌عنوان مثال یک نمونه بارز از کاربرد تحلیل‌های نوع سفر شناسایی مناطق وابسته و بهبود سامانه‌های حمل‌ونقل در راستای تشویق به استفاده از حمل‌ونقل عمومی است [۱۱]. اعمال محدودیت‌ها و قوانین ترافیکی مانند خطوط ویژه وسایل نقلیه عمومی<sup>۵</sup> در ساعات اوج شلوغی از دیگر کاربردهای مقدماتی تحلیل‌های نوع سفر است. پیش‌از این اطلاعات مربوط به نوع سفر از طریق مصاحبه‌های حضوری و تلفنی انجام

<sup>۳</sup> *GeoLife*<sup>۴</sup> *Transportation mode*<sup>۵</sup> *High Occupancy Vehicle (HOV) Lanes*<sup>۱</sup> *Global Positioning System*<sup>۲</sup> *Urban computing*

مورد مطالعه، جزئیات پیاده‌سازی روش تحقیق تشریح می‌شود. در بخش ششم نتایج تحقیق ارائه شده و مورد بحث قرار می‌گیرد. در نهایت، در بخش هفتم نتیجه‌گیری حاصل از تحقیق بیان شده و پیشنهادهایی برای تحقیقات آینده ارائه خواهد شد.

## ۲- مطالعات پیشین

در سال‌های اخیر داده‌کاوی بر روی داده‌های خام *GPS* که منجر به ایجاد خطوط سیر می‌شوند، به یک موضوع تحقیقاتی گسترده تبدیل شده است. این مطالعات شامل استخراج داده‌ها، پیش‌پردازش، مدیریت پایگاه داده و بازیابی اطلاعات، خوشه‌بندی و طبقه‌بندی، الگو کاوی خطوط سیر، عدم قطعیت<sup>۱</sup> و زمینه‌های دیگری است که در همکاری طی یک فرآیند داده‌کاوی به کاربردهای مختلفی منجر می‌شوند [۲ و ۱۴]. در برخی تحقیقات مانند [۱۵ و ۱۶] خطوط سیر به فرمت‌های دیگری مانند گراف، ماتریس و تنسور تبدیل می‌شوند تا تکنیک داده‌کاوی بیشتری بر روی آن‌ها قابل پیاده‌سازی باشد.

در تحقیقات داده‌کاوی خطوط سیر غالباً از داده‌های خام *GPS*، داده‌های شتاب سنج گوشی‌های همراه و داده‌های سامانه جهانی موبایل<sup>۲</sup> (*GSM*) استفاده می‌شود [۱۰]. همچنین افزودن برخی از داده‌های زمینه‌ای مانند دمای بدن، ضربان قلب، رطوبت و شدت نور در کنار داده‌های مکانی نظیر شبکه راه‌ها، نقاط جذاب<sup>۳</sup> و کاربری‌ها در گروهی از مطالعات مورد توجه قرار گرفته‌اند [۱۳]. در این زمینه مطالعات قابل توجهی با به‌کارگیری تکنیک‌های داده‌کاوی از جمله مدل‌های طبقه‌بندی، بر روی برآورد نوع حمل‌ونقل، برآورد نوع فعالیت و برآورد مقصد متمرکز شده‌اند. برخی از تحقیقات روش‌هایی برای برآورد نوع حمل‌ونقل ارائه

می‌پذیرفت که یک روش پرهزینه، کم‌دقت و ناقص بود، اما امروزه دسترسی ما به داده‌های خط سیر اطلاعات دقیق، کم‌هزینه، به‌روز و کاملی را در اختیار می‌گذارد [۱۰].

وجود منابع خطا داده‌های خام خطوط سیر را تحت تأثیر قرار می‌دهد و این خطا به ویژگی‌های خطوط سیر انتشار می‌یابد و دقت مدل‌ها را تحت تأثیر قرار می‌دهد [۱۲]. بنابراین پالایش داده‌ها و کاهش خطا جهت بهبود دقت شناسایی نوع‌های حمل‌ونقل ضروری است، در این راستا در این تحقیق یک فرآیند پیش‌پردازش چندمرحله‌ای جهت حذف داده‌های اشتباه و کاهش خطا پیاده‌سازی می‌شود و تأثیر آن بر دقت نتایج مورد بررسی قرار می‌گیرد.

برای شناسایی نوع حمل‌ونقل با استفاده از داده‌های خط سیر، استخراج ویژگی‌ها از اهمیت بالایی برخوردار است. در این تحقیق علاوه بر بهره‌گرفتن از مهم‌ترین و تأثیرگذارترین ویژگی‌هایی که در تحقیقات پیشین معرفی شده‌اند [۱۰، ۱۲ و ۱۳]، درصد حضور نقاط هر خط سیر در فاصله یک انحراف معیار از میانگین کل سرعت نوع‌های حمل‌ونقل به‌عنوان یک ویژگی جدید معرفی می‌شود. همچنین رفتار نوع‌های مختلف حمل‌ونقل در بازه‌های سرعت مختلف در جهت بهبود دقت شناسایی نوع‌های حمل‌ونقل مورد بررسی قرار می‌گیرد.

در سال‌های اخیر تحقیقات گسترده‌ای در راستای تشخیص نوع حمل‌ونقل بر مبنای داده‌های خط سیر *GPS* انجام شده است. این تحقیق تلاش می‌کند عملکرد روش‌های مختلف یادگیری ماشین را به ازای پارامترهای مختلف مورد بحث قرار داده و نتایج آن‌ها را مقایسه کند.

در ادامه این مقاله، در بخش دوم مطالعات انجام‌شده در زمینه تعیین نوع حمل‌ونقل در داده‌های خط سیر مرور می‌گردد. بخش سوم به تشریح مبانی نظری تحقیق و بخش چهارم به بیان روش پیاده‌سازی تحقیق می‌پردازد. در بخش پنجم ضمن معرفی داده‌های

<sup>۱</sup> Uncertainty

<sup>۲</sup> Global System for Mobile

<sup>۳</sup> Point Of Interest

عمرانی و همکاران (۲۰۱۵) از شبکه عصبی مصنوعی<sup>۳</sup> (*ANN*) از جمله شبکه عصبی پرسپترون چندلایه<sup>۴</sup> (*MLP*) و تابع پایه شعاعی<sup>۵</sup> (*RBF*) در راستای پیش‌بینی نوع حمل‌ونقل استفاده کردند. اندو و همکاران (۲۰۱۶) مدلی ارائه دادند که خطوط سیر را به یک تصویر دوبعدی تبدیل می‌کند و مقدار هر پیکسل برابر با مدت‌زمانی است که کاربر در آن مکان حضور داشته است. سپس با معرفی این تصویر به یک شبکه عصبی عمیق تعدادی ویژگی خودکار استخراج شده و پس از ادغام با ویژگی‌های سنتی مرسوم به مدل طبقه‌بندی معرفی می‌شود [۱۳]. آن‌ها روش‌های ماشین بردار پشتیبان<sup>۶</sup> (*SVM*)، درخت تصمیم و رگرسیون لجستیک<sup>۷</sup> پیاده‌سازی کردند و درخت تصمیم را به‌عنوان بهترین مدل معرفی کردند. این تحقیق علیرغم پیچیدگی زمانی بالایی که دارد موفق به کسب دقت قابل قبولی نشد. در تحقیق مشابه دیگری ونگ و همکاران (۲۰۱۷) با استخراج ویژگی‌هایی در سطح نقطه و معرفی این ویژگی به یک شبکه عصبی پیچشی<sup>۸</sup> ویژگی‌های بیشتری استخراج کرده و به مدل طبقه‌بندی معرفی کردند [۲۲].

در تحقیق دیگری خیائو و همکاران (۲۰۱۷) با تولید تعداد صد و یازده ویژگی آماری در سطح کل خطوط سیر و در سطح محلی برای قطعه‌های کوچک‌تر خطوط سیر، مدلی بر مبنای درخت تصمیم ارائه کردند [۲۳]. دبیری و همکاران در سال ۲۰۱۹، با به‌کارگیری شبکه عصبی پیچشی چهار ویژگی اصلی سرعت، شتاب، تغییرات شتاب و تغییر جهت را به لایه ورودی شبکه

کردند که از داده‌های ثابتی مانند نقشه راه‌ها استفاده می‌کند [۱۷ و ۱۸]. با توجه به ماهیت پویای ساختار شهری دسترسی آسان و به‌روز به این‌گونه اطلاعات خارجی میسر نیست [۱۳].

گروه تحقیقاتی شرکت مایکروسافت به سرپرستی ژنگ تحقیقات زیادی بر روی داده‌های خط سیر انجام داده‌اند. یکی از مهم‌ترین تحقیقات انجام‌شده توسط این گروه پروژه ژئولایف بود که هدف آن ارائه یک شبکه اجتماعی مکان‌مبنا بود که بر مبنای داده‌های خط سیر *GPS* که در شبکه به اشتراک گذاشته می‌شد عمل نماید [۱۹]. ژنگ و همکاران (۲۰۰۸) مدلی برای برآورد نوع حمل‌ونقل بر مبنای داده‌های پروژه ژئولایف که داده خام *GPS* محسوب می‌شوند، ارائه کردند. آن‌ها ابتدا با یک روش قطعه‌بندی نقطه مبنا تلاش کردند خطوط سیر با یک نوع حمل‌ونقل خاص ایجاد کنند. سپس به استخراج ویژگی‌های مهمی برای خطوط سیر پرداختند که پس از آن در تحقیقات دیگری نیز موردتوجه قرار گرفتند. از جمله مهم‌ترین ویژگی‌های استفاده شده در این تحقیق سرعت، انحراف معیار سرعت، شتاب، نرخ توقف، تغییر جهت و غیره بود. سپس از با پیاده‌سازی مدل‌های طبقه‌بندی به مدلی برای پیش‌بینی نوع حمل‌ونقل دست یافتند [۲۰]. سپس در تحقیق دیگری در همان سال با افزودن یک مرحله پس پردازشی گراف مبنا به بهبود مدل پیشنهادی خود پرداختند. این مرحله پس پردازش بر مبنای احتمال وقوع هر نوع حمل‌ونقل در یک دنباله از نوع‌های حمل‌ونقل عمل می‌کند [۲۱]. آن‌ها در تحقیقات خود از مدل‌های مختلف طبقه‌بندی استفاده کردند و بهترین نتایج را از مدل‌های رندم فارست<sup>۱</sup> و درخت تصمیم<sup>۲</sup> (*DT*) به دست آوردند.

<sup>۳</sup> Artificial Neural Network

<sup>۴</sup> Multi Layer Perceptron

<sup>۵</sup> Radial Basis Function

<sup>۶</sup> Support Vector Machine

<sup>۷</sup> Logistic Regression

<sup>۸</sup> Convolutional Neural Network

<sup>۱</sup> Random Forest

<sup>۲</sup> Decision Tree

نوع حمل‌ونقل داده‌های خط سیر استفاده می‌کند. با توجه به اینکه مناطق شهری دارای نقاط کور برای دریافت سیگنال توسط GPS های کاربران هستند و وجود منابع خطای دیگر، انجام یک فرآیند پیش‌پردازش چندمرحله‌ای در این تحقیق موردتوجه قرار گرفته است که در تحقیقات گذشته غالباً کمتر به آن توجه شده است. همچنین در این تحقیق ویژگی‌های جدیدی برای خطوط سیر ارائه می‌شود. در این تحقیق از چهار روش ماشین بردار پشتیبان، درخت تصمیم، شبکه عصبی پرسپترون چندلایه و بیز ساده جهت پیاده‌سازی مدل طبقه‌بندی استفاده می‌شود. روش ماشین بردار پشتیبان که در بسیاری از زمینه‌های تحقیقاتی عملکرد مناسبی از خود نشان داده است، در این تحقیق نیز مورد استفاده قرار گرفته است. در تحقیقات گذشته روش SVM برای طبقه‌بندی خطوط سیر بر اساس نوع حمل‌ونقل کمتر موردتوجه قرار گرفته و نتوانسته است دقت مناسبی از خود نشان دهد، درخت تصمیم نیز که پیش‌از این در تحقیقات مشابه با موفقیت به کار گرفته شده است، در این تحقیق با اعمال ضریب جینی، بهره‌های اطلاعاتی و نسبت بهره به‌عنوان معیارهای طبقه‌بندی، مورد استفاده قرار می‌گیرد. شبکه‌های عصبی امروزه در تحقیقات مختلفی توانایی خود را در یادگیری رفتار داده‌های مختلف و پیش‌بینی آن‌ها نشان داده‌اند. در این تحقیق از شبکه عصبی پرسپترون چندلایه استفاده شده است. روش طبقه‌بندی کننده بیز ساده نیز که پیش‌از این به‌ندرت در تحقیقات مشابه موردتوجه قرار گرفته است، در این تحقیق مورد ارزیابی قرار می‌گیرد. این تحقیق تلاش می‌کند با یک جست‌وجوی جامع بهترین پارامترهای ورودی برای روش‌های مختلف طبقه‌بندی را کشف کند و نتایج حاصل با نتایج تحقیقات پیشین مقایسه شود.

### ۳- مبانی نظری تحقیق

در این بخش به تشریح مبانی نظری تحقیق پرداخته می‌شود. در این تحقیق از چهار روش طبقه‌بندی

معرفی کردند و ویژگی‌های سطح بالا به شکل خودکار توسط شبکه استخراج شدند. با توجه به اینکه خطوط سیر با طول‌های متنوع قابل‌ارائه به ورودی شبکه عصبی پیچشی نبودند، آن‌ها هر خط سیر را به قطعات ۲۰۰ نقطه‌ای تقسیم کردند [۲۴].

نواز و همکاران (۲۰۲۰) نیز از شبکه عصبی پیچشی جهت استخراج ویژگی‌های سطح بالا استفاده کردند، سپس از مدل حافظه طولانی کوتاه‌مدت<sup>۱</sup> (LSTM) برای آموزش مدل شناسایی نوع حمل‌ونقل استفاده کردند. آن‌ها با افزودن داده‌های هواشناسی به لایه ویژگی‌ها موفق شدند دقت مدل پیشنهادی خود را هفت درصد نسبت به تحقیقات پیشین بهبود بخشند [۲۵]. جیمز و همکاران نیز در سال ۲۰۲۰، ضمن استفاده از شبکه عصبی عمیق از تبدیل موجک گسسته برای استخراج ویژگی‌های دامنه زمان بسامد برای کمک به شبکه عصبی استفاده کردند [۲۶].

در تحقیق دیگری هوانگ و همکاران (۲۰۲۰) با به‌کارگیری مدل جنگل ایزوله مشارکتی<sup>۲</sup> هر نوع حمل‌ونقل را به شکل جداگانه در یک مدل شناسایی کرده و نتایج را تلفیق کردند. این تحقیق نشان داد که مدل جنگل ایزوله مشارکتی در برخورد با داده‌های خط سیری که دارای خطای اتفاقی هستند انعطاف بیشتری نسبت به سایر مدل‌ها نشان می‌دهد [۲۷].

تحقیقات گذشته در راستای تشخیص نوع حمل‌ونقل داده‌های GPS، مدل‌های متنوعی را به کار گرفته‌اند و در این راستا ویژگی‌های مختلفی را از روش‌های متفاوت استخراج کرده‌اند. برخی از تحقیقات نیز از داده‌های کمکی دیگری در جهت بهبود عملکرد مدل پیشنهادی خود استفاده کرده‌اند. تحقیق حاضر بر مبنای داده‌های خام GPS و بدون کمک داده‌های مکانی خارجی، از چهار مدل طبقه‌بندی برای پیش‌بینی

<sup>۱</sup> Long Short Term Memory

<sup>۲</sup> Collaborative Isolation Forest

می‌کند به طوری که مشابه‌ترین داده‌ها کنار یکدیگر قرار گیرند. این روند تا زمانی که تمام داده‌ها از نظر یک ویژگی هدف در یک طبقه قرار گیرند، ادامه پیدا می‌کند. انتخاب بهترین ویژگی برای شروع طبقه‌بندی و همچنین یافتن یک حد آستانه مناسب برای اعمال طبقه‌بندی مهم‌ترین سؤال‌ها در این روش هستند. درخت تصمیم به طور کلی از مفهومی به نام آنتروپی<sup>۶</sup> برای بهینه‌سازی مسئله استفاده می‌کند [۳۰]. رابطه (۲) نحوه محاسبه آنتروپی را نمایش می‌دهد.

$$E(D) = -\sum_{i=1}^m P_i \log_2(P_i) \quad \text{رابطه (۲)}$$

که در رابطه (۲)،  $P_i$  احتمال تعلق داده‌ها به کلاس  $i$  ام و  $m$  تعداد کلاس‌ها هستند.

بهره اطلاعاتی<sup>۷</sup>، نسبت بهره<sup>۸</sup> و ضریب جینی<sup>۹</sup> مهم‌ترین معیارهایی هستند که در مدل‌های درخت تصمیم استفاده می‌شود. بهره اطلاعاتی حاصل از طبقه‌بندی داده‌ها بر مبنای ویژگی  $A$  از رابطه (۳) محاسبه می‌شود و در هر مرحله از تقسیم ویژگی با بیشترین بهره اطلاعاتی جهت طبقه‌بندی داده‌ها انتخاب می‌شود [۳۱].

$$\text{Gain}(A) = E(D) - \sum_{j=1}^n \frac{|D_j|}{|D|} \times E(D_j) \quad \text{رابطه (۳)}$$

در رابطه (۳)،  $E(D)$  آنتروپی و  $\frac{|D_j|}{|D|}$  وزن بخش  $j$  ام هستند.

معیار بهره اطلاعاتی به سمت داده‌ها با تنوع مقداری بالا دارای انحراف است [۳۱]. معیار نسبت بهره این مشکل را برطرف می‌کند و از رابطه (۴) محاسبه می‌شود.

مختلف برای ایجاد مدل پیش‌بینی کننده نوع حمل‌ونقل داده‌های خط سیر استفاده شده است.

### ۳-۱- ماشین بردار پشتیبان

روش ماشین بردار پشتیبان نخستین بار توسط وپنیک (۱۹۹۲) ارائه شد [۲۸]. هر ماشین بردار پشتیبان متشکل از یک یا چند ابرصفحه<sup>۱</sup> در فضای ابعاد بالا است که برای انجام وظیفه طبقه‌بندی یا رگرسیون کاربرد دارد. چنانچه کلاس‌ها به شکل خطی قابل جداسازی نباشند، اعمال یک تابع کرنل و انتقال داده‌های به فضای دیگر کارگشا خواهد بود. بدیهی است که هر چه فاصله نزدیک‌ترین داده از مجموعه داده‌های آموزشی از ابرصفحه بیشتر باشد، دقت روش در طبقه‌بندی افزایش می‌یابد. این فاصله حاشیه<sup>۲</sup> نامیده می‌شود. در اکثر مسائل طبقه‌بندی واقعی طبقات را نمی‌توان به شکل خطی از یکدیگر جدا کرد و پذیرش مقداری خطا غیرقابل اجتناب است [۲۹]. رابطه (۱) مسئله بهینه‌سازی نهایی را نشان می‌دهد.

رابطه (۱)

$$\min \frac{1}{2} W^2 + C \sum_{i=1}^n \xi_i, y_i (w \cdot x_i + b) - 1 + \xi_i \geq 0 \quad \forall_i$$

در رابطه (۱)،  $w$  و  $b$  پارامترهای ابرصفحه،  $\xi$  متغیر کمکی<sup>۳</sup> و  $C$  پارامتر تنظیم<sup>۴</sup> است. مقدار حاشیه<sup>۲</sup>  $\frac{2}{\|w\|}$  خواهد بود.

### ۳-۲- درخت تصمیم

درخت تصمیم یک الگوریتم یادگیری ماشین نظارت شده است که از درخت‌ها برای نمایش استفاده می‌کند. درخت تصمیم از گره ریشه شروع می‌کند و بر اساس یک ویژگی<sup>۵</sup> داده‌ها را به دو یا چند زیر درخت تقسیم

<sup>۱</sup> HyperPlane

<sup>۲</sup> Margin

<sup>۳</sup> Slack Variable

<sup>۴</sup> Regularization Parameter

<sup>۵</sup> Attribute

<sup>۶</sup> Entropy

<sup>۷</sup> Information Gain

<sup>۸</sup> Gain Ratio

<sup>۹</sup> Gini index

ساختار غیرخطی و تعمیم‌پذیری بالا از جمله ویژگی‌های مهم این شبکه‌ها هستند. شبکه عصبی پرسپترون چندلایه (*Multi Layer Perceptron (MLP)*) در زمره پرکاربردترین شبکه‌های عصبی در مسائل مختلف داده‌کاوی است. این شبکه حداقل دارای سه لایه است. آموزش در این شبکه بر اساس قانون پس انتشار (*Back Propagation*) خطا صورت می‌گیرد و مقدار خطا پس از محاسبه در مسیر برگشت از طریق لایه‌های شبکه عصبی در کل شبکه توزیع می‌شود [۳۲]. شکل (۱) یک شبکه *MLP* سه لایه را نمایش می‌دهد که در آن بایاس در لایه‌ها با نماد  $b$  و وزن میان نورون‌ها با نماد  $IW$  نمایش داده شده است. *Purelin* و *tansig* نیز توابع عملکرد در لایه خروجی و میانی هستند.

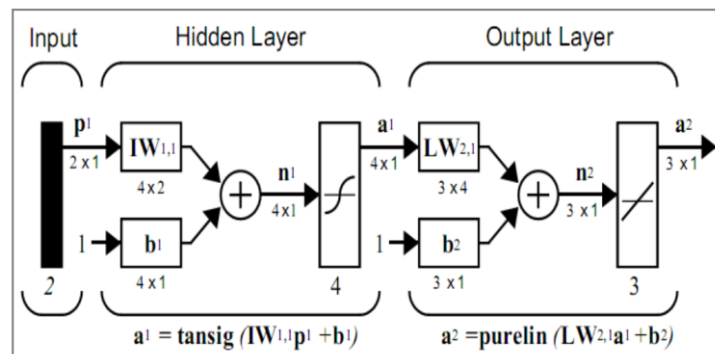
$$\text{رابطه (۴)} \quad \text{GainRatio}(A) = \frac{\text{Gain}(A)}{-\sum_{j=1}^n \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)}$$

معیار ضریب جینی نیز یکی پرکاربردترین معیارها در مدل‌های درخت تصمیم است. هنگام استفاده از این ضریب مراحل طبقه‌بندی داده‌ها به نحوی ادامه پیدا می‌کند که ضریب جینی حداقل شود. ضریب جینی طبق رابطه (۵)،  $P_i$  احتمال تعلق داده‌ها به کلاس  $i$  ام و  $m$  تعداد کلاس‌ها هستند [۳۱].

$$\text{رابطه (۵)} \quad \text{Gini}(D) = 1 - \sum_{i=1}^m P_i^2$$

### ۳-۳- شبکه عصبی پرسپترون چندلایه

شبکه‌های عصبی مصنوعی الهام گرفته از مغز انسان و متشکل از واحدهای پردازشگری به نام نورون هستند.



شکل ۱: شبکه عصبی پرسپترون چندلایه

$$\text{رابطه (۶)} \quad P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)}$$

با مستقل در نظر گرفتن متغیرهای تصادفی  $x_1, x_2, \dots, x_n$  کلاس هدف برای داده  $X$  مطابق رابطه (۷) قابل برآورد است که در آن  $K$  تعداد کلاس‌های هدف است.

$$\text{رابطه (۷)} \quad \hat{y} = \text{argmax} \left( P(C_k) \prod_{i=1}^n P(x_i|C_k) \right), \quad k = 1:K$$

### ۴- روش تحقیق

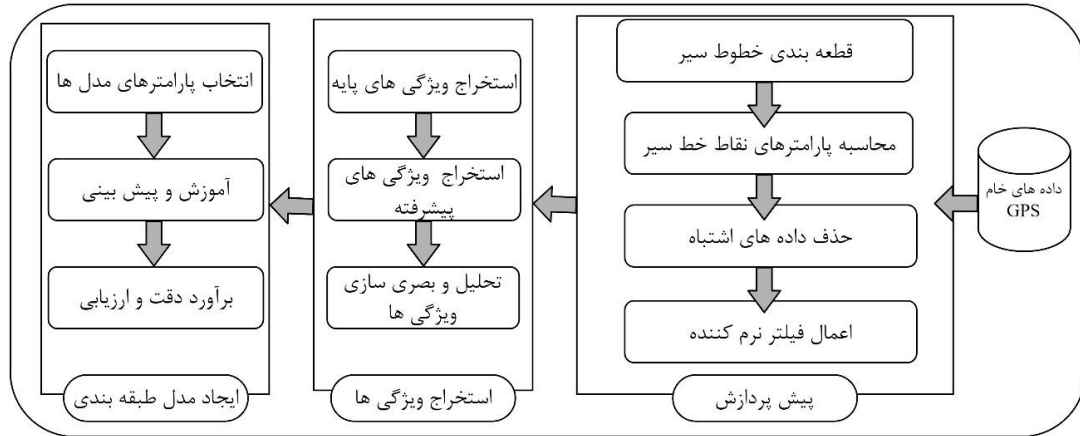
این تحقیق تلاش می‌کند با استخراج ویژگی از خطوط سیر داده‌های *GPS*، خطوط سیر را بر اساس نوع سفر

### ۳-۴- بیز ساده

روش بیز ساده<sup>۱</sup> (*NB*) با فرض استقلال متغیرهای تصادفی یک مدل احتمالاتی است که بر اساس قضیه بیز عمل می‌کند. این روش نظارت‌شده بر مبنای احتمال وقوع طبقه‌بندی‌ها را انجام می‌دهد [۳۱]. بر مبنای قضیه بیز احتمال وقوع کلاس  $C_i$  برای داده  $X$  طبق رابطه (۶) قابل محاسبه است.

<sup>۱</sup> Naive Bayes Classifier

طبقه‌بندی کرده و عملکرد مدل‌ها و ویژگی‌ها را در پیش‌بینی نوع سفر مورد بررسی قرار دهد. مراحل روش



شکل ۲: مراحل روش پیاده‌سازی تحقیق

$$B_{P_i} = |H_{i+1} - H_i| \quad \text{رابطه (۱۱)}$$

در روابط (۸)، (۹)، (۱۰) و (۱۱)،  $A_{P_i}$ ،  $S_{P_i}$ ،  $J_{P_i}$  و  $B_{P_i}$  و  $H_i$  به ترتیب سرعت، شتاب، تغییرات شتاب، تغییر جهت و زاویه با شمال برای نقطه  $P_i$  هستند. شکل (۳) پارامترهای رابطه (۱۱) را نمایش می‌دهد.

داده‌های خام  $GPS$ ، مجموعه‌ای متوالی از نقاط هستند که از مکان یک کاربر در زمان مشخصی ثبت شده‌اند. این داده‌ها به قسمت‌های کوچک‌تری به نام سفر<sup>۳</sup> تقسیم می‌شوند. هر سفر بخشی از داده‌ها است که فاصله زمانی بیشتر از یک حد آستانه مشخص در آن روی نداده باشد [۲۰]. سپس سفرها بر اساس نوع حمل‌ونقل قطع‌بندی (*Segmentation*) می‌شوند که هر قطعه دارای یک نوع حمل‌ونقل مشخص خواهد بود. از این پس در این تحقیق هر یک از این قطعه‌ها، خط سیر نامیده می‌شوند. داده‌های خام  $GPS$  ممکن است به دلیل وجود خطاهای مختلف نیازمند پیش‌پردازش‌هایی جهت حذف ثبت‌های اشتباه و کاهش خطاهای اتفاقی باشند.

#### ۴-۱- پیش‌پردازش و آماده‌سازی داده‌ها

داده‌های خام  $GPS$  شامل طول و عرض جغرافیایی و زمان برای هر نقطه هستند. در نخستین مرحله پنج پارامتر فاصله، سرعت، شتاب، تغییرات شتاب<sup>۱</sup> و تغییر جهت (زاویه) برای هر نقطه در داده‌های خام  $GPS$  محاسبه می‌گردد تا در ادامه تحقیق مورد استفاده قرار گیرند. به منظور محاسبه فاصله بین نقاط فرمول وینچنتی<sup>۲</sup> [۳۳] مورد بهره‌برداری قرار گرفت، سپس پارامترهای سرعت، شتاب، تغییرات شتاب و تغییر جهت از روابط (۸)، (۹)، (۱۰) و (۱۱) به دست می‌آیند.

$$S_{P_i} = \frac{\text{vincenty}(P_i, P_{i+1})}{\Delta t = t_{i+1} - t_i} \quad \text{رابطه (۸)}$$

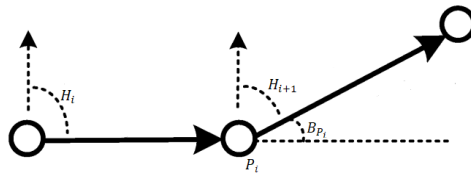
$$A_{P_i} = \frac{S_{i+1} - S_i}{\Delta t} \quad \text{رابطه (۹)}$$

$$J_{P_i} = \frac{A_{i+1} - A_i}{\Delta t} \quad \text{رابطه (۱۰)}$$

<sup>۱</sup> Jerk

<sup>۲</sup> Vincenty's formulae

<sup>۳</sup> Trip



شکل ۳: نمایش تغییر جهت (زاویه) برای یک نقطه از خط سیر

#### ۴-۲- استخراج ویژگی‌ها

پس از اینکه خطوط سیر پالایش شده از داده‌های خام به دست آمد، به منظور اعمال مدل‌های طبقه‌بندی ویژگی‌هایی از خطوط سیر که در میان نوع‌های مختلف حمل و نقل متفاوت هستند، استخراج می‌شود. به عنوان مثال ویژگی سرعت در نوع پیاده‌روی و قطار تفاوت چشم‌گیری خواهد داشت. در این بخش ویژگی‌های مختلفی که در این تحقیق مورد استفاده قرار گرفته‌اند، معرفی می‌شوند.

#### ۴-۲-۱- ویژگی‌های پایه

طول خط سیر و مدت زمان آن مهم‌ترین ویژگی هر خط سیر هستند. از آنجایی کاربران عموماً با توجه به همین ویژگی‌ها نوع حمل‌ونقل خود را انتخاب می‌کنند، بنابراین این ویژگی‌ها می‌تواند در بهبود عملکرد مدل‌های پیش‌بینی کننده مؤثر باشند.

سرعت یکی از ویژگی‌هایی است که در نوع‌های مختلف حمل‌ونقل مقدار متفاوتی دارد و می‌تواند در تمیز این نوع‌ها توسط مدل‌های پیش‌بینی کننده نقش ایفا کند. ویژگی‌هایی مانند سرعت میانگین، امید ریاضی سرعت، حداکثر سرعت و انحراف معیار سرعت از جمله ویژگی‌های مشتق شده از سرعت نقاط خط سیر هستند که در تحقیقات پیشین مورد توجه قرار گرفته‌اند و در این تحقیق نیز استفاده می‌شوند. از دیگر ویژگی‌های مهمی که می‌تواند در تعیین نوع حمل‌ونقل مفید باشد، شتاب است؛ زیرا شتاب وسایل مختلف حمل‌ونقل مانند اتوبوس، خودرو شخصی و قطار متفاوت است. میانگین شتاب، انحراف معیار شتاب و حداکثر شتاب ویژگی‌های مورد استفاده در این تحقیق هستند. کاربران زمانی که در حال پیاده‌روی هستند، معمولاً تغییر جهت‌های زیادی دارند، در حالی که در نوع حمل‌ونقلی قطار ممکن

در این تحقیق مراحل زیر برای پالایش داده‌ها انجام می‌گیرد:

- حذف قطعه‌های مورد مطالعه با طول کمتر از ۱۰ نقطه
- حذف نقاطی که لحظه زمانی آن‌ها بزرگ‌تر یا مساوی نقطه بعدی است.
- حذف نقاطی که سرعت یا شتاب آن‌ها بر اساس قوانین و قابلیت‌ها غیر معتبر باشد. مثلاً سرعت بیش از ۲/۵ متر بر ثانیه برای نوع پیاده‌روی انسان غیر معتبر است.
- اعمال یک کرنل نرم‌کننده برای حذف خطاهای اتفاقی

در این تحقیق از فیلتر ساویتزکی-گلی<sup>۱</sup> به منظور نرم کردن خطوط سیر استفاده شده است. به طور کلی نرم کردن فرآیندی است که طی آن هر نقطه از داده‌ها توسط همسایگان خود تحت یک تابع کرنل مشخص قرار می‌گیرد و به این ترتیب خطای اتفاقی کاهش پیدا می‌کند. از آنجایی هیچ اطلاعات قبلی در مورد شکل خطوط سیر وجود ندارد، باید از فیلتری استفاده شود که به شکل پیش‌فرض داده‌ها را به شکل خاصی متمایل نکند و روند خط سیر را دچار تغییر ننماید. فیلتر ساویتزکی-گلی ابتدا یک چندجمله‌ای درجه  $m$  به تعداد فرد  $n$  نقطه از داده‌ها که پنجره نامیده می‌شوند و نقطه مورد نظر در مرکز این پنجره قرار دارد، برازش می‌دهد. سپس مقدار جدید نقطه مرکزی توسط تابع چندجمله‌ای برآورد می‌شود [۳۴]؛ بنابراین با این روش شکل کلی و الگوی خط سیر حفظ می‌شود.

<sup>۱</sup> Savitzky-Golay Filter

است به ندرت جهت تغییر کند. همچنین میزان تغییر جهت برای سایر حمل و نقل‌ها نیز متغیر است؛ بنابراین میانگین زاویه تغییر جهت نیز در این تحقیق مورد استفاده قرار گرفته است.

#### ۴-۲-۲- ویژگی‌های پیشرفته

در نوع حمل و نقل پیاده روی یا دوچرخه ممکن است کاربران بارها جهت حرکت خود را به طور محسوس تغییر دهند، در حالی که انتظار می‌رود در سایر نوع‌ها تغییرات محسوس کمتر باشد، مثلاً در نوع قطار به ندرت جهت حرکت به شکل محسوس تغییر می‌کند و این تعداد احتمالاً برای اتوبوس و خودرو بیشتر خواهد بود. در نظر گرفتن درصد تغییر جهت‌های محسوس به عنوان یک ویژگی در بهبود عملکرد مدل‌های پیش‌بینی کننده می‌تواند مفید باشد. در این تحقیق با در نظر گرفتن یک حد آستانه  $\delta$  برای تشخیص تغییر جهت‌های محسوس این ویژگی برای خط سیری با  $n$  نقطه میانی طبق رابطه (۱۲) محاسبه می‌شود که در آن  $B_{P_i}$  تغییر جهت برای نقطه  $P_i$  است [۲۰].

$$\text{رابطه (۱۲)} \quad \frac{|B_{P_i} > \delta|}{n}, \quad i = 1:n$$

در نوع‌های حمل و نقلی مختلف ممکن است به دلایل متفاوتی بارها سرعت حرکت به صفر نزدیک شود و کاربر متوقف باشد. مثلاً بازدید از یک فروشگاه در نوع پیاده روی ممکن است کاربر را متوقف کند. در مقایسه نوع‌های اتوبوس و رانندگی نیز این ویژگی می‌تواند تأثیرگذار باشد، زیرا اتوبوس‌ها برای سوار کردن مسافران تعداد دفعات بیشتری متوقف می‌شوند. در این تحقیق با در نظر گرفتن یک حد آستانه برای سرعت  $\delta$ ، درصد نقاط متوقف محاسبه شده و به عنوان ویژگی خطوط سیر در روش‌های طبقه‌بندی مورد استفاده قرار می‌گیرد [۲۰]. سه ویژگی مشابه دیگر نیز بر اساس درصد نقاط خطوط سیر در بازه‌های سرعتی معین محاسبه و استفاده خواهد شد (رابطه (۱۳)).

$$\text{رابطه (۱۳)} \quad \frac{|V_i > \delta|}{n}, \quad i = 1:n$$

معمولاً با افزایش تعداد داده‌ها توزیع آن‌ها به منحنی توزیع نرمال میل پیدا می‌کند. در توزیع نرمال، ۶۸٪ داده‌ها در فاصله کمتر از یک انحراف معیار نسبت به میانگین قرار دارند. در این تحقیق این فاصله برای سرعت تمام داده‌ها در نوع‌های حمل و نقل متفاوت محاسبه می‌شود، سپس درصد قرار گرفتن نقاط هر خط سیر در هریک از این فواصل به عنوان ویژگی به مدل معرفی می‌گردد.

#### ۴-۳- مدل‌های طبقه‌بندی

پس از محاسبه ویژگی‌ها برای خطوط سیر، ۸۰ درصد داده‌ها به عنوان داده‌های آموزشی انتخاب شده و برای آموزش مدل‌ها مورد استفاده قرار می‌گیرند. ۲۰ درصد داده‌ها نیز به عنوان داده آزمایشی برای اندازه‌گیری دقت مدل‌ها در نظر گرفته خواهند شد. در این تحقیق از چهار روش طبقه‌بندی مختلف برای ایجاد مدل پیش‌بینی کننده نوع حمل و نقل داده‌های خط سیر استفاده شده است.

در این تحقیق از روش *SVM* با استفاده از رویکرد یک در مقابل همه<sup>۱</sup> و انتخاب کرنل‌های مختلف برای طبقه‌بندی چند کلاسه استفاده می‌شود. در رویکرد یک در مقابل همه، به ازای هر کلاس یک *SVM* استفاده می‌شود که وظیفه آن جدا کردن داده‌های یک کلاس از بقیه کلاس‌ها است. انتخاب مقدار مناسب برای پارامتر تنظیم نیز در این روش حائز اهمیت است [۳۱]. درخت تصمیم نیز در این تحقیق با اعمال ضریب جینی، بهره اطلاعاتی و نسبت بهره به عنوان معیارهای طبقه‌بندی، مورد استفاده قرار می‌گیرد. انتخاب عمق درخت و حداقل اعضای برگ‌ها پارامترهایی هستند که انتخاب صحیح آن‌ها می‌تواند در دقت مدل تأثیرگذار باشد [۳۱]. شبکه عصبی پرسپترون چندلایه روش دیگری است که در این تحقیق پیاده‌سازی می‌شود. انتخاب تعداد لایه‌ها و تعداد گره‌ها در هر لایه مهم‌ترین

<sup>۱</sup> One vs all

سیر در این مجموعه داده با استفاده از مجموعه‌ای از نقاط با فاصله زمانی مشخص، نمایش داده می‌شود و شامل طول و عرض جغرافیایی، ارتفاع و زمان نمونه‌برداری هستند [۳۵ و ۳۶].

این داده‌ها شامل ۱۷۶۲۱ خط سیر هستند که در مجموع در طی ۵۰۱۷۶ ساعت و به طول ۱۲۹۲۹۵۱ کیلومتر جمع‌آوری شده است. داده‌ها شامل انواع گسترده‌ای از فعالیت‌های حرکتی کاربران هستند و تنها به سفرهای معمول بین خانه و محل کار محدود نمی‌شوند و فعالیت‌های سرگرمی و ورزش مانند خرید، غذا، هاکی، گردش و غیره را شامل می‌شوند. الگو کاوی حرکتی، شبکه‌های اجتماعی مکان‌مبنا، تشخیص فعالیت کاربر، حریم خصوصی و توصیه مکان از جمله کاربردهایی است که می‌توان از این داده‌ها در آن بهره گرفت [۲۱ و ۳۷].

اگرچه این داده‌ها از چندین شهر کشور چین و حتی شهرهای امریکا و اروپا جمع‌آوری شده است ولی بخش عمده این داده‌ها مربوط به شهر پکن هستند، شکل (۴) نقشه حرارتی از پراکندگی این داده‌ها را در شهر پکن نشان می‌دهد.

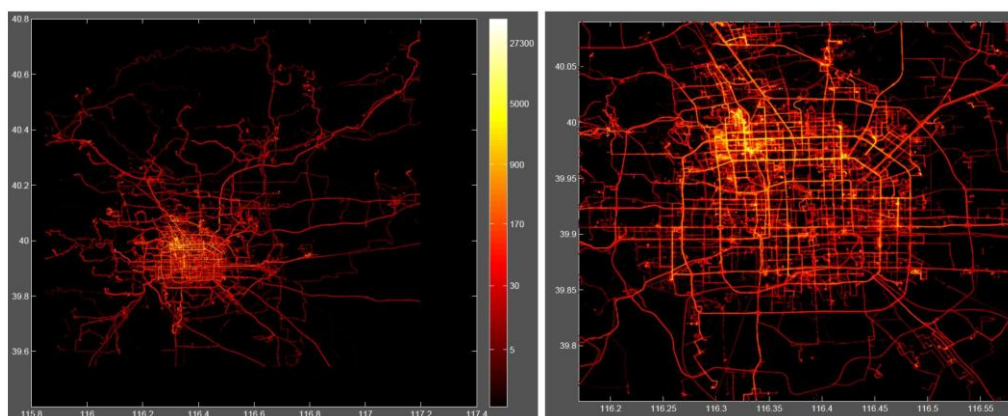
پارامترهایی هستند که دقت این روش را تحت تأثیر قرار می‌دهند [۳۱]. در روش طبقه‌بندی کننده بیز ساده دقت با برآورد چگالی کرنل افزایش پیدا می‌کند. در این تحقیق از رویکرد برآورد حریصانه استفاده شده است که پارامترهای تعداد کرنل‌ها و حداقل عرض باند باید تنظیم شوند [۳۱].

## ۵- پیاده‌سازی

در این بخش به ارائه نتایج حاصل از پیاده‌سازی روش پیاده‌سازی تحقیق پرداخته خواهد شد. ابتدا داده‌های مورد مطالعه معرفی شده و سپس توزیع آماری هر یک از ویژگی‌ها و دقت مدل‌های مختلف مورد بحث قرار می‌گیرند.

### ۵-۱- داده‌های مورد مطالعه

در این تحقیق، از داده‌های پروژه ژئولایف استفاده شده است. این داده‌ها مربوط به شهر پکن است و توسط ۱۸۲ کاربر در مدت ۵ سال (آوریل ۲۰۰۷ تا اوت ۲۰۱۲) جمع‌آوری شده است. نرخ نمونه‌برداری این داده‌ها متفاوت است اما ۹۱٫۵ درصد از داده‌ها از چگالی مناسبی برخوردار هستند. بیشتر این داده با فاصله ۱ تا ۵ ثانیه و یا ۵ تا ۱۰ متر نمونه‌برداری شده‌اند. هر خط



ب) شهر پکن و حومه

الف) شهر پکن

شکل ۴: نقشه حرارتی داده‌های مورد مطالعه [۳۵].

موجود پیاده، خودرو، تاکسی، اتوبوس، مترو، قطار، دوچرخه و غیره (دویدن، هواپیما، موتور و قایق) هستند و این بخش از داده‌ها مورد توجه این تحقیق هستند.

### ۵-۲- پیش‌پردازش و آماده‌سازی داده‌ها

از بین ۱۸۲ کاربر در پروژه ژئولایف، ۶۹ کاربر نوع حمل‌ونقل خود را نیز مشخص کرده‌اند. از جمله نوع‌های

جدول (۱) خلاصه‌ای از آمار کلی نوع‌های حمل‌ونقل را ارائه می‌کند.

جدول ۱: آمار نوع حمل‌ونقل در داده‌های ژئولایف

نوع سفر	مسافت (km)	زمان (h)
پیاده	۱۰۱۲۳	۵۴۶۰
دوچرخه	۶۴۹۵	۲۴۱۰
اتوبوس	۲۰۲۸۱	۱۵۰۷
رانندگی (ماشین و تاکسی)	۳۲۸۶۶	۲۳۸۴
قطار	۳۶۲۵۳	۷۴۵
هواپیما	۲۴۷۸۹	۴۰
سایر (قایق، موتور و دویدن)	۹۴۹۳	۴۰۴
کل	۱۴۰۳۰۴	۱۲۹۵۳

سیر تعیین کرده‌اند. همچنین هر قسمت پیوسته از سفرها با نوع حمل‌ونقل خاص به‌عنوان واحد کاری انتخاب شدند و از این‌پس خط سیر نامیده می‌شوند. سپس بر مبنای قوانین ترافیکی و محدودیت‌های فیزیکی، حداکثر سرعت و حداکثر شتاب برای هر نوع حمل‌ونقل تعیین شد [۱۰]. با اعمال کردن این شرایط بر روی نقاط خطوط سیر، نقطه‌های اشتباه از خط سیر حذف شده و پارامترهای نقاط همسایه به‌روزرسانی می‌شوند. شرایط ذکر شده در جدول (۲) شرح داده شده‌اند.

در این تحقیق نوع‌های سفر کم‌اهمیت نادیده گرفته شده و پنج نوع سفر رانندگی (خودرو شخصی و تاکسی)، پیاده‌روی، دوچرخه، اتوبوس و قطار (قطار و مترو) مورد بررسی قرار گرفتند.

در اولین قدم هرگاه فاصله بیشتر از ۲۰ دقیقه در طی یک خط سیر وجود داشت، خط سیر به واحدهای کوچک‌تری به نام سفر تقسیم شد [۲۰]. به‌منظور آماده‌سازی داده‌ها برای انجام تحلیل‌های مکانی بعدی داده‌های فاقد نوع حمل‌ونقل حذف شدند زیرا برخی از کاربران نوع حمل‌ونقل را فقط برای قسمتی از خطوط

جدول ۲: حداکثر سرعت و شتاب مجاز برای نوع‌های حمل‌ونقل

نوع حمل‌ونقل	حداکثر سرعت (m/s)	حداکثر شتاب ( $m/s^2$ )
پیاده	۲/۵	۳
دوچرخه	۱۱	۳
خودرو و تاکسی	۳۳	۱۰
اتوبوس	۲۷	۲
قطار	۹۷	۳

پارامترهای اصلی نقاط به‌روزرسانی گردید. جدول (۳) تعداد خطوط سیر و میانگین سرعت مربوط به آن‌ها را پیش و پس از مرحله پالایش شرح می‌دهد.

در آخرین مرحله پیش‌پردازش فیلتر ساویتزکی-گلی با تنظیم کردن پنجره با ابعاد ۷ نقطه و با تابع چندجمله‌ای درجه ۳ بر روی داده‌ها اعمال شد و

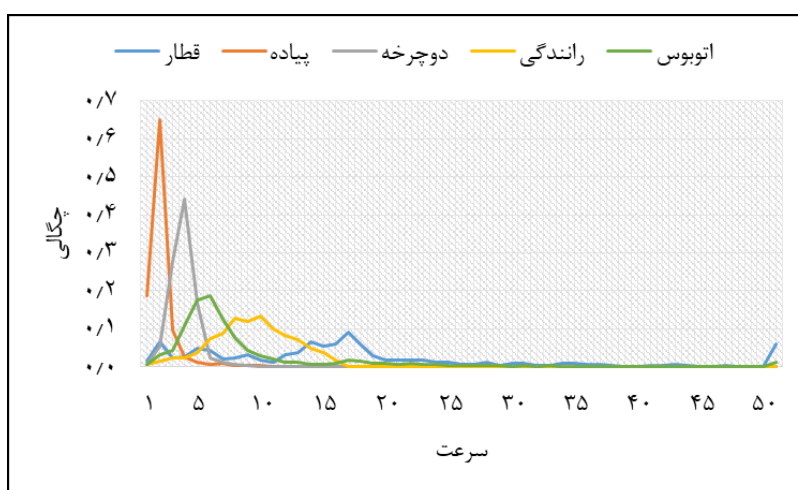
جدول ۳: آمار خطوط سیر پس از پالایش

نوع حمل‌ونقل	پیش از پالایش		پس از پالایش	
	تعداد خطوط سیر	میانگین سرعت (m/s)	تعداد خطوط سیر	میانگین سرعت (m/s)
پیاده	۴۸۲۲	۲	۴۱۷۶	۱٫۷
دوچرخه	۱۶۰۱	۳٫۶	۱۵۷۰	۳٫۶
رانندگی	۱۴۰۳	۱۱٫۲	۱۳۹۳	۱۱٫۳
اتوبوس	۱۹۹۶	۶٫۷	۱۸۹۲	۶٫۶
قطار	۷۹۵	۲۱٫۸	۷۹۲	۲۱٫۹

### ۵-۳- استخراج و تحلیل ویژگی‌ها

در این بخش ویژگی‌های خطوط سیر محاسبه و از طریق برخی آمار و نمودارها مورد بررسی قرار گرفته و مقایسه می‌شوند. سرعت مهم‌ترین ویژگی خطوط سیر در تعیین نوع حمل‌ونقل است، شکل (۵) نمودار توزیع

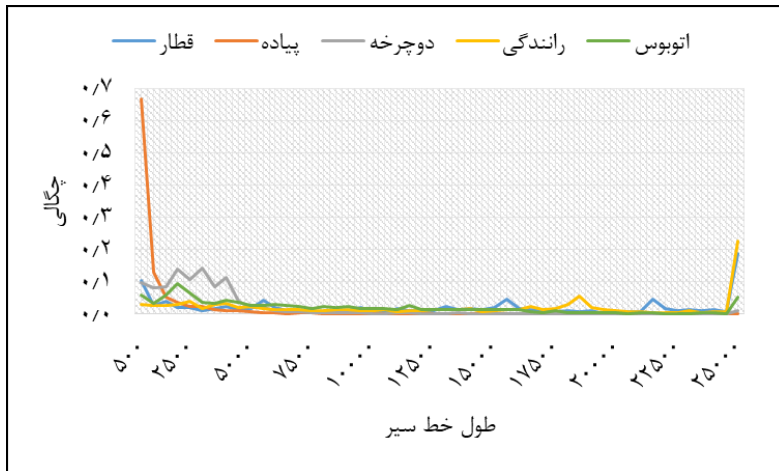
میانگین سرعت برای نوع‌های مختلف حمل‌ونقل را نشان می‌دهد. مطابق شکل نوع‌های حمل‌ونقل مختلف از نظر ویژگی سرعت رفتار متفاوت و تعیین‌کننده‌ای دارند. فقط نوع حمل‌ونقل قطار در سرعت‌های بالای ۵۰ متر بر ثانیه فراوانی قابل‌توجهی دارد.



شکل ۵: نمودار توزیع میانگین سرعت بر حسب m/s

شکل (۶) طول خطوط سیر را برای نوع‌های حمل‌ونقل مختلف نمایش می‌دهد. همان‌طور که انتظار می‌رود طول بخش عمده‌ای از خطوط سیر پیاده‌روی کمتر از

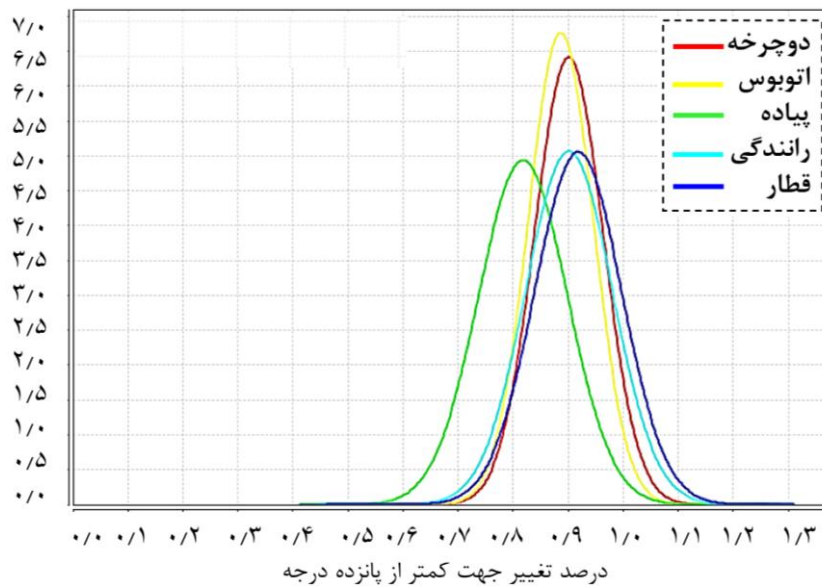
دو کیلومتر است. همچنین طول کم خطوط سیر دوچرخه نسبت به سایر نوع‌ها مشهود است.



شکل ۶: نمودار توزیع طول خطوط سیر برحسب متر

محاسبه درصد نقاط توقف انتخاب شد [۲۱]. این نمودار نشان می‌دهد که نوع پیاده‌روی دارای بیشترین تغییر جهت بیشتر از ۱۵ درجه است.

به‌منظور کشف تغییر جهت‌های محسوس از حد آستانه ۱۵ درجه استفاده شد [۲۰]. شکل (۷) توزیع آماری این ویژگی را برای نوع‌های حمل‌ونقل مختلف نشان می‌دهد. همچنین حد آستانه ۳ متر بر ثانیه برای



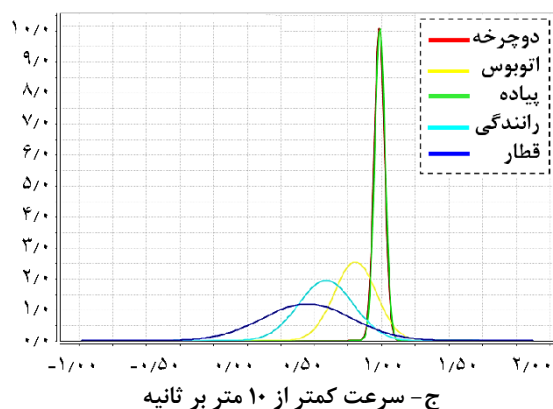
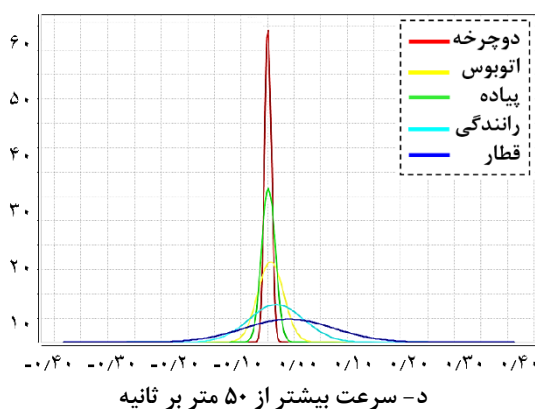
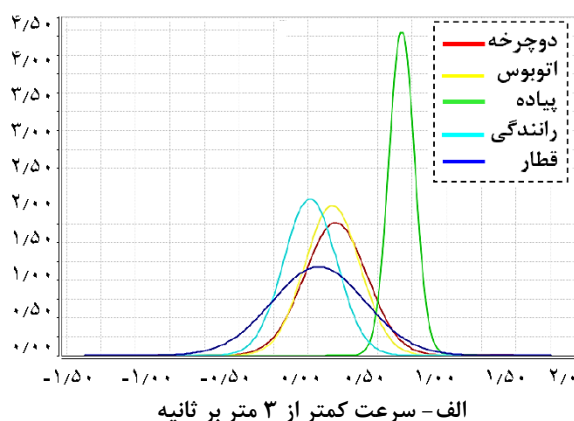
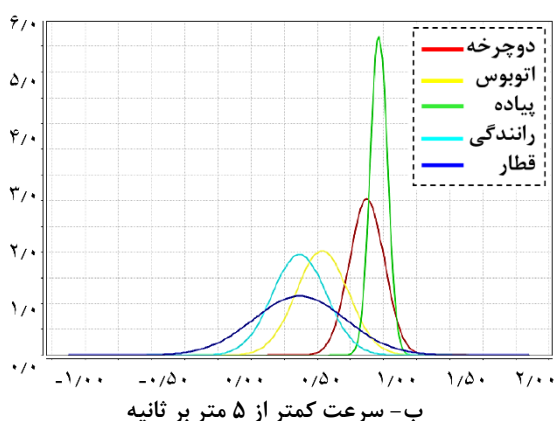
شکل ۷: نمودار توزیع چگالی تغییر جهت

قرار می‌گیرند. رفتار متفاوت هر نوع حمل‌ونقل در بازه‌های سرعتی مختلف در شکل (۸) مشهود است. مطابق شکل (۸-ج) نوع‌های پیاده‌روی و دوچرخه دارای

شکل (۸) توزیع داده‌های خط سیر در نوع‌های مختلف حمل‌ونقل را برای چهار بازه سرعت نمایش می‌دهد که به‌عنوان ویژگی‌های جدید در این تحقیق مورد استفاده

درصد بالایی از نقاط دارای سرعت کمتر از ۵ ثانیه هستند. در سایر شکل‌ها نیز تفاوت رفتار نمودار برای نوع‌های مختلف حمل‌ونقل مشهود است.

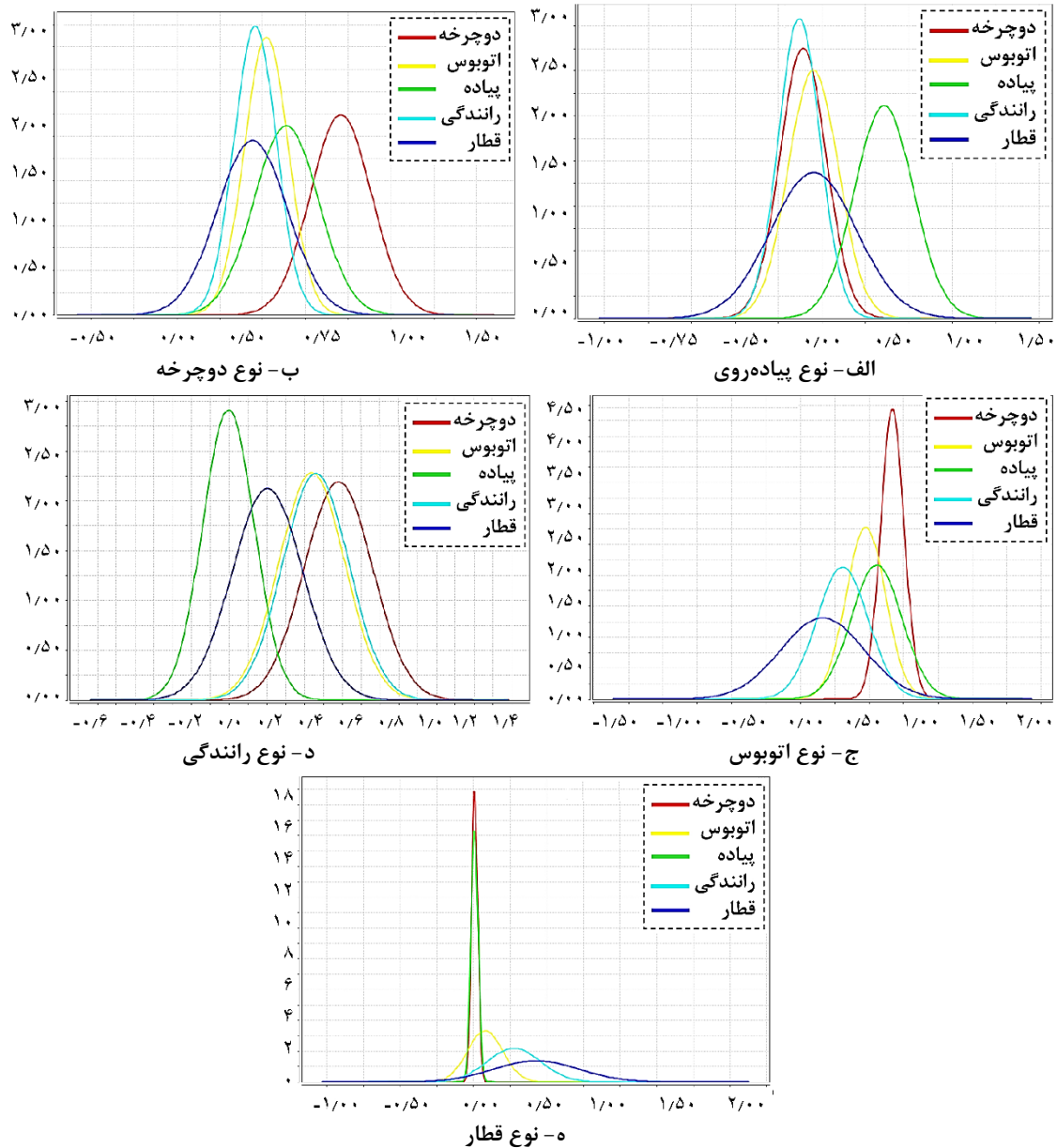
توزیع یکسانی برای سرعت‌های کمتر از ۱۰ متر بر ثانیه هستند، اما شکل (۸-ب) نشان می‌دهد که دو نوع پیاده‌روی و دوچرخه در سرعت‌های کمتر از ۵ متر بر ثانیه دارای توزیع متفاوتی هستند و در نوع پیاده‌روی



شکل ۸: نمودار توزیع چگالی درصد نقاط خطوط سیر در بازه‌های سرعتی مختلف

می‌توانند در بهبود عملکرد مدل پیاده‌سازی تحقیق تأثیرگذار باشند. به‌عنوان نمونه شکل (۹-ه) نشان می‌دهد که سرعت نقاط نوع‌های دوچرخه و پیاده‌روی به‌ندرت در فاصله یک انحراف معیار از میانگین سرعت نوع قطار قرار می‌گیرند. درحالی‌که مطابق انتظار نوع قطار در این نمودار دارای توزیع نرمال است و نوع‌های رانندگی و اتوبوس نیز به ترتیب نزدیک‌ترین توزیع را به توزیع نرمال دارند.

به‌عنوان یک ویژگی جدید، پس از محاسبه درصد حضور نقاط هر خط سیر در فاصله یک انحراف معیار از میانگین کل سرعت نوع‌های حمل‌ونقل، توزیع این ویژگی برای همه حالات در شکل (۹) نمایش داده شده است. نمودار زردرنگ شکل (۹-الف) نحوه توزیع سرعت خطوط سیر اتوبوس نسبت به میانگین سرعت پیاده‌روی را نشان می‌دهد. به‌طورکلی رفتار متفاوت داده‌های هر نوع حمل‌ونقل در ویژگی‌های محاسبه‌شده در این تحقیق، در شکل‌های (۹) مشهود است؛ بنابراین



شکل ۹: نمودار توزیع چگالی درصد نقاط خطوط سیر در بازه یک انحراف معیار از میانگین کل سرعت در نوع‌های مختلف

مناسبی بودند. شکل (۱۰) توزیع مکانی داده‌های آزمایشی را نشان می‌دهد. تمام روش‌ها به ازای پارامترهای ورودی مختلف پیاده‌سازی شد و بهینه‌ترین پارامترها به کمک جستجوی شبکه‌ای (*Grid Search*) تنظیم شد و نتایج زیر حاصل گردید. در جستجوی شبکه‌ای بر اساس بهترین مقادیر ارائه شده در جستجوی تصادفی، شبکه‌ای از مقادیر مختلف برای پارامترهای مدل‌ها ایجاد می‌گردد و جستجو آغاز می‌شود [۳۸].

#### ۵-۴- پیاده‌سازی مدل‌های طبقه‌بندی

در این مرحله چهار روش طبقه‌بندی *SVM*، *MLP*، درخت تصمیم و بیز ساده بر روی داده‌ها پیاده‌سازی شد. داده‌ها به دو بخش داده آموزشی و آزمایشی تقسیم شدند. به‌منظور مقایسه نتایج ۲۰ درصد ثابت از داده‌ها به‌عنوان داده آزمایشی در نظر گرفته شدند. داده‌های آزمایشی با توجه به حجم بالای مجموعه داده ژئولایف، دارای توزیع مکانی



شکل ۱۰: نقشه توزیع مکانی داده‌های آزمایشی

خطی به ازای مقدار پارامتر تنظیم برابر ۰/۰۲ شناسایی شد. کرنل خطی یا ضرب داخلی به شکل  $K(x, y)$  تعریف می‌شود. جدول (۴) دقت کلی برآورد نوع حمل‌ونقل را به ازای پارامترهای متفاوت روش SVM نشان می‌دهد.

در پیاده‌سازی روش SVM چهار کرنل خطی، چندجمله‌ای، گوسین و تابع پایه شعاعی (RBF) به ازای پارامتر تنظیم در  $CE \{0.001, 0.01, 0.02, 0.03, 0.05, 0.1\}$  مورد جستجو قرار گرفتند. بهترین عملکرد توسط کرنل

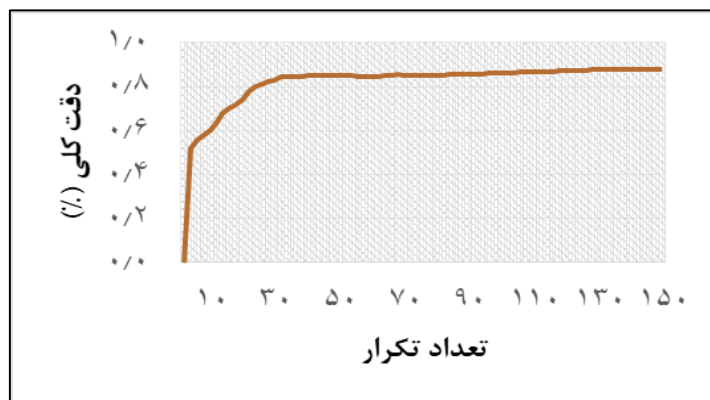
جدول ۴: تأثیر پارامترهای روش SVM بر دقت

پارامتر تنظیم						تابع کرنل
۰/۱	۰/۰۵	۰/۰۳	۰/۰۲	۰/۰۱	۰/۰۰۱	
۰/۸۴۵	۰/۸۵۱	۰/۸۵۷	۰/۸۵۹	۰/۸۵۸	۰/۸۳۱	کرنل خطی
۰/۸۴۱	۰/۸۴۷	۰/۸۵۵	۰/۸۵۵	۰/۸۵۵	۰/۸۳۹	کرنل چندجمله‌ای
۰/۸۰۶	۰/۸۱۹	۰/۸۲۵	۰/۸۲۵	۰/۸۲۲	۰/۸۰۱	تابع گوسین
۰/۸۴۰	۰/۸۴۲	۰/۸۵۰	۰/۸۵۲	۰/۸۵۰	۰/۸۳۷	تابع پایه شعاعی

عملکرد را از خود نشان داد. همچنین حداکثر عمق درخت برابر ۱۰ و حداقل تعداد اعضای برگ دو انتخاب شد. جدول (۶) دقت کلی برآورد نوع حمل‌ونقل را به ازای پارامترهای متفاوت درخت تصمیم نشان می‌دهد. در پیاده‌سازی روش بیز ساده با رویکرد برآورد حریمانه از چهار کرنل با حداقل باند ۰/۰۸ استفاده شد. در این روش تعداد کرنل‌ها بین ۱ تا ۱۰ و حداقل باند در  $\{0.5\}$  قرار گرفتند. جدول (۷) دقت کلی برآورد نوع حمل‌ونقل را به ازای پارامترهای متفاوت روش بیز ساده نشان می‌دهد.

برای روش MLP یک لایه میانی با ۲۵ گره انتخاب شد و نرخ یادگیری برابر ۰/۳ تنظیم شد. این روش به ازای یک تا پنج لایه میانی و تعداد گره‌ها ۴۰، ۱۱، ...،  $n=10$  پیاده‌سازی شد. جدول (۵) دقت کلی برآورد نوع حمل‌ونقل را به ازای معماری متفاوت در روش MLP نشان می‌دهد. شکل (۱۱) نمودار همگرایی روش MLP را نشان می‌دهد. درخت تصمیم با معیارهای بهره اطلاعاتی<sup>۱</sup>، نسبت بهره<sup>۲</sup> و شاخص جینی پیاده‌سازی شد. شاخص جینی بهترین

<sup>۱</sup> Information Gain<sup>۲</sup> Gain Ratio



شکل ۱۱: نمودار همگرایی روش MLP

جدول ۵: تأثیر پارامترهای روش MLP بر دقت

تعداد گره‌ها						تعداد لایه‌ها
۳۵	۳۰	۲۵	۲۰	۱۵	۱۰	
۰.۸۷۸	۰.۸۷۸	۰.۸۷۸	۰.۸۷۰	۰.۸۵۲	۰.۸۴۱	۱
۰.۸۷۷	۰.۸۷۷	۰.۸۷۷	۰.۸۷۰	۰.۸۵۵	۰.۸۴۴	۲
۰.۸۷۳	۰.۸۷۳	۰.۸۷۳	۰.۸۷۰	۰.۸۵۸	۰.۸۴۵	۳

جدول ۶: تأثیر پارامترهای درخت تصمیم بر دقت

حداکثر عمق درخت						معیار
۱۴	۱۲	۱۰	۸	۶	۴	
۰.۸۳۷	۰.۸۴۰	۰.۸۴۲	۰.۸۴۶	۰.۸۲۹	۰.۷۹۵	بهره اطلاعاتی
۰.۷۵۶	۰.۷۵۱	۰.۷۴۸	۰.۷۴۳	۰.۷۴۳	۰.۶۹۸	نسبت بهره
۰.۸۵۰	۰.۸۵۱	۰.۸۵۳	۰.۸۴۹	۰.۸۳۸	۰.۷۸۷	شاخص جینی

جدول ۷: تأثیر پارامترهای بیز ساده بر دقت

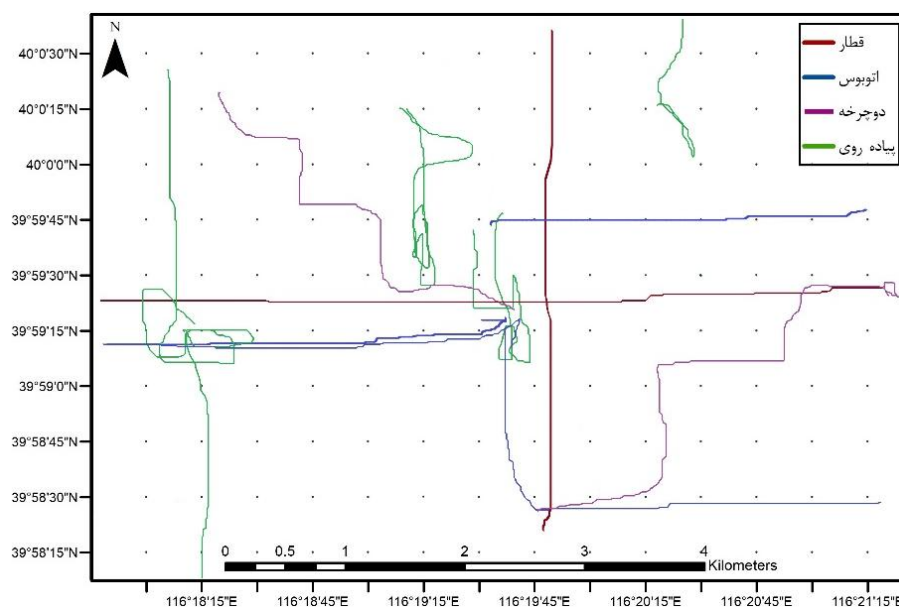
حداقل باند						تعداد کرنل‌ها
۰.۵	۰.۲	۰.۱	۰.۰۸	۰.۰۵	۰.۰۱	
۰.۷۳۰	۰.۷۳۳	۰.۷۳۴	۰.۷۳۵	۰.۷۴۰	۰.۷۳۳	۱
۰.۷۵۸	۰.۷۶۱	۰.۷۶۵	۰.۷۶۵	۰.۷۶۴	۰.۷۶۱	۳
۰.۷۶۸	۰.۷۷۱	۰.۷۸۱	۰.۷۸۲	۰.۷۷۷	۰.۷۶۹	۴
۰.۷۶۸	۰.۷۷۰	۰.۷۸۰	۰.۷۸۲	۰.۷۷۴	۰.۷۶۱	۵
۰.۷۶۰	۰.۷۶۳	۰.۷۶۹	۰.۷۷۳	۰.۷۶۴	۰.۷۵۱	۸

این تحقیق با دقت مشابه‌ترین تحقیقات گذشته مقایسه می‌شود.

شکل (۱۲) نتایج حاصل از روش *MLP* بر روی بخش بسیار کوچکی از داده‌های مورد مطالعه را نشان می‌دهد. در این شکل تمام نوع‌های حمل‌ونقل به درستی شناسایی شده‌اند.

## ۶- نتایج و بحث

در این بخش نتایج حاصل از پیاده‌سازی تحقیق بیان می‌شود و به مقایسه دقت عملکرد روش‌های مختلف طبقه‌بندی در تعیین نوع حمل‌ونقل پرداخته خواهد شد. همچنین دقت طبقه‌بندی نوع‌های حمل‌ونقل در



شکل ۱۲: بخشی از نتایج تحقیق به روش *MLP*

همچنین روش *SVM* عملکرد بهتری نسبت به روش‌های بیز ساده و درخت تصمیم داشت. جدول (۱۲) نیز مقدار معیارهای دقت را برای روش‌های مختلف شرح می‌دهد.

## ۶-۱- مقایسه نتایج روش‌های طبقه‌بندی

به منظور مقایسه نتایج حاصل از پیاده‌سازی روش‌های مختلف طبقه‌بندی دو معیار دقت کلی<sup>۱</sup> و شاخص کاپا<sup>۲</sup> محاسبه گردید. جدول‌های (۸) الی (۱۱) ماتریس ابهام<sup>۳</sup> را برای چهار روش طبقه‌بندی نشان می‌دهند. در میان روش‌های پیاده‌سازی این تحقیق روش *MLP* در هر دو معیار دقت کلی و شاخص کاپا توانست دقت بهتری نسبت به سایر روش‌ها از خود نشان دهد.

<sup>۱</sup> Overall accuracy

<sup>۲</sup> Kappa index

<sup>۳</sup> Confusion matrix

جدول ۹: ماتریس ابهام روش MLP

واقعی					تجزیه و تحلیل
قطار	رانندگی	پیاده	اتوبوس	دوچرخه	
قطار	۱۲۵	۶	۰	۹	۰
رانندگی	۲	۲۱۶	۳	۵۱	۱
پیاده	۲۵	۷	۸۲۴	۲۲	۱۸
اتوبوس	۴	۴۶	۳	۲۸۵	۲۱
دوچرخه	۲	۴	۵	۱۱	۲۷۴

جدول ۸: ماتریس ابهام روش درخت تصمیم

واقعی					تجزیه و تحلیل
قطار	رانندگی	پیاده	اتوبوس	دوچرخه	
قطار	۱۰۷	۱۷	۴	۶	۰
رانندگی	۱۱	۱۹۴	۱۰	۳۷	۳
پیاده	۳۲	۱۰	۸۱۸	۲۹	۳۳
اتوبوس	۵	۵۴	۱	۳۰۱	۲۲
دوچرخه	۳	۴	۲	۵	۲۵۶

جدول ۱۱: ماتریس ابهام روش بیز ساده

واقعی					تجزیه و تحلیل
قطار	رانندگی	پیاده	اتوبوس	دوچرخه	
قطار	۷۳	۱۴	۱۱	۱۳	۱
رانندگی	۴۲	۲۰۲	۷	۶۶	۱
پیاده	۲۹	۱۲	۷۸۶	۲۸	۳۱
اتوبوس	۱۱	۴۱	۱۴	۲۱۸	۲۴
دوچرخه	۳	۱۰	۱۷	۵۳	۲۵۷

جدول ۱۰: ماتریس ابهام روش SVM

واقعی					تجزیه و تحلیل
قطار	رانندگی	پیاده	اتوبوس	دوچرخه	
قطار	۱۰۶	۱۰	۱	۶	۰
رانندگی	۱۶	۱۸۸	۲	۳۱	۱
پیاده	۳۱	۱۲	۸۲۴	۲۹	۳۱
اتوبوس	۳	۶۳	۳	۳۰۶	۱۹
دوچرخه	۲	۶	۵	۶	۲۶۳

جدول ۱۲: مقایسه دقت روش‌های طبقه‌بندی

مدل	دقت کلی	شاخص کاپا
درخت تصمیم	۰٫۸۵۳	۰٫۷۹۵
MLP	۰٫۸۷۸	۰٫۸۵۹
SVM	۰٫۸۵۹	۰٫۸۰۳
بیز ساده	۰٫۷۸۲	۰٫۶۹۸

## ۶-۲- بررسی تأثیر پیش‌پردازش داده‌ها و ویژگی

### جدید بر دقت نتایج

به‌منظور بررسی میزان تأثیرگذاری مراحل پیش‌پردازش بر دقت نتایج تحقیق، چهار روش مورد استفاده بر روی داده‌های خام نیز پیاده‌سازی شد. مطابق جدول (۱۳) معیار دقت کلی برای مقایسه نتایج روش‌ها با پیاده‌سازی بر روی داده‌های خام و داده‌های پالایش‌شده محاسبه شده است. مقایسه دقت کلی نتایج نشان‌دهنده تأثیر قابل توجه پیش‌پردازش داده‌ها بر روی

دقت نتایج است، به‌طوری‌که در همه روش‌ها دقت نتایج پس از انجام مراحل پیش‌پردازش بیش از ۴ درصد بهبود یافته است. همچنین در این تحقیق از درصد حضور نقاط هر خط سیر در فاصله یک انحراف معیار از میانگین کل سرعت نوع‌های حمل‌ونقل به‌عنوان یک ویژگی جدید استفاده شده است. معیار دقت کلی برای چهار روش مورد استفاده در تحقیق با حضور و عدم حضور ویژگی جدید نیز در جدول (۱۴) مورد مقایسه قرار گرفته است. نتایج نشان می‌دهد

افزودن این ویژگی به مدل‌های طبقه‌بندی موجب بهبود دقت کلی نتایج شده است.

جدول ۱۳: مقایسه دقت کلی نتایج پس از پیش‌پردازش داده‌ها

مدل	بدون پیش‌پردازش	با پیش‌پردازش
درخت تصمیم	۰٫۸۰۸	۰٫۸۵۳
MLP	۰٫۸۳۷	۰٫۸۷۸
SVM	۰٫۸۱۲	۰٫۸۵۹
بیز ساده	۰٫۷۴۲	۰٫۷۸۲

جدول ۱۴: مقایسه دقت کلی نتایج با افزودن ویژگی جدید

مدل	بدون حضور ویژگی جدید	با حضور ویژگی جدید
درخت تصمیم	۰٫۸۴۲	۰٫۸۵۳
MLP	۰٫۸۶۹	۰٫۸۷۸
SVM	۰٫۸۵۱	۰٫۸۵۹
بیز ساده	۰٫۷۷۲	۰٫۷۸۲

### ۳-۶- مقایسه با تحقیقات پیشین

به‌منظور سنجیدن دقت روش‌های پیاده‌سازی این تحقیق در پیش‌بینی نوع حمل‌ونقل خطوط سیر، دقت این تحقیق با مطالعاتی که پیش‌ازین بر روی داده‌های پروژه ژئولایف انجام شده بود، مقایسه شد. مطابق با آنچه در جدول (۱۵) آورده شده است، نتایج این تحقیق با استفاده از مدل MLP دقت بهتری نسبت به مطالعات مشابه از خود نشان داده است. اگرچه مدل‌های دیگر این تحقیق نیز توانسته‌اند دقت بهتری نسبت به اکثر

مطالعات قبلی از خود نشان دهند. این بهبود دقت می‌تواند ناشی از پیش‌پردازش مناسب داده‌ها، استخراج ویژگی‌های جدید و تنظیم دقیق پارامترهای ورودی مدل‌ها باشد. در این تحقیق برخلاف تحقیق ژنگ و همکاران [۲۰] از مسئله قطعه‌بندی خودکار خطوط سیر توسط مدل چشم‌پوشی شده است و خطوط سیر به شکل قطعه‌بندی شده به مدل معرفی شده است و همین مسئله بخشی از این اختلاف دقت را به وجود آورده است.

جدول ۱۵: مقایسه نتایج تحقیق با تحقیقات مشابه

تحقیق	دقت کلی	تحقیق	دقت کلی
هوانگ و همکاران (۲۰۲۰) [۲۷]	۰٫۸۱	ونگ و همکاران (۲۰۱۷) [۲۲]	۰٫۷۴
جیمز و همکاران (۲۰۲۰) [۲۶]	۰٫۹۳	اندو و همکاران (۲۰۱۶) [۱۳]	۰٫۶۹
نواز و همکاران (۲۰۲۰) [۲۵]	۰٫۸۴	ژنگ و همکاران (۲۰۱۰) [۲۱، ۱۲]	۰٫۷۶
دبیری و همکاران (۲۰۱۸) [۱۰]	۰٫۸۵	بهترین مدل این تحقیق (MLP)	۰٫۸۸

## ۷- نتیجه‌گیری

در این تحقیق، برای ایجاد مدل شناسایی نوع حمل‌ونقل کاربران بر مبنای داده‌های خطوط سیر روش‌های مختلف طبقه‌بندی به کار گرفته شد. مطالعه این تحقیق بر روی داده‌های خام GPS صورت گرفت و هیچ داده خارجی بر آن اضافه نشد. در مرحله پیش‌پردازش داده‌های غیر معتبر حذف شدند و به‌منظور کاهش خطای اتفاقی یک کرنل نرم‌کننده بر روی خطوط سیر اعمال شد. سپس یک مجموعه از ویژگی‌های خطوط سیر شامل سرعت، طول، شتاب و غیره استخراج شد و تعدادی ویژگی آماری نیز بر آن‌ها افزوده شد. سپس مدل‌های طبقه‌بندی با چهار روش درخت تصمیم، بیز ساده، SVM و MLP ایجاد شد و عملکرد آن‌ها مورد بررسی و مقایسه قرار گرفت.

نتایج حاصل از این تحقیق نشان داد که شبکه عصبی پرسپترون چندلایه در بین روش‌های مختلف بهترین کارایی را در شناسایی نوع حمل‌ونقل کاربران دارد. ارزیابی نتایج نشان داد که مدل‌های پیاده‌سازی این تحقیق نسبت به تحقیقات پیشین عملکرد بهتری از خود نشان می‌دهد. عملکرد مناسب مدل‌های

## مراجع

- [1] Y. Zheng, L. Capra, O. Wolfson, and H. Yang, "Urban computing: concepts, methodologies, and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 5, pp. 1-55, 2014.
- [2] Y. Zheng, "Trajectory data mining: an overview," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 6, pp. 1-41, 2015.
- [3] E. Camossi, P. Villa, and L. Mazzola, "Semantic-based anomalous pattern discovery in moving object trajectories," *arXiv preprint arXiv:1305.1946*, 2013.
- [4] O. Wolfson, B. Xu, S. Chamberlain, and L. Jiang, "Moving objects databases: Issues and solutions," in *Proceedings. Tenth International Conference on Scientific and Statistical Database Management (Cat. No. 98TB100243)*, 1998, pp. 111-122.
- [5] A. Elragal and N. El-Gendy, "Trajectory data mining: integrating semantics," *Journal of Enterprise Information Management*, vol. 26, pp. 516-535, 2013.
- [6] S. Chen, C. S. Jensen, and D. Lin, "A benchmark for evaluating moving object indexes," *Proceedings of the VLDB Endowment*, vol. 1, pp. 1574-1585, 2008.
- [7] R. H. Güting, T. Behr, and C. Düntgen, "SECONDO: A Platform for Moving Objects Database Research and for Publishing and Integrating Research Implementations," *IEEE Data Eng. Bull.*, vol. 33, pp. 56-63, 2010.
- [8] N. Pelekis, Y. Theodoridis, S. Vosinakis, and T. Panayiotopoulos, "Hermes-a

پیاده‌سازی این تحقیق می‌تواند ناشی از پیش‌پردازش دقیق و جامع داده‌ها و استخراج ویژگی‌ها مناسب باشد. همچنین جستجوی گسترده صورت گرفته برای تنظیم پارامترهای ورودی مدل‌ها منجر به نتایج بهتر شد. خطوط سیر در این تحقیق به شکل بخش‌بندی شده بر مبنای نوع حمل‌ونقل به مدل‌های طبقه‌بندی معرفی شدند. درحالی‌که داده‌های خام ممکن است به شکل دنباله پیوسته‌ای از نوع‌های مختلف حمل‌ونقل ایجاد شوند؛ بنابراین بخش‌بندی خطوط سیر بر مبنای تغییر نوع حمل‌ونقل می‌تواند به‌عنوان یک موضوع تحقیقاتی در آینده مورد توجه قرار گیرد. همچنین برای تحقیقات آینده می‌توان از منابع داده دیگری مانند داده‌های شبکه‌های اجتماعی مکان‌مبنا استفاده کرد و نتایج حاصل را مقایسه نمود. همچنین تعمیم روش پیاده‌سازی این تحقیق به داده‌هایی که شامل برچسب نوع حمل‌ونقل نیستند و به‌کارگیری روش‌های نظارت نشده می‌تواند یک زمینه تحقیقاتی در آینده باشد. استفاده از بعد زمانی داده‌ها و استخراج ویژگی‌های دیگر و مطالعه دقیق آن‌ها نیز می‌تواند در تحقیقات آینده موجب بهبود دقت شوند.

- framework for location-based data management," in *International Conference on Extending Database Technology*, 2006, pp. 1130-1134.
- [9] O. Wolfson, P. Sistla, B. Xu, J. Zhou, and S. Chamberlain, "DOMINO: Databases for moving objects tracking," *ACM SIGMOD Record*, vol. 28, pp. 547-549, 1999.
- [10] S. Dabiri and K. Heaslip, "Inferring transportation modes from GPS trajectories using a convolutional neural network," *Transportation research part C: emerging technologies*, vol. 86, pp. 360-371, 2018.
- [11] N. Eluru, V. Chakour, and A. M. El-Geneidy, "Travel mode choice and transit route choice behavior in Montreal: insights from McGill University members commute patterns," *Public Transport*, vol. 4, pp. 129-149, 2012.
- [12] Y. Zheng, Y. Chen, Q. Li, X. Xie, and W.-Y. Ma, "Understanding transportation modes based on GPS data for web applications," *ACM Transactions on the Web (TWEB)*, vol. 4, pp. 1-36, 2010.
- [13] Y. Endo, H. Toda, K. Nishida, and A. Kawanobe, "Deep feature extraction from trajectories for transportation mode estimation," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 2016, pp. 54-66.
- [14] X. Kong, M. Li, K. Ma, K. Tian, M. Wang, Z. Ning, et al., "Big trajectory data: A survey of applications and services," *IEEE Access*, vol. 6, pp. 58295-58306, 2018.
- [15] J. Shang, Y. Zheng, W. Tong, E. Chang, and Y. Yu, "Inferring gas consumption and pollution emission of vehicles throughout a city," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2014, pp. 1027-1036.
- [16] X. Xiao, Y. Zheng, Q. Luo, and X. Xie, "Inferring social ties between users with human location history," *Journal of Ambient Intelligence and Humanized Computing*, vol. 5, pp. 3-19, 2014.
- [17] R. C. Shah, C.-y. Wan, H. Lu, and L. Nachman, "Classifying the mode of transportation on mobile phones using GIS information," in *Proceedings of the 2014 ACM international joint conference on pervasive and ubiquitous computing*, 2014, pp. 225-229.
- [18] D. J. Patterson, L. Liao, D. Fox, and H. Kautz, "Inferring high-level behavior from low-level sensors," in *International Conference on Ubiquitous Computing*, 2003, pp. 73-89.
- [19] Y. Zheng, Y. Chen, X. Xie, and W.-Y. Ma, "GeoLife2. 0: a location-based social networking service," in *2009 tenth international conference on mobile data management: systems, services and middleware*, 2009, pp. 357-358.
- [20] Y. Zheng, L. Liu, L. Wang, and X. Xie, "Learning transportation mode from raw gps data for geographic applications on the web," in *Proceedings of the 17th international conference on World Wide Web*, 2008, pp. 247-256.
- [21] Y. Zheng, Q. Li, Y. Chen, X. Xie, and W.-Y. Ma, "Understanding mobility based on GPS data," in *Proceedings of the 10th international conference on Ubiquitous computing*, 2008, pp. 312-321.
- [22] H. Wang, G. Liu, J. Duan, and L. Zhang, "Detecting transportation modes using deep neural network," *IEICE TRANSACTIONS on Information and Systems*, vol. 100, pp. 1132-1135, 2017.
- [23] Z. Xiao, Y. Wang, K. Fu, and F. Wu, "Identifying different transportation modes from trajectory data using tree-based ensemble classifiers," *ISPRS International Journal of Geo-Information*, vol. 6, p. 57, 2017.
- [24] S. Dabiri, C.-T. Lu, K. Heaslip, and C. K. Reddy, "Semi-supervised deep learning approach for transportation mode identification using GPS trajectory data," *IEEE Transactions on Knowledge and Data*

- Engineering*, vol. 32, pp. 1010-1023, 2019.
- [25] A. Nawaz, H. Zhiqiu, W. Senzhang, Y. Hussain, I. Khan, and Z. Khan, "Convolutional LSTM based transportation mode learning from raw GPS trajectories," *IET Intelligent Transport Systems*, vol. 14, pp. 570-577, 2020.
- [26] J. James, "Travel Mode Identification With GPS Trajectories Using Wavelet Transform and Deep Learning," *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [27] Z. Huang, P. Wang, and Y. Liu, "Statistical characteristics and transportation mode identification of individual trajectories," *International Journal of Modern Physics B*, vol. 34, p. 2050092, 2020.
- [28] B. E. Boser, I. M. Guyon, and V. N. Vapnik, "A training algorithm for optimal margin classifiers," in *Proceedings of the fifth annual workshop on Computational learning theory*, 1992, pp. 144-152.
- [29] Y. Ma and G. Guo, *Support vector machines applications* vol. 649: Springer, 2014.
- [30] S. R. Safavian and D. Landgrebe, "A survey of decision tree classifier methodology," *IEEE transactions on systems, man, and cybernetics*, vol. 21, pp. 660-674, 1991.
- [31] J. Han, J. Pei, and M. Kamber, *Data mining: concepts and techniques*: Elsevier, 2011.
- [32] R. Hecht-Nielsen, "Theory of the backpropagation neural network," in *Neural networks for perception*, ed: Elsevier, 1992, pp. 65-93.
- [33] T. Vincenty, "Direct and inverse solutions of geodesics on the ellipsoid with application of nested equations," *Survey review*, vol. 23, pp. 88-93, 1975.
- [34] R. W. Schafer, "What is a Savitzky-Golay filter?[lecture notes]," *IEEE Signal processing magazine*, vol. 28, pp. 111-117, 2011.
- [35] Y. Zheng, H. Fu, X. Xie, W.-Y. Ma, and Q. Li, "Geolife GPS trajectory dataset-user guide," *Microsoft Research*, 2011.
- [36] Y. Zheng, L. Zhang, X. Xie, and W.-Y. Ma, "Mining interesting locations and travel sequences from GPS trajectories," in *Proceedings of the 18th international conference on World wide web*, 2009, pp. 791-800.
- [37] Y. Zheng, X. Xie, and W.-Y. Ma, "GeoLife: A collaborative social networking service among user, location and trajectory," *IEEE Data Eng. Bull.*, vol. 33, pp. 32-39, 2010.
- [38] I. Syarif, A. Prugel-Bennett, and G. Wills, "SVM parameter optimization using grid search and genetic algorithm to improve classification performance," *Telkomnika*, vol. 14, p. 1502, 2016.



## *Analyzing the performance of different machine learning methods in determining the transportation mode using trajectory data*

Morteza Tayebi <sup>1</sup>, Parham Pahlavani <sup>2\*</sup>

1- GIS PhD student, School of Surveying and Geospatial Engineering, College of Engineering, University of Tehran

2- Associate Professor, School of Surveying and Geospatial Engineering, College of Engineering, University of Tehran

### **Abstract**

With the widespread advent of the smart phones equipping with Global Positioning System (GPS), a huge volume of users' trajectory data was generated. To facilitate the urban management and present appropriate services to users, studying these data was raised as a widespread research field and has been developing since then. In this research, the transportation mode of users' trajectories was identified based on their raw GPS data. These data are often associated with errors, it was attempted to minimize them by applying a comprehensive pre-processing procedure in this research. Accordingly, various features were extracted to identify the transportation modes including walk, bike, train, bus, and driving. In this regard, four classification methods including decision tree, multilayer perceptron neural network, Naïve Bayes, and support vector machine were used to build a predictive model. In order to improve the performance of the implementation methods, the percentage of the points of each trajectory on the distance of one standard deviation from the total speed average of transportation modes has been used as a new feature. The above-mentioned four models were implemented with different regularization parameters and their values were set to the optimal values by applying a comprehensive grid search. Then, Kappa and the overall accuracy indices were employed to evaluate different methods. The results of this study show that the multilayer perceptron neural network with overall accuracy of 0.88 has the best results compared to the other models.

**Key words:** Trajectory data, Determining the transportation mode, Classification, Machine learning.