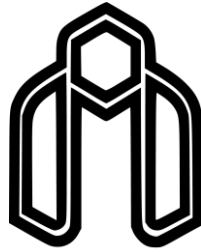


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی شاهرود

دانشکده مهندسی عمران

پایان نامه کارشناسی ارشد مهندسی راه و ترابری

ارائه‌ی روشی برای خوشه‌بندی تصادفات با استفاده از داده‌های مکانی

نگارنده

امین گنجعلی خسروشاهی

استاد راهنما

دکتر ایمان آقایان

استاد مشاور

دکتر عبدالاحد چوپانی

بهمن ۱۳۹۷

تقدیم به

پدر و مادر عزیزم

که همواره در تمام مراحل زندگی ام برایم تکیه‌گاهی محکم بودند.

تقدیر و تشکر

این جانب لازم می‌دانم که بدین وسیله از زحمات بی‌دریغ و نیز همراهی و راهنمایی‌های مستمر اساتید راهنمای این پایان‌نامه جناب آقای دکتر ایمان آقایان و جناب آقای دکتر عبدالاحد چوپانی تشکر و قدردانی نمایم.

تعهدنامه

این جانب امین گنجعلی خسروشاهی دانشجوی دوره کارشناسی ارشد رشته عمران- راه و ترابری دانشکده مهندسی عمران دانشگاه صنعتی شاهرود نویسنده پایان نامه ارائه‌ی روشی برای خوشه‌بندی و شناسایی عوامل مؤثر در تصادفات با استفاده از داده‌های مکانی، تحت راهنمایی دکتر ایمان آقایان و دکتر عبدالاحد چوپانی متعهد می‌شوم:

- تحقیقات در این پایان نامه توسط این جانب انجام شده است و از صحت و اصالت برخوردار است .
- در استفاده از نتایج پژوهش‌های محققان دیگر به مرجع مورداستفاده استناد شده است .
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است .
- کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود می‌باشد و مقالات مستخرج با نام « دانشگاه صنعتی شاهرود » و یا « Shahrood University of Technology » به چاپ خواهد رسید
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده‌اند در مقالات مستخرج از پایان نامه رعایت می‌گردد.
- در کلیه مراحل انجام این پایان نامه ، در مواردی که از موجود زنده (یا بافتهای آن‌ها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است .
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری ، ضوابط و اصول اخلاق انسانی رعایت شده است.

تاریخ

امضای دانشجو

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه‌های رایانه‌ای، نرم‌افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می‌باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی‌باشد.

چکیده

با مشاهده آمار تصادفات به وضوح قابل درک است که تصادفات ترافیکی آثار زیان باری را به سلامت جامعه وارد می کند و هزینه های اقتصادی نامطلوبی را به افراد تحمیل می کند، از این رو محققین را به مطالعه ای تمام جوانب و روش های ممکن برای شناسایی تصادفات وادار می سازد. امروزه روش های داده کاوی بالأخص خوشه بندی برای تحلیل داده های تصادف و پیدا کردن الگوهای مکانی آنها مورد استفاده قرار می گیرند. در این تحقیق از روشی ترکیبی برای خوشه بندی داده های تصادفات با عنوان الگوریتم GridBSCAN استفاده کردیم. خوشه بندی DBSCAN به علت در نظر گرفتن تراکم داده های مکانی و GRID با شبکه بندی به علت بالا بردن دقت و سرعت محاسبات در داده های انبوه همچون داده های تصادفات قابل توجیه است. همچنین از روش های دیگر خوشه بندی همچون K-Means ، Nnh ، KDE و خودهمبستگی مکانی استفاده شد تا با روش GridBSCAN مقایسه شود. نتایج تحقیق نشان داد که الگوریتم Nnh دقیق ترین روش برای شناسایی مکانی نقاط حادثه خیز می باشد و الگوریتم GridBSCAN نیز برای تفکیک مناطق پرتراکم تصادفات در داده های با حجم بالا کاربردی می باشد. خوشه های به دست آمده از این الگوریتم، تصادفات مناطق مختلف شهری را با توجه به تراکم آنها تفکیک کرد. با بررسی پارامترهای موجود در این خوشه ها می توان در جهت کاهش تصادفات و درک عوامل مؤثر بر آنها قدم برداشت.

کلمات کلیدی: تصادفات درون شهری، الگوریتم های داده کاوی، تحلیل خوشه بندی مکانی،

الگوریتم GridBSCAN

فهرست مطالب

چکیده	۵
فصل ۱ مقدمه	۱
۱-۱- مقدمه	۲
۲-۱- تعریف مسئله و ضرورت انجام تحقیق	۴
۳-۱- روش اجرای تحقیق	۴
۴-۱- اهداف تحقیق	۶
۵-۱- نوآوری تحقیق	۶
۶-۱- ساختار تحقیق	۷
فصل ۲ مبانی نظری و پیشینه تحقیق	۹
۱-۲- مقدمه	۱۰
۲-۲- عوامل مؤثر بر تصادفات	۱۰
۱-۲-۲- عامل محیط	۱۱
۲-۲-۲- وسیله نقلیه	۱۳
۳-۲-۲- مشخصات رانندگان	۱۴
۴-۲-۲- پویایی ترافیک	۱۵
۳-۲- روش‌های خوشه‌بندی تصادفات	۱۵

۱۶ ۱-۳-۲- داده کاوی
۲۱ ۲-۳-۲- خوشه بندی
۲۲ ۳-۳-۲- الگوریتم های خوشه بندی غیر آماری
۳۹ ساختار ترکیبی الگوریتم GriDBSCAN
۴۱ ۴-۳-۲- شاخص های ارزیابی خوشه بندی
۵۳ ۵-۳-۲- انواع الگوریتم های خوشه بندی آماری
۵۷ ۶-۳-۲- استفاده از GIS در مطالعات پیشین
۶۳ فصل ۳ روش تحقیق
۶۴ ۱-۳- مقدمه
۶۵ ۲-۳- محدوده مورد مطالعه
۶۹ ۲-۳- معرفی داده ها
۷۲ ۴-۳- چارچوب تعیین نقاط پرتراکم با روش های خوشه بندی
۷۲ ۱-۴-۳- آماده سازی داده ها
۷۷ ۵-۳- خوشه بندی Nnh و K-Means در نرم افزار CrimeStat
۷۹ ۱-۵-۳- نحوه کار در نرم افزار CrimeStat
۸۲ ۶-۳- خوشه بندی GriDBSCAN
۸۳ ۱-۶-۳- حذف نویزها:
۸۵ ۸-۳- مقایسه روش های خوشه بندی
۸۵ ۱-۸-۳- شاخص دقت پیش بینی (PAI)

۸۶ANOVA روش ۲-۸-۳
۸۷ضریب همبستگی اتا ۳-۸-۳
۸۹ فصل ۴ نتایج
۹۰ ۱-۴-۱ مقدمه
۹۰ ۲-۴-۲ نتایج حاصل از خوشه بندی GridBSCAN
۹۴ ۳-۴-۳ روش خوشه بندی Nnh
۱۰۴ ۵-۴-۵ روش خوشه بندی K-Means
۱۰۶ ۶-۴-۶ روش تخمین چگالی کرنل
۱۰۹ ۷-۴-۷ نتایج روش‌های خوشه‌بندی و بررسی عوامل موثر بر تصادفات در نقاط حادثه‌خیز
۱۱۰ ۱-۷-۱ توصیف خوشه‌های حادثه‌خیز شهر گیزه
۱۱۶ ۲-۷-۲ توصیف خوشه‌های حادثه‌خیز شهر ایزمیت
۱۲۳ ۸-۴-۸ مقایسه روش‌های خوشه‌بندی و شاخص‌های ارزیابی خوشه‌ها
۱۲۵ ۹-۴-۹ نتایج حاصل از توابع موجود در GIS
۱۲۷ فصل ۵ بحث و نتیجه‌گیری
۱۲۸ ۱-۵-۱ مقدمه
۱۲۸ ۲-۵-۲ خلاصه نتایج و جمع‌بندی
۱۳۰ ۳-۵-۳ پیشنهادات کارهای آینده
۱۳۱ مراجع:

فهرست شکل‌ها

- شکل ۱-۱: روند آمار افراد فوتی ناشی از تصادفات در کشور ترکیه (WHO,2015)..... ۳
- شکل ۱-۲: روند کلی تحقیق ۵
- شکل ۲-۱: زنجیره عوامل مؤثر بر تصادف (اوه و همکاران، ۲۰۰۱)..... ۱۰
- شکل ۲-۲: پارامترهای انسانی مؤثر بر تصادفات سال ۲۰۰۵ تا ۲۰۱۴ در کشور نروژ (ساگبرگ، ۲۰۱۸) ۱۵
- شکل ۲-۳: (الف) تراکم قابل دسترس q و p (ب) اتصال تراکم p و q ۳۱
- شکل ۲-۴: تعیین Eps بهینه با استفاده از نمودار Knn-distance ۳۳
- شکل ۲-۵: شبکه همسایگی در فضای دو بعدی ۳۸
- شکل ۲-۶: پارتیشن بندی فضایی با استفاده از شبکه ۳۹
- شکل ۲-۷: پیرامون سلول ϵ ۴۰
- شکل ۲-۸: نقطه p درون سلول راست قرار دارد، اما در سلول چپ نیز شامل می شود ۴۰
- شکل ۲-۹: روند ادغام دو روش Grid و DBSCAN ۴۱
- شکل ۲-۱۰: نمای شماتیک تعریف precision و recall ۴۹
- شکل ۳-۱: استان کوجائلی ترکیه و توزیع نقاط تصادف درون شهری و برون شهری برای سال ۲۰۱۳-۲۰۱۴ ۶۶
- شکل ۳-۲: عکس ماهواره‌های شهرهای اطراف دریای مرمره ۶۷
- شکل ۳-۳: تعداد تصادفات کل به تفکیک شهرهای استان کوجائلی در سال ۲۰۱۳-۲۰۱۴ ۶۸

- شکل ۳-۴: تعداد تصادفات درون شهری به تفکیک شهرها در سال ۲۰۱۳ و ۲۰۱۴ پس از پاکسازی داده‌ها ۷۰
- شکل ۳-۵: درصد تصادفات (الف) فوتی و جراحی در داده های (ب) درون شهری و برون شهری در کل استان کوجائلی ۷۰
- شکل ۳-۶: نقاط تصادف فوتی (الف) و جراحی (ب) در کل استان کوجائلی ۷۱
- شکل ۳-۷: انواع تصادف در مناطق درون شهری به تفکیک شهر ایزمیت و گبزه در سالهای ۲۰۱۳-۲۰۱۴ ۷۵
- شکل ۳-۸: شکل ۳-۹: توزیع زمانی تصادفات جراحی شهر گبزه (۶۴۷) در سال ۲۰۱۳-۲۰۱۴ ۷۵
- شکل ۳-۱۰: توزیع زمانی تصادفات جراحی شهر ایزمیت (۶۱۱) در سال ۲۰۱۳-۲۰۱۴ ۷۶
- شکل ۳-۱۱: تعداد تصادفات به تفکیک دوره‌های زمانی در ۷۶
- شکل ۳-۱۲: مرحله اولیه وارد کردن فایل shp داده‌های تصادف ۷۹
- شکل ۳-۱۳: وارد کردن مشخصات جغرافیایی ۸۰
- شکل ۳-۱۴: پارامترهای اندازه‌گیری ۸۱
- شکل ۳-۱۵: محیط نرم‌افزار CrimeStat IV برای محاسبه خوشه‌بندی Nnh ۸۲
- شکل ۳-۱۶: جمع‌بندی سطوح کرنل نرمال برای ۵ نقطه (لویین، ۲۰۱۰) ۸۵
- شکل ۴-۱: نمودار Knn با مقدار $K=4$ و مقدار بهینه $Eps(m)$ برای شهر (الف) ایزمیت (ب) گبزه .. ۹۱
- شکل ۴-۲: خوشه بندی GridBSCAN با مقدار $Eps = 350m$ و $minPts = 15$ و $d = 700m$ برای شهر گبزه ۹۲
- شکل ۴-۳: خوشه بندی GridBSCAN با مقدار $Eps = 250m$ و $minPts = 10$ و $d = 500m$ برای شهر ایزمیت ۹۳

- شکل ۴-۴: نقشه موضوعی تصادفات جراحی با شدت نقاط حادثه‌خیز به روش N_{nh} با مقادیر $n_{min} = 20$ و فاصله جستجو ۱۰۰۰ متر ۹۵
- شکل ۴-۵: فواصل جستجو (d) مختلف با $n_{min} = 10$ برای محاسبه خوشه‌بندی N_{nh} برای تصادفات شهر ایزمیت ۹۶
- شکل ۴-۶: یکی از تقاطعات شهر ایزمیت که با معیار $d=100m$ و $n_{min} = 10$ به عنوان یک از خوشه های حادثه‌خیز بدست آمد ۹۷
- شکل ۴-۷: مقدار $n_{min} = 10$ و $d = 100$ برای خوشه‌بندی N_{nh} شهر ایزمیت ۹۸
- شکل ۴-۸: مقدار $n_{min} = 5$ و $d = 100$ برای خوشه‌بندی N_{nh} شهر ایزمیت ۹۹
- شکل ۴-۹: مقدار $n_{min} = 10$ و $d = 200$ برای خوشه‌بندی N_{nh} شهر ایزمیت ۱۰۰
- شکل ۴-۱۰: مقدار $n_{min} = 5$ و $d = 200$ برای خوشه‌بندی N_{nh} شهر ایزمیت ۱۰۱
- شکل ۴-۱۱: مقدار $n_{min} = 5$ و $d = 100$ برای خوشه‌بندی N_{nh} شهر ایزمیت ۱۰۲
- شکل ۴-۱۲: مقدار $n_{min} = 5$ و $d = 100$ برای خوشه‌بندی N_{nh} شهر ایزمیت ۱۰۳
- شکل ۴-۱۳: نمودار شاخص اعتبارسنجی دیویس-بولدین برای یافتن تعداد بهینه خوشه‌ها در روش K-Means برای شهر گبزه ۱۰۴
- شکل ۴-۱۴: نمودار شاخص اعتبارسنجی دیویس-بولدین برای یافتن تعداد بهینه خوشه‌ها در روش K-Means برای شهر ایزمیت ۱۰۴
- شکل ۴-۱۵: خوشه‌بندی K-Means شهر گبزه ۱۰۵
- شکل ۴-۱۶: خوشه‌بندی K-Means شهر ایزمیت ۱۰۵
- شکل ۴-۱۷: نقشه تخمین ترام کرنل با شعاع جستجوی ۲۰۰ متر برای شهر گبزه ۱۰۷
- شکل ۴-۱۸: نقشه تخمین ترام کرنل با شعاع جستجوی ۲۰۰ متر برای شهر ایزمیت ۱۰۸

- شکل ۴-۱۹: تقاطعی در منطقه مسکونی که در همهی روش‌های خوشه‌بندی به عنوان نقطه حادثه‌خیز شناخته شده است. ۱۰۹.....
- شکل ۴-۲۰: ترکیب خوشه CL3 از روش GriDBSCAN و خوشه‌های حاصل از Nnh ۱۱۲.....
- شکل ۴-۲۱: تصاویر هوایی و تصاویر دید خیابان در نرم‌افزار گوگل ارث ۱۱۳.....
- شکل ۴-۲۲: میدان پرتصادف در شهر گبزه..... ۱۱۴.....
- شکل ۴-۲۳: ادغام خوشه بندی Nnh، KDE و GriDBSCAN ۱۱۴.....
- شکل ۴-۲۴: توزیع تصادفات در خوشه CL5 در روش GriDBSCAN ۱۱۵.....
- شکل ۴-۲۵: تصادفات خوشه CL2 و خوشه‌های روش Nnh در شهر گبزه ۱۱۶.....
- شکل ۴-۲۶: تصویر ماهواره‌ای محور D100 عبوری از شهر ایزمیت و توزیع نقاط تصادف ۱۱۷.....
- شکل ۴-۲۷: نقطه حادثه‌خیز در نزدیکی آتشنشانی ایزمیت ۱۱۹.....
- شکل ۴-۲۸: محدوده خوشه CL6 شهر ایزمیت..... ۱۲۰.....
- شکل ۴-۲۹: تقاطع غضنفریبلگه شهر ایزمیت خوشه CL6 ۱۲۱.....
- شکل ۴-۳۰: خوشه CL2 شهر ایزمیت ۱۲۲.....
- شکل ۴-۳۱: نتیجه تابع خود همبستگی موران برای شهر ایزمیت ۱۲۵.....
- شکل ۴-۳۲: نتیجه تابع خود همبستگی موران برای شهر ایزمیت ۱۲۶.....

فهرست جدول‌ها

- جدول ۱-۱: تعداد تصادفات شامل کشته‌ها و زخمی‌ها بر اساس نوع منطقه در کوجائلی در سال ۲۰۱۳-۲۰۱۴..... ۳
- جدول ۱-۲: طبقه بندی الگوریتم‌های خوشه‌بندی غیر آماری (هان و همکاران، ۲۰۱۱)..... ۲۳
- جدول ۱-۳: جمعیت شهرهای استان کوجائلی (سال ۲۰۱۳-۲۰۱۴) (Turkstat, ۲۰۱۸) ۶۵
- جدول ۲-۳: خلاصه داده‌های عددی..... ۷۳
- جدول ۳-۳: داده‌های توصیفی تصادفات جرحی شهر گبزه و ایزمیت ۷۳
- جدول ۴-۱: نتایج آزمون ANOVA برای خوشه بندی GridBSACN خوشه‌های شهر گبزه ۱۱۱
- جدول ۴-۲: توزیع نقاط در تقاطعات مختلف در خوشه CL3 شهر گبزه ۱۱۲
- جدول ۴-۳: نتایج آزمون ANOVA برای خوشه‌های شهر ایزمیت ۱۱۸
- جدول ۴-۴: آمار توصیفی تعداد خودروهای درگیر تصادف در خوشه CL11 ۱۱۹
- جدول ۴-۵: جدول شاخص‌های خوشه‌بندی GriDBSCAN ۱۲۳
- جدول ۴-۶: مقادیر مورد استفاده برای الگوریتم GriDBSCAN ۱۲۴
- جدول ۴-۷: مقادیر نرخ ضربه و شاخص دقت پیش‌بینی روش‌های مختلف..... ۱۲۴
- جدول ۴-۸: نتایج تابع خودهمبستگی موران برای شهر ایزمیت ۱۲۶

فصل ۱

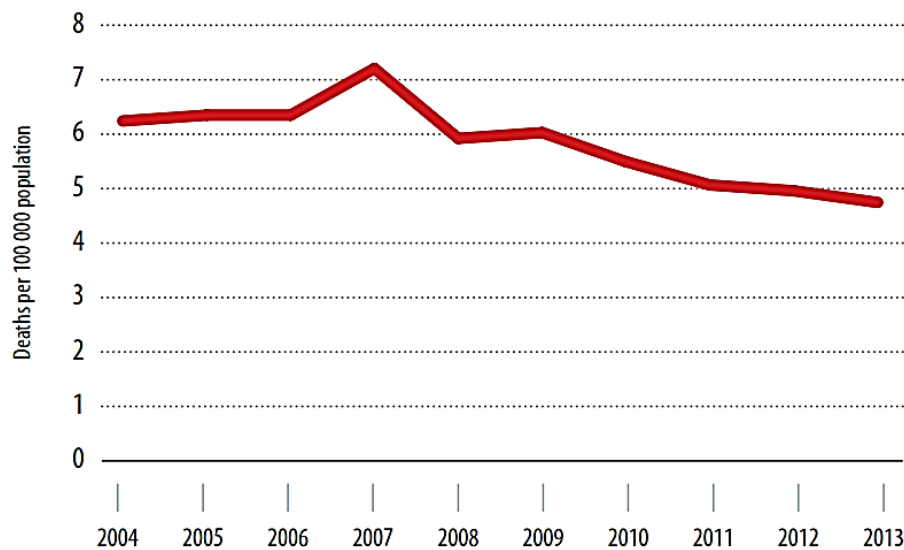
مقدمه

۱-۱- مقدمه

در دنیای امروز به علت افزایش جمعیت و نیاز به استفاده از راه‌ها مسئله‌ی حمل‌ونقل دارای اهمیت بالایی است. در هر کشور بر اساس زیرساخت‌های توسعه‌یافته در آن، یکی از شیوه‌های حمل‌ونقل صورت متداول‌تری به خود می‌گیرد. در کشور ترکیه نیز مانند بسیاری از کشورها، سفرهای جاده‌ای کاربرد بیشتری نسبت به سایر شیوه‌های حمل‌ونقل دارند؛ به همین دلیل جاده‌ها شریان‌های حیاتی برقرارکننده ارتباط بین مناطق مختلف هستند. از طرفی صدمات جانی و مالی ناشی از تصادفات رانندگی در جاده‌ها جوامع را با مشکلات عدیده‌ای روبه‌رو ساخته‌اند. حوادث رانندگی و به دنبال آن مرگ‌ومیر در حال حاضر به‌عنوان یکی از مشکلات اصلی ایمنی و سلامت عمومی در سراسر جهان مطرح است.

طبق آمار جهانی از سازمان جهانی بهداشت^۱ در سال ۲۰۱۵، در سال‌های ۲۰۰۰ تا ۲۰۱۳ میزان مرگ‌ومیر ناشی از تصادفات در سراسر کشورها ۱۳ درصد افزایش یافته است. در سال ۲۰۱۳، به دلیل تصادفات جاده‌ای در سراسر جهان، حدود ۱,۲۵ میلیون نفر کشته شده‌اند و همچنین ۵۰ میلیون نفر دچار جراحت شده‌اند. علاوه بر این، تصادفات یکی از دلایل اصلی مرگ‌ومیر به‌ویژه در بین جوانان بین ۱۵ تا ۲۹ سال بوده است (سازمان بهداشت جهانی، ۲۰۱۵). با توجه به آخرین گزارش سازمان بهداشت جهانی در سال ۲۰۱۵ برای کشور ترکیه ۳۶۸۵ فوتی (۷۷٪ مردان و ۲۳٪ زنان) گزارش شده است که در چند سال اخیر روند کاهشی داشته است. شکل شماره ۱-۱ در گزارش سازمان بهداشت جهانی آمده است.

¹ World Health Organization



شکل ۱-۱: روند آمار افراد فوتی ناشی از تصادفات در کشور ترکیه (WHO,2015)

گزارش کامل تر افراد فوتی و جراحت دیده در سامانه رسمی اینترنتی سازمان آمار ملی پلیس ترکیه^۱ در سال ۲۰۱۳ و ۲۰۱۴ آمار تصادفات مربوط به مناطق درون شهری و برون شهری در جدول ۱-۱ آمده است. (سایت آمار تصادفات ترکیه، ۲۰۱۸)

جدول ۱-۱: تعداد تصادفات شامل کشته‌ها و زخمی‌ها بر اساس نوع منطقه در کوجائلی در سال ۲۰۱۳-۲۰۱۴

تعداد افراد جراحت دیده		تعداد افراد کشته شده			تعداد کل تصادفات منجر به فوت و جراحت			سال	
مجموع	منطقه درون شهری	منطقه برون شهری	مجموع	منطقه درون شهری	منطقه برون شهری	مجموع	منطقه درون شهری		منطقه برون شهری
۲۷۴ ۸۲۹	۱۸۳ ۳۰۷	۹۱ ۵۲۲	۳ ۶۸۵	۱ ۳۷۲	۲ ۳۱۳	۱۶۱ ۳۰۶	۱۲۰ ۰۹۵	۴۱ ۲۱۱	۲۰۱۳
۲۸۵ ۰۵۹	۱۹۱ ۶۵۳	۹۳ ۴۰۶	۳ ۵۲۴	۱ ۲۴۳	۲ ۲۸۱	۱۶۸ ۵۱۲	۱۲۶ ۵۲۷	۴۱ ۹۷۵	۲۰۱۴

آمار تصادفات ترافیکی راه‌های ترکیه، TurkStat

¹ TURKISH NATIONAL POLICE

با مشاهده آمار تصادفات به وضوح قابل درک است که تصادفات ترافیکی آثار زیان باری را به سلامت جامعه وارد می کند و هزینه های اقتصادی به افراد تحمیل می کند. در تحقیقات مربوط به تصادفات روش های مختلفی برای تحلیل تصادفات وجود دارد. امروزه روش های داده کاوی^۱ بالأخص خوشه بندی و طبقه بندی برای کاهش عدم تجانس داده های تصادف و پیدا کردن الگوهای مبهم مورد استفاده قرار می گیرند. در این تحقیق از روشی ترکیبی و قابل توجه برای خوشه بندی داده های تصادفات با عنوان الگوریتم GridBSCAN استفاده کردیم.

۲-۱- تعریف مسئله و ضرورت انجام تحقیق

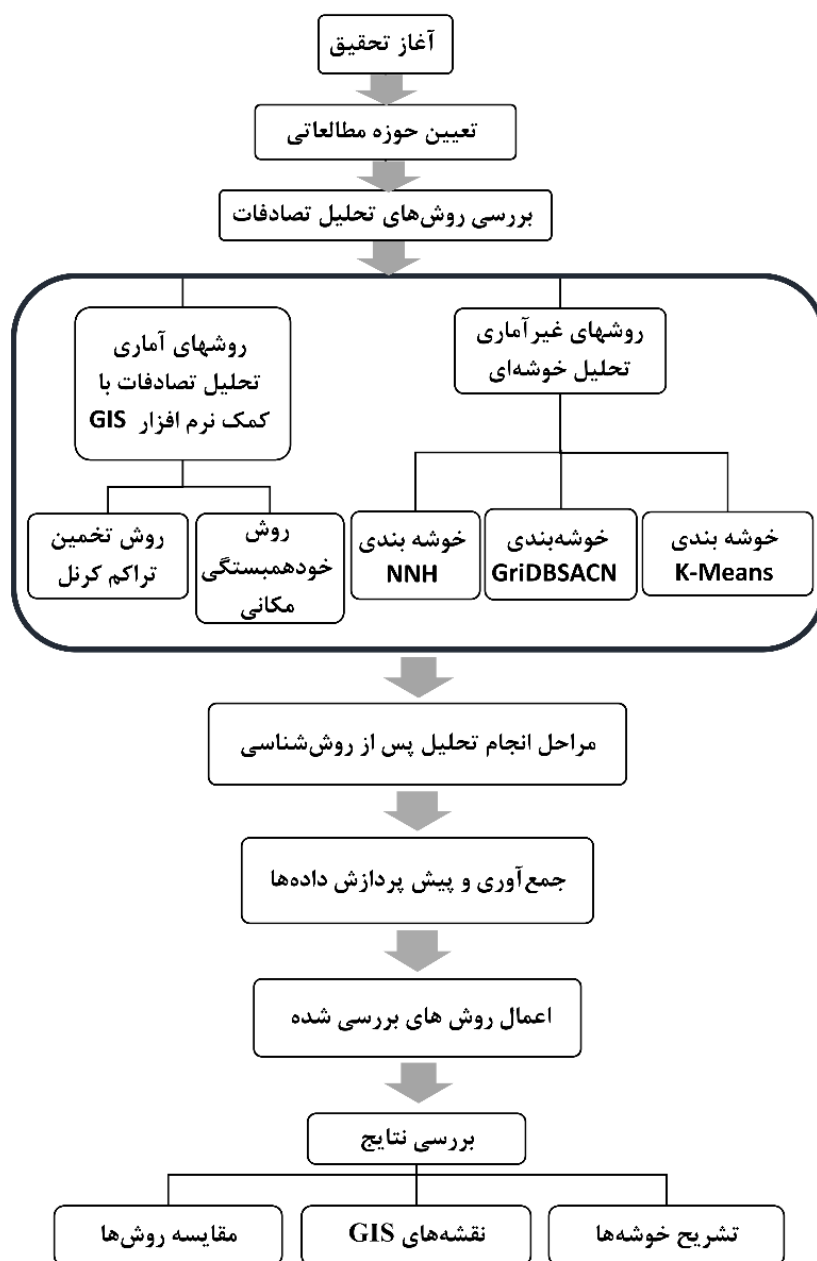
ایمنی، یک پارامتر بسیار مهم مورد بررسی در بسیاری از کشورهای پیشرفته و در حال توسعه است که نشان دهنده اهمیت بالای این موضوع است. با مطالعه تصادفات می توان عوامل مؤثر در رخداد تصادفات ترافیکی را شناسایی کرد. تعریف و توصیف عوامل مؤثر در تصادفات گامی اساسی در جهت اقدامات عملی برای کاهش و به حداقل رساندن عوامل ایجاد تصادف است که باعث ایمنی بیشتر سیستم حمل و نقل و بهبود آن می شود. نظر به اینکه تصادف یک پدیده مکانی است، در این تحقیق رویکرد مبتنی بر مکان مورد تأکید قرار گرفته است. با استفاده از نرم افزارهایی مکان-مبنا همچون GIS، مهندسان و دیگر متخصصین می توانند تحلیل های متعددی را در رابطه با منطقه جغرافیایی خاص و اتفاقات مربوط به آن منطقه مانند تصادفات را نمایش دهند.

۳-۱- روش اجرای تحقیق

در شکل ۱-۲ روند تحقیق نشان داده شده است. همانند دیگر تحقیقات ابتدا مطالعات در حوزه تصادفات و روش های تحلیل تصادف، انجام می شود. سپس به آماده سازی داده ها پرداختیم. سپس از

¹ Data mining

الگوریتم‌های خوشه‌بندی آماری و غیر آماری برای انجام تحلیل‌های موردنظر استفاده می‌شود. سپس به بررسی معیارهای مؤثر در خوشه‌های تصادفات پرداخته می‌شود. به عبارتی تصادفات در این مرحله بررسی می‌شوند و می‌توان نقاط پرخطر تصادفات را شناسایی و ویژگی‌های آن‌ها را به دست آورد.



شکل ۱-۲: روند کلی تحقیق

۱-۴- اهداف تحقیق

هدف اصلی این تحقیق ارائه‌ی روشی برای شناسایی نقاط حادثه‌خیز تصادفات در مناطق درون‌شهری با استفاده از روش‌های تحلیل داده‌های مکانی با استفاده از GIS و ترکیب با روش‌های داده‌کاوی همچون خوشه‌بندی است. از دیگر اهداف این پژوهش می‌توان به موارد زیر اشاره کرد:

- تعیین تصادفات غالب در مناطق شهری
 - دسته‌بندی تصادفات به گروه‌های همگن یا خوشه‌ها
 - تعیین روشی مؤثر و کارآمد برای خوشه‌بندی
 - تعیین خودهمبستگی مکانی عوامل مؤثر شناسایی‌شده در دوره زمانی داده‌ها
 - بررسی ویژگی خوشه‌های حادثه‌خیز شناسایی‌شده در نقاط مختلف شهر
- در این تحقیق به شناسایی پارامترهای تأثیرگذار در نقاط پر تصادف در مناطق مختلف شهری می‌پردازیم. با تشخیص این پارامترها و سیاست‌گذاری در جهت رفع عوامل ایجاد تصادف می‌توان به میزان مطلوبی از ایمنی در راه‌ها برسیم.

۱-۵- نوآوری تحقیق

مطالعات زیادی در زمینه تحلیل تصادفات انجام شده است. پژوهش‌های اخیر هر کدام بر روی پارامترهای خاص در بحث تصادفات تمرکز دارند؛ مانند نوع تصادف، نوع خودرو، نوع مسیر و تأثیرات آب‌وهوایی و... داده‌های در دسترس برای این تحقیق دارای پارامترهای مختلفی همچون شرایط فیزیکی مسیر (وجود گارد ریل، شانه مسیر، شرایط روشنایی مسیر، مسیر شیب‌دار یا مسطح، مسیر قوس یا مستقیم) و برخی ویژگی‌های تصادف همچون نوع تصادف یا تعداد خودروهای درگیر تصادف هستند. ما در این تحقیق با در اختیار داشتن داده‌های جغرافیایی با مختصات جغرافیایی دقیق، با استفاده از ابزار

و روش‌های تحلیل مکانی تصادفات درون‌شهری برای دو شهر ایزمیت و گبزه (در کشور ترکیه) با ویژگی‌های جغرافیایی و شبکه شهری متفاوت را مورد بررسی قرار دادیم. در گذشته تصادفات در شهرهای دیگر ترکیه همچون افیون قره‌حصار^۱ با روش‌های طبقه‌بندی مدل‌سازی شدند (پلات و دوردوران، ۲۰۱۱). در این تحقیق از روش خوشه‌بندی با الگوریتم ترکیبی GridBSCAN استفاده شده است. خوشه‌بندی DBSCAN به علت در نظر گرفتن تراکم^۲ داده‌های مکانی و GRID با شبکه‌بندی به علت بالا بردن دقت و سرعت محاسبات در داده‌های انبوه همچون داده‌های تصادفات قابل توجه است. همچنین از روش‌های دیگر خوشه‌بندی همچون K-Means^۳، NNH^۴، KDE^۵ و خودهمبستگی مکانی^۶ استفاده شد تا با روش GridBSCAN مقایسه شود. به صورت کلی استفاده از روش‌های خوشه‌بندی در تصادفات اولاً استفاده از داده‌های دارای زیرشاخه را در خوشه‌بندی را ممکن می‌سازد. ثانیاً، به جای استفاده از سنجش اکتشافی^۶ فاصله‌ی بین تصادفات، از تشابه مبتنی بر تراکم برای تخصیص تصادفات به خوشه‌های مختلف استفاده می‌کنند.

۱-۶- ساختار تحقیق

در فصل اول به معرفی و ضرورت انجام تحقیق و سپس به جنبه نوآوری و اهداف تحقیق پرداخته شده است. در فصل دوم تحقیقات و پژوهش‌های پیشین در سه بخش عوامل مؤثر بر تصادفات، روش‌های مختلف خوشه‌بندی تصادفات و استفاده از GIS در تحلیل تصادفات که در گذشته انجام شده، تشریح گردیده است.

^۱ Afyonkarahisar

^۲ Density

^۳ Nearest Neighborhood Hierarchical

^۴ Kernel Density Estimation

^۵ Spatial Autocorrelation

^۶ Heuristic measuring

در فصل سوم ضمن معرفی داده‌های مورد استفاده در این تحقیق به تشریح روش‌های متداول خوشه‌بندی نقاط تصادفات و انتخاب مناسب‌ترین و کارآمدترین آن‌ها برای این تحقیق پرداخته شده است.

در فصل چهارم نتایج خوشه‌بندی‌های مختلف بررسی شده و پارامترهای معنی‌دار در هر خوشه بررسی شده و به صورت آماری و تحلیل ترافیکی مورد تحلیل و بررسی قرار گرفته‌اند. در نهایت شاخص دقت پیش‌بینی برای هر نوع خوشه‌بندی محاسبه شد تا بهترین مدل مشخص شود.

در فصل پنجم خلاصه نتایج حاصل از این پژوهش ارائه شده است. در انتها با توجه به نتایج گرفته شده پیشنهادهایی جهت ادامه تحقیقات ارائه گردیده است.

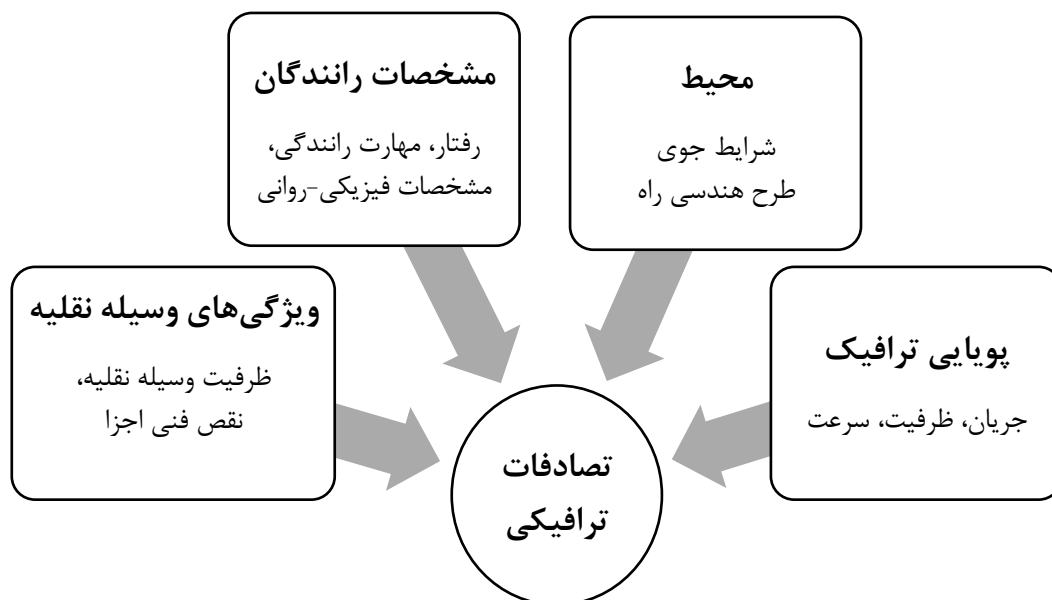
فصل ۲

مبانی نظری و پیشینه تحقیق

در این فصل ابتدا به عوامل مؤثر بر ایجاد تصادفات می‌پردازیم و سپس به روش‌های خوشه‌بندی در تصادفات در مطالعات پیشین و در انتها به کاربرد GIS در کارهای گذشته می‌پردازیم.

۲-۲- عوامل مؤثر بر تصادفات

عوامل مؤثر بر تصادفات به چهار دسته کلی تقسیم‌بندی می‌شوند. (اوه^۱ و همکاران، ۲۰۰۱). با توجه به شکل ۱-۲، این عوامل با تأثیرات متقابل بر یکدیگر بر رخداد تصادفات مؤثر هستند و با وجود نقص در یکی از این عوامل، این چرخه دچار تزلزل شده و این ناپایداری باعث وقوع تصادف می‌شود. نقاط حادثه‌خیز تصادفات از عوامل انسانی و نقص خودرو، طرح هندسی راه، کاربری‌ها و عوارض اطراف، فراوانی و نوع تصادفات تحت تأثیر قرار می‌گیرند.



شکل ۱-۲: زنجیره عوامل مؤثر بر تصادف (اوه و همکاران، ۲۰۰۱)

¹ Oh et al

طرح هندسی راه

از عوامل مهم محیطی در تصادفات می‌توان به طرح هندسی راه‌ها اشاره کرد. عوامل مؤثر بر هندسه راه را می‌توان به صورت زیر طبقه‌بندی کرد (معتمد، ۱۳۸۸):

- عوامل مربوط به طرح هندسی و نواقص فنی مربوط به راه
- عوامل مربوط به علائم و تجهیزات کنترلی جاده و حریم راه‌ها

در گزارشی از کمیته فنی ۳.۲ طراحی و عملکرد ایمن‌تر زیرساخت راه‌ها، انجمن جهانی راه^۱ عوامل مهم در رخداد تصادفات به ویژگی‌های مسیر، طرح هندسی آن، علائم، تابلوهای هشداردهنده و کنترل ترافیک اشاره کردند. طرح هندسی که شامل ویژگی‌هایی از قبیل مسیر قوس، مسیر مستقیم، سربالایی، سرپایینی و مسیر مسطح است، از عوامل مهم در وقوع تصادفات رانندگی هستند. (Association, 2016)

چن و همکارانش اخیراً تأثیر طرح هندسی راه را در سه بزرگراه در ایالات متحده بررسی کردند. نتایج تحقیق نشان داد که تصادفات خسارتی در بزرگراه‌ها حساسیت بیشتری نسبت به تغییر حجم ترافیک، میانگین شیب طولی راه، عرض روفوژ میانی، عرض داخل شانه و وضعیت پیاده‌رو دارد؛ اما تصادفات فوتی و جرحی به تغییرات عرض خط، وضعیت پیاده‌رو و عرض شانه خارجی حساس است (چن و همکاران، ۲۰۱۹).

¹ World Road Association

شرایط جوی

اندرسون و چپمن^۱ (۲۰۱۱) در مطالعه‌ای تصادفات جاده‌ای در وست میدلند^۲ انگلستان را با استفاده از سناریوهای تغییرات آب و هوایی UKCIP و شبیه‌سازی زمانی ارتباط بین دما و تصادفات شدید جاده‌ای بررسی کردند. نتایج آن‌ها نشان داد که کاهش پیش‌بینی‌شده روزهای یخبندان به علت لغزندگی مسیر، منجر به کاهش پنجاه درصدی تعداد تصادفات شد.

زائرزاده (۱۳۸۴) در بررسی نقاط حادثه‌خیز و مطالعه موردی چند نقطه حادثه‌خیز استان خراسان از دیگر عوامل مؤثر در بروز تصادفات جاده‌ای را عوامل محیطی بیان کرد. از جمله مهم‌ترین عوامل محیطی، عوامل طبیعی و شرایط جوی است.

سجادی (۱۳۸۹) در پایان‌نامه‌ای با موضوع تحلیل تصادفات با تأکید بر خصوصیات جاده و محیط، شرایط جوی مانند بارش باران یا برف را که از عوامل تأثیرگذار بر میزان تصادفات راه‌ها است، علاوه بر تأثیر در میدان دید، سبب کاهش اصطکاک سطح جاده و به دنبال آن کاهش کارایی سیستم ترمز در وسایل نقلیه بیان کردند. نسبت تصادف در زمان بارانی نسبت به زمان غیر بارانی ۱/۷۸ برابر می‌شود.

عفتی و همکارانش (عفتی و همکاران، ۱۳۹۳) با تحلیل عوامل مکانی مؤثر بر تمرکز تصادفات در راه‌های برون‌شهری با استفاده از GIS و داده‌کاوی بیان کردند انعکاس نور خورشید نیز از عوامل محیطی است که در تسلط راننده بر محیط اطراف خود مؤثر است. در ساعات اولیه و یا پایانی روز که زاویه ارتفاعی خورشید در کمترین مقدار خود است، نور مستقیم خورشید قابلیت دید راننده را کاهش می‌دهد.

¹ Andersson & Chapman

² West Midland

پراسانکومار و همکارانش^۱ (۲۰۱۱) با خوشه‌بندی مکانی-زمانی تصادفات جاده‌ای، اثر بادهای و باران‌های موسمی را در توزیع مکانی تصادفات بررسی کردند. نتایج به‌دست‌آمده از پژوهش آن‌ها حاکی از آن است که الگوی توزیع مکانی تصادفات در هر نوع وضعیت جوی متغیر است.

۲-۲-۲- وسیله نقلیه

در این قسمت به بررسی وسیله نقلیه، به‌عنوان یکی از عوامل مؤثر بر وقوع تصادفات جاده‌ای، پرداخته می‌شود. تیان و همکارانش^۲ (۲۰۱۰) همچنین مهم‌ترین عوامل مربوط به خودرو که در وقوع تصادفات مؤثر هستند را به‌صورت زیر دسته‌بندی کرد:

- نقص در سیستم برقی (خرابی چراغ‌ها، زاویه اشتباه نور)
- نقص در سیستم فرمان
- نقص در سیستم ترمز (ترمز گیری با انحراف و با تأخیر)
- نقص در سیستم سفر (خرابی لاستیک‌ها، باد لاستیک‌ها)
- دیگر موارد (کیسه هوا، کمربند ایمنی و ...)

ایوان و همکارانش^۳ (۲۰۰۰) نشان دادند که در تصادفاتی که تنها شامل یک خودرو می‌شود رابطه تعداد خودروها با تصادف منفی است؛ اما در مورد تصادفات شامل چند خودرو در صورتی که جاده تقسیم‌نشده باشد با افزایش تعداد خودرو، تعداد تصادفات نیز زیاد می‌شود، اما در جاده تقسیم‌شده (جاده‌ای که به‌وسیله نرده، چمن و غیره از هم جدا شده‌اند) تصادف کاهش می‌یابد.

¹ Prasannakumar et al.

² Tian et al.

³ Ivan et al.

۳-۲-۲- مشخصات رانندگان

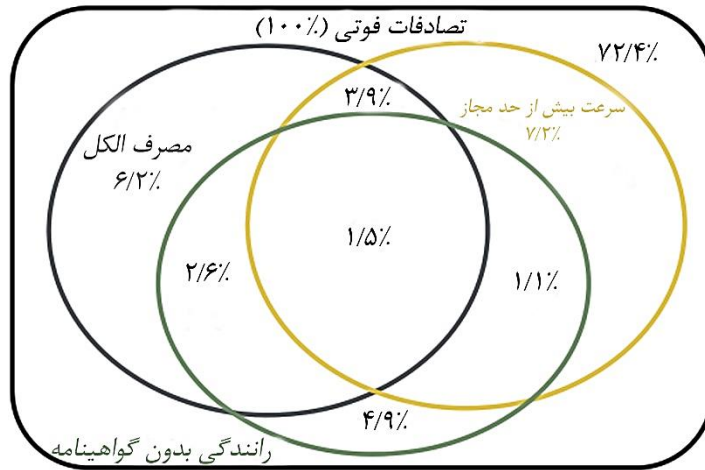
کارشناسان ایمنی بر این عقیده‌اند که تصادفات ترافیکی عمدتاً به دلیل خطاهای انسانی است. بر اساس تحقیقات سلمانی و همکارانش علل انسانی و نیز علل انسانی-محیطی () به ترتیب ۵۷٪ و ۲۶٪ از تصادفات را شامل می‌شوند. لذا، در نظر گرفتن تقابل عامل انسانی با عوامل محیطی و طرح هندسی راه در تصادفات را نمی‌توان انکار کرد (سلمانی و همکاران، ۱۳۸۷).

از عوامل انسانی مؤثر بر تصادفات می‌توان به خصوصیات فردی راننده مانند سن راننده، خستگی و خواب‌آلودگی، مصرف مواد مخدر، تخلف عمدی در هنگام رانندگی اشاره کرد. سلمانی و همکاران در سال ۱۳۸۷ مهم‌ترین عوامل انسانی تأثیرگذار بر تصادفات را به چهار گروه زیر تقسیم‌بندی کردند:

- قوه بینایی و تشخیص: عدم تمرکز ذهنی، تشخیص نادرست فاصله یا سرعت و یا ضعف دید.
- آسیب‌های جسمی: مصرف الکل، خستگی مفرط، مصرف دارو و مسائل روحی و روانی.
- شیوه رانندگی: سرعت‌های غیرمجاز و زیاد و رانندگی تهاجمی.
- عدم مهارت: تجربه ناکافی، ضعف تشخیص و واکنش نادرست رانندگان.

در تحقیق جدیدی ساگبرگ^۱ (۲۰۱۸) در کشور نروژ تأثیر رانندگی افراد بدون گواهینامه در تصادفات فوتی سال ۲۰۰۴-۲۰۱۵ را بررسی کرد. این افراد ۱۰٪ کل تصادفات فوتی را شامل بودند؛ که اکثر این افراد به دو علت رایج مصرف مشروبات الکلی یا سرعت بیش از حد مجاز منجر به تصادف فوتی شدند (شکل ۲-۲).

¹ Sagberg



شکل ۲-۲: پارامترهای انسانی مؤثر بر تصادفات سال ۲۰۰۵ تا ۲۰۱۴ در کشور نروژ (ساگبرگ^۱، ۲۰۱۸)

۲-۲-۴- پویایی ترافیک

شرایط ترافیکی راه (حجم ترافیک، سرعت مجاز، ظرفیت راه و دیگر مشخصات ترافیکی) از عوامل مهم بر تصادفات است. اوه و همکارانش (اوه و همکاران، ۲۰۰۵) در مطالعه‌ای عوامل تأثیرگذار بر تصادفات در آزادراه را با استفاده از داده‌های زمان واقعی^۲ (داده‌های برخط) از شرایط ترافیکی به‌عنوان مثال کاهش ظرفیت راه و تغییرات سرعت در مسیر تونل یا عوارضی و همچنین شرایط ترافیک بلندمدت مانند میانگین ترافیک روزانه برای هر سال را در مدل‌سازی و موردبررسی قراردادند. نتایج کار آن‌ها نشان داد که کاهش تغییرات سرعت در مسیر باعث کاهش تصادفات می‌شود.

۲-۳- روش‌های خوشه‌بندی تصادفات

در این قسمت ابتدا به اهمیت داده‌کاوی^۳ در مطالعات ترافیکی می‌پردازیم. سپس به بررسی خوشه‌بندی که زیرشاخه‌ی داده‌کاوی است می‌پردازیم و روش‌های استفاده‌شده در کارهای گذشته را بررسی می‌کنیم.

¹ Sagberg

² Real-time

³ Data mining

۲-۳-۱- داده کاوی

روش داده کاوی به بررسی و تحلیل داده‌ها به منظور پیدا کردن الگوها و قوانین پنهان و معنادار درون داده‌ها اطلاق می‌گردد. انتخاب مجموعه‌ای از تکنیک‌ها برای به کارگیری در موقعیت خاص بستگی به سه عامل ماهیت عمل داده کاوی، ماهیت داده‌های موجود و مهارت‌ها و ترجیحات داده کاوان دارد (پرابها^۱ و شانواز^۲، ۲۰۱۴). روش‌های داده کاوی به طور گسترده‌ای در زمینه‌های مختلف از جمله ایمنی ترافیک مورد استفاده قرار می‌گیرند. کشف دانش^۳ در پایگاه داده موجب دستیابی آسان‌تر به اطلاعات مفید و الگوهای موجود در داده‌ها می‌شود (دونهام، ۲۰۰۶).

داده کاوی در دو نوع هدایت شده^۴ و غیر هدایت شده^۵ ظاهر می‌شود. داده کاوی هدایت شده، دارای متغیر هدفی خاص و از پیش تعیین شده است که به دنبال الگویی خاص می‌گردد. در حالی که هدف داده کاوی غیر هدایت شده، یافتن الگوها یا تشابهات بین گروه‌هایی از اطلاعات، بدون داشتن متغیر هدفی خاص و یا مجموعه‌ای از دسته‌ها و الگوهای از پیش تعیین شده است. برای تبدیل مسئله مورد نظر به مسئله داده کاوی باید آن را به یکی از شش فعالیت داده کاوی (دسته بندی^۶، تخمین^۷، پیش بینی^۸، گروه بندی شباهت^۹، خوشه بندی^{۱۰} و توصیف و نمایه سازی^{۱۱})، تبدیل نمود (شهرابی، ۱۳۹۲).

تعدادی از روش‌های داده کاوی که به طور مؤثر در زمینه ایمنی ترافیک استفاده شده‌اند، عبارت‌اند

از:

¹ Prabha

² Shanavas

³ Knowledge Discovery

⁴ Supervised

⁵ Unsupervised

⁶ Classification

⁷ Estimation

⁸ Prediction

⁹ Affinity Grouping

¹⁰ Clustering

¹¹ Profiling

شبکه عصبی مصنوعی

در پژوهشی کنت^۱ و آقایان در سال ۲۰۱۱ با استفاده از الگوریتم ژنتیک (GA)، روش جستجوی مدل و شبکه‌های عصبی مصنوعی (ANN) شدت تصادفات آژادراه تهران-قم با دوازده پارامتر مربوط به تصادفات در سال ۲۰۰۷ را مدل‌سازی کردند. این مدل‌ها با استفاده از پارامترهای ورودی رانندگی سن و جنس، استفاده از کمر بند ایمنی، نوع و ایمنی وسیله نقلیه، شرایط آب و هوایی، سطح جاده، نسبت سرعت، زمان سقوط، نوع سقوط، نوع برخورد و جریان ترافیکی توسعه داده شد. نتایج نشان داد GA برای ارزیابی الگوریتم‌ها بهترین گزینه را به دست آورد. سپس روش‌های GA و PS با استفاده از بهترین معادله GA ترکیب شدند. شبکه عصبی با معماری چندلایه پیش پردازنده (MLP) استفاده شده است که شامل یک شبکه چندلایه‌ای خوراکی برای نورون‌های خروجی سیگموئید^۲ و خروجی خطی است که می‌تواند به راحتی با نقشه‌های چندبعدی هم سازگار باشد. ANN در طول آموزش، تست و اعتبار سنجی اعمال شد و ۱۲ ورودی، ۲۵ نورون‌های لایه پنهان و ۳ نورون در خروجی وی داشت. مدل مناسب با توجه به R-value، RMSE، میانگین خطاهای مطلق (MAE) و مجموع خطای مربع (SSE) انتخاب شد. بالاترین R-value برای ANN در حدود ۰٫۸۷ به دست آمد که نشان می‌دهد ANN بهترین پیش‌بینی را ارائه می‌دهد. ترکیبی از روش‌های GA و PS برای رتبه‌بندی پیش‌بینی‌های مختلف از روابط خطی به معادلات پیچیده را اجازه می‌دهد. مزیت این مدل‌ها این است که با اضافه کردن داده‌های جدید می‌تواند بهبود یابد (کنت و همکاران، ۲۰۱۱).

در مطالعه‌ای دلن^۳ و همکارانش با استفاده از یک سری از شبکه‌های عصبی ریاضی به مدل‌سازی روابط بالقوه غیرخطی بین سطوح شدت آسیب‌دیدگی و عوامل مربوط به تصادف پرداختند. سپس آنالیز حساسیت بر روی مدل شبکه‌های عصبی آموزش دیده را برای شناسایی اولویت اهمیت عوامل مرتبط با

¹ Kunt

² Sigmoid

³ Delen

تصادف انجام دادند. در این فرآیند، مشکل پیش‌بینی پنج کلاسه به مجموعه‌ای از مدل‌های پیش‌بینی باینری تجزیه می‌شود (با استفاده از یک نمونه آماری ملی از ۳۰۳۵۸ گزارش تصادف ثبت‌شده پلیس آمریکا) تا جزئیات اطلاعات موردنیاز برای شناسایی علت و معلول واقعی روابط بین عوامل مرتبط با تصادف و سطوح مختلف شدت جراحت را به دست آید. نتایج حاصل تغییر فاکتورهای مؤثر بر تصادفات را با تغییر شدت‌های مختلف جراحت نشان داد. (دلن و همکاران، ۲۰۰۶)

چیمبا^۱ و ساندو (۲۰۰۹) با استفاده از شبکه عصبی مصنوعی از روش برگشتی^۲ برای پیش‌بینی شدت جراحت در تصادفات استفاده کردند. همچنین دقت این روش را با روش پروبیت ترتیبی مقایسه کردند. نتایج نشان داد که روش شبکه عصبی مصنوعی با نوروها و گره‌های بهینه‌شده بهتر از روش پروبیت ترتیبی رایج می‌باشد.

در مطالعه‌ای دیگر با استفاده از یک سری از شبکه‌های عصبی مصنوعی مدل برآورد شدت تصادف‌ها و شناسایی عوامل مربوط به تصادف در بزرگراه‌های شهری پرداخته شده است. اعمال شبکه‌های عصبی مصنوعی در علوم مهندسی در سال‌های اخیر ثابت شده است. نتایج به دست آمده نشان می‌دهد که متغیرهایی مانند عرض بزرگراه، برخورد از روبرو، نوع وسیله نقلیه در تصادف، نادیده گرفتن ترمز جانبی، فاصله پس از آن، عدم کنترل خودرو، نقض سرعت مجاز و انحراف در سمت چپ توسط رانندگان، مهم‌ترین عوامل افزایش تصادفات در بزرگراه‌های شهری است. (مقدم و همکاران، ۱۳۸۹)

درخت‌های رگرسیونی و طبقه‌بندی یا درخت‌های تصمیم^۳

طبقه‌بندی و رگرسیون دو موضوع مهم در آمار و داده‌کاوی است. هر یک از این دو روش یک متغیر پاسخ را با توجه به مقادیر بردار متغیرهای پیش‌بینی x به دست می‌آورد. فرض کنید X دامنه x

¹ Chimba & Sando

² Backpropagation

³ Decision Tree

و Y را دامنه y را مشخص کند. اگر y یک متغیر گسسته باشد، مانند مقادیر واقعی (مثلا وزن یک ماشین یا تعداد تصادفات)، مسئله رگرسیون است. در غیر این صورت، اگر Y یک مجموعه نامحدود از مقادیر غیر ارادی (مثلا نوع ماشین یا کشور مبدأ آن) باشد، مشکل طبقه‌بندی نامیده می‌شود (لو، ۲۰۰۸).

درختان تصمیم‌گیری، یا طبقه‌بندی و یا رگرسیون درختی، برای سه دلیل اصلی مورد توجه است. اولاً آنها نمایشی بصری دارند، در نتیجه مدل برای افراد به آسانی قابل فهم و درک است. دوم، درخت تصمیم‌گیری، مدل‌های غیر پارامتری است، هیچ مداخله‌ای از کاربر لازم نیست و به همین دلیل آنها بسیار مناسب برای کشف دانش اکتشافی هستند. سوم، الگوریتم‌های مقیاس‌پذیر، به این معنا که با توجه به افزایش اندازه داده‌های آموزشی^۱ عملکرد ضعیف است، برای مدل‌های ساخت درخت تصمیم‌گیری وجود دارد. آخر، دقت روش درخت تصمیم‌گیری قابل مقایسه یا برتر از سایر مدل‌ها است (دوبرا، ۲۰۰۲).

چانگ و چن^۲ از طبقه‌بندی و رگرسیون درخت^۳، به‌عنوان یکی از گسترده‌ترین تکنیک‌های داده‌کاوی، در مدیریت کسب‌وکار، صنعت و مهندسی استفاده کردند. در این مطالعه اطلاعات مربوط به تصادفات ۲۰۰۱-۲۰۰۲ آزادراه ملی تایوان را جمع‌آوری کردند. مدل CART و مدل رگرسیون دوجانبه منفی برای ایجاد روابط تجربی بین حوادث ترافیکی و متغیرهای هندسی بزرگراه، خصوصیات ترافیک و عوامل محیطی تهیه شد. نتایج CART نشان داد که حجم ترافیک روزانه و متغیرهای بارش باران در محدوده زمانی معین، عوامل کلیدی برای فراوانی حوادث آزادراه بودند. با مقایسه عملکرد پیش‌بینی بین CART و مدل‌های رگرسیون منفی دوگانه منفی، این مطالعه نشان می‌دهد که CART یک روش جایگزین مناسب برای تجزیه و تحلیل فراوانی تصادفات آزادراه است. با مقایسه عملکرد پیش‌بینی بین

¹ Training data

² Chang and Chen

³ Classification And Regression Tree (CART)

CART و مدل‌های رگرسیون دوجانبه منفی، این مطالعه نشان می‌دهد که CART یک روش مناسب برای تجزیه و تحلیل فراوانی تصادفات آزادراه است (چانگ و چن، ۲۰۰۵).

- یادگیری قوانین وابستگی^۱

گرتز^۲ و همکارانش در سال ۲۰۰۳ با استفاده از الگوریتم یادگیری قوانین وابستگی تصادفات مربوط به سال ۱۹۹۹ مربوط در بلژیک را مورد بررسی قرار دادند. این الگوها برای تحلیل و مقایسه با ویژگی‌های تصادفات کم‌رخداد نیز استفاده شد. از مزایای این الگوریتم، شناسایی متغیرهای مرتبط است که سهم قوی در درک بهتر شرایط حادثه و تشخیص الگوهای تصادفات توصیفی را از شرایط تصادفات بیشتر به نقاط سیاه و مناطق سیاه نشان می‌دهد. استفاده از این الگوریتم داده کاوی به ویژه در زمینه مجموعه داده‌های انبوه همچون تصادفات جاده‌ای مفید است، زیرا داده کاوی می‌تواند با استخراج اطلاعات از مقادیر زیادی از داده‌ها، توصیف شود. نتایج نشان داد که جنبه‌های رفتاری رانندگان در تحلیل الگوهای تصادف نمایان می‌شود. به طور کلی این عوامل نقش مهمی را در شناسایی مشکلات ایمنی ترافیک بازی می‌کنند (گرتز و همکاران، ۲۰۰۳).

- شبکه‌های کوهن^۳

تی سنگ و همکارانش در پژوهشی تکنیک‌های داده کاوی را برای پیدا کردن رابطه‌ی بین عدم توجه راننده و وسایل نقلیه موتوری به کار بردند. داده‌های مورد استفاده در این تحقیق از سیستم گزارش تحلیلی فوتی‌ها در اداره ایمنی حمل و نقل بزرگراه‌های ملی، در سالهای ۲۰۰۰ تا ۲۰۰۳ در مریلند^۴ و واشنگتن^۵، به دست آمده است. داده‌ها با استفاده از شبکه‌های کوهنون جمع شدند. سپس الگوهای قواعد داده‌ها توسط درخت تصمیم‌گیری و مدل‌های شبکه عصبی بررسی شد. نتایج نشان داد، زمانی

¹ Association Rule Learning

² Geurts

³ Kohonen networks

⁴ Meriland

⁵ Washington

که عدم توجه و شرایط فیزیکی-ذهنی راننده در یک زمان اتفاق بیافتد، راننده گرایش بالایی به تصادف با اشیاء ثابت دارد. علاوه بر این، با توجه به نوع برخورد تصادف، اهمیت نسبی تصادف با یک وسیله نقلیه متحرک نسبت به برخورد با یک جسم ثابت دو برابر بیشتر در هر تصادف است(تی سنگ و همکاران، ۲۰۰۵).

۲-۳-۲- خوشه‌بندی

خوشه‌بندی به عمل تقسیم جمعیت ناهمگن به تعدادی از زیرمجموعه‌ها یا خوشه‌های همگن گفته می‌شود. در خوشه‌بندی داده‌ها صرفاً بر اساس تشابه گروه‌بندی می‌شوند و هیچ دسته از پیش تعیین‌شده‌ای وجود ندارد و هر گروه نیز توسط کاربر نام‌گذاری می‌شود. معمولاً از خوشه‌بندی به‌عنوان پیش‌زمینه‌ای در تحلیل‌های داده‌کاوی یا مدل‌سازی استفاده می‌شود.

در این قسمت به روش‌های خوشه‌بندی نقاط تصادف و تحلیل‌های خوشه‌ای تصادفات در کارهای گذشته می‌پردازیم. علاوه بر آن، دو روش رایج خوشه‌بندی، روش‌های غیر آماری و آماری، مورد مطالعه قرار گرفته‌اند. روش غیر آماری شامل خوشه‌بندی‌های تفکیکی^۱، مبتنی بر چگالی^۲ و مبتنی بر شبکه بندی^۳ و در روش خوشه‌بندی آماری نیز چگالی کرنل^۴ در دو سطح شبکه^۵ و صفحه‌ای^۶ و توابع خودهمبستگی مکانی^۷ بررسی شده‌اند. به‌منظور یافتن خوشه‌هایی که بهترین تناسب را با داده‌های موردنظر داشته باشند، درنهایت شاخص‌های ارزیابی ارائه می‌گردند.

در روش تحلیل خوشه‌ای که یک روش تحلیل چندمتغیری است که شامل اندازه‌های هم‌زمان روی چند متغیرند. کاهش داده‌ها یا آسان‌سازی ساختاری و دسته‌بندی کردن داده‌ها، از جمله اهدافی

¹ Partitioning

² Density_based

³ Grid-based

⁴ Kernel Density

⁵ Network

⁶ Planar

⁷ Spatial Autocorrelation

می‌باشند که در روش‌های چندمتغیری و خوشه‌بندی دنبال می‌شوند. تحلیل خوشه‌ای هم‌چنین به‌عنوان طبقه‌بندی اشیاء مشابه تعریف شده است به‌طوری‌که در آن، شکل تعدادی از گروه‌ها ناشناخته هستند (روستو و کافمن^۱، ۱۹۹۰). به‌طور کلی تحلیل خوشه‌ای از دو جزء اصلی تشکیل شده است که عبارت‌اند از اندازه‌گیری شباهت (فاصله) و الگوریتم خوشه‌بندی. هدف اصلی از تحلیل خوشه‌ای سازمان‌دهی اقلام به گروه‌هایی است که درجه شباهت برای اقلام در یک گروه حداکثر و برای آن‌هایی که بین دو گروه هستند حداقل است. علاوه بر این، روش فوق‌زمانی استفاده می‌شود که محققان به‌جای آزمون فرضیه-های علت و معلولی در مجموعه داده، به ویژگی‌های منحصر به فرد اقلام، علاقه‌مند باشند.

هدف از آنالیز خوشه‌ای تصادف‌ها، شناسایی گروه‌های همسان تصادف در مطالعاتی است که به‌منظور شناسایی راه‌کارها و اقدامات طبقه‌بندی شده کاهش تصادف در سطح کلان انجام می‌پذیرد. با توجه به اینکه این گروه‌ها خصوصیات مشترک بیشتری دارند، می‌توان ایده‌های کاهش تصادف را در آن‌ها به وجود آورد و راه‌کارهای رسیدن به این هدف را تدوین کرد. آنالیز خوشه‌ای عموماً روشی است که پیش از آنکه هدف‌گرا باشد، موضوع‌گرا است و می‌تواند به‌عنوان یک روش مناسب به کار رود.

۲-۳-۳- الگوریتم‌های خوشه‌بندی غیر آماری

تفاوت خوشه‌بندی آماری و خوشه‌بندی غیر آماری در این است که روش‌های خوشه‌بندی آماری خودهمبستگی بین پارامترهای مؤثر در تصادفات را مشخص می‌کند، درحالی‌که روش‌های غیر آماری بر اساس تابع فاصله هستند به عبارتی شباهت بین داده‌ها بر اساس فاصله آن‌ها تعیین می‌شود. چهار نوع طبقه‌بندی رایج در روش‌های خوشه‌بندی غیر آماری وجود دارد: (هان و همکاران، ۲۰۱۱)

¹ Rousseeuw and Kaufman

- الگوریتم‌های خوشه‌بندی افزایشی^۱
- خوشه بندی سلسله مراتبی^۲
- خوشه‌بندی مبتنی بر تراکم^۳
- الگوریتم‌های مبتنی بر شبکه^۴

جدول ۱-۲: طبقه بندی الگوریتم‌های خوشه‌بندی غیر آماری (هان و همکاران، ۲۰۱۱)

روش خوشه‌بندی	خصوصیات عمومی
افزایی	<ul style="list-style-type: none"> - یافتن خوشه های منحصر به فرد از شکل کروی - مبتنی بر فاصله - استفاده از مدل Means یا Medoide (و غیره) برای نشان دادن مرکز خوشه - مؤثر برای مجموعه داده های کوچک و متوسط - تعداد خوشه‌ی از پیش تعیین شده
سلسله مراتبی	<ul style="list-style-type: none"> - خوشه بندی در این نوع یک تجزیه سلسله مراتبی است. - نمی توان اشتباهات ادغام شده و یا تقسیم شده را تصحیح کرد. - ممکن است تکنیک های دیگری مانند میکرو خوشه بندی را در نظر بگیریم یا ارتباطات شی را در نظر بگیریم.
مبتنی بر تراکم	<ul style="list-style-type: none"> - می تواند خوشه های با شکل دلخواه پیدا کند - خوشه ها، مناطق متراکم از اشیاء در فضا هستند که با مناطق کم تراکم از هم جدا می شوند - چگالی خوشه: هر نقطه باید حداقل تعداد نقاط در داخل آن همسایگی داشته باشد. - داده‌هایی که معیارهای خوشه‌بندی را ندارند داده‌های خارج از محدوده^۵ به حساب می‌آیند.
مبتنی بر شبکه	<ul style="list-style-type: none"> - از یک ساختار شبکه چندتایی داده استفاده می‌کنند. - زمان پردازش سریع (به طور معمول مستقل از تعداد اشیاء داده، در عین حال وابسته به اندازه شبکه)

¹ Partitioning
² Hierarchical
³ Density-Based
⁴ Grid-Based
⁵ Outlier

خوشه‌بندی افرازی

الگوریتم k -means و الگوریتم پیشینه‌سازی امید ریاضی^۱ (EM)، نمایندگان اصلی برای الگوریتم‌های خوشه بندی افرازی می‌باشند. الگوریتم k -means نشان‌دهنده یک خوشه با استفاده از یک نقطه واحد است، در حالی که EM از توزیع احتمالی برای نشان دادن یک خوشه استفاده می‌کند. الگوریتم‌های تفکیکی نیاز به ارائه پارامتر، k ، تعداد خوشه‌ها و انجام مکانیابی مجدد عضو بودن داده را دارند تا زمانی که عضویت دیگر تغییر نکند یا تغییر در محدوده قابل تحمل باشد. کیفیت خوشه بندی به مقدار k وابسته است و به طور کلی، الگوریتم‌های تفکیکی نمی‌توانند به صورت خوشه ای با شکل‌های دلخواه درآیند.

در الگوریتم افرازی، همه خوشه‌ها به یک‌باره مشخص می‌شوند؛ به عبارتی روش افرازی، داده‌ها را به تعداد خوشه‌های از پیش تعیین‌شده (k) از گروه‌های دوبه‌دو ناسازگار و جامع، تقسیم‌بندی می‌کند. (Han et al., 2011) در خوشه‌بندی افرازی داده‌ها به گروه‌هایی تقسیم‌بندی می‌شوند که هر گروه نشان‌دهنده یک خوشه است؛ به عبارت دیگر این روش تقسیم مجموعه‌ای از داده‌ها به زیرمجموعه غیر همپوشان^۲ است به گونه‌ای که هر شیء داده دقیقاً در یک زیرمجموعه است؛ بنابراین تعداد خوشه‌ها از قبل تعیین شده است (Kaufman and Rousseeuw, 1990).

الگوریتم K-Means

این روش خوشه‌بندی یکی از مشهورترین و ساده‌ترین رویکردهای خوشه‌بندی است. مک کوئین^۳ در سال ۱۹۶۷ اولین شخصی بود که از عبارت K-Means استفاده کرد (مک کوئین، ۱۹۶۷). هدف، تقسیم بندی n مشاهده (اشیاء) به k خوشه است به طوری که هر مشاهده متعلق به خوشه با نزدیک‌ترین میانگین

¹ Expectation maximization

² Non-overlapping

³ MacQueen

است. به‌طور کلی الگوریتم فوق تکرارشونده است و معمولاً پس از چند بار تکرار برای بهبود محلی، همگرا می‌شود. در این خوشه‌بندی، تعداد خوشه‌ها در ابتدا انتخاب می‌شوند و الگوریتم صرفاً هر شیء را به خوشه‌های از پیش تعیین‌شده، اختصاص می‌دهد. الگوریتم با اولین انتخاب تصادفی k موقعیت برای مرکز ثقل خوشه، کار می‌کند. هر مشاهده به خوشه‌ای اختصاص داده می‌شود که نزدیک‌ترین مرکز ثقل را دارد و ثقل با استفاده از مقدار میانگین مقادیر اختصاص داده‌شده، محاسبه می‌شود. سپس الگوریتم این فرآیند را تا زمانی که مرکز جرم خوشه‌ها تغییر نکند یا تغییراتی کمتر از حد آستانه داده‌شده داشته باشند، تکرار می‌کند. بدین ترتیب روش K-Means، عددی، نظارت‌نشده، غیرقطعی و بسیار تکرارشونده است. در زیر مراحل مختلف الگوریتم K-Means آورده شده است (Bradley et al., 1997)

- برای تعریف مرکز k خوشه، با k نقطه انتخابی تصادفی شروع می‌شود
- هر آیتم به نزدیک‌ترین نقطه اختصاص داده می‌شود
- مرکز (میانگین) هر خوشه محاسبه می‌شود
- از k میانگین برای تعریف مرکز k خوشه جدید استفاده می‌کنیم و هر آیتم در خوشه با نزدیک‌ترین مرکز جابه‌جا می‌شود.
- دو مرحله قبل تا زمانی که هیچ تغییری در ماهیت خوشه‌ها بین مراحل وجود نداشته باشد، تکرار می‌شود.

در الگوریتم K-Means همواره k خوشه وجود دارد و همواره حداقل یک مقدار در هر خوشه وجود دارد. خوشه‌ها افزای هستند و همپوشانی ندارند. هر عضوی از یک خوشه به خوشه خود نزدیک‌تر است تا به خوشه دیگر، چراکه نزدیکی همواره مرکز خوشه‌ها را درگیر نمی‌کند. از مزایای این روش می‌توان به خوشه‌بندی نسبتاً ساده اشاره کرد. زمانی که مقادیر زیاد و تعداد خوشه‌های کوچکی وجود دارند، خوشه‌بندی K-Means ممکن است از لحاظ محاسباتی سریع‌تر باشد. علاوه بر این خوشه‌بندی K-Means تمایل به تولید خوشه‌های فشرده‌تر دارد به‌خصوص اگر از لحاظ شکل گروهی باشند. با این

معنی که هر خطی که بین دو عضو یا اشیاء خوشه کشیده می‌شود درون مرزهای خوشه باقی می‌ماند. از معایب این الگوریتم نیز اینست که برای بررسی خوشه‌های غیر کروی و خوشه‌ها با اندازه‌ها و تراکم‌های مختلف، مناسب نیست. خوشه‌های کروی معمولاً مرز بسیار پیچیده دارند. این روش خوشه‌بندی اشیاء را به k کلاس انحصاری از قبل تعیین شده اختصاص می‌دهد. در حقیقت اشیاء نمی‌توانند به بیش از یک خوشه با درجات مختلفی از احتمال عضویت، تعلق داشته باشند. به‌علاوه، این روش در خوشه‌بندی داده‌هایی که شامل نقاط نویز است، خیلی مؤثر نیست و از این رو تشخیص نقاط نویز و حذف آن‌ها الزامی است.

کاربرد الگوریتم K-Means در تصادفات

(۱) اندرسون^۱ در سال ۲۰۰۹ با استفاده از سیستم اطلاعات جغرافیایی^۲ (GIS) و تخمین تراکم کرنل^۳ به مطالعه الگوهای مکانی تصادفات جراحی سال ۱۹۹۹ تا ۲۰۰۳ در لندن پرداخت. همچنین با استفاده از داده‌های محیطی تصادف (شرایط جاده، شرایط راننده و...) و نتایج تخمین تراکم کرنل به منظور طبقه‌بندی نقاط حادثه‌خیز با الگوریتم K-Means خوشه‌بندی نمود. پس از خوشه‌بندی داده‌ها به ۵ گروه و ۱۵ خوشه تقسیم بندی شدند و هر کدام از خوشه‌ها با توجه به داده‌های محیطی توصیف شدند (اندرسون، ۲۰۰۹).

(۲) در تحقیق سینگ چانگ و لی یین^۴ (۲۰۱۱) با استفاده از تکنیک مدل‌سازی K-Means به کشف تأثیرات عوامل غیررفتاری که شامل ویژگی‌های هندسی بزرگراه، عوامل ترافیکی که شامل حجم ترافیک در روز و شرایط محیطی که شامل بارش سالیانه می‌باشد، پرداخته شده

¹ Anderson

² Geographic Information System

³ Kernel density estimation

⁴ Hsing-chung & Li-yen Chang

است. نتایج پژوهش نشان می‌دهد عوامل فوق که عوامل خطر نامیده می‌شوند، در بروز شدت تصادفات نقش مؤثری دارند.

(۳) طی تحقیقی در سال ۲۰۱۲ توسط رساک و دیوید^۱ با استفاده از الگوریتم K-Means به منظور تشخیص نقش عوامل انسانی در بروز و شدت برخوردها و تصادفات جاده‌ای پرداخته شده است. نتایج پژوهش نشان می‌دهد که بستن کمربند ایمنی، داشتن گواهینامه، سن و جنس در بروز تصادفات نقش مؤثری دارند.

خوشه بندی سلسله مراتبی^۲

برعکس خوشه‌بندی افرازی که اشیاء را در گروه‌های مجزا تقسیم می‌کند، خوشه‌بندی سلسله مراتبی، در هر سطح از فاصله، نتیجه خوشه‌بندی را نشان می‌دهد. این سطوح به صورت سلسله مراتبی هستند. خوشه بندی سلسله مراتبی تکنیک یادگیری هدایت‌نشده^۳ محبوبی است که به دنبال ایجاد سلسله مراتب در خوشه‌ها است. این خوشه‌بندی به طور معمول به دو دسته زنجیره‌ای^۴ و منفصل^۵ تقسیم می‌شود. خوشه‌بندی زنجیره‌ای به دنبال یک رویکرد پایین به بالا، یعنی هر شیء داده در خوشه خود شروع می‌شود و اشیاء بسته بیشتر با هم ادغام می‌شوند و خوشه جدید را تشکیل می‌دهند. این فرایند تکرار می‌شود تا خوشه‌ای برای ادغام باقی نگذاشته شود. برخلاف خوشه بندی زنجیره‌ای، خوشه بندی منفصل رویکرد بالا به پایین است که در آن تمام اشیاء داده در یک خوشه شروع می‌شوند و تقسیم در میان اشیاء داده ادامه می‌یابد تا هر اشیاء داده به یک خوشه واحد تعلق داشته باشند. با این حال، خوشه

¹ Rasak & david

² Hierarchical

³ Unsupervised

⁴ Agglomerative

⁵ Divisive

بندی زنجیره‌ای از نظر محاسباتی کارآمدتر از الگوریتم خوشه بندی تقسیم شده است (کومار و توشینوال^۱، ۲۰۱۶) (تان و همکاران^۲، ۲۰۰۵)

خوشه‌بندی سلسله مراتبی نزدیک‌ترین همسایگی^۳ Nnh

بر اساس این روش، تفاوت بین دو خوشه حداقل فاصله‌ی بین اجسام درون هر خوشه است. فاصله در اینجا به معنی حداقل فاصله هر شیء تا شیء $d(x_i, y_j)$ محاسبه می‌شود که x_i به خوشه اول تعلق دارد، و y_j متعلق به خوشه دوم است. به عبارت دیگر، این روش حداکثر شباهت بین دو گروه از اشیاء یا خوشه را محاسبه می‌کند. این روش زنجیره‌های طولانی را ایجاد می‌کند که خوشه‌های ضعیف را تشکیل می‌دهند. این تکنیک در برخورد با اشکال غیر بیضوی مفید است، اما بسیار حساس به داده‌های نویز و ناخوشه است. فاصله ریاضی، $D(X, Y)$ بین خوشه X و Y به صورت زیر بیان می‌شود:

$$D(X, Y) = \min d(x, y), \quad x \in X, y \in Y \quad (1-2)$$

که در آن $d(x, y)$ فاصله بین دو نقطه x و y است و X و Y دو خوشه هستند.

روش خوشه‌بندی سلسله مراتبی نزدیک‌ترین همسایگی^۴ (Nnh)، گروهی از رخدادها را که از نظر مکانی نزدیک هستند، شناسایی می‌کند. این خوشه‌ها بر مبنای یک معیار برآورده می‌شوند. خوشه بندی تکرار می‌شود تا زمانی که تمام نقاط به یک خوشه تقسیم می‌شوند و یا معیار خوشه ناکام می‌ماند. روش تحلیل خوشه‌ای Nnh برای تشخیص نقاط حادثه‌خیز تصادف استفاده می‌شود که بر شناسایی گروه‌هایی از داده‌ها که به صورت مکانی نزدیک هستند، متمرکز می‌شود. در خوشه بندی Nnh، فاصله اقلیدسی بین هر جفت نقطه داده بررسی و به عنوان یک معیار برای خوشه بندی استفاده

¹ Kumar & Toshniwal

² Tan et al.

³ Nerast neighbourhood Hierarchical

⁴ Nearest Neighborhood Hierarchical

می شود. معیار اول فاصله جستجو (d) است که پس از تعیین آن، نقاط با فاصله‌های کوچکتر در یک خوشه قرار می‌گیرند. معیار دوم حداقل تعداد نقاط n_{min} در خوشه می‌باشد. سپس نقاط منطبق با معیارهای فاصله جستجو (d) و n_{min} به عنوان یک خوشه شناسایی می‌شود. اگر تعداد نقاط تعریف شده نزدیکتر از فاصله جستجو باشد، یک خوشه جدید تولید می‌شود. پس از محاسبه خوشه‌های مرتبه اول، خوشه‌های مرتبه دوم و مرتبه بالا با همان شیوه تشکیل می‌شوند تا زمانی که فقط یک خوشه باقی بماند یا با معیارهای مذکور تطابق نداشته باشد. به همین دلیل، الگوریتم خوشه بندی N_{nh} خوشه‌ای از تمام نقاط در منطقه مورد مطالعه ایجاد نمی‌کند. از آنجایی که این روش به معیارهای مربوط به فاصله جستجو (d) و n_{min} بستگی دارد، تنها نقاطی که معیارهای مورد نیاز را برآورده می‌کنند، خوشه بندی می‌شوند (لوین، ۲۰۱۰).

خوشه‌بندی مبتنی بر تراکم^۱

استر^۲ و همکاران در سال ۱۹۹۶ برای اولین بار روش خوشه‌بندی مبتنی بر تراکم را پیشنهاد دادند. این روش‌های خوشه‌بندی بر این اصل استوارند که خوشه‌ها، ناحیه‌هایی از فضای داده با چگالی زیادی هستند که توسط نواحی با چگالی کمتر از همدیگر جدا شده‌اند. این الگوریتم‌ها هر شی را با مقدار چگالی تعیین شده توسط تعداد اجسام همسایه خود در یک شعاع داده شده مرتبط می‌کند. یک جسم که تراکم آن بیشتر از آستانه تعریف شده باشد، به عنوان یک شیء متراکم تعریف می‌شود و در ابتدا خوشه‌ای شکل می‌گیرد. الگوریتم‌های خوشه بندی مبتنی بر چگالی عبارت‌اند از DBSCAN، OPTICS، HOP و DENCLUE. این روش می‌تواند نویز (خروجی‌ها) را جدا کند، خوشه‌های شکل دلخواه را پیدا کند و هیچ فرضیه‌ای در مورد توزیع داده‌های پایه ایجاد نمی‌کند. با این حال، از لحاظ

¹ Density-Based clustering

² Ester et al.

محاسباتی بسیار سخت هستند؛ به خصوص در مراحل تولید تراکم داده ها و جستجو برای همسایگان متراکم (استر و همکاران، ۱۹۹۶).

الگوریتم DBSCAN

از جمله الگوریتم‌های مبتنی بر تراکم، DBSCAN است. این روش خوشه‌بندی بر این اصل استوار است که خوشه‌ها، ناحیه‌هایی از دو فضای داده با تراکم زیاد هستند که توسط نواحی با تراکم کمتر از همدیگر جدا شده‌اند. مفهوم اصلی، تعریف خوشه جدید یا توسعه خوشه موجود بر اساس همسایگی (Eps) است. پارامتر Eps، بیشترین شعاع همسایگی و پارامتر MinPts، کمترین تعداد نقاط در شعاع همسایگی نقطه هستند. با توجه به مجموعه داده‌ی D، تابع توزیع فاصله و پارامترهای Eps و MinPts، تعاریف زیر برای تعریف DBSCAN استفاده می‌شود:

- هسته اصلی: در همسایگی شعاع ϵ حاوی حداقل اجزای minpts باشد.
- نقطه مرزی: همسایگی یک جسم خاص کمتر از minpts باشد
- نقاط نویزی: همسایگی یک جسم خاص شامل ۰ است.
- تراکم قابل دسترس مستقیم: اگر شی p یک شی هسته‌ای است، آنگاه شی q در همسایگی p قرار دارد، سپس q و p را تراکم قابل دسترس مستقیم می‌نامند.
- تراکم قابل دسترس: پیوند داده‌ای بین p_1, p_2, \dots, p_n وجود دارد، جایی که $p_n = q, p_1 = p_2$ ، اگر شی p_i و شی p_{i+1} ، چگالی قابل دسترس مستقیم باشد، سپس ما q و p را تراکم قابل دسترس می‌نامیم (شکل ۲-۳ الف).
- اتصال تراکم: اگر جسم o در مجموعه داده وجود داشته باشد که جسم p و q را تراکم قابل دسترس برای o می‌سازد، سپس شیء p و q تراکم قابل دسترس است (شکل ۲-۳ ب)

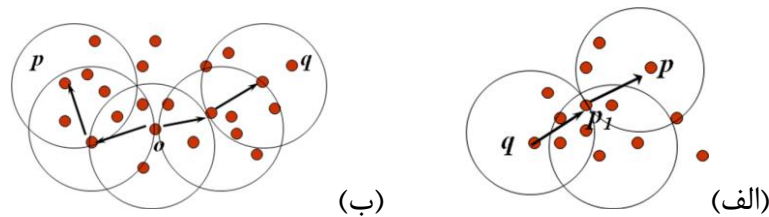
- فاصله هسته: حداقل شعاع همسایگی که یک شی را به یک هسته اصلی به نام فاصله هسته تبدیل می کند. اگر شی تبدیل هسته اصلی نیست، هیچ فاصله هسته ای نیز وجود ندارد.

- فاصله قابل دسترس: فاصله اقلیدس از شی q به هسته p به عنوان فاصله قابل دسترس نامیده می شود، اگر شی p یک هسته اصلی نیست، فاصله ای بین p و q وجود ندارد.

نقطه دلخواه $p(p \in D)$ همسایگی p به صورت رابطه ۲-۲ تعریف می شود:

$$N_{Eps}(p) = \{q \in D | dist(p, q) \leq Eps\} \quad (2-2)$$

اگر $|N_{Eps}(p)| \geq MinPts$ باشد.



شکل ۲-۳: (الف) تراکم قابل دسترس q و p (ب) اتصال تراکم p و q

اگر p نقطه مرکزی و q همسایه p باشد؛ q متعلق به این خوشه است. همه همسایه های q برای اضافه شدن به خوشه بررسی می شوند و در غیراینصورت q با عنوان نویز برچسب دهی می شود. این روند برای هر نقطه در همسایگی تکرار می شود. اگر خوشه ای نتواند گسترش یابد، DBSCAN نقطه لخواه بدون برچسب دیگری را انتخاب کرده، فرایند تکرار می شود. به عبارت دیگر الگوریتم DBSCAN به این صورت عمل می کند که ابتدا نقطه دلخواه x انتخاب می شود سپس همه نقاطی که از سوی x در دسترس چگالی هستند جستجو می شوند؛ اگر x یک نقطه هسته باشد خوشه تشکیل می شود و اگر x یک نقطه مرزی باشد هیچ نقطه ای از طرف x در دسترس چگالی نخواهد بود؛ این مراحل برای نقاط بعدی تا زمانی

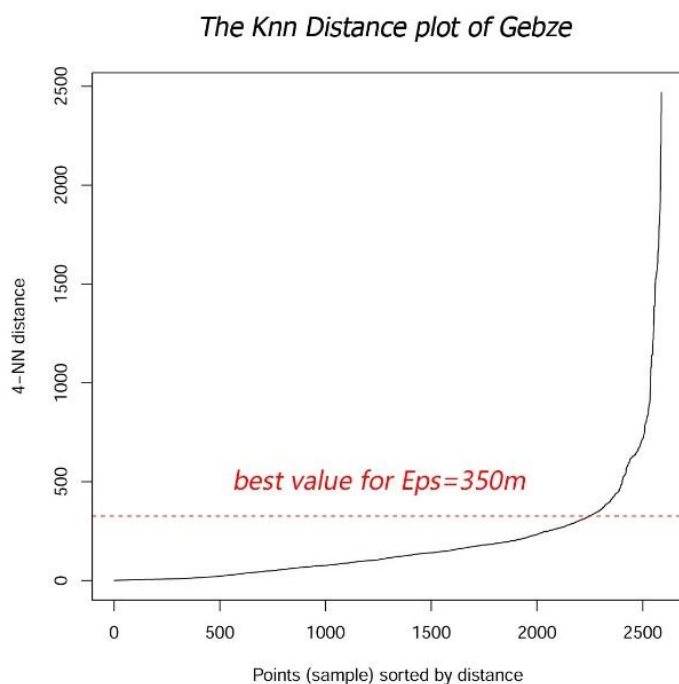
که همه نقاط مجموعه داده در خوشه‌ها قرار گیرند یا به‌عنوان نویز برچسب دهی شوند، ادامه می‌یابد. به‌عبارت‌دیگر روش‌های خوشه‌بندی مبتنی بر چگالی تلاش می‌کنند خوشه‌هایی را بر اساس تراکم نقاط در منطقه پیدا کنند. مناطق متراکمی که از یکدیگر قابل‌دسترس هستند برای تشکیل خوشه باهم ادغام می‌شوند. در این روش خوشه‌بندی، هر داده متعلق به یک خوشه، در دسترس چگالی سایر داده‌های متعلق به آن خوشه است و در دسترس چگالی هیچ داده دیگری قرار ندارد. روش‌های خوشه‌بندی مبتنی بر تراکم خوشه‌ها را با اشکال دلخواه پیدا می‌کنند (استر و همکاران^۱، ۱۹۹۶).

شاخص البو برای تعیین مقدار شعاع *Eps*

تحقیقات استر و همکارانش^۲ در ۱۹۹۶ حاکی از آن بود که روش البو برای تعیین پارامترهای بهینه خوشه‌ها در خوشه‌بندی DBSCAN استفاده می‌شود. این روش بر اساس k نزدیک‌ترین همسایگی، مقدار $MinPts$ کمترین تعداد نقاط در شعاع همسایگی نقطه و Eps را محاسبه می‌کند. پارامتر $MinPts$ برای تشخیص داده‌های اصلی مورد استفاده قرار می‌گیرد. داده اصلی، داده‌ای است که در شعاع همسایگی آن حداقل به تعداد $MinPts$ داده وجود داشته باشد.

¹ Ester et al.

² Ester et al.



شکل ۲-۴: تعیین Eps بهینه با استفاده از نمودار^۱ Knn-distance

شکل ۲-۴ روش البو را نشان می‌دهد که در آن محور افقی تعداد نمونه‌های مرتب شده بر اساس K نقطه همسایگی است و محور عمودی، فاصله K نزدیک‌ترین همسایگی است. بر اساس تحقیقات استر در سال ۱۹۹۶ میلادی برای داده‌های دوبعدی (x,y) می‌توان k را برابر با ۴ در نظر گرفت؛ بنابراین ۴ همسایگی حساب می‌شود و فاصله از تمامی نقاط دیگر به دست می‌آید. در واقع ۴ نقطه‌ای که نزدیک به نقطه مرکز هستند، بررسی می‌شوند. بیشترین فاصله تا ۴ نقطه حساب می‌شود و این کار برای همه نقاط مجموعه داده انجام می‌شود و نموداری رسم می‌شود که در شکل ۲-۴ نشان داده شده است. در نمودار به دست آمده نقطه‌ای که تغییراتش شدید می‌شود (نقطه عطف)، نقطه مورد نظر است که با عمود کردن نقطه عطف به محور عمودی پارامتر Eps دست می‌آید. برای پارامتر Minpts نیز حداقل مقدار در این روش عدد K و بیشتر می‌باشد و به صورت سعی و خطا باید مقداری را در نظر گرفت که کمترین مقدار نویز با حفظ شاخص‌های خوشه بندی بدست آید.

^۱ K nearest neighbourhood distance

مزایا و معایب روش DBSCAN

این الگوریتم نیز همانند دیگر الگوریتم‌ها دارای مزایا و معایبی می‌باشد. از مزایای این الگوریتم اینست که نیازی به مشخص بودن تعداد خوشه‌ها، به صورت اطلاع قبلی ندارد (برخلاف روش‌های خوشه-بندی مبتنی بر مرکز مانند مثل K-Means) می‌تواند خوشه‌های با اشکال مختلف را بیابد. همچنین می‌تواند خوشه‌ای را که کاملاً توسط خوشه‌ی دیگر احاطه شده است (ولی به آن متصل نیست)، شناسایی کند. مشکل تک‌یالی (خوشه‌های مختلف با مسیر نازکی از نقاط به هم متصل باشند) نیز با پارامتر $MinPts$ کاهش یافته است. از مفهوم نویز پشتیبانی می‌کند. به دو پارامتر نیاز دارد Eps و $MinPts$ و نسبت به ترتیب قرار گرفتن نقاط در پایگاه داده حساس نیست. (البته در مورد برخی از نقاط مرزی، این امر صادق نیست. نقاطی که در مرز دو خوشه‌ی متفاوت باشند، با توجه به ترتیب در نظر گرفتن نقاط در خوشه‌بندی، ممکن است در خوشه‌ها جابه‌جا شوند).

از معایب این الگوریتم این است که نمی‌تواند مجموعه داده‌هایی با اختلاف چگالی زیاد را خوشه-بندی کند، به دلیل اینکه مقادیر $MinPts$ و Eps نمی‌توانند برای تمام خوشه‌ها مناسب باشند (برای رفع این مشکل، توسعه DBSCAN-DD برای آن ارائه شده است). برای نمونه، در نواحی تراکم، همسایگی نقاط مختلف اشتراک زیادی باهم دارند. در نتیجه نقاط تکراری زیادی را در محاسبات خود لحاظ می‌کند و باعث افت کارایی این الگوریتم می‌شود برای رفع این مشکل توسعه‌ی IDBSCAN برای آن ارائه شده است.

کیفیت DBSCAN به نوع اندازه‌گیری فاصله‌ی نقاط بستگی دارد. فاصله‌ی اقلیدسی رایج‌ترین نوعی است که استفاده می‌شود. در داده‌های دارای ابعاد بالا، این نوع اندازه‌گیری فاصله بی‌فایده می‌شود. در نتیجه یافتن مقداری برای Eps دشوار می‌شود. در مقایسه با روش‌های سنتی تشخیص تراکم، روش خوشه‌بندی مبتنی بر تراکم به‌طور ذاتی می‌تواند تراکم مکانی را کشف کند. تقسیم‌بندی مجموعه داده

در آغاز این فرایند موردنیاز نیست. همچنین می‌تواند اشکال با تراکم دلخواه را از مجموعه داده پیدا کند. به‌رحال بیشتر روش‌های خوشه‌بندی مبتنی بر تراکم، محدودیت DBSCAN اصلی را به ارث می‌برند که ویژگی‌های غیرمکانی را به حساب نمی‌آورند و نمی‌توانند به‌طور مستقیم برای مجموعه داده تصادف روی شبکه مکانی به کار گرفته شوند (موریرا و همکاران^۱، ۲۰۰۵).

- کاربرد الگوریتم DBSCAN در تصادفات

(۱) دوئرو و سوباشی در سال ۲۰۱۵ در مقاله‌ای چگونگی استفاده از تکنیک‌های داده کاوی در ITS برای تشخیص و پیش‌گیری تصادفات در بزرگراه‌ها را نشان دادند. در زمان ترافیک، خودروها رفتار مشابهی با خودروهای واقع در نواحی محله‌ای بسته را نشان می‌دهند. در مواردی مانند تصادفات، تغییر خط نامناسب و سرعت می‌تواند به عنوان ناهنجاری‌ها تشخیص داده شود. آن‌ها یک تصادف رانندگی شبیه‌سازی کرده و اثربخشی تکنیک‌های خوشه‌بندی متفاوت برای تشخیص تصادف رانندگی موردبررسی قرار دادند. آن‌ها نشان دادند که با دو الگوریتم DBSCAN و خوشه بندی سلسله مراتبی زنجیره‌ای^۲ AHC، اگر مقدار سرعت و موقعیت هر وسیله نقلیه داده شود، رفتار وسایل نقلیه را می‌توان تجزیه و تحلیل کرد و حوادث را می‌توان به راحتی شناسایی کرد. موفقیت الگوریتم‌های پیشنهادی با استفاده از شبیه‌سازی در یک سناریو بزرگراه نشان داده شده است. شبیه‌سازی نشان داد که ابزارهای داده کاوی به طور موفقیت‌آمیزی تصادفات را با نرخ تشخیص متوسط خطر ۱۰۰٪ و میزان اشتباه ۰٪ با استفاده از DBSCAN و خوشه بندی سلسله مراتبی نشان داد.

(۲) العتیبی^۳ (۲۰۱۸) با استفاده از الگوریتم‌های داده کاوی DBSCAN و فرکانس

موازی^۴ به بررسی پارامترهای موثر بر تصادفات پرداخت. ابتدا مناطق تصادف را به تعداد K

¹ Moreira et al.

² Agglomerative hierarchical clustering

³ Alotaibi

⁴ Parallel Frequent

خوشه بر اساس فرکانس تصادف آنها با الگوریتم DBSCAN تقسیم کرد. سپس الگوریتم فرکانس موازی را روی این خوشه ها اعمال کرد، تا ارتباط بین ویژگی های متمایز در داده های تصادفات رانندگی را برای تحلیل آنها در جهت پیش بینی عوامل مختلف که بر تصادفات جاده ای در مکان های مختلف تأثیر می گذارد، نمایان گردد. هدف اصلی برداشت داده های تصادف، شناسایی مسائل کلیدی در زمینه ایمنی جاده ها است. داده های تصادف جاده ای با استفاده از ابزار Weka تحلیل شدند و نتایج نشان داد با استفاده از الگوریتم داده کاوی DBSCAN و PF نشان می دهد که این روش می تواند در داده های تصادفات جدید با ویژگی های اضافی برای شناسایی عوامل مختلف مرتبط با تصادفات جاده ای استفاده شود.

الگوریتم های مبتنی بر شبکه (GRID)

الگوریتم های خوشه بندی مبتنی بر شبکه ابتدا دامنه فضای مسئله را با شبکه مش بندی یکنواخت پوشش می دهد. صفات آماری برای تمام اشیاء داده شده در هر سلول شبکه مش جمع آوری شده و سپس خوشه بندی در شبکه به جای اشیاء داده ها انجام می شود. این الگوریتم ها به طور معمول زمان پردازش سریع تری نسبت به بقیه الگوریتم ها دارند، زیرا از یک «مجموعه داده» برای محاسبه مقادیر آماری شبکه استفاده می کنند و عملکرد خوشه بندی فقط به اندازه شبکه ها بستگی دارد که معمولاً بسیار کمتر از اشیاء داده است. الگوریتم خوشه بندی مبتنی بر شبکه بر پایه STING، WaveCluster و CLIQUE است. تمام این روشها یک شبکه یکپارچه را برای پوشش دادن کل مسئله استفاده می کنند. برای داده های با توزیع بسیار نامنظم، ابعاد شبکه مش باید به اندازه کافی برای به دست آوردن کیفیت خوشه نرم تر تنظیم شود یعنی تعداد شبکه بیشتری در نظر گرفته شود. یک مش نرم تر¹ می تواند موجب

¹ Finer mesh

اندازه مش نزدیک یا حتی بزرگ‌تر از اندازه اشیاء داده برسد که می‌تواند بار محاسباتی را برای خوشه سازی افزایش دهد (مهران و مهرا^۱، ۲۰۰۸).

رویکردهای مبتنی بر تراکم و شبکه برای خوشه‌بندی در یک فضای چندبعدی هستند که در آن خوشه‌ها به‌عنوان مناطق متراکم‌تر از محیط اطراف آن‌ها محسوب می‌شوند. پیچیدگی محاسباتی اکثر الگوریتم‌های خوشه‌بندی حداقل به‌صورت خطی متناسب با اندازه مجموعه داده است. مزیت مهم خوشه‌بندی مبتنی بر شبکه کاهش قابل توجهی در پیچیدگی محاسباتی است، به‌ویژه برای خوشه‌بندی مجموعه داده‌های بسیار بزرگ است. رویکرد خوشه‌بندی مبتنی بر شبکه متفاوت از الگوریتم‌های خوشه‌بندی متعارف است که در آن نه تنها با نقاط داده بلکه با ارزش فضای اطراف نقاط داده را احاطه می‌کند. به‌طور کلی، یک الگوریتم خوشه‌بندی مبتنی بر شبکه معمولی شامل پنج مرحله پایه است: (گرباستز^۲ و بوریسوف^۳، ۲۰۰۲)

۱. ایجاد ساختار شبکه، یعنی تقسیم فضای داده به تعداد محدودی از سلول‌ها.

۲. محاسبه تراکم سلول برای هر سلول.

۳. مرتب‌سازی سلول‌ها با توجه به تراکم آن‌ها.

۴. شناسایی مراکز خوشه‌ای

۵. پیروی از سلول‌های همسایه.

- تقسیم مجموعه داده بر اساس شبکه

تعریف سلول شبکه در این پایان‌نامه به این صورت است که دو برابر شعاع همسایگی را به عنوان طول سلول شبکه انتخاب شده است (Mahran and Mahar, 2008). سپس مجموعه داده‌ها را در هر ابعاد با طول 2ϵ یکسان تقسیم می‌شود. بنابراین، فضای داده به سلول‌های شبکه‌های مختلف تقسیم

¹ Mahran & Mahar

² Grabusts

³ Borisov

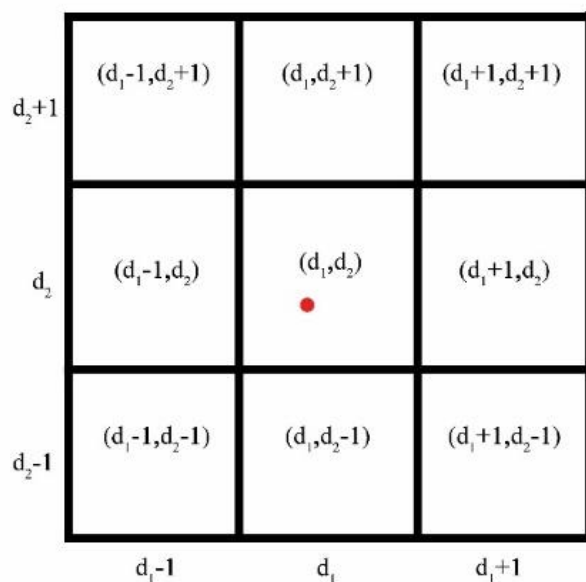
شده و با d_1, d_2, \dots, d_n مشخص می‌شود، که n ابعاد داده است، d_1 تعداد توالی سلول شبکه در ابعاد مختلف $1 \leq i \leq n$ را نشان می‌دهد (Yi-ling et al., 2008).

تعریف سلول شبکه p متعلق: x_1, x_2, \dots, x_n مختصات n از شیء p را در هر ابعاد نشان می‌دهد، شماره d_1 شبکه i ام اید دارای شرایط زیر باشد:

$$x_i/2\varepsilon \leq d_i \leq (x_i + 2\varepsilon)/2\varepsilon \quad (3-2)$$

d_i یک عدد صحیح است.

ساختار الگوریتم شبکه فقط نیاز به عبور از ۹ سلول شبکه را دارد که در هنگام محاسبه Eps همسایگی اشیاء پس از سلول شبکه تنظیم شده است. داده های دو بعدی را به عنوان مثال در نظر بگیرید، جسم $p(x, y)$ متعلق به شبکه (d_1, d_2) ، دامنه عبور در شکل ۲-۵ نشان داده شده است:

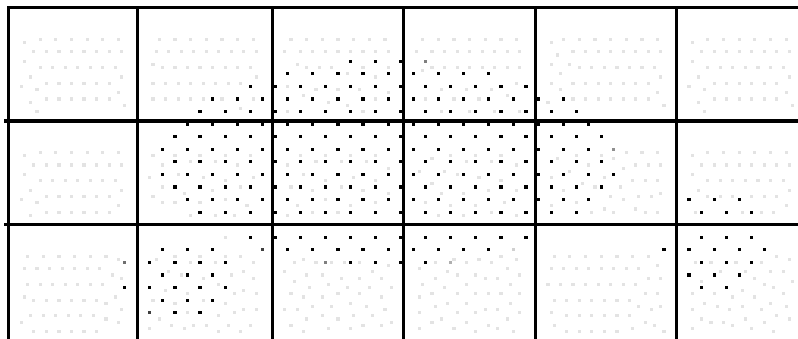


شکل ۲-۵: شبکه همسایگی در فضای دو بعدی

ساختار ترکیبی الگوریتم GridBSCAN

الگوریتم پیشنهادی بر اساس تکنیک تقسیم کن و غلبه کن^۱ است. ما یک شبکه را ایجاد می کنیم که فضای اطراف آن را پارتیشن بندی می کند و به همین ترتیب داده ها به سلول ها منتقل می شود. هر پارتیشن به طور جداگانه DBSCAN است و سپس خوشه های حاصل از همه پارتیشن ها برای ادغام خوشه بندی داده ها به یکدیگر متصل می شوند. این روند باید همان نتیجه ای را تولید کند که اگر DBSCAN در کل مجموعه داده اعمال می شد.

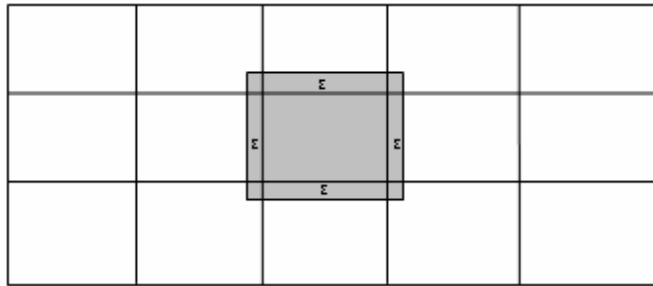
با استفاده از صفحات هم محور و هم فاصله با ابعاد "d-1" با فرض ابعاد داده ای "d" (شکل ۲-۶)، یک شبکه ساخته می شود. تعداد سلول ها به انتخاب کاربر است، اما نقطه ای بودن داده ها تأثیر زیادی بر عملکرد الگوریتم دارد. جزئیات دقیق تر شبکه بندی عملکرد آن را کاهش می دهد بنابراین ما پیشنهاد می کنیم که عرض سلول باید در هر ابعادی ($\epsilon > 2$) باشد (مهران و مهرا، ۲۰۰۸).



شکل ۲-۶: پارتیشن بندی فضایی با استفاده از شبکه

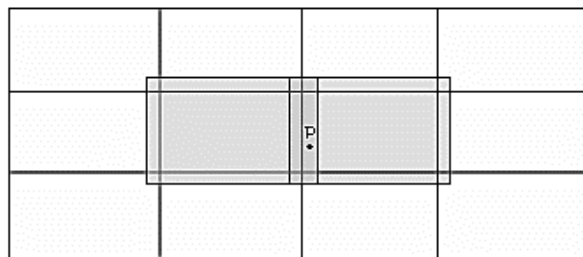
هر پارتیشن شامل تمام نقاط در داخل سلول (نقاط داخلی) و نقاط اطراف سلول در محوطه سلول (نقاط بیرونی) است (شکل ۲-۷).

¹ Divide-And-Conquer



شکل ۲-۷: پیرامون سلول ϵ

این به این معنی است که بعضی از نقاط در پارتیشن های مختلف گنجانده می شوند، اما یک نقطه تنها در یک پارتیشن درونی خواهد بود و نقطه دیگر در پارتیشن های دیگر نیز خواهد بود (شکل ۲-۸). به این ترتیب ما مطمئن هستیم که تمام پرسش های محدوده برای نقاط داخلی، دقیق هستند، زیرا دایره ϵ -شعاع در اطراف یک نقطه داخلی، به طور کامل در پارتیشن قرار می گیرد، هرچند از معیار فاصله استفاده می شود، به طوری که تمام نقاط هسته به درستی برای نقاط داخلی تعیین می شود. محدوده درخواست های بیرونی دقیق نیست زیرا حلقه های ϵ -شعاعی اطراف آنها ممکن است قطع شوند.



شکل ۲-۸: نقطه p درون سلول راست قرار دارد، اما در سلول چپ نیز شامل می شود

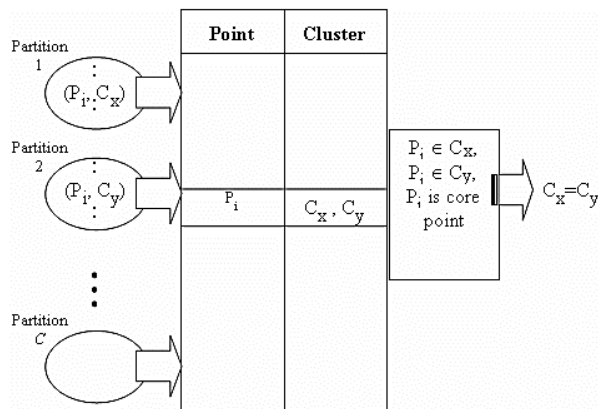
DBSCAN در هر بخش از داده ها به طور جداگانه با توجه به تمایز نقاط اصلی و نقاط مرزی

اعمال می شود و مجموعه ای متمایز از خوشه-ID برای هر پارتیشن را رزرو می کند.

خوشه بندی حاصل از همه پارتیشن ها به صورت پیوسته برای ساخت یک مجموعه شامل تمام

نقاط و شناسه خوشه خود استفاده می شود. از آنجا که یک نقطه می تواند به چندین پارتیشن تعلق

دارد، ممکن است به خوشه های چندگانه اختصاص داده شود، در این صورت تمام این خوشه ها باید تحت شرایطی که هر نقطه نقطه هسته است (شکل ۲-۹) ادغام شود. به عبارت دیگر یک نقطه اصلی برای جمع کردن تمام خوشه هایی که به آن اختصاص داده شده کافی است.



شکل ۲-۹: روند ادغام دو روش Grid و DBSCAN

۲-۳-۴- شاخص های ارزیابی خوشه بندی^۱

نتایج حاصل از اعمال الگوریتم های خوشه بندی روی یک مجموعه داده با توجه به انتخاب های پارامترهای الگوریتم ها می توانند بسیار متفاوت از یکدیگر باشند. به عبارتی تخمین تعداد بهینه خوشه ها در دسته بندی بدون نظارت^۲ داده ها، از زمینه های چالش برانگیز برای محققان در سال های اخیر بوده که منجر به ارائه شاخص های اعتبار خوشه بندی متعددی شده است. در این تحقیق شاخص های اعتبارسنجی رایج برای محاسبه تعداد بهینه خوشه ها برای ورود به الگوریتم های خوشه بندی بدون نظارت بررسی می شوند. هدف از اعتبارسنجی خوشه ها یافتن خوشه هایی است؛ که بهترین تناسب را با داده های موردنظر داشته باشند (کواچز و همکاران^۳، ۲۰۰۵).

¹ Cluster Evaluation

² Unsupervised

³ Kovács et al.

روش‌های ارزیابی خوشه‌های حاصل از خوشه‌بندی به سه دسته کلی تقسیم می‌شوند: معیارهای نسبی، معیارهای درونی، معیارهای خارجی. هم معیارهای خارجی و هم معیارهای درونی بر مبنای روش‌های آماری عمل می‌کنند و پیچیدگی محاسباتی بالایی را نیز دارا هستند. معیارهای خارجی عمل ارزیابی خوشه‌ها را با استفاده از بینش خاص کاربران انجام می‌دهند. معیارهای درونی عمل ارزیابی خوشه‌ها را با استفاده از مقادیری که از خوشه‌ها و نمای آن‌ها محاسبه می‌شود، انجام می‌دهند. پایه معیارهای نسبی، مقایسه بین خوشه‌بندی (الگوریتم به‌علاوه پارامترهای آن) مختلف است. یک و یا چندین روش مختلف خوشه‌بندی چندین بار با پارامترهای مختلف روی یک مجموعه داده اجرا می‌شوند و بهترین خوشه‌بندی از بین تمام خوشه‌ها انتخاب می‌شود. در این روش مبنای مقایسه، شاخص‌های اعتبارسنجی هستند؛ که این تحقیق در دسته معیارهای نسبی قرار دارد (کواچ و همکاران، ۲۰۰۵).

شاخص‌های اعتبارسنجی برای سنجش میزان صحت نتایج خوشه‌بندی به‌منظور مقایسه بین روش‌های خوشه‌بندی مختلف یا مقایسه نتایج حاصل از یک روش با پارامترهای مختلف مورد استفاده قرار می‌گیرند. شاخص‌های ارزیابی نسبی بسیار متنوعی پیشنهاد شده‌اند؛ از جمله، شاخص F-measure، شاخص precision و recall، شاخص Rand و AIR، شاخص MF که در این بخش به تشریح این شاخص‌ها می‌پردازیم.

شاخص رند^۱

برای نشان دادن میزان شباهت بین دو شیوه برچسب‌گذاری می‌توان از شاخص رند استفاده کرد. این شاخص توسط دانشمند آمار ویلیام رند (William Rand) در سال ۱۹۷۱ در مقاله‌ای با عنوان معیارهای هدف برای ارزیابی روش‌های خوشه‌بندی^۲ معرفی شد (رند، ۱۹۷۱). برای محاسبه آن باید دو پارامتر را اندازه‌گیری کنیم. A تعداد زوج‌هایی که هم در خوشه‌ها و هم در دسته‌ها در کنار هم هستند.

^۱ Rand Index

^۲ Objective criteria for the evaluation of clustering methods

به بیان دیگر هم در خوشه‌ها دارای برچسب یکسانی هستند و هم برچسب دسته‌ها برای آن‌ها یکسان است. B تعداد زوج‌هایی که هم در خوشه‌ها و هم در دسته‌ها از یکدیگر جدا هستند. یعنی برچسب خوشه‌هایشان متفاوت است و البته برچسب دسته‌های متفاوتی نیز دارند. حال برای محاسبه شاخص رند کافی است که حاصل جمع A و B را به تعداد کل زوج‌ها تقسیم کنیم.

$$\text{Rand}(S, C) = \frac{A+B}{\binom{N}{2}} = \frac{A+B}{\frac{N(N-1)}{2}} \quad (4-2)$$

اگر خوشه‌ها مطابق با دسته‌ها ایجاد شده باشند، شاخص رند برابر با ۱ خواهد بود. ولی اگر خوشه‌بندی به صورت تصادفی ایجاد شده باشد، دلیلی ندارد که مقدار این شاخص برابر با صفر باشد. برای رفع این مشکل از شاخص رند اصلاح شده^۱ استفاده می‌شود. این شاخص را به صورت ARI نشان داده و به شکل زیر محاسبه می‌شود:

$$\text{ARI}(S, C) = \frac{\text{Rand}(S,C) - E(\text{Rand}(S,C))}{\max(\text{Rand}(S,C)) - E(\text{Rand}(S,C))} \quad (5-2)$$

منظور از E نیز امید-ریاضی شاخص رند است.

خصوصیات شاخص رند اصلاح شده

حدود مقادارها ARI: مقداری بین ۱ و -۱ خواهد بود. در حالتی که ARI=1 باشد، مطابقت کامل بین برچسب‌های واقعی و خوشه‌ای وجود دارد و در مقابل اگر مقدار این شاخص برابر با -۱ باشد نشانگر برچسب‌گذاری تصادفی در حین خوشه‌بندی است.

¹ Adjusted Rand Index

- شاخص کارا برای مقایسه چندین روش: از آنجایی که این شاخص به توافق بین برچسب‌ها تکیه دارد، می‌توان از آن برای مقایسه دو روش خوشه‌بندی نیز استفاده کرد. برای مثال می‌توان مطابقت بین شیوه خوشه‌بندی و برچسب‌گذاری در الگوریتم K-Means را با روش فازی-C-Means بررسی کرد.

- بدون وابستگی به تعداد خوشه‌ها: با توجه به شیوه محاسبه این شاخص، تفاوت بین تعداد خوشه‌ها و دسته‌ها نیز در آن لحاظ شده است.

- تقارن در ARI: اگر جای دسته‌ها و خوشه‌ها عوض شود، شاخص ARI تغییری نمی‌کند. به این معنی که $ARI(S, C) = ARI(C, S)$.

- عدم حساسیت به تغییر برچسب‌ها: اگر برچسب‌های خوشه‌ها تغییر کند در نتیجه ARI تغییری بوجود نمی‌آید. برای مثال اگر همه برچسب‌های ۱ در خوشه‌بندی به برچسب ۲ تبدیل شوند، حاصل ARI تغییری نخواهد کرد.

فرض کنید در اینجا n_{ij} تعداد عناصر مشترک در دسته i و خوشه j باشد. همچنین a_i نیز مجموع مقدارهای n_{ij} برای سطر i ام و b_j نیز جمع مقدارهای n_{ij} در ستون j ام باشد. جدول متقاطع زیر تعداد عناصر مشترک در بین دسته‌ها و خوشه‌ها را نشان می‌دهد.

$C \setminus S$	S_1	S_2	...	S_l	Sums
C_1	n_{11}	n_{12}	...	n_{1l}	a_1
C_2	n_{21}	n_{22}	...	n_{2l}	a_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
C_k	n_{k1}	n_{k2}	...	n_{kl}	a_k
Sums	b_1	b_2	...	b_l	

حال براساس محاسبه زیر می توان شاخص رند اصلاح شده را بدست آورد:

$$ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}}{\frac{1}{2} [\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}] - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}} \quad (۶-۲)$$

شاخص فولکز - مالوز^۱

برای استفاده از این شاخص باید در مورد چند اصطلاح، آگاهی داشته باشید:

- **مثبت صحیح:** اگر زوجی از مشاهدات که در یک دسته هستند، در یک خوشه نیز قرار بگیرند، نتیجه خوشه‌بندی را برای این زوج مثبت صحیح^۲ می‌نامیم. مثبت صحیح را با TP نیز نشان می‌دهند.
- **منفی صحیح:** اگر زوجی از مشاهدات که در دو دسته مجزا قرار دارند، در دو خوشه مجزا نیز جای گیرند، نتیجه خوشه‌بندی را برای این زوج منفی صحیح^۳ می‌نامیم. منفی صحیح را با TN نیز نشان می‌دهند.
- **مثبت کاذب:** اگر زوجی از مشاهدات که در دو دسته مجزا قرار دارند، در یک خوشه جای گیرند، نتیجه خوشه‌بندی را برای این زوج مثبت کاذب^۴ است. مثبت کاذب را با FP نیز نشان می‌دهند.
- **منفی کاذب:** اگر زوجی از مشاهدات که در یک دسته قرار دارند به اشتباه در دو خوشه قرار گیرند، نتیجه خوشه‌بندی را برای این زوج منفی کاذب^۵ است. منفی کاذب را با FN نیز نشان می‌دهند.

¹ Fowlkes-Mallows Index

² True Positive

³ True Negative

⁴ False Positive

⁵ False Negative

جدول زیر که وظیفه مقایسه بین تعداد نتایج خوشه‌بندی توسط الگوریتم و دسته‌بندی اصلی را به عهده دارد، به ماتریس درهم‌ریختگی^۱ شهرت دارد که توجه به آن به درک و شیوه محاسبه شاخص فولکر-مالوز و اصطلاحات گفته شده، کمک می‌کند.

جدول ۲-۲: تعاریف زوج‌های مثبت و منفی

دسته‌بندی واقعی		نتیجه خوشه‌بندی برای زوج $(x,y)(x,y)$	
FP	TP	در یک خوشه	خوشه‌بندی
TN	FN	در دو خوشه مجزا	

حال اگر تعداد زوج‌های مثبت صحیح را با $|TP|$ ، تعداد مثبت کاذب را با $|FP|$ و تعداد منفی کاذب را با $|FN|$ نشان بدهیم، شاخص فولکر-مالوز که با FMI نشان داده می‌شود به صورت زیر محاسبه می‌شود:

$$FM = \sqrt{\frac{TP}{TP+FP} \cdot \frac{TP}{TP+FN}} \quad (۷-۲)$$

شاخص فولکر-مالوز یک روش ارزیابی خارجی است که برای تعیین شباهت بین دو خوشه بندی (خوشه‌های به دست آمده پس از الگوریتم خوشه بندی) استفاده می‌شود. این اندازه گیری شباهت می‌تواند بین دو خوشه بندی سلسله مراتبی یا دسته بندی و طبقه بندی معیار باشد. مقادیر بیشتری برای شاخص فولکر-مالوز نشان می‌دهد که شباهت بیشتری بین خوشه‌ها و طبقه بندی‌های معیار وجود

¹ Confusion Matrix

دارد. شاخص فولکز- مالوز، هنگامی که نتایج حاصل از دو الگوریتم خوشه بندی برای ارزیابی نتایج استفاده می شود، به صورت زیر تعریف می شود: (Fowlkes and Mallows, 1983)

که TP تعداد مثبت واقعی است، FP تعداد مثبت کاذب است و FN تعداد منفی کاذب است. از آنجا که این شاخص مستقیماً با تعداد مثبت های واقعی متناسب است، یک شاخص بالاتر، به معنی همبستگی بیشتر بین دو خوشه بندی مورد استفاده برای تعیین شاخص استفاده شده است. یکی از روش های اصلی برای بررسی اعتبار این شاخص، مقایسه دو دسته بندی است که به یکدیگر وابسته نیستند. فولکز و مالوز نشان دادند که برای دو خوشه بندی غیر مرتبط، ارزش این شاخص زمانی که تعدادی داده های مورد نظر برای خوشه بندی افزایش می یابد به صفر می رسد؛ در حالی که ارزش شاخص رند برای داده های مشابه به سرعت به ۱ تبدیل می شود؛ بنابراین فولکز- مالوز شاخص بسیار دقیق تری برای داده های غیر مرتبط است. این شاخص همچنین زمانی که نویز به یک مجموعه داده اضافه شود به خوبی عمل می کند. فولکز و مالوز نشان داد که مقدار شاخص با افزایش مولفه نویز کاهش می یابد. این شاخص همچنین حتی زمانی که مجموعه داده های نویز دارای تعداد خوشه ای متفاوت از خوشه های مجموعه داده اصلی بودند، شباهت را نشان داد؛ بنابراین این شاخص را یک ابزار قابل اعتماد برای اندازه گیری شباهت بین دو خوشه می سازد.

شاخص اندازه^۱ F

همانطور که دیده شد، استفاده از مفاهیم مثبت کاذب، مثبت صحیح، منفی کاذب و منفی صحیح در محاسبه شاخص ارزیابی نتایج خوشه بندی اهمیت زیادی دارد. با توجه به این مفاهیم می توان شاخص رند را به صورت زیر نیز محاسبه کرد:

^۱ F-measure Index

$$\text{Rand}(S, C) = \frac{|TP|+|TN|}{|TP|+|FP|+|TN|+|FN|} \quad (۸-۲)$$

در این حالت وزن یا اهمیت مربوط به منفی کاذب (FN) یا مثبت کاذب (FP) در محاسبه شاخص یکسان است. ولی در اندازه F می توان وزن آن ها را با پارامتر β تغییر داد. اگر P و R را به صورت زیر تعریف کنیم:

$$P = \frac{|TP|}{|TP|+|FP|}, R = \frac{|TP|}{|TP|+|FN|} \quad (۹-۲)$$

آنگاه شاخص F را به ازاء $\beta \neq 0$ به شکل زیر محاسبه می شود:

$$F_{\beta} = \frac{(\beta^2+1)PR}{\beta^2P+R} \quad (۱۰-۲)$$

در تجزیه و تحلیل آماری طبقه بندی دوتایی، نمره F_1 (همچنین F-score یا F-measure) یک نوع اندازه گیری از دقت آزمون است (ساسکی^۱، ۲۰۰۷).

از هر دو آزمون تشخیص^۲ p و حساسیت^۳ r برای محاسبه نمره F_1 استفاده می شود. P تعداد نتایج مثبت صحیح تقسیم بر تعداد تمام نتایج مثبت توسط طبقه بندی شده است و r تعداد نتایج منفی صحیح تقسیم بر تعداد تمام نمونه های مربوطه (تمام نمونه هایی که باید مثبت باشند). نمره F_1 میانگین هارمونیک دقت و حساسیت است، جایی که نمره F_1 به بهترین وجه باشد به ۱ (precision و recall کامل) و در بدترین حالت به ۰ می رسد. اندازه گیری رایج F یا F-score متعادل (نمره F_1) میانگین هارمونی precision و recall است:

¹ Sasaki

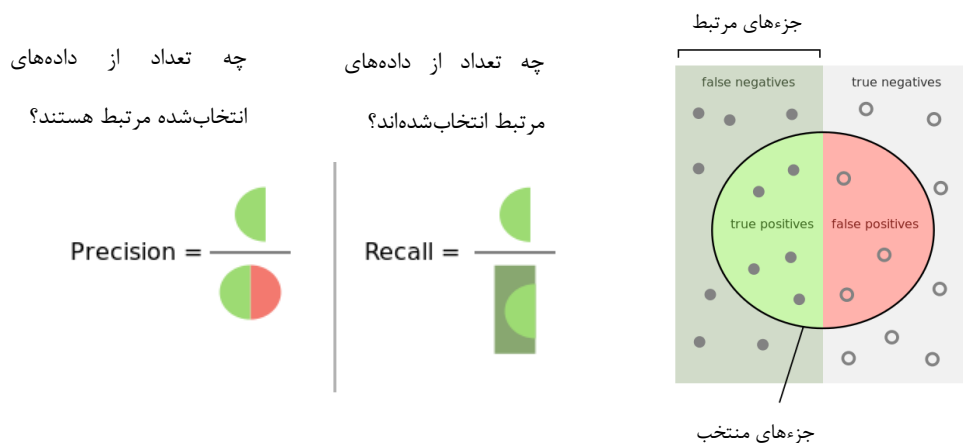
² Precision

³ Recall

$$F_1 = \left(\frac{recall^{-1} + precision^{-1}}{2} \right)^{-1} = 2 \cdot \frac{precision \cdot recall}{precision + recall} \quad (11-2)$$

شاخص دقت و حساسیت

در تشخیص الگو، بازیابی اطلاعات و طبقه‌بندی دوتایی^۱، precision (ارزش پیش‌بینی مثبت) کسری از نمونه‌های مربوطه در میان نمونه‌های بازیابی شده است، در حالی که recall (همچنین به عنوان حساسیت شناخته می‌شود) کسری از نمونه‌های مربوط است که در کل مقدار نمونه‌های مربوطه بازیابی شده است؛ بنابراین precision و recall بر درک و اندازه‌گیری ارتباط استوار است. به عبارت ساده، precision بیشتر به این معنی است که یک الگوریتم نتایج مرتبط تری نسبت به نتایج نامناسب ارائه می‌دهد، در حالی که recall بیشتر به این معنی است که از یک الگوریتم، نتایج کاملاً مرتبط بدست می‌آید.



شکل ۲-۱۰: نمای شماتیک تعریف recall و precision

¹ Binary

شاخص جاکارد^۱

شاخص جاکارد، همچنین به عنوان اشتراک در اتحاد یک مجموعه و ضریب تشابه جاکارد (ضریب اصل ضوابط de communauté توسط پل جاکارد)، آماری برای مقایسه شباهت و تنوع مجموعه های نمونه مورد استفاده قرار می گیرد. ضریب تشابه جاکارد اندازه گیری شباهت بین مجموعه نمونه های محدود است و اندازه اشتراک تقسیم بر اندازه اتحاد از مجموعه های نمونه تعریف شده است. (جاکارد، ۱۹۰۱) از رابطه ۲-۵ این شاخص محاسبه می شود:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|} \quad (12-2)$$

اگر A و B هر دو تهی باشند، $J(A, B) = 1$ تعریف می کنیم و $0 \leq J(A, B) \leq 1$.

شاخص NMI^۲

نرمال اطلاعات متقابل (NMI) یعنی امتیاز اطلاعات متقابل (MI^۳) نرمال شده به مقیاس نتایج خوشه، بین ۰ (نبود اطلاعات متقابل) و ۱ (همبستگی کامل) است. اطلاعات متقابل یعنی با کاهش انترپی برچسب های کلاس بندی که ما دریافت می کنیم اگر برچسب خوشه را بدانیم نشان می دهد. از رابطه ۲-۷ این مقدار بدست می آید.

$$NMI(Y, C) = \frac{2 \times I(Y; C)}{[H(Y) + H(C)]} \quad (13-2)$$

که در آن

Y = برچسب کلاس

¹ Jaccard

² Normalized Mutual Information

³ Mutual information

برچسب های خوشه ای $C =$

آنتروپی $H(.) =$

$I(Y; C) =$ و C اطلاعات متقابل $b/w Y$

NMI اندازه گیری خوبی برای تعیین کیفیت خوشه بندی است. یک اندازه گیری خارجی است، زیرا ما به برچسب کلاس نمونه ها نیاز داریم تا NMI را تعیین کنیم. از آنجا که نرمال است، ما می توانیم NMI را در بین خوشه بندی های مختلف با تعداد مختلف خوشه اندازه گیری و مقایسه کنیم.

شاخص جینی^۱

در اقتصاد، ضریب جینی که گاهی اوقات شاخص جینی یا نسبت جینی نامیده می شود، اندازه گیری پراکندگی آماری است که برای نشان دادن توزیع درآمد یا ثروت، توزیع ساکنان یک ملت و مهمترین آن، اندازه گیری نابرابری استفاده می شود. این شاخص توسط جامعه شناس و محقق ایتالیایی Corrado Gini و در سال ۱۹۱۲ ارائه شد. ضریب جینی نابرابری را در بین مقادیر توزیع فراوانی (به عنوان مثال، سطوح درآمد) اندازه گیری می کند. ضریب جینی صفر برابری مطلق را نشان می دهد، در حالیکه تمام مقادیر یکسان هستند (به عنوان مثال که هر کس درآمد مشابهی دارد). ضریب جینی ۱ (یا ۱۰۰٪) بیانگر نابرابری حداکثر بین ارزش ها می باشد (به عنوان مثال برای تعداد زیادی از مردم که تنها یک نفر تمام درآمد و مصرف و همه دیگران را دارد، ضریب جینی تقریباً نزدیک به یک است) با این وجود، اگر یک نفر به سهم منفی در مجموع (به عنوان مثال، داشتن درآمد منفی یا ثروت)، یک مقدار بیشتر از یک ممکن است رخ دهد. برای گروه های بزرگتر، مقادیر نزدیک به ۱ یا بالاتر از ۱ در عمل بسیار بعید است.

¹ Gini index

شاخص خلوص

یکی از ساده‌ترین شاخص‌های ارزیابی بیرونی در خوشه‌بندی، شاخص خلوص^۱ است که درصد مطابقت بین برچسب‌های خوشه‌بندی و برچسب‌های واقعی را می‌سنجد. در این حالت برچسب هر خوشه با برچسب واقعی دسته‌ای که بیشترین اشتراک را دارد مطابقت پیدا کرده و تعداد نقاطی از خوشه که در دسته صحیح طبقه‌بندی شده‌اند شمارش می‌شوند. نسبت این تعداد به تعداد کل نقاط شاخص خلوص را می‌سازد و شکل محاسباتی آن به صورت زیر است:

$$Purity(S, C) = \frac{\sum_m \max_n |S_m \cap C_n|}{N} \quad (۱۴-۲)$$

باید توجه داشت که منظور از $|S_m \cap C_n|$ تعداد نقاط مشترک از خوشه C_n با دسته S_m و N نیز تعداد کل نقاط است. با توجه به شیوه محاسبه شاخص خلوص، مشخص است که حداکثر مقدار برای آن ۱ خواهد بود و این در زمانی اتفاق می‌افتد که برچسب‌های حاصل از خوشه‌بندی کاملاً با برچسب‌های واقعی مطابقت داشته باشند. همین‌طور اگر هیچ برچسب خوشه‌ای با برچسب واقعی مطابقت نداشته باشد، این شاخص صفر می‌شود.

خصوصیات شاخص خلوص

سادگی در محاسبات: با توجه به شیوه محاسبه این شاخص، محاسبات طولانی برای بدست آوردن میزان مطابقت لازم نیست.

مستقل از تعداد خوشه‌ها: شاخص خلوص به تعداد خوشه‌ها توجه ندارد. در نتیجه نمی‌توان این شاخص را به عنوان معیاری برای سنجش مطابقت تعداد خوشه‌ها نیز در نظر گرفت.

^۱ Purity Index

کاهش کارایی با افزایش تعداد خوشه‌ها: اگر تعداد خوشه‌ها زیاد باشد و هیچ هماهنگی نیز بین برجسب‌های واقعی و خوشه‌ای وجود نداشته باشد ممکن است شاخص خلوص به ۱ نزدیک شود که یک عیب برای چنین شاخصی است.

۲-۳-۵- انواع الگوریتم‌های خوشه‌بندی آماری

در سال ۱۹۹۰ با تمرکز بر شناسایی نقاط پرحادثه، خوشه‌بندی آماری، حرکتی از آماره‌های سراسری به وجود آمد. همچنین روش‌هایی مانند برآورد تراکم کرنل و خودهمبستگی مکانی برای شناسایی موقعیت خوشه‌ها توسعه یافتند. دو روش اصلی که ما در این تحقیق مورد بررسی قرار دادیم عبارتند از روش تراکم کرنل در سطح صفحه‌ای و آماره خودهمبستگی مکانی که با استفاده از نرم افزار GIS امکان‌پذیر است.

روش تخمین تراکم کرنل صفحه‌ای KDE^۱

این روش در سال ۱۹۵۶ توسط محققى به نام رزنالت در جهت تخمین ناپارامتری تراکم معرفی شد. این ابزار برای تخمین تراکم عوارض در همسایگی یک نقطه را محاسبه می‌کند. تخمین تراکم کرنل روشی رایج برای تحلیل داده‌های یک‌بعدی و دوبعدی و همچنین شناسایی نقاط متمرکز یا حادثه‌خیز است و از طرفی اطلاعات مفیدی را در مورد داده‌ها ارائه می‌دهد (دنگ^۲، ۲۰۰۷). این روش می‌تواند به سرعت نقاط حادثه‌خیز را از پایگاه‌های داده‌ای بزرگ مانند تصادفات را شناسایی کرده و در نتیجه خروجی آماری قابل قبولی را ارائه دهد.

از مزایای این روش می‌توان به نمایش نقاط متمرکز تصادفات در سطح (منطقه‌ای) و نقطه و شبکه راه‌ها به خصوص در تصادفات اشاره کرد. این روش برخلاف روش‌های خودهمبستگی مکانی که

^۱ Planar Kernel Density Estimation

^۲ Duong

نقاط متمرکز را شناسایی می کنند می توانند مناطق حادثه خیز را شناسایی کنند. مزیت اصلی استفاده از تخمین تراکم کرنل، تشخیص محدوده ی گسترش خطر تصادفات است درحالی که در خوشه بندی، صرفاً دسته بندی بر اساس تراکم داده ها خوشه بندی انجام می شود. در خوشه بندی، اطراف هر خوشه صرفاً در نهایت دارای درجه ای خطر تصادف برای آن هایی است که وارد آن خوشه می شوند (سابل و همکاران^۱، ۲۰۰۵).

هر نقطه داده x_i و $i = 1, 2, \dots, n$ نمونه مستقل و با توزیع یکسانی با چگالی نامعلوم f هستند. تابع برآورد تراکم کرنل $(\hat{f}_h(x))$ چگالی نامعلوم $f(x)$ را در نقطه x برآورد می نماید. تابع برآورد تراکم کرنل $\hat{f}_h(x)$ به صورت رابطه شماره ۲-۸ محاسبه می شود:

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) \quad (15-2)$$

که در آن $\hat{f}_h(x)$ تابع برآورد تراکم کرنل، $k\left(\frac{x-x_i}{h}\right)$ تابع کرنل مثبت متقارن با انتگرال واحد است و h دلالت بر پهنای باند یا شعاع همسایگی دارد. به ازای کاهش پهنای باند (h)، منحنی حاصل از برآورد کرنل ناهموارتر شده و جزئیات جعلی بیشتری را از تراکم واقعی به نمایش می گذارد و به ازای افزایش پهنای باند، منحنی هموار و باعث محو شدن جزئیات واقعی تابع تراکم می گردد. تخمین کرنل مجموع «برآمدگی»هایی است که مشاهدات در مرکز آن ها واقع شده اند. تابع کرنل k ، شکل برآمدگی ها و پهنای باند (h)، طول آن ها را مشخص می کند. از تابع کرنلی که به طور معمول استفاده می شود می توان به توزیع نرمال اشاره کرد.

¹ Sabel et al.

تحلیل خودهمبستگی مکانی^۱

زمانی که مقادیر یک متغیر که از نظر جغرافیایی به هم نزدیک هستند باهم مرتبط می‌باشند که به این تعریف خودهمبستگی مکانی می‌گویند. ارزیابی روش‌های خودهمبستگی مکانی عمدتاً به‌وسیله تحقیقات انسلین^۲ ۱۹۹۵ و موران^۳ ۱۹۴۸ شروع شد. هنگامی که حوادث وابسته هستند گفته می‌شود که همبستگی مکانی مثبت است. اگر عوارض و یا مقادیر متغیرهای مربوط به آن‌ها به‌طور تصادفی در فضا توزیع شده باشند ظاهراً نباید بین آن‌ها ارتباطی وجود داشته باشد (گیریفیس^۴، ۲۰۰۸). چندین روش برای آزمایش همبستگی مکانی وجود دارد نظیر Moran's I، * Gi Ord-Getis و Anselin Local Moran's I که در ادامه به بررسی آن‌ها پرداخته می‌شود (Moran 1948).

روش خودهمبستگی مکانی Moran's I

Moran's I که مقدار خروجی آن بین ۱ و ۱- تعریف شده است، یکی از روش‌های خودهمبستگی مکانی است که در زمینه تحلیل الگوهای توزیع عوارض و پدیده‌ها در مکان استفاده می‌شود. این ابزار یکی از کاربردی‌ترین و مهم‌ترین ابزارهای تحلیلی برای تحقیق در مورد داده‌های مکانی است. به‌طور مثال در روش‌هایی که نیاز به داده‌های مستقل است با کمک ابزار خودهمبستگی مکانی می‌توان دریافت که آیا داده‌ها از نظر مکانی مستقل از یکدیگر هستند و یا خیر. در صورتی که نتایج تحلیل خودهمبستگی Moran's I نشان‌دهنده وجود خودهمبستگی مکانی بین داده‌ها باشند، در آن صورت نمی‌توان از روش‌هایی که نیازمند داده‌های مستقل هستند استفاده نمود. ابزار تحلیل خودهمبستگی مکانی Moran's I به بررسی خودهمبستگی مکانی بر اساس مکان دو مقدار ویژگی مورد نظر عوارض جغرافیایی می‌پردازد. اگر تعدادی عوارض جغرافیایی با ویژگی مشخصی موجود باشند

¹ Spatial Autocorrelation

² Anselin

³ Moran

⁴ Griffith

این ابزار نشان می‌دهد که الگوی پراکنش این عوارض با در نظر گرفتن مقادیر ویژگی مورد مطالعه از الگوی خوشه‌ای یا پراکنده برخوردار است. نتایج حاصل از این تحلیل نشان می‌دهد که آیا عوارض به صورت تصادفی، پراکنده یا خوشه‌ای در مکان توزیع شده‌اند. این ابزار در حقیقت آماره و یا شاخص Moran را محاسبه می‌کند و با استفاده از امتیاز استاندارد Z (Z-Score) و P-Value به ارزیابی و معنادار بودن شاخص محاسبه شده می‌پردازد رابطه، نشان‌دهنده شاخص Moran است.

$$I = \frac{n \sum_{i=1}^n \sum_{j=1}^n w_{ij} (x_i - \bar{x})(x_j - \bar{x})}{(\sum_{i \neq j}^n \sum_{j=1}^n w_{ij})(\sum_{i=1}^n (x_i - \bar{x})^2)} \quad (16-2)$$

در این رابطه n نشان‌دهنده تعداد رخدادهای \bar{x} مقدار متوسط کل رخدادهای x_j تعداد رخدادهای در موقعیت i ام و j ام است. $(x_i - \bar{x})$ نیز تفاضل بین مقدار ویژگی عارضه i با میانگین آن است. w_{ij} ماتریس وزن مکانی تعریف شده برای تعیین درجه تأثیرات محلی رویدادها است. فرمول ریاضی Z-Score، اهمیت نتایج به دست آمده را به وسیله رابطه ۲-۱۰ آزمایش می‌کند:

$$Z = \frac{I_0 - I_E}{SD_{IE}} \quad (17-2)$$

در رابطه ۲-۱۰، I_0 ، I_E و SD_{IE} به ترتیب نشان‌دهنده مشاهدات، انحراف مورد انتظار (میانگین) و انحراف استاندارد هستند.

Z-Score بزرگ و P-Value کوچک برای هر نقطه نشان‌دهنده خوشه‌بندی مکانی مقادیر بالا است. Z-Score پائین و منفی و P-Value کوچک نشان‌دهنده خوشه‌بندی مکانی مقادیر پایین است. برای ابزارهای تحلیل الگو، مقدار P-Value احتمال تولید الگوی مکانی مشاهده شده به وسیله فرآیندهای تصادفی است. P-Value خیلی کوچک (احتمال پایین) به این معنا است که الگوی مکانی مشاهده شده نتیجه فرایندهای تصادفی بوده است، بنابراین فرض صفر رد می‌شود. به طور کلی اگر مقدار شاخص Moran نزدیک به عدد +۱ باشد داده‌ها دارای خودهمبستگی مکانی و دارای الگوی خوشه‌ای بوده و اگر

مقدار شاخص Moran نزدیک به عدد ۱- باشد، آنگاه داده‌ها از هم گسسته و پراکنده هستند. در مورد این ابزار فرضیه صفر آن است که هیچ نوع خوشه‌بندی مکانی بین مقادیر ویژگی مرتبط با عوارض مکانی مورد نظر وجود ندارد. حال زمانی که مقدار P-Value بسیار کوچک و مقدار Z-Score محاسبه‌شده (قدر مطلق آن) بسیار بزرگ باشد (خارج از محدوده اطمینان قرار گیرد)، آنگاه می‌توان فرضیه صفر را رد کرد. اگر مقدار شاخص Moran بزرگتر از صفر باشد، داده‌ها نوعی خوشه‌بندی مکانی را نشان می‌دهند. اگر مقدار شاخص کمتر از صفر باشد عوارض مورد مطالعه دارای الگوی پراکنده هستند.

۲-۳-۶- استفاده از GIS در مطالعات پیشین

رایج‌ترین کاربرد نرم‌افزار GIS در سطح ماکروسکوپی تعیین نقطه پرحادثه می‌باشد. نقاط حادثه‌خیز مناطقی هستند که در آن‌ها یک واقعه با تمرکز بالا رخ می‌دهد. در پژوهش حاضر منظور از نقاط حادثه‌خیز مناطقی می‌باشد که دارای تمرکز بالای تصادفات منجر به جراحت می‌باشند. نتایج تحلیل نقاط حادثه‌خیز می‌توانند به صورت تصویری در غالب نقشه‌هایی که تحت عنوان «سیستم اطلاعات جغرافیایی» GIS نامیده می‌شوند به نمایش در آیند. این قابلیت باعث افزایش امکان شناخت وضع موجود و مقایسه نقاط حادثه‌خیز با شرایط کلی منطقه می‌گردد. این نتایج تصویری می‌توانند توسط مدیران جامعه و طراحان شهری مورد بررسی قرار گیرند تا دلایل احتمالی شکل‌گیری چنین الگویی‌هایی در وقوع حوادث خاص مشخص شده و راهکارهای مؤثر در جهت حل مشکلات موجود اتخاذ گردد. درحالی‌که متمرکز بودن وقایعی نظیر جرایم، حوادث طبیعی و همینطور تصادفات رانندگی منطقاً دور از انتظار نیست اما به دلیل دخیل بودن عوامل مختلف در تحلیل نقاط حادثه‌خیز، لازم است حداکثر دقت در انتخاب روش تحلیل و پارامترهای مؤثر در هر روش صورت گیرد تا بهترین و قابل اعتمادترین نتایج حاصل گردند. با اینکه امروزه روند شناسایی نقاط حادثه‌خیز با استفاده از برنامه‌های

رایانه ای خیلی پیچیده نمی‌باشد در هر حال جهت کسب بهترین نتایج، آگاهی از زمینه های تئوری و منطق عملکردی این ابزارها ضروری می‌باشد.

یکی از مزایای تحلیل‌های GIS محور امکان روی هم گذاری نتایج این تحلیل ها با لایه های مختلف مشخصات شبکه های حمل و نقل می‌باشد که قابلیت بررسی ارتباط بین تراکم تصادفات با هر یک از مشخصات را ممکن می‌سازد.

نقاط حادثه‌خیز توسط روش‌هایی که بمنظور شناسایی تمرکز وقایع طراحی شده اند تعیین می‌گردند. روش‌های بسیاری وجود دارند که برای تعیین نقاط حادثه‌خیز مورد استفاده قرار می‌گیرند و در برخی مواقع نتایج متفاوتی خواهند داشت. درحالی‌که اکثر مطالعات GIS نه تنها به منظور نمایش تصویری نتایج بلکه به عنوان ابزاری جهت جمع بندی اطلاعات از منابع مختلف نیز به کار گرفته شده هنوز تعدادی از ابزارهای تحلیل آمار مکانی آن طور که باید و شاید از طرف ارگان های مربوطه مورد استفاده قرار نمی‌گیرد. همانطور که قبلاً نیز اشاره شد باوجود پیشرفت های سریع GIS در سال های اخیر همچنان بسیاری از نهاد های مرتبط به دلیل نداشتن پایگاه داده GIS محور، امکان استفاده از GIS و ابزار های سودمند و در عین حال ساده و صریح این نرم افزار ها را ندارند، یکی دیگر از مضرات عدم استفاده از GIS در تعیین نقاط حادثه‌خیز تصادفات، احتمال نادیده گرفته شدن برخی اطلاعات مربوط به موقعیت مکانی می‌باشد که ممکن است تأثیر بسزایی نیز در وقوع تصادفات داشته باشند. تحلیل نقاط حادثه‌خیز تصادفات نیازمند درک هم جانبه روند وقوع تصادف، شدت صدمات بوجود آمده و محیط اطراف محل حادثه می‌باشد. محیط GIS به دلیل ارائه سیستم کارآمد مرتبط نمودن پایگاه های داده متعدد مناسب استفاده در این گونه موارد می‌باشد. این پلت فرم یک سیستم مرجع مکانی برای گزارش گیری در سطوح مختلف ارائه می‌کند که ضمن قابلیت وارد کردن داده‌های آماری و موقعیتی تصادفات امکان برآورد خطر و شدت صدمات وارده را فراهم می‌نماید. طبقه بندی نقاط حادثه‌خیز تصادفات معمولاً بر اساس اطلاعات موجود درباره همان تصادفات مانند نوع تصادف، زمان وقوع، نوع مصدوم، نوع

وسیله نقلیه و... صورت می گیرد. همین موضوع اغلب باعث محدود شدن جوانب تحلیل مسئله می گردد و منجر به از قلم افتادن برخی از عوامل محیطی، اجتماعی و اقتصادی به دلیل نبود داده‌های اولیه می گردد. با این همه، نکته مثبت، رویکرد پژوهشگران عرصه حمل و نقل به سوی استفاده روزافزون از این ابزارها و مطالعه و بررسی در بهبود و گسترش کاربرد بیش از پیش این روش‌ها در تحلیل مشکلات موجود می‌باشد.

پاراسانا کومار و همکارانش (پاراسانا کومار و همکاران، ۲۰۱۱)، تحلیل خوشه‌ای را با توجه به نوع حادثه و زمان وقوع با استفاده از Index Moran روش خودهمبستگی مکانی مطالعه کردند. تحلیل خوشه‌ای با آمار KDE و Gtis-Ord GI انجام شد. در نتیجه تجزیه و تحلیل‌ها، لکه‌های گرم و سرد نقاط تصادف تعیین شد. موریرا و همکاران (۲۰۱۲) سه روش مختلف برای شناسایی مکان‌های خطرناک (HRL) در شهر ویارئال پرتغال را مورد بررسی قرار دادند: الگوریتم خوشه بندی KDE، Nnh و تراکم نقطه‌ای. سرعت بالا به عنوان یک عامل شاخص تصادفات تک خودرویی (تصادف خروج از جاده) پیدا شد. اکی و کرتی (۲۰۱۲) خوشه‌ای از تصادفات در شبکه را با استفاده از نرم افزار Sanet جستجو کردند. خوشه مکانی- زمانی با روش ناکس انجام شد.

تسا اندرسون (۲۰۰۹) با ارائه دو روش به مطالعه تصادفات جاده‌ای پرداختند. در روش اول آن‌ها با استفاده از سیستم اطلاعات جغرافیایی و تخمین تراکم کرنل، الگوهای مکانی مرتبط با تصادفات جراحی را در لندن مورد مطالعه قرار دادند. هم‌چنین آن‌ها از روش خوشه‌ای، با بهره‌گیری از داده‌های محیطی و نیز نتایج حاصل از روش تراکم کرنل، جهت دسته‌بندی مناطق حادثه‌خیز جاده‌ای استفاده کردند. پردازش خوشه‌ای در مطالعه آن‌ها ۵ گروه و ۱۵ خوشه ایجاد کرد و تفاوت اصلی بین گروه‌ها و خوشه‌ها، واریانس متغیرهای ساختاری می‌باشد. نتایج آن‌ها نشان داد که واریانس متغیرها در هر گروه به علت ساختار سلسله مراتبی به‌طور قابل ملاحظه‌ای کاهش می‌یابد. مزیت اصلی روش تخمین تراکم

کرنل، تشخیص توزیع ریسک تصادف در محدوده خوشه‌ای است که به دلیل استقلال مکانی، احتمال رخداد تصادف در آن بیشتر می‌باشد.

از مزایای استفاده از GIS می‌توان به امکان مدیریت و پردازش داده‌های مکانی و اطلاعات مرتبط با این داده‌ها اشاره کرد (Longley et al., 2005). تشخیص نقاط حادثه‌خیز با استفاده از GIS امکان‌پذیر شده است و از آنجایی که در سال‌های اخیر ثبت داده‌ها به صورت دقیق و توسط GPS انجام می‌شود، تشخیص مناطقی که از جنبه‌های مختلف به هم مشابهت دارند، می‌تواند درک دقیق‌تر و جامع‌تری را از ویژگی تصادفات در اختیار قرار دهد. به این منظور، در کنار استفاده از GIS و روش خوشه‌بندی می‌توان به این هدف دست‌یافت.

باران اولاک و همکارانش در سال ۲۰۱۷ تصادفات را با در نظر گرفتن سن افراد مورد بررسی قراردادند. سه مرحله برای تحلیل مکانی انجام شد. ابتدا خوشه‌بندی بر اساس تخمین تراکم کرنل در شبکه (با استفاده از افزونه^۱ SANET در GIS) استفاده شد. سپس ضریب تفاوت نرخ تراکم (DRD)^۲ برای دو گروه سنی برای مقایسه تراکم بیشینه و نرمال شده به دست آمد. در مرحله سوم ضریب جمعیت (PF)^۳ برای بررسی تأثیر وابستگی مکانی محاسبه شد. با استفاده از این ضرایب می‌توان مناطق مهم برای سیاست‌گذاری در جهت تصمیم‌گیری برای ایمنی با محوریت سن افراد را مشخص نمود (Ulak, Ozguven et al. ۲۰۱۷).

سلطانی و عسکری در سال ۲۰۱۷ برای خیابان‌های شریانی درجه اول شیراز توسط GIS خوشه‌بندی مکانی را با روش‌هایی همچون انجام دادند. این تحقیق به منظور یافتن الگوهای زمانی و مکانی تصادفات در محیط شهری انجام شد. نقشه‌های موازی برای مقایسه سه فاکتور زمان رخداد، شدت و موقعیت مکانی تصادف استفاده شد. تحلیل فراوانی سالانه تصادفات در ۱۵۶ منطقه تصادفات

¹ Spatial analysis along network

² Density ratio difference

³ population factor

ترافیکی در شیراز از سال ۲۰۱۰ تا ۲۰۱۴ نشان داد که دو روش شاخص I موران^۱ و گتیس اورد Gi*^۲ خوشه‌بندی معناداری از الگوهای تصادف ایجاد می‌کند. این تحقیق نشان داد که تصادفات دارای الگوی خوشه‌ای است اما در تحلیل مکانی-زمانی الگوی تصادفی از خود نشان داد.

شفابخش و همکاران^۳ در سال (۲۰۱۷) برای تصادفات در نقاط درون‌شهری مشهد از GIS برای تحلیل مکانی از روش تخمین چگالی کرنل و از افزونه SANET برای تخمین چگالی کرنل در شبکه راه‌ها استفاده کردند. سپس با کمک نتایج به‌دست آمده، سطح مورد مطالعه را در پنج منطقه با اولویت‌های متفاوت جهت اقدامات ایمنی طبقه‌بندی کردند.

^۱ Moran's I index

^۲ Getis ord Gi*

^۳ Shafabakhsh et al.

فصل ۳

روش تحقیق

۳-۱- مقدمه

در این بخش ابتدا داده‌های مورد استفاده و محدوده مطالعه را بررسی کرده و سپس به شرح روش‌های ذکر شده در فصل دوم می‌پردازیم و روند دقیق انجام شده در این پایان‌نامه آورده شده است. در تحلیل نقاط حادثه‌خیز تصادفات شهری، همانند نقاط حادثه‌خیز بزرگراه‌ها برای شناسایی نقاط حادثه‌خیز تعریف یا چارچوب معینی در کارهای گذشته وجود ندارد. این عمدتاً به این دلیل است که پیچیدگی و ترکیب بودن مدهای مختلف حمل و نقل در شبکه شهری می‌باشد و همچنین مبداءها و مقصدهای فراوانی برای مناطق درون شهری وجود دارد. با توجه به اینکه حجم ترافیک در مناطق شهری به صورت مکانی تصادفی نیستند، تصادفات نیز از لحاظ مکانی تصادفی نیستند. علاوه بر این وجود جابجایی عابرین پیاده در مناطق شهری بر پیچیدگی تحلیل شبکه شهری می‌افزاید که بیشتر به کابری اطراف مسیر مرتبط می‌باشد. به همین دلیل در مناطق شهری نمی‌توان مناطق را به بخش‌های یکنواخت و نرمال تقسیم کرد.

پس از آشنایی و تحلیل پارامترهای مختلف تأثیرگذار بر وقوع تصادفات، همچنین بررسی روش‌های خوشه‌بندی مکانی که در فصل دوم ارائه شد، شرایط لازم جهت ارائه روشی برای خوشه‌بندی بر مبنای اطلاعات مکانی به منظور تحلیل مکانی تصادفات مهیا است. در این فصل ابتدا به تشریح منطقه مورد مطالعه، جمع‌آوری و آماده‌سازی داده‌ها پرداخته می‌شود. روش‌های خوشه‌بندی غیر آماری K-Means و Nnh و GriDBSCAN به طور جداگانه برای دو شهر ایزمیت و گبزه به منظور تعیین خوشه‌ها با بیشترین تعداد تصادفات انجام می‌شود. بعد از آن روش‌های آماری همچون توابع همبستگی Anselin Local Moran's I و Gi ord-Getis، تابع برآورد تراکم کرنل بر روی داده‌ها اعمال شد. سپس مقایسه‌ای بین خروجی حاصل از روش‌های خوشه‌بندی غیر آماری و آماری صورت پذیرفت تا مشخص شود که کدام روش خوشه‌بندی غیر آماری مناطق حادثه‌خیز را بهتر مدل می‌کنند. سپس از بین روش‌های مذکور بر اساس داده‌های اولیه، روش خوشه‌بندی غیر آماری که مناطق

حادثه‌خیز را بهتر مدل می‌کند و همچنین بهترین روش از بین روش‌های آماری انتخاب می‌شوند. بعد از آن، روش‌های برگزیده روی داده‌های واقعی پیاده‌سازی می‌شوند و در نهایت نقاط حادثه‌خیز با داده‌های واقعی شناسایی می‌شوند.

۳-۲- محدوده مورد مطالعه

در این پایان‌نامه داده‌های مربوط به استان کوجائلی^۱ استفاده شده است. داده‌های مورد استفاده در این تحقیق شامل موارد زیر است: نقشه راه‌های استان کوجائلی در فرمت shp. و در قالب یک لایه خطی و فاقد اطلاعات توصیفی راه، موقعیت نقاط تصادفات موجود در فرمت xlsx. با مشخصات مکانی، شرایط جوی، شرایط هندسی مسیر، و... که به تشریح در ادامه بررسی شده است. با توجه به دقت بالای ثبت داده‌های تصادف در این استان که با GPS به صورت دقیق مختصات مکانی نقاط تصادف ثبت شده است لذا در تحقیقات دیگر نیز با داده‌های مکانی با مختصات دقیق می‌توان مطالعاتی از این دست انجام داد. ما ۲ شهر ایزمیت^۲، گبزه^۳ از این استان را به تفکیک مورد بررسی قرار دادیم. این شهرها در مسیر راه‌های شریانی درجه یک ترکیه و در کنار دریای مرمره می‌باشند. بر اساس آمار نفوس رسمی ترکیه سال ۲۰۱۸ جمعیت هر کدام از شهرها برای سال‌های ۲۰۱۳-۲۰۱۴ در جدول شماره ۱-۳ آمده است. در شکل ۳-۲ موقعیت جغرافیایی استان و شهرهای انتخابی قابل مشاهده است.

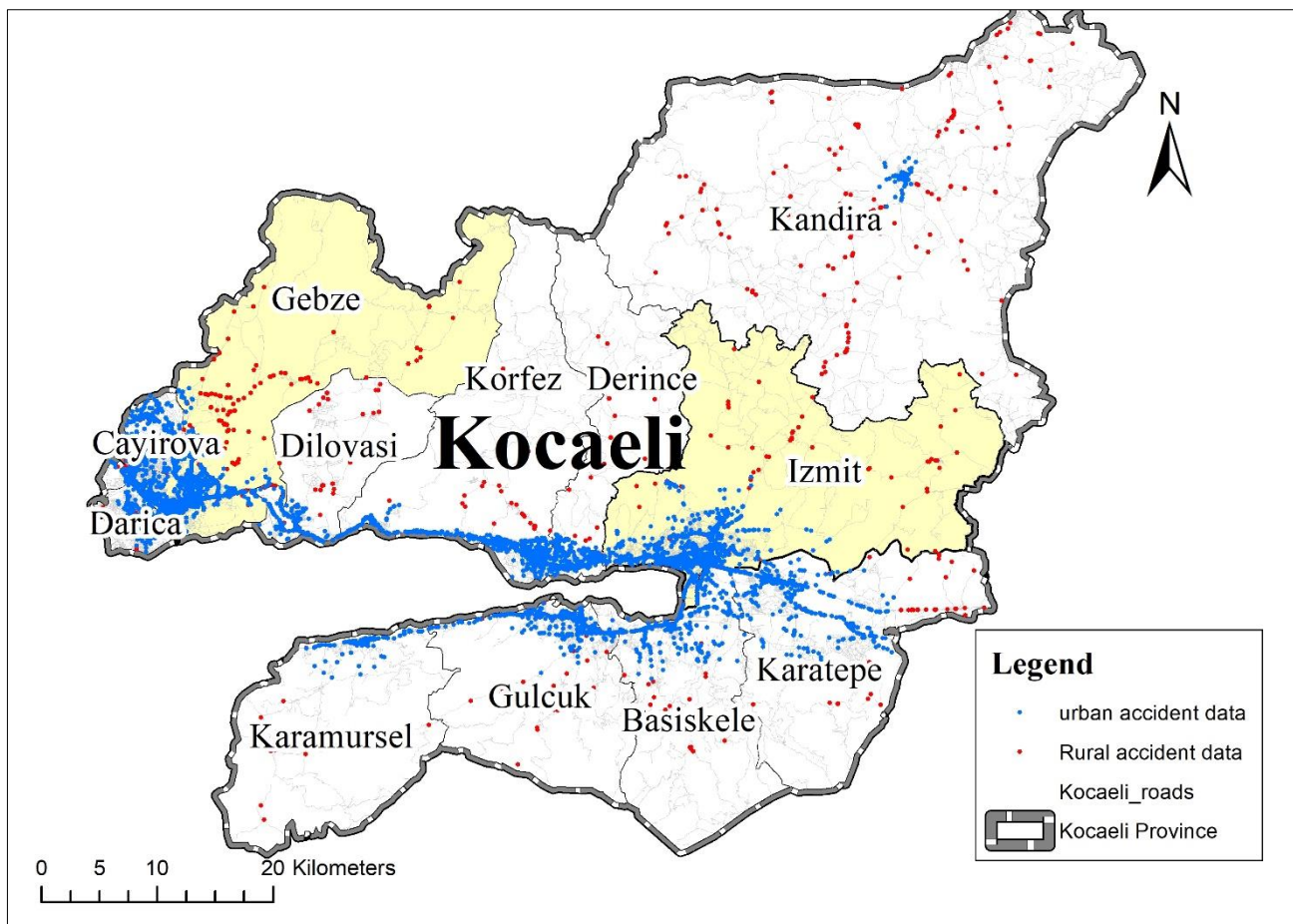
جدول ۱-۳: جمعیت شهرهای استان کوجائلی (سال ۲۰۱۳-۲۰۱۴) (Turkstat, ۲۰۱۸)

منطقه	۲۰۱۳	۲۰۱۴
کوجائلی	۱,۶۷۶,۲۰۲	۱,۷۲۲,۷۹۵
ایزمیت	۳۳۲,۷۵۴	۳۳۸,۷۱۰
گبزه	۳۲۹,۱۹۵	۳۳۸,۴۱۲

^۱ Kocaeli

^۲ İzmit

^۳ Gebze



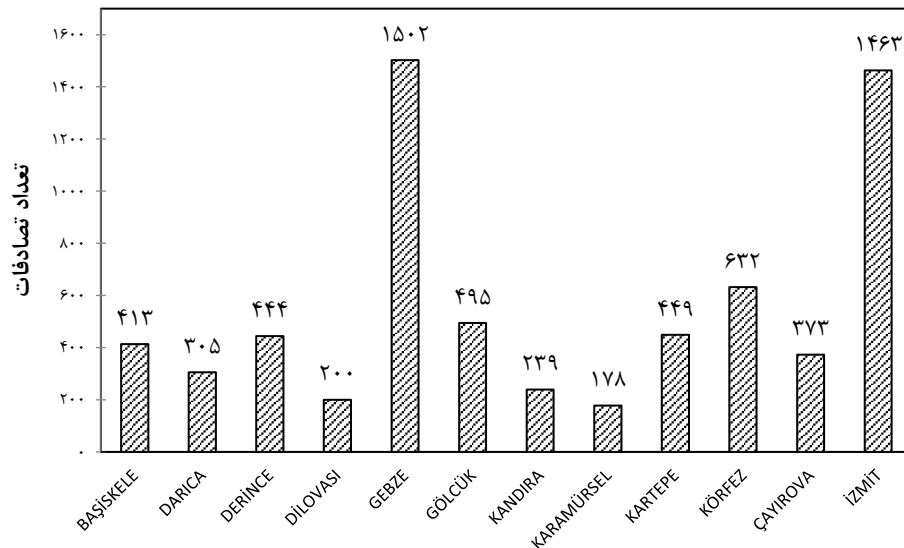
شکل ۳-۱: استان کوجائلی ترکیه و توزیع نقاط تصادف درون شهری و برون شهری برای سال ۲۰۱۳-۲۰۱۴



شکل ۳-۲: عکس ماهواره‌ای شهرهای اطراف دریای مرمره

در شکل ۳-۱ توزیع نقاط در کل استان قابل مشاهده است و نقاط آبی نشان‌دهنده نقاط درون-شهری است که تعداد تصادفات درون‌شهری به مراتب بیشتر است. در شکل ۳-۲ موقعیت جغرافیایی محدوده مورد مطالعه قابل مشاهده است.

با توجه به آمار گزارش شده از سازمان پلیس ملی ترکیه شکل ۳-۳ تعداد کل تصادفات در دو شهر گبزه و ایزمیت بیشترین تعداد تصادفات را در بر داشته است که به ترتیب ۱۵۰۲ و ۱۴۶۳ تصادف از ۶۶۸۹ تصادف این استان مربوط به این دو شهر است.



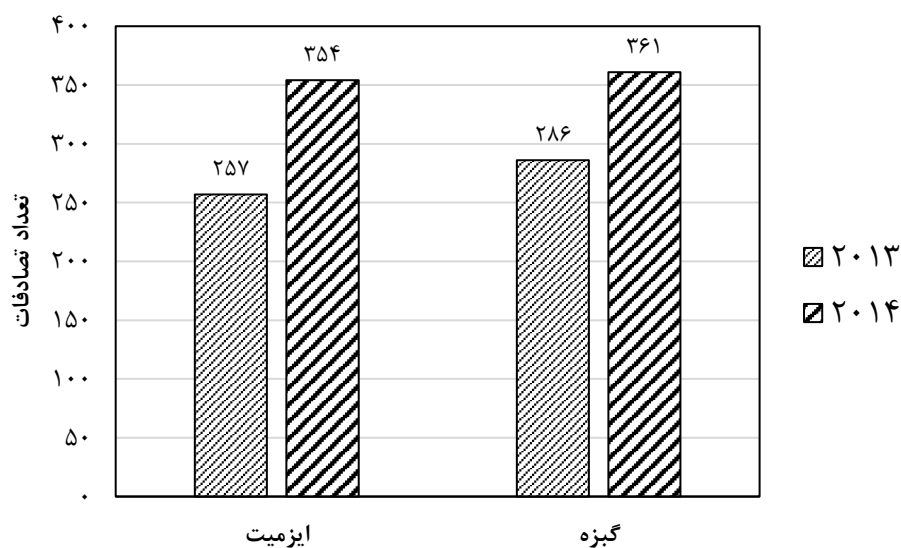
شکل ۳-۳: تعداد تصادفات کل به تفکیک شهرهای استان کوجائلی در سال ۲۰۱۳-۲۰۱۴

ابتدا شهر گبزه که واقع در ۶۵ کیلومتری (۳۰ متری) جنوب شرقی استانبول، در خلیج ایزمیت، بازوی شرقی دریای مرمره واقع شده است، بررسی کردیم. گبزه دومین منطقه بزرگ در کوجائلی پس از ایزمیت، مرکز استانی از نظر جمعیت است. گبزه در سال های اخیر رشد سریع جمعیتی را تجربه کرده است، از ۱۵۹۱۱۶ نفر در سال ۱۹۹۰ به ۲۹۹۰۴۷ نفر در سال ۲۰۱۱. منطقه گبزه در قسمت غربی استان کوجائلی واقع شده است، با همسایگان توزلا در غرب؛ پندیک، شمال غربی؛ شیل، شمال؛ و کورفز، به شرق. این ناحیه شامل ۶ شهرداری (گبزه، داریچا، چایروا، دیلووآسی، شکرپینار و توسانچیل) و ۲۲ روستا است. انتهای شمالی پل عثمان Gazi در این منطقه قرار دارد؛ این پل با طول ۴ کیلومتر دریای مرمره را از کابابورون به دیلبورونو را وصل می کند.

سپس شهر ایزمیت واقع در قسمت شرقی این استان که در مسیر محور اصلی D100 (جاده متصل کننده شرق به غرب) و در کنار دریای مرمه قرار دارد، مورد بررسی قرار دادیم. ساختار توپوگرافی ایزمیت از دو بخش تشکیل شده است. ابتدا در دشت‌های مسطح ساخته شده بود که مرکز شهر می‌باشد که شبکه‌های راه آهن و بزرگراه از این منطقه عبور می‌کنند که نزدیک مرز دریایی است. بخش دوم بر روی تپه‌ها ساخته شد که بسیاری از خانه‌های تاریخی از دوران عثمانی در این بخش ساخته شده است.

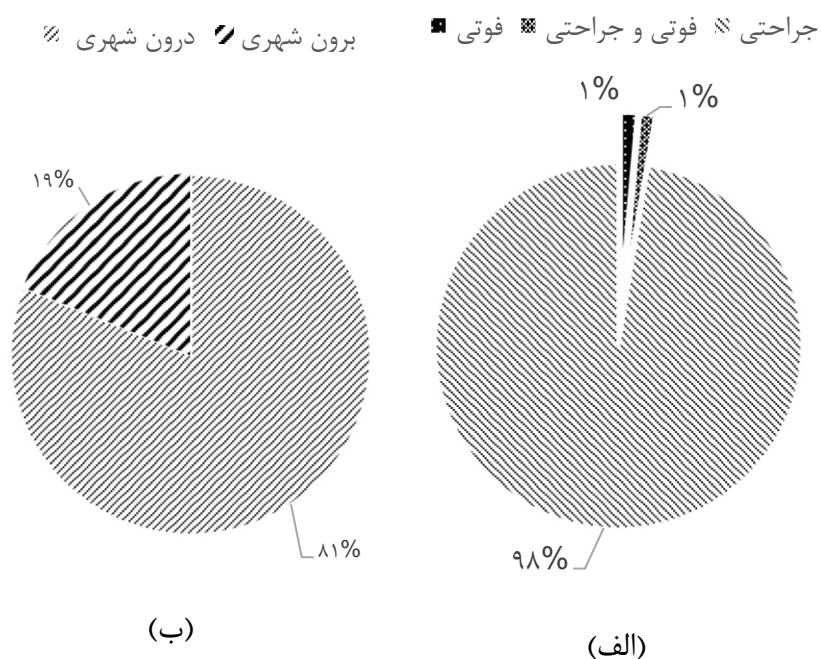
۲-۳- معرفی داده‌ها

تعداد کل تصادفات در این استان برای سال ۲۰۱۳-۲۰۱۴ به تعداد ۶۶۸۹ تصادف بود. آمار تصادفات برای مناطق درون‌شهری ۵۶۸۹ تصادف و برون‌شهری ۱۰۰۰ تصادف بود. پاکسازی داده‌ها در دو جهت سطری و ستونی (اصلاح سطری یعنی اصلاح تک به تک رکوردها از لحاظ کامل بودن پارامترها برای هر رکورد تصادف و اصلاح ستونی یعنی اصلاح زیر شاخه‌های هر پارامتر در هر ستون) انجام شد. برای مثال تعدادی از ستونها به دلیل ثبت ناقص یا تعداد کم رکوردهای تصادفات حذف شد. پس از پاکسازی داده‌ها ۶۴۷ و ۶۱۱ تصادف درون‌شهری به ترتیب برای شهر گبزه و ایزمیت برای تحلیل و اعمال روشهای خوشه‌بندی مورد استفاده قرار گرفت. در شکل ۳-۴ تعداد تصادفات درون‌شهری (پس از پاکسازی) به تفکیک برای سال ۲۰۱۳-۲۰۱۴ قابل مشاهده است.

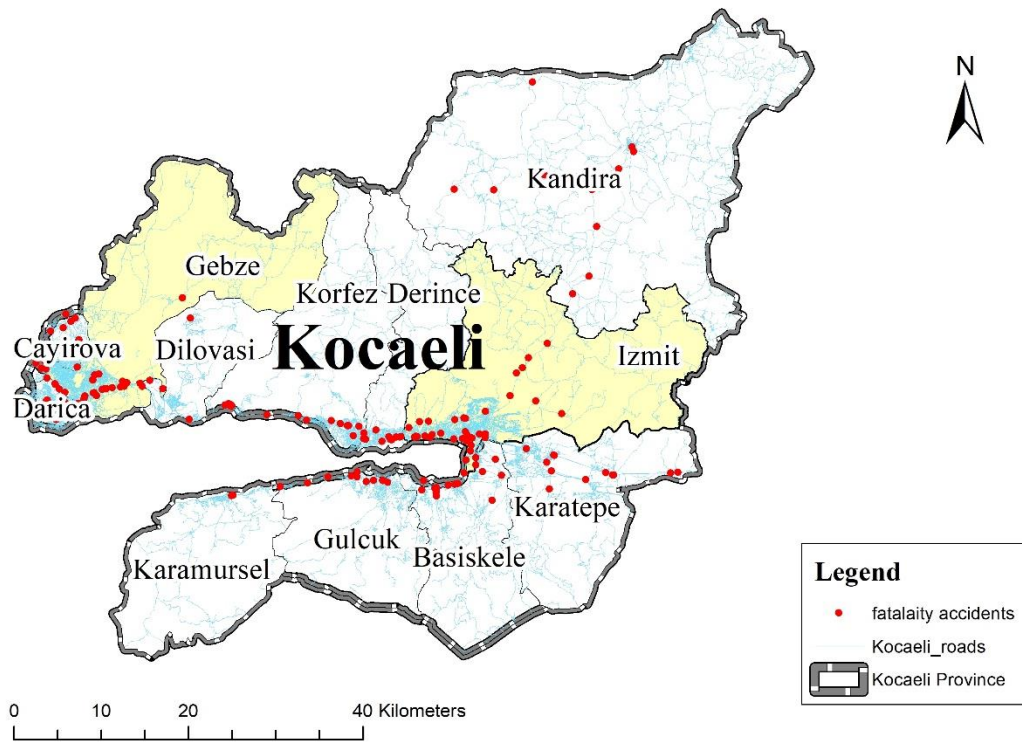


شکل ۳-۴- تعداد تصادفات درون شهری به تفکیک شهرها در سال ۲۰۱۳ و ۲۰۱۴ پس از پاکسازی داده‌ها

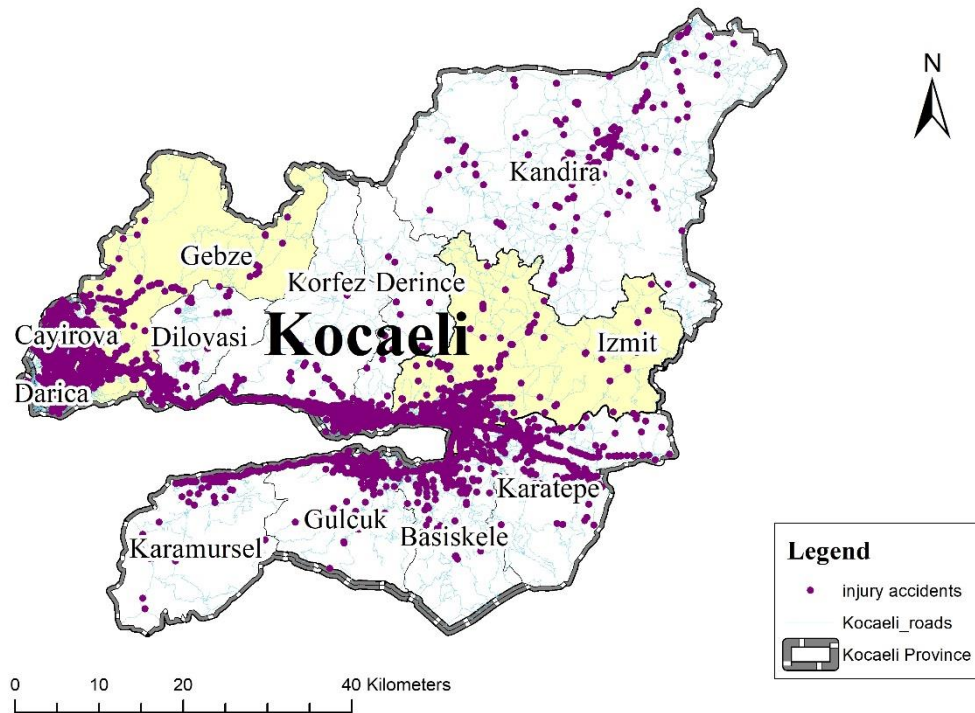
با توجه به فراوانی داده‌های درون شهری شکل ۳-۵ صرفاً از داده‌های درون شهری استفاده شد و همینطور به دلیل ناچیز بودن تعداد تصادفات فوتی گزارش شده شکل ۳-۵ در داده‌ها صرفاً از داده‌های جراحی استفاده شد. در شکل ۳-۶ نیز بر روی نقشه استان کوجائلی قابل تشخیص می‌باشد.



شکل ۳-۵: درصد تصادفات (الف) فوتی و جراحی در داده‌های (ب) درون شهری و برون شهری در کل استان



(الف)



(ب)

شکل ۳-۶: نقاط تصادف فوتی (الف) و جراحی (ب) در کل استان کوجاالی

۳-۴- چارچوب تعیین نقاط پرتراکم با روش‌های خوشه‌بندی

۳-۴-۱- آماده سازی داده ها

برای انجام صحیح تحلیل های GIS و همینطور خوشه‌بندی مناسب به داده‌هایی با کیفیت بالا نیاز است. بنابراین اطلاعات ورودی باید آماده‌سازی گردد. در این تحقیق نیز با توجه به دقت داده‌های تصادف و نیز انتخاب روش‌های مکانی برای تحلیل تصادفات در مناطق و شبکه درون‌شهری، جهت آماده‌سازی داده‌ها انجام شد. پس از آن، داده‌های تصادف را هم در خوشه بندی هم در تحلیل های GIS مورد استفاده قرار گرفت. آماده سازی داده ها بدین ترتیب بود که ابتدا نقاط کل تصادف برای این استان را که شامل همه ی تصادفات رخ داده در سال ۲۰۱۳-۲۰۱۴ بود را از لحاظ جامعه آماری برای نقاط تصادفات درون‌شهری و برون‌شهری مقایسه کردیم و با توجه به تعداد قابل قبول تصادفات در نقاط درون‌شهری این نقاط برای مطالعه انتخاب گردید. سپس تک تک رکوردهای تصادف را از لحاظ کامل بودن داده ها در تمام پارامترهای مرتبط با تصادف بررسی و پاکسازی انجام شد.

در این نوع تحلیل ها هر چقدر منطقه مورد مطالعه مخصوصا برای هر شهر مقیاس کوچکتری باشد، تحلیل های دقیق تری انجام خواهد شد. به همین علت، در انتها با توجه به بررسی مطالعات پیشین و توصیه های محققین، داده های هر شهر را تفکیک نموده و جداگانه برای خوشه بندی و تحلیل های GIS مورد استفاده قرار دادیم.

عوامل ثبت شده در داده‌های موجود به تعداد ۴۰ پارامتر بود اما به دلیل کامل نبودن عوامل تنها ۲۰ پارامتر از توصیفی و عددی از کل گزارش سازمان پلیس ترکیه مورد استفاده قرار گرفت. در جدول ۳-۳ سه داده عددی (سرعت مجاز و تعداد خودروهای درگیر تصادف و تعداد افراد مجروح) و جدول ۴-۳ داده‌های توصیفی را به همراه فراوانی داده‌های تصادف، نشان می‌دهد که شامل ۱۷ داده کیفی می‌باشد.

جدول ۳-۲- خلاصه داده‌های عددی به تفکیک شهر برای سال ۲۰۱۳-۲۰۱۴

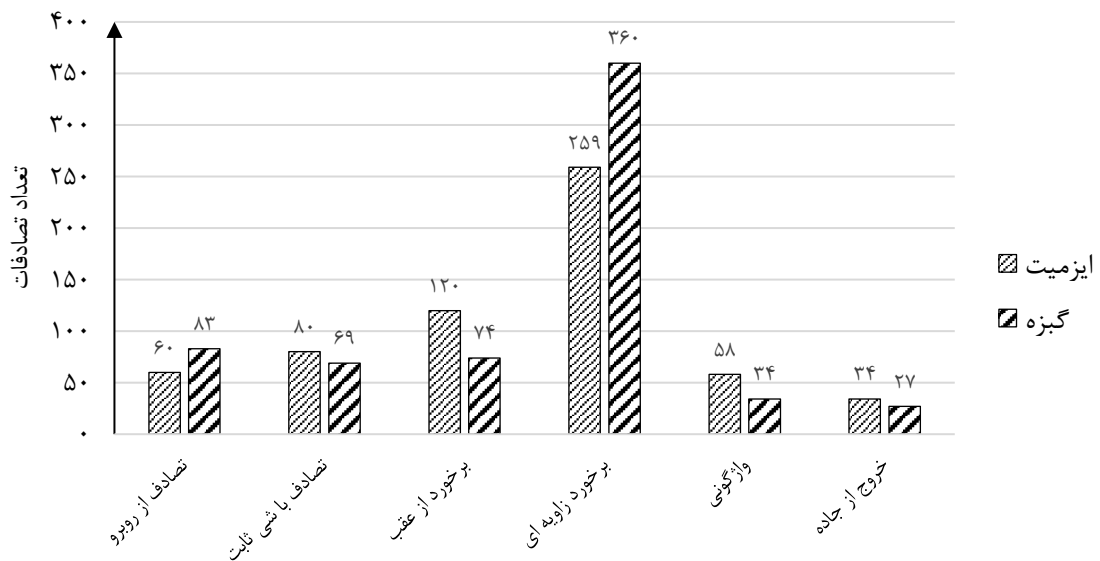
متغیرهای عددی		کمینه	بیشینه	میانگین	انحراف معیار
تعداد مجروحین در هر تصادف	گبزه	۱	۱۷	۱/۷۵۸	۱/۴۱۱
	ایزمیت	۱	۹	۱/۶۹۹	۱/۱۷۷
تعداد خودروی درگیر در هر تصادف	گبزه	۱	۴	۱/۹۴۶	۰/۵۴۲
	ایزمیت	۱	۶	۱/۹۵۳	۰/۷۴۲

جدول ۳-۳: داده‌های توصیفی تصادفات جرحی درون‌شهری شهر گبزه و ایزمیت برای سال ۲۰۱۳-۲۰۱۴

ردیف	متغیرهای توصیفی	زیرشاخه	گبزه	ایزمیت
			فراوانی نسبی / (مطلق)	فراوانی نسبی / (مطلق)
۱	نوع راه	جداشده	۴۸/۷ (۳۱۵)	۵۵/۱ (۳۳۷)
		یک‌طرفه	۶/۵ (۴۲)	۹/۳ (۵۷)
		دوطرفه جدا نشده	۴۴/۸ (۲۹۰)	۳۵/۶ (۲۱۷)
۲	طبقه راه	بزرگراه	۷/۵ (۴۹)	۶۹/۶ (۴۲۵)
		خیابان	۹۲/۵ (۵۹۸)	۳۰/۴ (۱۸۶)
۳	وجود گارد ریل	خیر	۸۳/۱ (۵۳۸)	۷۱/۷ (۴۳۸)
		بله	۱۶/۹ (۱۰۹)	۲۸/۳ (۱۸۳)
۴	وجود شانه مسیر	خیر	۸۵/۵ (۵۵۳)	۷۵/۱ (۴۵۲)
		بله	۱۴/۵ (۹۴)	۲۴/۹ (۱۵۹)
۵	وجود خط‌کشی	خیر	۳۳/۴ (۲۱۶)	۲۵ (۱۵۳)
		بله	۶۶/۶ (۴۳۱)	۷۵ (۴۵۸)
۶	علائم راهنمایی رانندگی	خیر	۵۵,۷۳ (۳۲۹)	۴۸/۵ (۲۹۶)
		بله	۴۴,۲۶ (۳۱۸)	۵۱/۵ (۳۱۵)
۷	علائم چراغ‌دار	خیر	۹۲/۵ (۵۹۹)	۸۵ (۵۲۵)
		بله	۷/۵ (۴۸)	۱۵ (۸۶)
۸	تقاطع	چهارراه	۲۵/۶ (۱۶۶)	۱۸/۲ (۱۱۱)
		دور از تقاطع	۳۵/۱ (۲۲۷)	۵۷/۸ (۳۵۳)
		میدان	۸/۳ (۵۴)	۳ (۱۸)
		تقاطع T شکل	۲۷/۲ (۱۷۶)	۱۵/۸ (۹۷)
		تقاطع Y شکل	۴,۸۸ (۲۴)	۵/۲ (۳۲)
۹	محل تصادف	میان راه	۲/۹ (۱۹)	۶/۲ (۳۸)

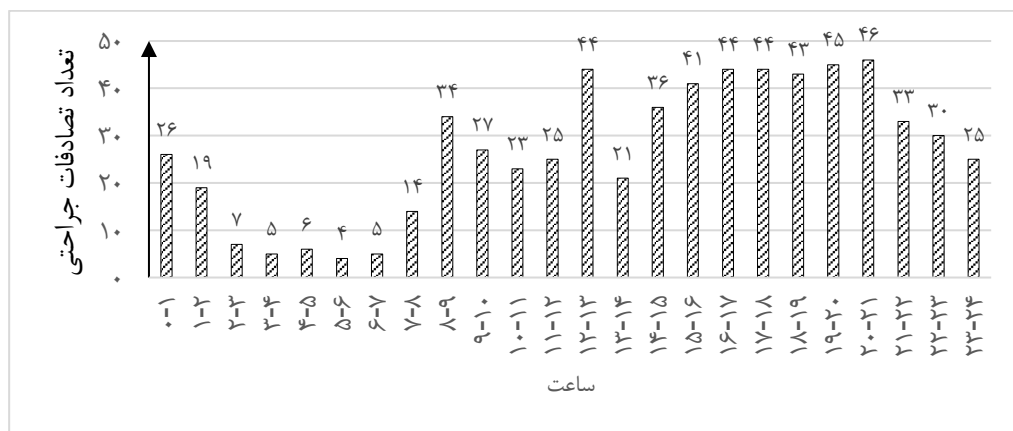
		روی مسیر	%۸۷/۹ (۵۶۹)	%۸۱/۸ (۵۰۰)
		در شانه راه	%۱/۶ (۱۰)	%۲/۴ (۱۵)
		در کناره راه	%۴/۳ (۲۸)	%۵/۵ (۳۳)
		در پیاده‌رو	%۳/۳ (۲۱)	%۴/۱ (۲۵)
۱۰	روز و شب	روز	%۴۸/۷ (۳۱۵)	%۵۵/۳ (۳۳۸)
		شب	%۵۱/۳ (۳۳۲)	%۴۴/۷ (۲۷۳)
۱۱	آب‌وهوا	هوای صاف	%۸۶/۲ (۵۵۸)	%۸۳/۷ (۵۱۱)
		بارانی / برفی	%۱۳/۸ (۸۹)	%۱۶/۳ (۱۰۰)
۱۲	وضعیت سطح جاده	خشک	%۸۱/۶ (۵۲۸)	%۷۸/۵ (۴۸۰)
		خیس	%۱۸/۴ (۱۱۹)	%۲۱/۵ (۱۳۱)
۱۳	هندسه افقی مسیر	قوس	%۱۴/۷ (۹۵)	%۱۶/۲ (۹۹)
		مسیر مستقیم	%۸۵/۳ (۵۵۲)	%۸۳/۸ (۵۱۲)
۱۴	هندسه قائم مسیر	شیب‌دار	%۳۵/۵ (۲۳۰)	%۲۵/۳ (۱۵۵)
		شیب صفر	%۶۴/۵ (۴۱۷)	%۷۴/۷ (۴۵۶)
۱۵	روشنایی مسیر	خیر	%۴۳/۴ (۲۸۱)	%۲۰/۶ (۱۲۶)
		بله	%۵۶/۶ (۳۶۶)	%۷۹/۴ (۴۸۵)
۱۶	نوع تصادف	تصادف از روبرو	%۱۲/۸ (۸۳)	%۹/۸ (۶۰)
		برخورد با شیء ثابت	%۱۰/۷ (۶۹)	%۱۳/۱ (۸۰)
		برخورد از عقب	%۱۱/۴ (۷۴)	%۱۹/۹ (۱۲۰)
		برخورد زاویه‌ای	%۵۵/۶ (۳۶۰)	%۴۲/۴ (۲۵۹)
		واژگونی	%۵/۳ (۳۴)	%۹/۵ (۵۸)
		خروج از جاده	%۴/۲ (۲۷)	%۵/۵ (۳۴)
۱۷	تعداد خودرو	چند خودرویی	%۸۲/۷ (۵۳۵)	%۷۶ (۴۶۴)
		تک خودرویی	%۱۷/۳ (۱۱۲)	%۲۴ (۱۴۷)

انواع تصادف در داده های موجود را می توان در شکل ۳-۷ مشاهده کرد. همانطور که قابل ملاحظه است تصادف از نوع برخورد زاویه‌ای از بیشترین تصادفات است چرا که در تقاطعات بیشترین نوع برخورد، برخورد زاویه‌ای است. برای در نظر داشتن معناداری داده ها آزمون خی دو بر روی هر کدام از پارامترها انجام شد. که در جدول ۳-۴ آمده است.

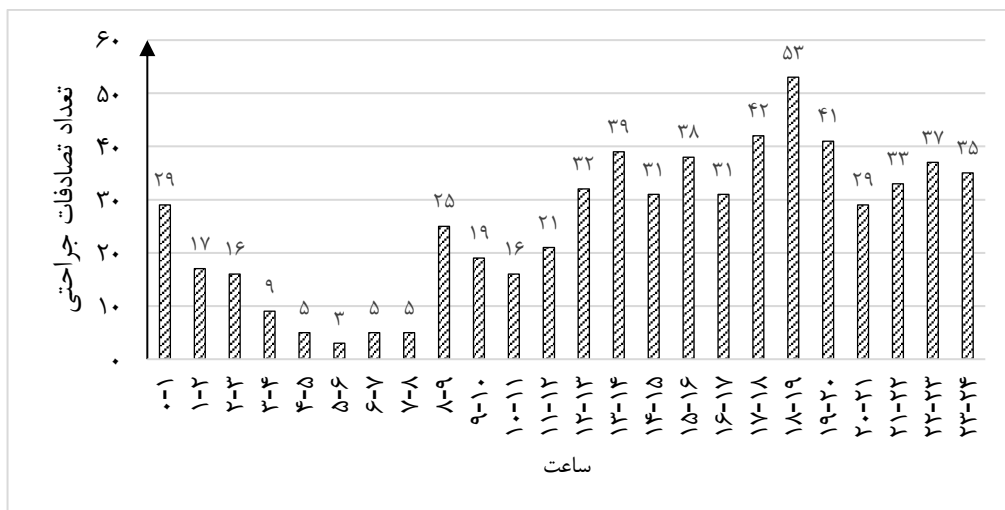


شکل ۳-۷: انواع تصادف در مناطق درون شهری به تفکیک شهر ایزمیت و گبزه در سالهای ۲۰۱۳-۲۰۱۴

علاوه بر نوع تصادف داده ها، زمان وقوع تصادف نیز مورد توجه قرار گرفته است. در شکل ۳-۸ و ۳-۹ توزیع ساعتی اطلاعات مربوط به تصادفات جراحی دو شهر نشان داده شده است. طبق این گزارش، تعداد حوادث به خصوص در ساعات شب بسیار زیاد است. از آنجایی که مطالعه در حوادث ساعت ها دشوار است، تصادفات محاسبه می شود و با توجه به ساعت های پیک و غیر اوج، دوره های زمانی محاسبه می شوند.

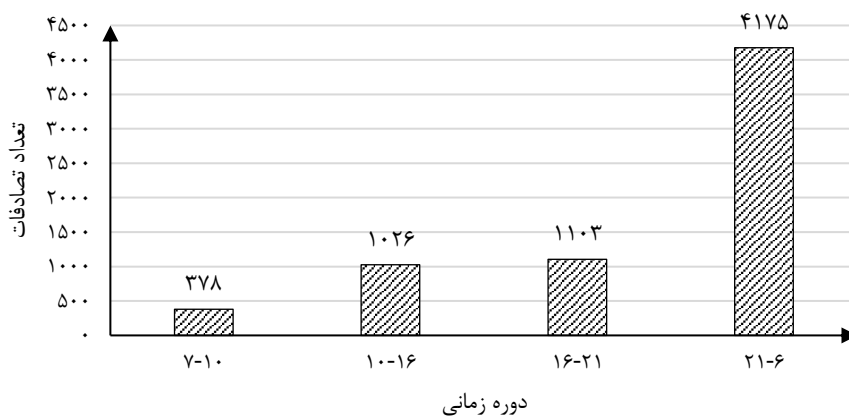


شکل ۳-۸: شکل ۳-۹: توزیع زمانی تصادفات جراحی شهر گبزه (۶۴۷) در سال ۲۰۱۳-۲۰۱۴



شکل ۳-۱۰: توزیع زمانی تصادفات جراحی شهر ایزمیت (۶۱۱) در سال ۲۰۱۳-۲۰۱۴

با توجه به زمان وقوع تصادفات جراحی در دوره های زمانی مختلف در شکل ۳-۱۰، از تعداد کل تصادفات ۶۶۹۸؛ بین ساعت ۷ تا ۹ صبح ۸۶۱ (۱۲/۸ درصد) تصادف مشاهده شد که کمترین مقدار تصادفات در ساعت ترافیک اوج صبح بوده است. بین ۱۰ صبح تا ۳ بعد از ظهر ۲۰۹۶ (۳۱/۱ درصد)، بین ساعت ۴ بعد از ظهر تا ۸ شب ۲۲۶۱ (۳۳/۵ درصد) و بین ساعت ۹ شب تا ۶ صبح ۱۵۳۲ (۲۲/۷ درصد) دیده شد. تعداد تصادفات مربوط به هر کدام از دوره های زمانی از یکدیگر متفاوت هستند. این تفاوت با کمک آزمون Chi-Square مورد آزمایش قرار گرفت. براساس نتایج، توزیع حوادث مربوط به دوره های زمانی از لحاظ آماری معنی دار هستند.



شکل ۳-۱۱: تعداد تصادفات به تفکیک دوره های زمانی در

در طول دو سال ۲۰۱۳-۲۰۱۴ از ۶۶۹۸ تصادفات در این استان تعداد ۱۱۴۷۴ نفر مجروح شده اند. از این تعداد مجروح ۱۲۷۳ نفر تصادف با عابر پیاده بود که به علت کم بودن جامعه آماری نسبت به تعداد کل مقدار قابل توجهی برای بررسی تصادفات عابرین پیاده نبود.

شناسایی نقاط حادثه‌خیز می‌تواند هم با کل داده‌های تصادف، هم با داده‌های تفکیک شده بر حسب نوع تصادف یا تعداد افراد کشته یا مجروح انجام شود. نقاط حادثه‌خیز شناسایی شده به صورت نقشه موضوعی با تعداد افراد کشته یا مجروح که خروجی خوشه‌بندی های مختلف به همراه نقطه حادثه‌خیز می‌باشد. از آنجایی که روش‌های تحلیل موجود در GIS از روش‌های پرکاربرد مورد استفاده در پژوهش‌های اخیر مربوط به شناسایی نقاط حادثه‌خیز می‌باشد نتایج دیگر خوشه‌بندی‌ها باید با این روش مقایسه شود. روش‌های مورد استفاده از نرم افزارها برای بررسی تراکم و خوشه‌بندی تصادفات به ترتیب زیر است:

- خوشه‌بندی Nnh و Kmeans در نرم افزار CrimeStat

- خوشه بندی GridBSACN در نرم‌افزار Elki

- روش تراکم کرنل صفحه‌ای و شبکه‌ای، در نرم افزار GIS

بعد از به دست آمدن خوشه‌ها، برای ارزیابی عوامل مؤثر بر تصادفات می‌توان با آزمون‌های آماری همچون ANOVA معنی‌داری داده‌های خوشه بندی را بررسی نمود.

۳-۵- خوشه‌بندی Nnh و K-Means در نرم افزار CrimeStat

همانطور که در فصل دوم ذکر شده است، تکنیک‌های بسیار زیادی وجود دارد که در آن ایده اصلی محاسبه تعدادی نقاط یا داده مانند تصادفات در یک منطقه محدود که توسط محقق تعیین شده است، می‌باشد (لویین، ۲۰۱۰). در این مطالعه، روش تحلیل خوشه‌ای Nnh برای تشخیص نقاط حادثه‌خیز

تصادف استفاده می شود که بر شناسایی گروه هایی از داده ها که به صورت مکانی نزدیک هستند، متمرکز می شود. در خوشه بندی N_{nh} ، فاصله اقلیدسی بین هر جفت نقطه داده بررسی و به عنوان یک معیار برای خوشه بندی استفاده می شود. معیار اول فاصله جستجو (d) است که پس از تعیین آن، نقاط با فاصله های کوچکتر در یک خوشه قرار می گیرند. معیار دوم حداقل تعداد نقاط n_{min} در خوشه می باشد. سپس نقاط منطبق با معیارهای فاصله جستجو (d) و n_{min} به عنوان یک خوشه شناسایی می شود. اگر تعداد نقاط تعریف شده نزدیکتر از فاصله جستجو باشد، یک خوشه جدید تولید می شود. پس از محاسبه خوشه های مرتبه اول، خوشه های مرتبه دوم و مرتبه بالا با همان شیوه تشکیل می شوند تا زمانی که فقط یک خوشه باقی بماند یا با معیارهای مذکور تطابق نداشته باشد. به همین دلیل، الگوریتم خوشه بندی N_{nh} خوشه ای از تمام نقاط در منطقه مورد مطالعه ایجاد نمی کند. از آنجایی که این روش به معیارهای مربوط به فاصله جستجو (d) و n_{min} بستگی دارد، تنها نقاطی که معیارهای مورد نیاز را برآورده می کنند، خوشه بندی می شوند.

همانطور که در بالا ذکر شد، چندین روش تحلیل مکانی وجود دارد که می تواند مورد استفاده قرار گیرد. یکی دیگر از روش های شناخته شده تحلیل خوشه های نقاط، خوشه بندی K -Means است که داده ها را به گروه های K که توسط کاربر تعریف شده است، تقسیم می کند. این روش K تعداد خوشه از تمام داده ها ایجاد می کند. تمام نقاط به یک گروه اختصاص داده می شوند و تمام نقاط در این روش خوشه بندی می شوند و در یک خوشه قرار می گیرند. این روش زمانی می تواند مفید باشد که می خواهیم گروه بندی تمام داده ها را کنترل کنیم. برای مثال، اگر قرار است ۵ منطقه از بهترین مناطق ایستگاه آمبولانس را پیدا کنیم، روش خوشه بندی K -means با تعریف ۵ خوشه یک ایده در مورد متراکم ترین مناطق می دهد. برخلاف روش خوشه بندی K -Means، الگوریتم خوشه بندی N_{nh} امکان کنترل اندازه گروه بندی با تعیین فاصله آستانه و n_{min} در خوشه را فراهم می کند. خوشه بندی N_{nh} فرصتی را

برای شناسایی محدوده جغرافیایی متراکم داده‌ها فراهم می‌آورد. این مزیت الگوریتم Nnh، شناسایی نقاط حادثه‌خیز را برای هدف این مطالعه فراهم می‌کند.

۳-۵-۱- نحوه کار در نرم‌افزار CrimeStat

اجرای الگوریتم Nnh در نرم‌افزاری به نام CrimeStat IV قابل اجرا است. با توجه به شکل ۳-۱۲ مراحل اجرای نرم‌افزار بدین ترتیب است:

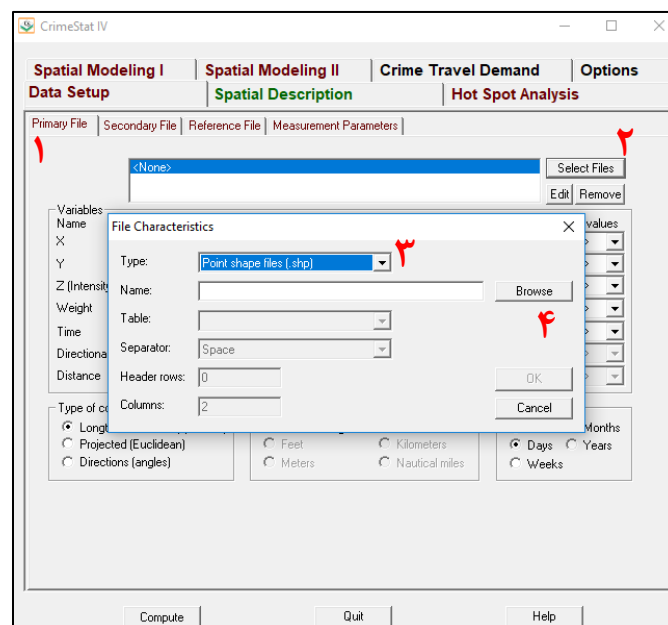
۱- کلیک بر روی سربرگ Primary File

۲- سپس بر روی Select Files کلیک می‌کنیم.

۳- سپس در پنجره باز شده File Characteristics در قسمت Type نوع Point shape

files (.shp) را انتخاب می‌کنیم

۴- سپس از Browse فایل داده‌ها به فرمت shp. وارد نرم‌افزار می‌شود.



شکل ۳-۱۲: مرحله اولیه وارد کردن فایل shp داده‌های تصادف

در ادامه باتوجه به شکل ۳-۱۳ به ترتیب:

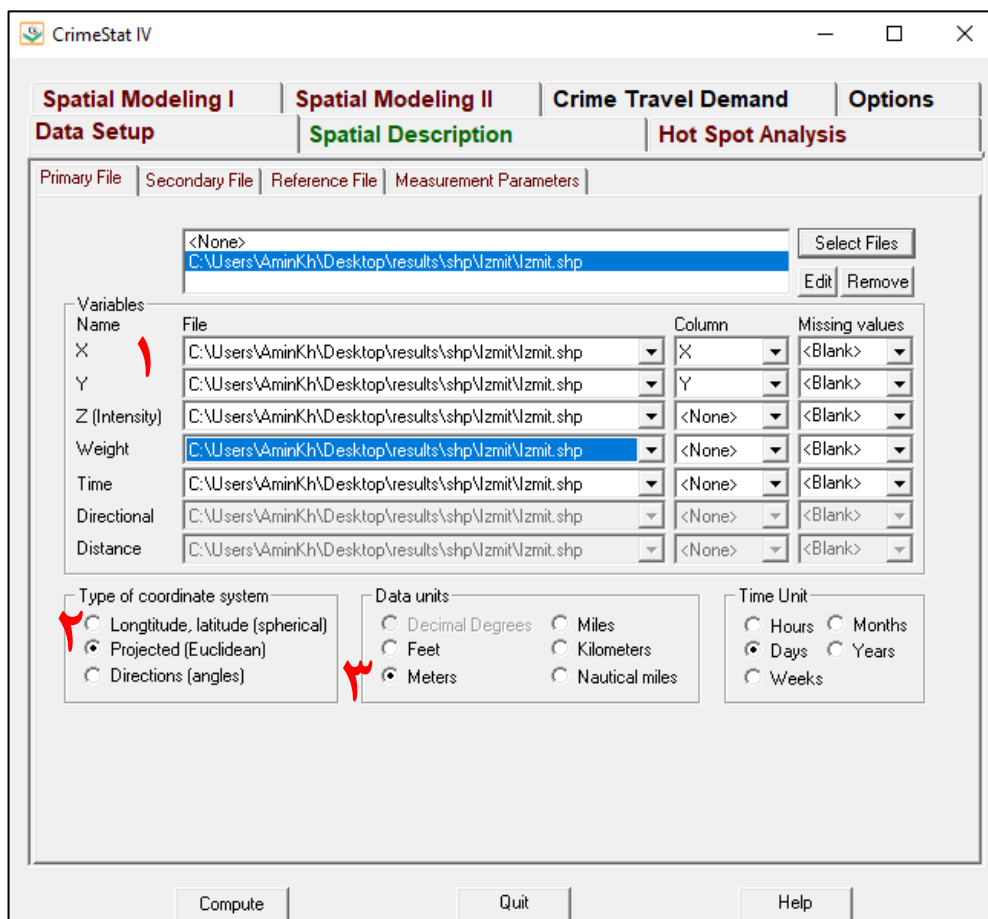
۱- مشخصات جغرافیایی X و Y را در قسمت Column انتخاب می‌کنیم.

۲- سپس در قسمت Type of coordinate system نوع سیستم مختصات

projected(Euclidean) یا همان مختصات اقلیدسی که به از نوع UTM می‌باشد را

انتخاب می‌کنیم.

۳- در قسمت Data units واحد طول متر را انتخاب می‌کنیم.

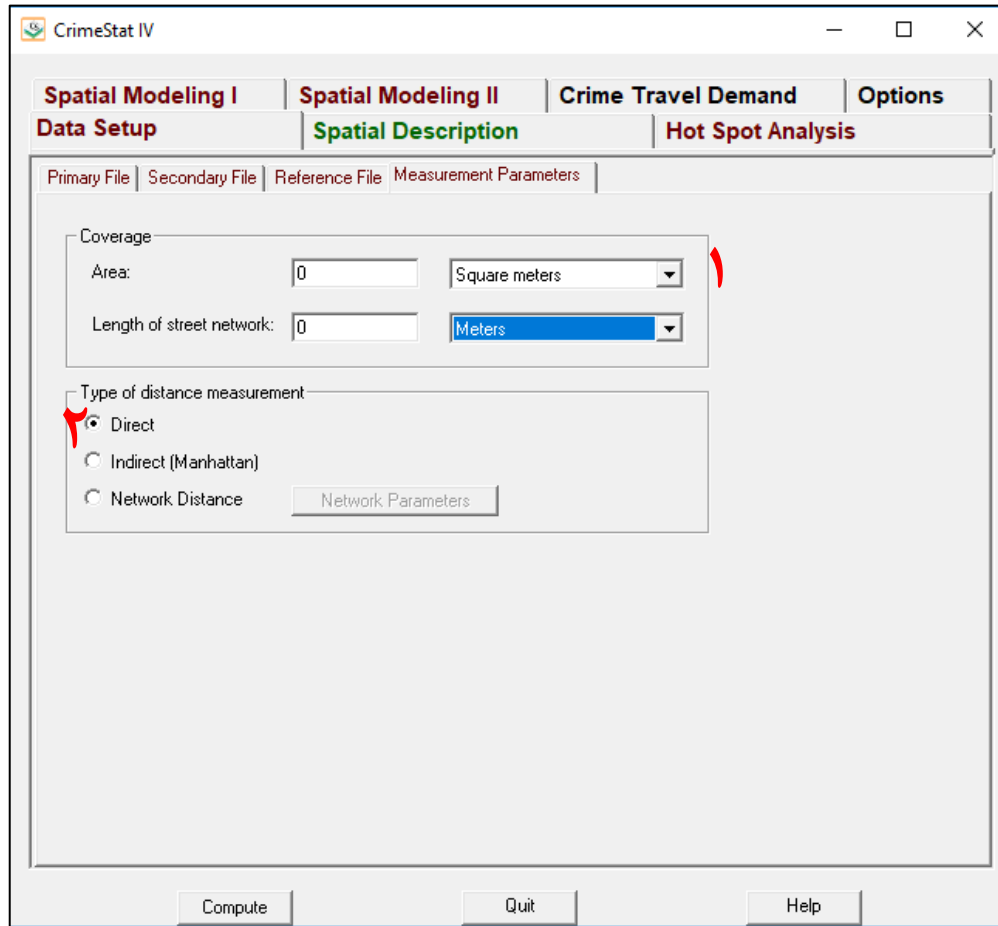


شکل ۳-۱۳: وارد کردن مشخصات جغرافیایی

با توج به شکل ۳-۱۴ در سربرگ Measurement Parameters داریم:

۱- واحد مساحت را متر مربع انتخاب می‌کنیم.

۲- نوع اندازه‌گیری فاصله را Direct که همان نوع اقلیدسی است انتخاب می‌کنیم.



شکل ۳-۱۴: پارامترهای اندازه‌گیری

در سربرگ Hot Spot Analysis در سربرگ اول Hot Spot Analysis I خوشه‌بندی Nnh را فعال و با

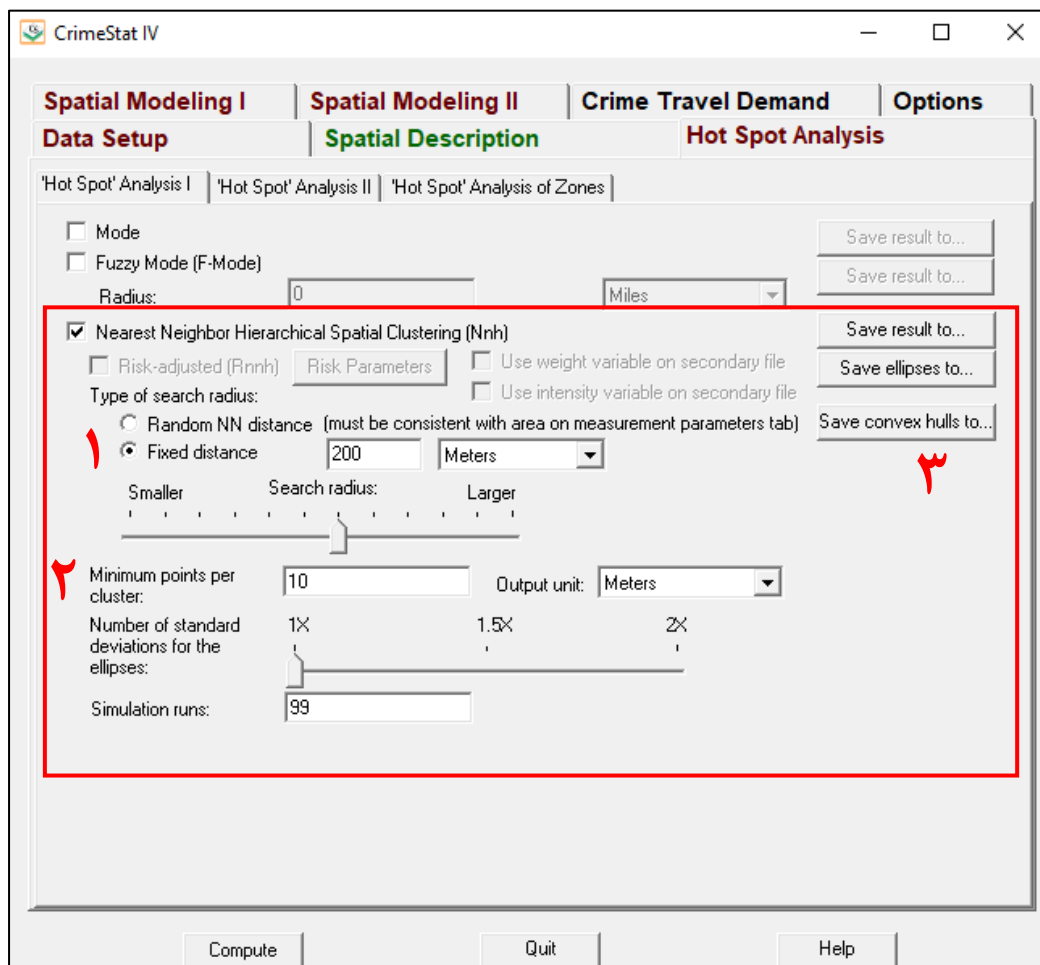
ترتیب پارامترها را وارد می‌کنیم. با توجه به شکل ۳-۱۵ داریم:

۱- انتخاب Fixed distance که همان پارامتر d یا همان شعاع جستجو خوشه است و در قسمت

کناری واحد متر را برای شعاع در نظر می‌گیریم.

۲- در قسمت Minimum points per cluster که همان n_{min} می‌باشد را وارد می‌کنیم.

۳- در سمت راست سه نوع ذخیره خروجی نتایج وجود دارد. گزینه اول به صورت جدولی است. مختصات نقاط مرکز خوشه بدست آمده از الگوریتم را می‌توان در محیط GIS مشاهده نمود. گزینه دوم Save Ellipse to... برای خروجی خوشه‌ها به صورت بیضی شکل و گزینه سوم Save convex hull to... به صورت چندضلعی حاصل از وصل نقاط مرزی هر خوشه است.



شکل ۳-۱۵: محیط نرم‌افزار CrimeStat IV برای محاسبه خوشه‌بندی NNH

۳-۶- خوشه بندی GridBSCAN

همانطور که در فصل دوم اشاره شد، در این مرحله ما باید مقادیر مختلف Eps را برای هر سطح چگالی پیدا کنیم. سپس ما از این مقادیر Eps برای خوشه بندی برای پیدا کردن تعداد صحیح خوشه

ها استفاده خواهیم کرد. اولاً، برای تعیین پارامترهای Eps و $MinPts$ ، باید رفتار فاصله را از نقطه به نزدیکترین همسایه $k-k$ ، یعنی $k-dist$ ببینیم. هر کدام از $k-dist$ برای تمام نقاط داده به ازای برخی از مقادیر (k) محاسبه می شود. پس از آن انتظار می رود که تغییر ناگهانی در گراف به دست آید. این تغییر شدید در مقدار $k-dist$ با مقدار مناسب Eps برای هر سطح چگالی مجموعه داده مطابقت دارد. که خط افقی عمود بر محور عمودی مقدار Eps را به ما می دهد. نمودار $K-dist$ برای انتخاب نه تنها پارامترهای Eps ، بلکه تجزیه و تحلیل سطح چگالی مجموعه داده نیز کشیده می شود. توجه داشته باشید که مقدار Eps که در این روش تعیین می شود بستگی به k دارد، اما با تغییر k مقدار Eps به طور چشمگیری تغییر نمی کند (گائونکار و ساوانت، ۲۰۱۳).

اگر مجموعه داده ها با تراکم متنوعی باشد، بستگی به تراکم خوشه و توزیع تصادفی نقاط، نقاط دارای سطح چگالی یکسان، محدوده‌ی تنوع بزرگی نخواهد داشت، در حالی که تغییرات شدیدی انتظار می رود بین دو سطح چگالی دیده شود. به این ترتیب، منحنی های صاف و متفاوتی وجود خواهد داشت.

۳-۶-۱- حذف نویزها:

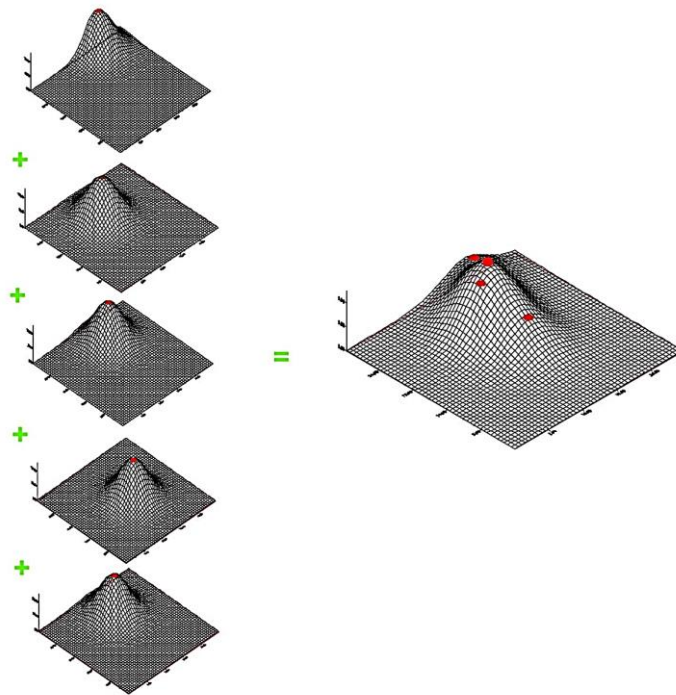
هر مجموعه داده ای تقریباً همیشه حاوی نویز است. این داده‌ها به هیچ یک از خوشه ها تعلق ندارند. به این معنی که محدوده های ناخوشه‌ها به طور کلی نسبت به نقاط خوشه ها پراکنده هستند، و فاصله یک ناخوشه به نزدیکترین خوشه فاصله زیادی نسبت به فاصله نقاط خوشه ها در هر خوشه است. هر روش خوشه بندی نیاز به مکانیسمی برای از بین بردن خوشه‌ها دارد.

در الگوریتم DBSCAN ۸۸، ناخوشه‌ها با توجه به فاصله‌ی بیشترشان از نقاط دیگر، کمتر تمایل دارند با نقاط دیگر ادغام شوند و رشد کمتری نسبت به خوشه های واقعی داشته باشند. بنابراین خوشه هایی که به آرامی در حال رشد هستند شناسایی و به عنوان ناخوشه از خوشه‌بندی حذف می‌شوند. همچنین، از آنجایی که تعداد نقاط در مجموعه ای از نویزها معمولاً کمتر از تعداد در یک

خوشه است و این نویزها خوشه های بسیار کوچکی هستند، ما می توانیم به راحتی چنین گروه های کوچکی را شناسایی کرده و آنها را حذف کنیم. به این ترتیب، مرحله نهایی، از بین رفتن خروجی، گام لازم برای خوشه بندی خوب است. در الگوریتم DBSCAN این اتفاق می افتد و نقاط نویز به راحتی مشخص می شوند.

۳-۷- روش تراکم کرنل صفحه ای و شبکه ای، در نرم افزار GIS

KDE یک تکنیک درونیایی است که در شناسایی نقاط حادثه خیز مورد استفاده قرار می گیرد. این روش براساس یک تابع چگالی نقطه ای است که برای تمام سلول های یک شبکه (محدوده مطالعه) محاسبه می شود. محاسبه KDE در GIS بدین ترتیب است که ابتدا محدوده مورد مطالعه باید تعریف شود. سپس تفکیک فضایی سلول شبکه ای (یا با فاصله سلول یا تعداد ستون ها) تعریف می شود. تابع کرنل دقیقاً بر روی هر نقطه تصادف قرار گرفته و سطح متقارن را تشکیل می دهد. سطح تابع چگالی با توجه به مقدار شعاع جستجو توزیع کاربرتعریف می شود. برای هر سلول در شبکه مشخص شده (محدوده مورد مطالعه)، این روش ارزش تمام سطوح را برای این سلول مرجع، جمع می کند (لویین، ۲۰۱۰). در شکل ۳-۱۲، خلاصه ای از روش محاسبه KDE را می توان مشاهده کرد.



شکل ۳-۱۶: جمع‌بندی سطوح کرنل نرمال برای ۵ نقطه (لوین، ۲۰۱۰)

۳-۸-۸- مقایسه روش های خوشه بندی

۳-۸-۱- شاخص دقت پیش‌بینی (PAI)

همانطور که در پیشینه تحقیق توضیح داده شد، تکنیک های بسیاری برای شناسایی ناحیه های حادثه‌خیز تصادف وجود دارد. با این وجود، وقتی تصمیم به استفاده از این تکنیک ها می‌شود، انتخاب روش بهتر مشکل می‌شود. تکنیک های مختلف شناسایی نقاط حادثه‌خیز نتیجه ی مختلفی از لحاظ اندازه، شکل و مکان را ارائه می‌دهند. به همین دلیل نتایج تکنیک‌های خوشه‌بندی داده‌های نقطه‌ای باید با یکدیگر مقایسه شوند. با این حال، هیچ دانش دقیق و یا تعریف دقیق در مورد تکنیک شناسایی نقطه حادثه‌خیز برای تشخیص حوادث با اهمیت بالا وجود ندارد. با توجه به این موضوع، توسط چی نی و همکارانش^۱ (۲۰۰۸) برای تعیین اینکه آیا تفاوتی برای لحاظ کردن یا پیش‌بینی نقاط حادثه‌خیز وجود

¹ Chainey et al.

دارد یا نه PAI^1 را مطرح کرد. این روش در نقشه برداری شناسایی نقاط جرم برای مقایسه دو تکنیک در نظر گرفتن و پیش‌بینی این نقاط استفاده شده است. روش شناسی بستگی به نسبت درصد امتیازات در نقاط حادثه‌خیز در تمام نقاط در کل منطقه مورد مطالعه به درصد از نقاط داغ مشخص شده در تمام منطقه مورد مطالعه است. درصد امتیاز در نقاط داغ در تمام نقاط منطقه مورد مطالعه به عنوان نرخ ضربه^۲ (HR) نامیده می‌شود. به طور مثال، پاترن و همکارانش (۲۰۰۹) PAI را به عنوان یک نرخ درصدی، یعنی درصد ضربه به درصد منطقه برای ۲۰ نقطه حادثه‌خیز شناخته شده توصیف کردند.

توابع ریاضی درصد ضربه (معادله) و PAI در زیر آمده است :

$$HR = \left(\frac{n}{N}\right) * 100 \quad (۱-۳)$$

$$PAI = \left(\frac{\left(\frac{n}{N}\right)*100}{\left(\frac{a}{A}\right)*100}\right) \quad (۲-۳)$$

که در آن n تعداد حوادث در نقطه داغ و N تعداد حوادث ناحیه مطالعه است، A مساحت نقطه داغ است و A منطقه ناحیه مطالعه است. PAI اندازه‌گیری شده براساس داده‌های در طول زمان است و معیاری است که توسط آن تکنیک‌های شناسایی نقاط حادثه‌خیز می‌توانند با دیگر مقایسه شوند.

۳-۸-۲- روش ANOVA

تجزیه و تحلیل واریانس یک طرفه (ANOVA) یک روش آماری پارامتری است که یک تحلیل یکنواخت واریانس را برای یک متغیر وابسته توسط یک عامل واحد (مستقل) تولید می‌کند و به نسبت F بستگی دارد. به طور عمده برای مقایسه اختلافات مهم بین دو یا چند گروه وابسته انتخاب شده استفاده می‌شود. بر خلاف فرض صفر که کل جمعیت به معنی برابر است، می‌توان ادعا کرد که حداقل

¹ Prediction Accuracy Index

² Hit Rate

میانگین یک گروه متفاوت است. اگر این مقدار به طور قابل توجهی متفاوت باشد، این گروه ها می توانند از طریق تست دامنه post-hoc و مقایسه چندگانه عددی مرتب شوند. در نرم افزار SPSS، یک ابزار آماری که معمولاً مورد استفاده قرار می گیرد، تست های دامنه ها، زیرمجموعه های همگن را که از یکدیگر متفاوت نیست، تشخیص می دهند و در سطح آلفا ۰,۰۵ به طور معناداری متفاوت هستند.

۳-۸-۳- ضریب همبستگی اتا

هنگامی که یک متغیر دارای مقیاس اسمی و رتبه ای (مثل جنس، نژاد، میزان رضایت و ...) و متغیر دیگر مقیاس فاصله ای یا رتبه ای داشته باشد مثل درآمد، معدل، اندازه، تعداد و ... باید شاخصی انتخاب شود که بتوان از روی یک متغیر، متغیر دیگر را پیش بینی نمود. از جمله شاخص ها در این وضعیت شاخص "نسبت همبستگی" است که با نماد η^2 (مجذور اتا) نشان داده می شود. برای فهم بهتر این ضریب، فرضیه زیر را در نظر داشته باشید:

"بین تعداد افراد مجروح و نوع تقاطع ارتباط وجود دارد."

در این فرضیه، نوع تقاطع یک متغیر کیفی با مقیاس اسمی و تعداد افراد مجروح متغیری کمی با مقیاس فاصله ای (نسبی) است. برای توضیح این شاخص باید گفت که این شاخص بر دو مفهوم کلی بنیان گذارده شده است، نخست آن که با فرض ندانستن نوع تقاطع ممکن است چه میزان خطا در پیش بینی تعداد افراد مجروح رخ دهد (e_1) و دوم آن که با فرض دانستن نوع تقاطع، به چه میزان خطای مرحله قبل (e_1) کاهش خواهد یافت. چنان چه ارتباطی بین این دو متغیر وجود داشته باشد می بایست تعداد خطاهای مرحله اول (e_1) کاهش یافته باشد. میزان این ضریب از رابطه زیر به دست می آید:

$$\eta^2 = \frac{e_1 - e_2}{e_1} \quad (3-3)$$

برخلاف متغیرهای اسمی، که در پیش بینی آن‌ها درست و غلط وجود دارد، در پیش بینی متغیرهای کمی (تعداد افراد مجروح) درست و غلط مطرح نمی‌شود. چنانچه تعداد واقعی افراد مجروح ۳ نفر باشد، تعداد ۱ و ۴ هر دو اشتباه می‌باشند ولی پیش بینی تعداد ۲ نزدیک به واقعتر می‌باشد، به بیان دیگر هرچه $|x_i - \bar{x}|$ کوچک تر باشد میزان خطای کمتری رخ داده است. در آمار به جای $|x_i - \bar{x}|$ معمولاً از مجذور خطاها استفاده می‌شود $((x_i - \bar{x})^2)$ ، پس می‌توان گفت که اگر از متغیر نوع تصادفات در مثال فوق هیچ اطلاعی وجود نداشته باشد مجموع مجذور خطاها برابر است با:

$$\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \hat{x})^2 + \sum_{i=1}^n (\hat{x} - \bar{x})^2 \quad (4-3)$$

\bar{x} مقادیر پیش بینی شده این متغیر و \hat{x} مقدار واقعی متغیر تعداد افراد مجروح، x_i در این فرمول میانگین آن است. پس می‌توان گفت که کل مجذور خطا برابر است با:

$$\text{میزان تغییر خطا} + \text{میزان تغییر پیش بینی شده} = \text{کل مجذور خطا} \quad (5-3)$$

فصل ۴

نتایج

۴-۱- مقدمه

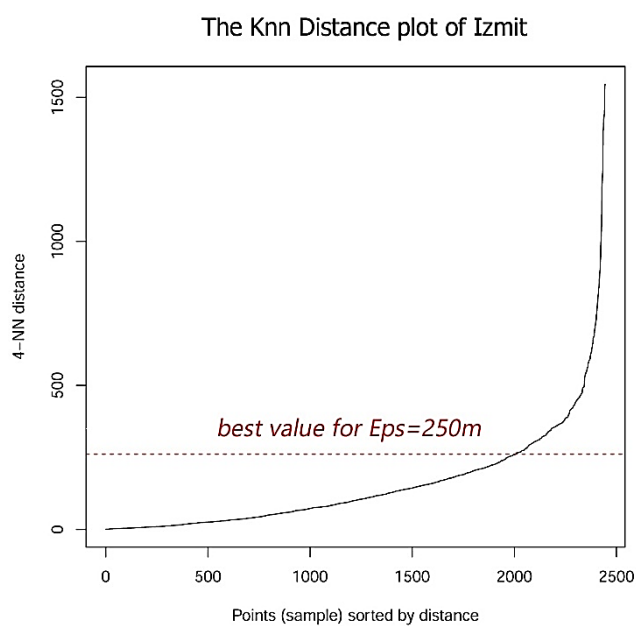
در این فصل به بررسی و تشریح نتایج حاصل از اعمال روش‌های ذکر شده در فصل سوم پرداخته شده است. نتایج برای دو شهر گبزه و ایزمیت به صورت جداگانه مورد تحلیل قرار می‌گیرند. روند تحلیل بدین گونه است که ابتدا نتایج انواع روش‌های ذکر شده برای خوشه بندی را بررسی می‌کنیم و سپس با استفاده از شاخص PAI^1 یا شاخص دقت پیش‌بینی و خوشه‌ها را از لحاظ پارامترهای موجود در آن خوشه بررسی می‌کنیم. هر خوشه را به صورت جداگانه از لحاظ آماری بررسی کرده و داده‌های معنا دار برای هر خوشه را مورد بررسی و تحلیل قرار می‌دهیم.

۴-۲- نتایج حاصل از خوشه بندی GridBSCAN

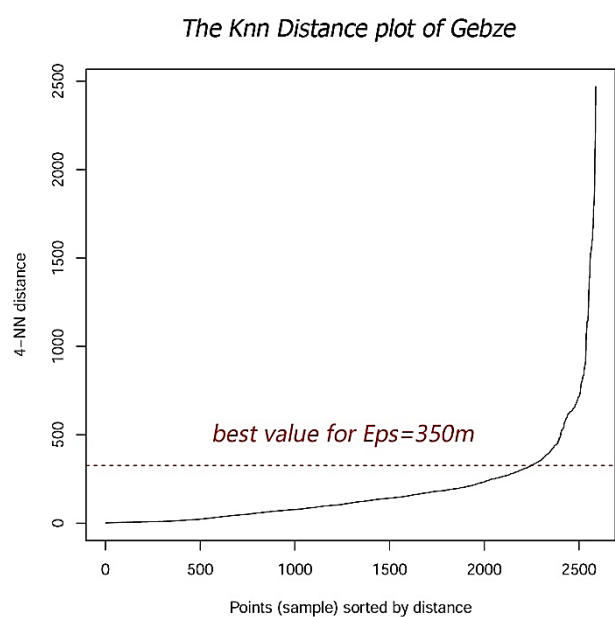
با توجه به روش ذکر شده در فصل سوم برای هر کدام از شهرها با توجه به نمودار $KnnDistance$ در نرم‌افزار Rstudio 1.0.136، برای هر شهر مقادیر بهینه شعاع حداکثر جستجو Eps برای خوشه بندی در نظر گرفته شد (شکل ۴-۱). بر اساس تحقیقات استر در سال ۱۹۹۶ میلادی برای داده‌های دوبعدی (x,y) می‌توان k را برابر با ۴ در نظر گرفت چرا که استر ثابت کرد برای داده‌های دو بعدی مقادیر بیشتر از ۴ نمودار تغییری نخواهد کرد. مقدار $d > 2Eps$ نیز برای شبکه‌بندی با توجه به مطالعه مهرا و مهرا همانطور که در فصل پیشین به تشریح ذکر شد، در نظر گرفته می‌شود (مهرا و مهرا^۲، ۲۰۰۸).

¹ Prediction Accuracy Index

² Mahran & Mahar



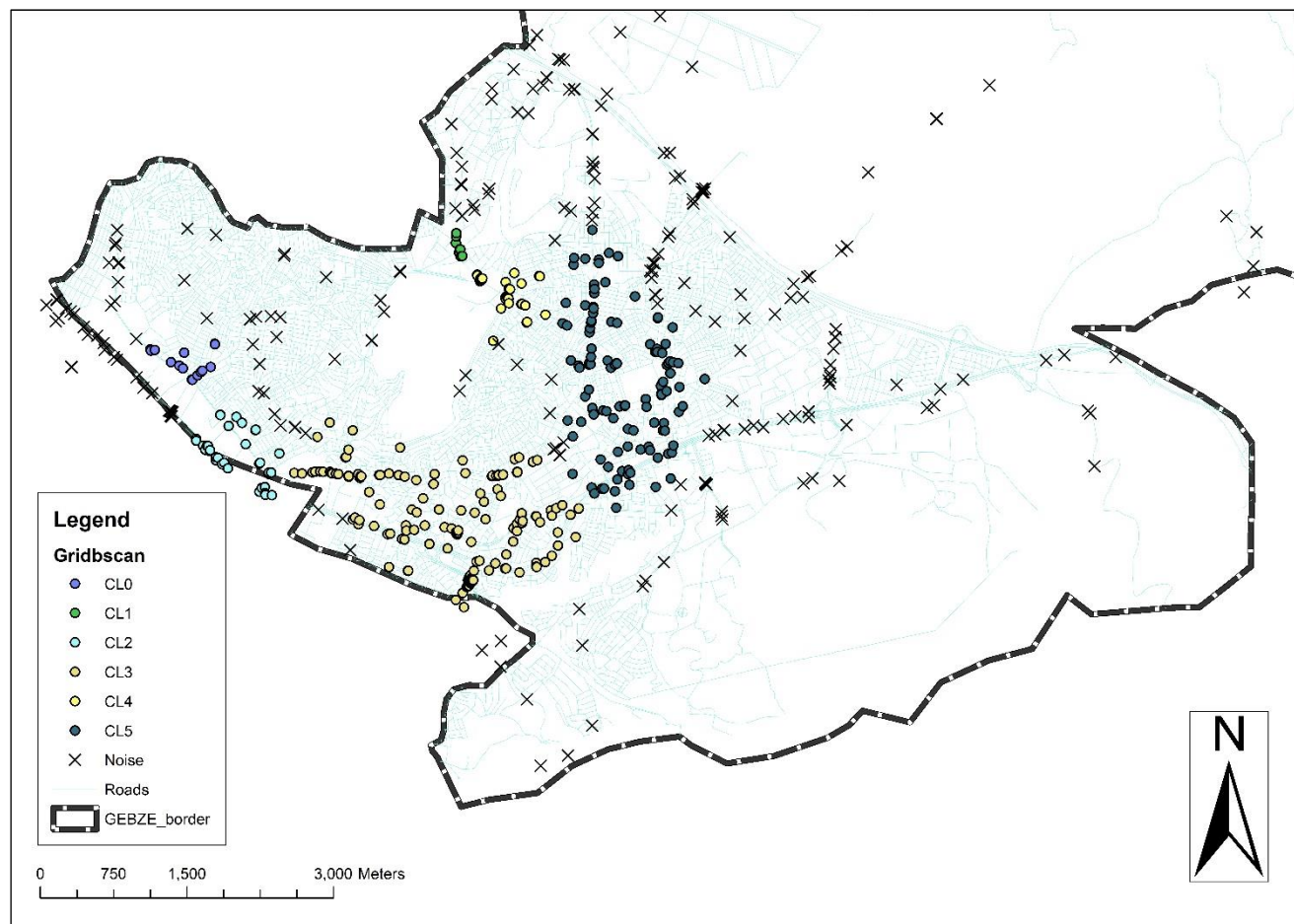
(الف)



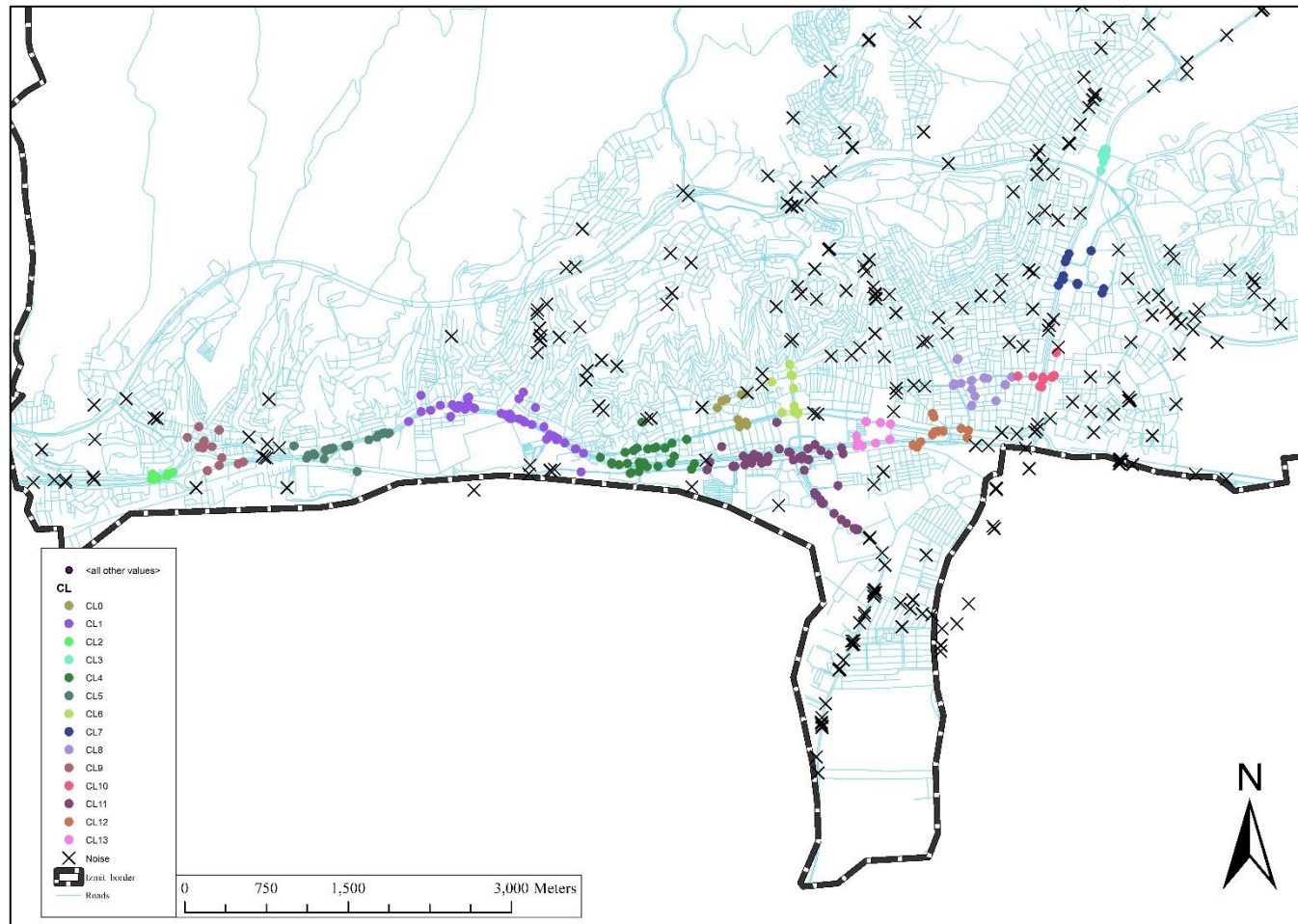
(ب)

شکل ۴-۱: نمودار Knn با مقدار $K=4$ و مقدار بهینه $Eps(m)$ برای شهر (الف) ایزمیت (ب) گبزه همانطور که در شکل ۴-۲ و ۴-۳ قابل مشاهده است الگوریتم GridBSCAN مکان‌هایی که نقاط تصادف در آن تراکم بیشتری نسبت به مکان‌های دیگر دارد را خوشه‌بندی می‌کند و نقاطی که پراکنده هستند به عنوان نقاط نویز یا ناخوشه می‌شناسد.

در ادامه خوشه بندی شهر ایزمیت و گبزه به روش GridBSCAN در شکل های ۳-۴ و ۴-۴ قابل مشاهده است.



شکل ۴-۲: خوشه بندی GridBSCAN با مقدار $d=700m$ و $minPts=15$ و $Eps=350m$ برای شهر گبزه

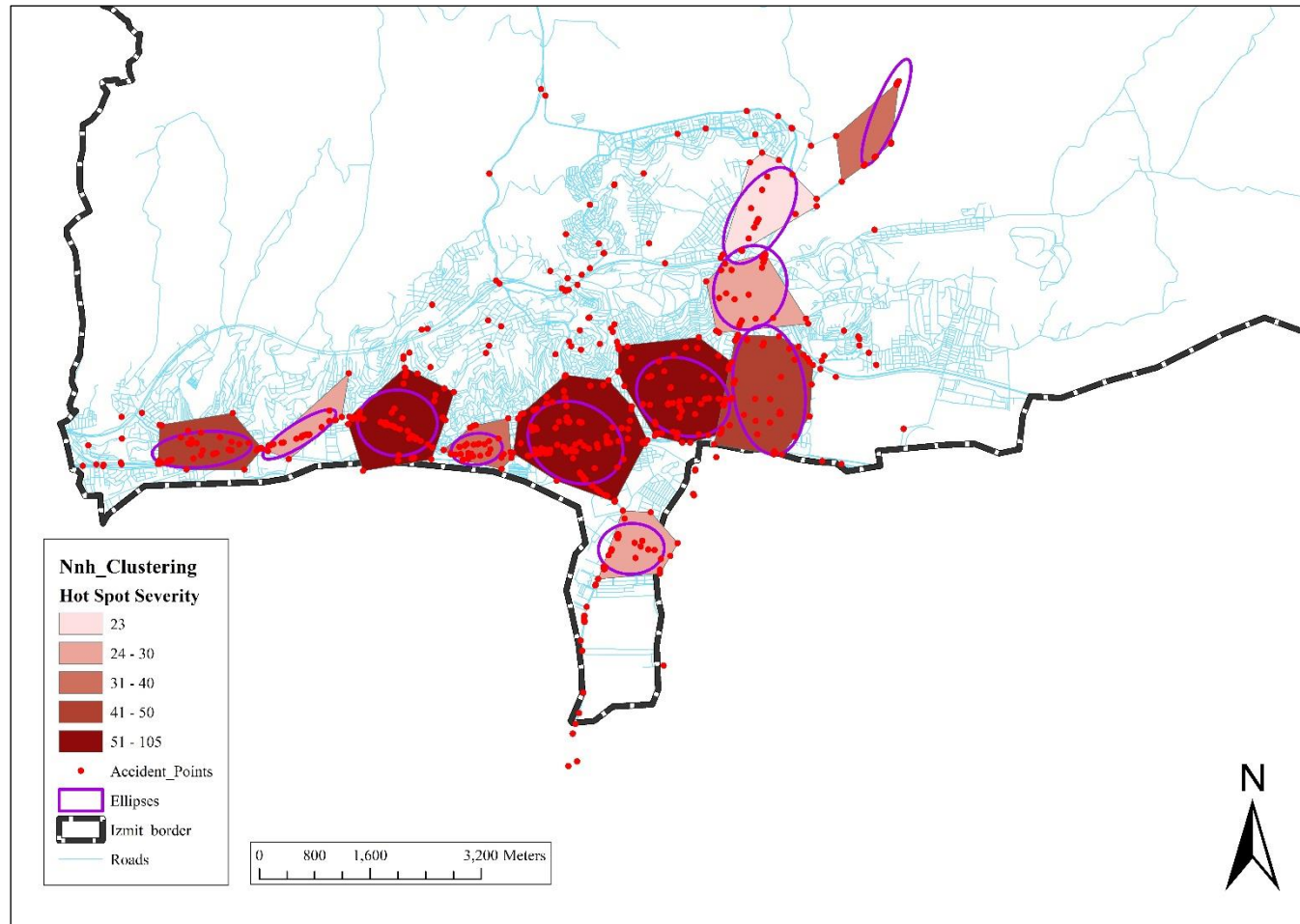


شکل ۴-۳: خوشه بندی GridDBSCAN با مقدار $d = 500m$ و $\minPts = 10$ و $Eps = 250m$ برای شهر ایزمیت

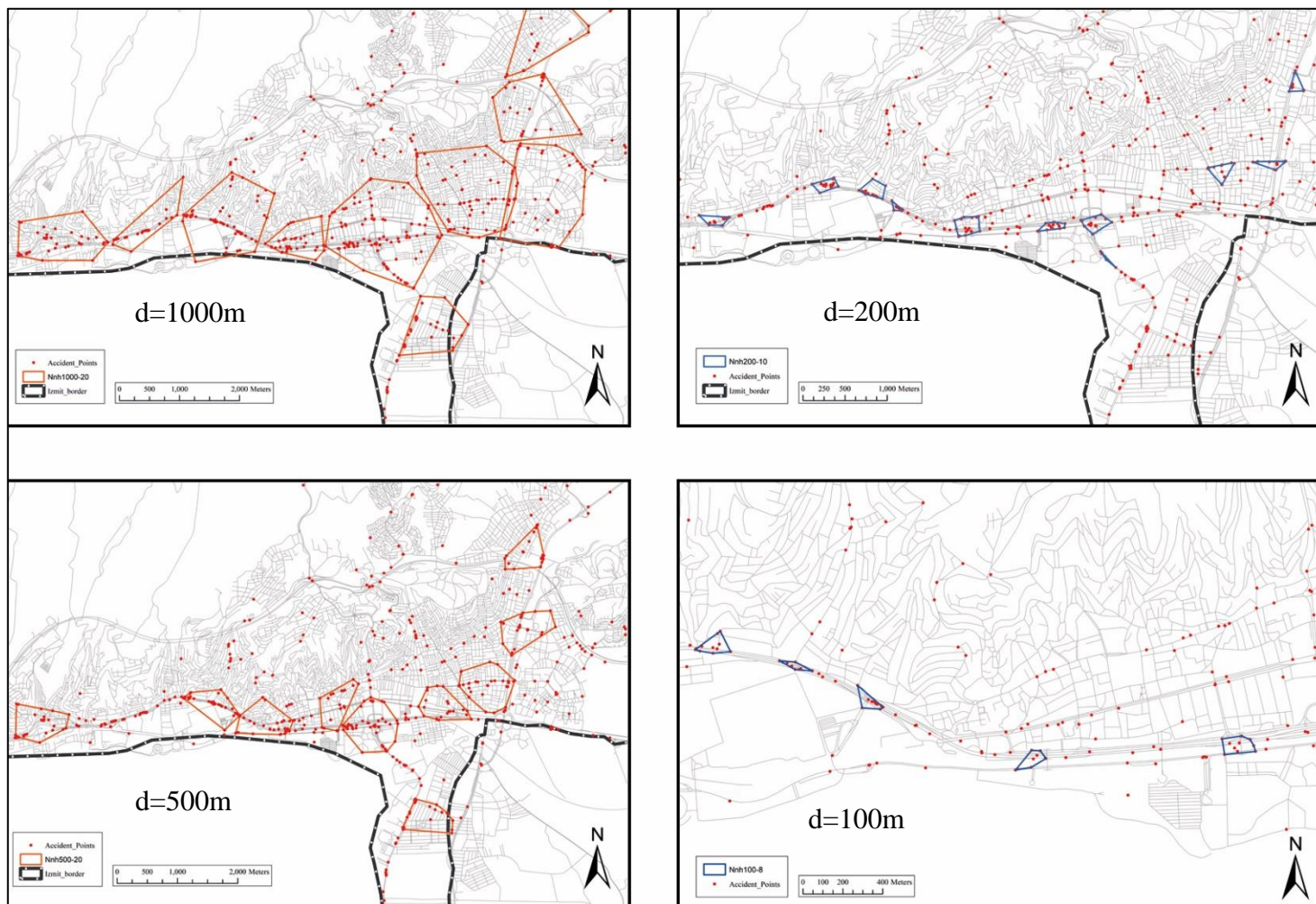
۴-۳- روش خوشه بندی Nnh

همانطور که در روش تحقیق ذکر شد، در این روش تعیین شعاع جستجو و تعداد حداقل نقاط به عهده کاربر می‌باشد. تشخیص نقاط حادثه خیز در این روش با فاصله آغازین ۱ کیلومتر آغاز شده است که در روش های تحلیل نقاط حادثه‌خیز تصادفات فوتی نیز استفاده می‌شود. برای اینکه تعداد زیادی از تصادفات جراحتی را در نظر بگیریم، ابتدا مقدار n_{min} عدد ۲۰ انتخاب شد. در هنگام ایجاد نقشه برای تصویر سازی نتیجه روش خوشه بندی، از روش نقشه موضوعی با مقادیر شدت نقاط حادثه خیز استفاده می‌شود. همانطور که در شکل ۴-۴ دیده می‌شود، روش خوشه بندی Nnh با معیارهای داده شده، خوشه های بسیار بزرگی ایجاد شده است که بیش از یک جاده و تقاطع را پوشش می‌دهد. به منظور درک بهتر ارتباط تصادفات با اجزای محیط شهری، شناسایی مناطق کوچکتر مناسب‌تر و معنی دار است. در شکل ۴-۴ هر دو نوع خروجی بیضی و چندضلعی قابل مشاهده است که تفاوتشان صرفاً در محصور کردن هر خوشه است.

در حین بررسی تصادفات درون شهری برای تشخیص نقاط حادثه‌خیز، فاصله جستجو نقاط حادثه-خیز امری مهم به شمار می‌آید. از آنجا که هدف اصلی این پژوهش این است که ویژگی های نقاط تصادف منجر به جراحات در ارتباط با محیط شهری را درک کنیم، فاصله جستجو نقاط حادثه‌خیز بسیار مهم است. همانطور که در شکل ۴-۴ قابل مشاهده است، خوشه های به دست آمده برای نقاط حادثه‌خیز تصادفات با محدوده ۱۰۰۰ متر و مقدار $n_{min} = 20$ بیش از یک جاده را پوشش می‌دهند. به همین دلیل، فاصله را تا ۵۰۰ متر کاهش دادیم. نقاط حادثه‌خیز شناسایی شده توسط فاصله آستانه ۵۰۰ متر نیز منطقه‌ی بزرگی را پوشش داد. بر این اساس، بار دیگر با ۲۰۰ متر محاسبه کردیم که نتایج بهتری را نسبت به مقادیر قبلی نشان داد.



شکل ۴-۴: نقشه موضوعی تصادفات جراحاتی با شدت نقاط حادثه خیز به روش Nnh با مقادیر $n_{min} = 20$ و فاصله جستجو ۱۰۰۰ متر



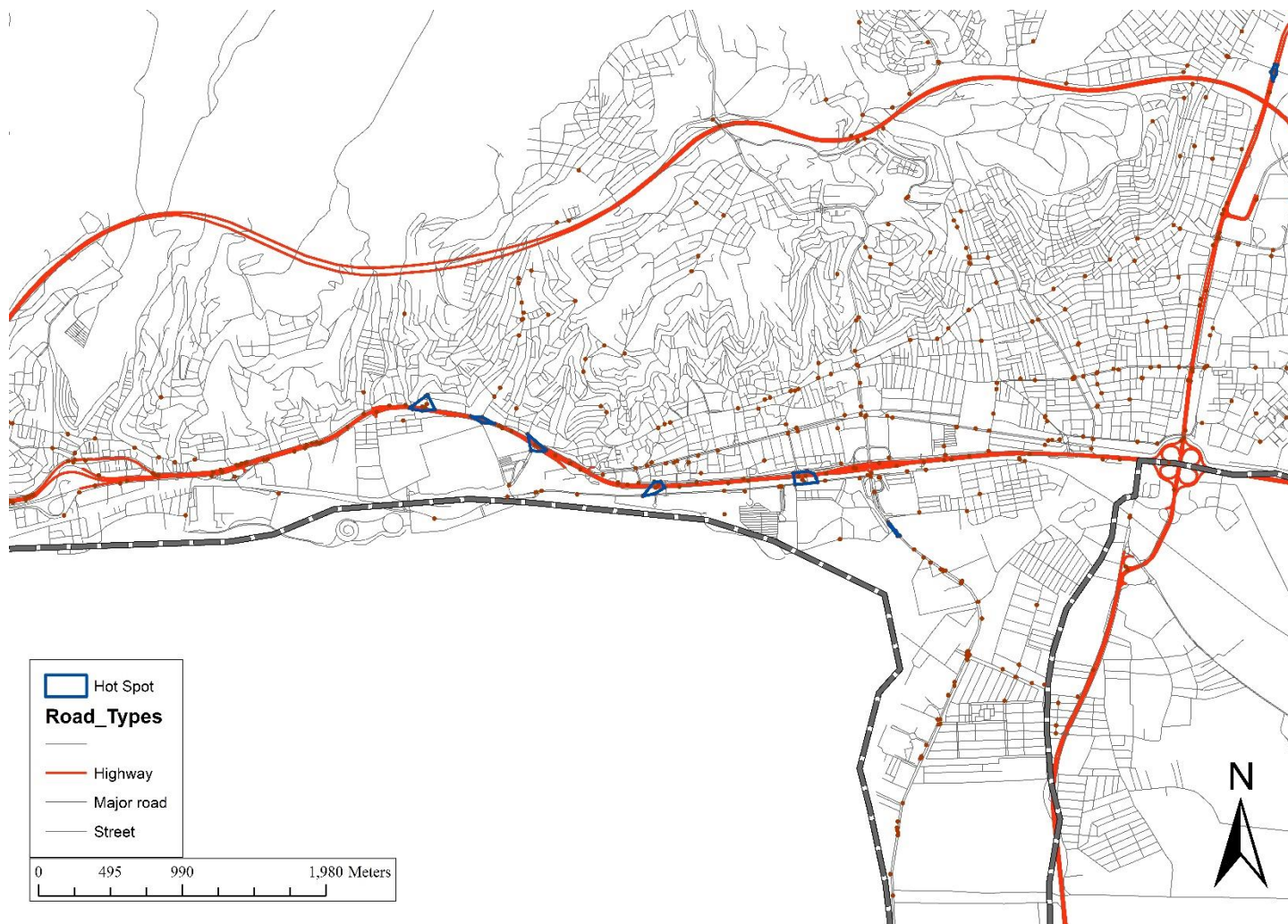
شکل ۴-۵: فواصل جستجو (d) مختلف با $n_{min} = 10$ برای محاسبه خوشه‌بندی Nnh برای تصادفات شهر ایزمیت

با این حال، با مقدار ۲۰۰ متر نیز در بعضی از مناطق شهری، بیش از یک جاده را پوشش می‌دهند و بزرگتر از تقاطع‌ها هستند. در نتیجه، پس از یک دوره تکرار برای پیدا کردن مقدار آستانه مناسب برای خوشه بندی (شکل ۴-۵)، فاصله ۱۰۰ متر به اندازه کافی برای نشان دادن نقاط حادثه‌خیز ترافیکی در مناطق شهری به دست آمد. انتخاب فاصله در روش N_{nh} امری مهم است، چرا که در اصل، دامنه تحلیل خوشه‌ای را تعریف می‌کند. برای بررسی صرف بزرگراه‌ها با محدوده‌ی سرعت مجاز بالا، ممکن است فاصله ۱ کیلومتر قابل قبول باشد، اما برای مکان‌های شهری؛ محدوده ۱۰۰ متر، به طور کلی، برای دستیابی به نقاط حادثه‌خیز واقع در تقاطعات اصلی شهری مناسب است. در نظر گرفتن فاصله‌های کوچکتر از قبیل ۵۰ متر یا ۲۰ متر برای تشخیص نقاط حادثه‌خیز در نزدیکی تقاطعات، بر اساس معیار n_{min} قابل انجام است اما امکان نبودن خوشه، هنگام نبودن نقاط کافی برای جستجو وجود دارد.

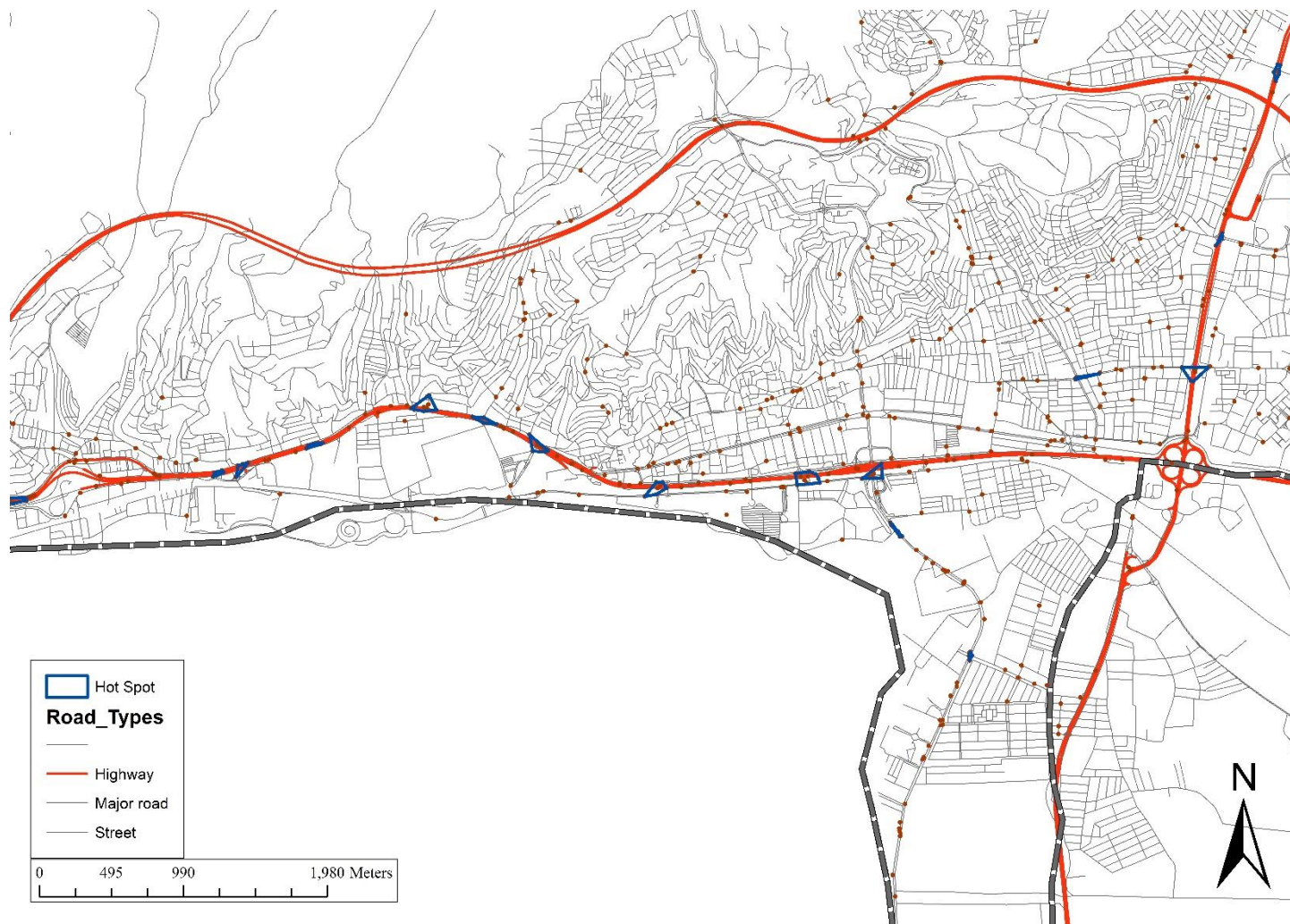


شکل ۴-۶: یکی از تقاطعات شهر ایزمیت که با معیار $d=100m$ و $n_{min} = 10$ به عنوان یک از خوشه‌های حادثه‌خیز بدست آمد.

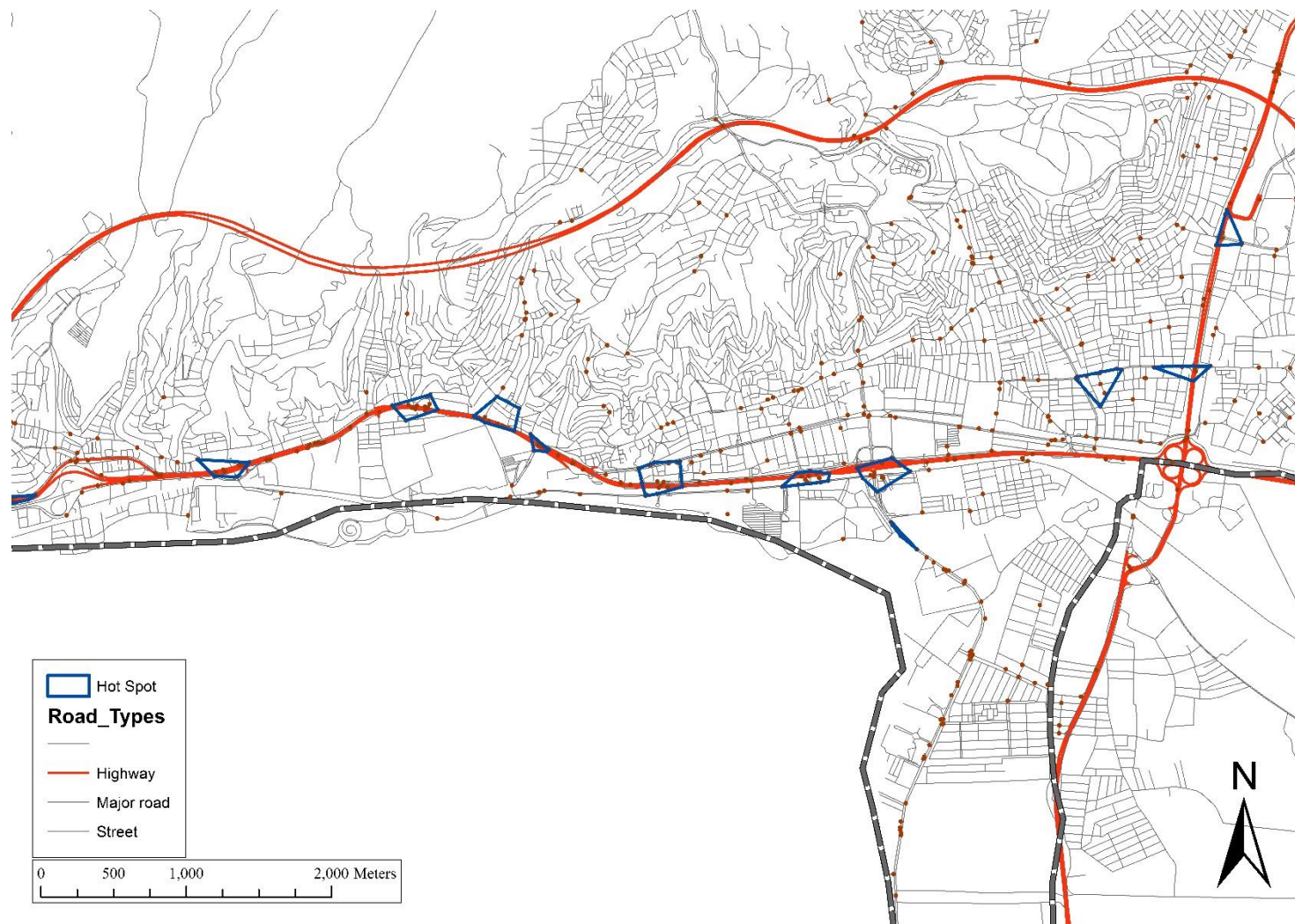
همانطور که در شکل ۴-۶ ملاحظه می‌کنید مقدار انتخابی d و n_{min} به خوبی نقاط تصادف در یک تقاطع را در بر گرفته است. در ادامه برای مقایسه‌ی مقادیر مختلف n_{min} و d فاصله آستانه‌ی جستجو را برای درک تفاوت این دو پارامتر در خوشه بندی N_{nh} بررسی می‌کنیم.



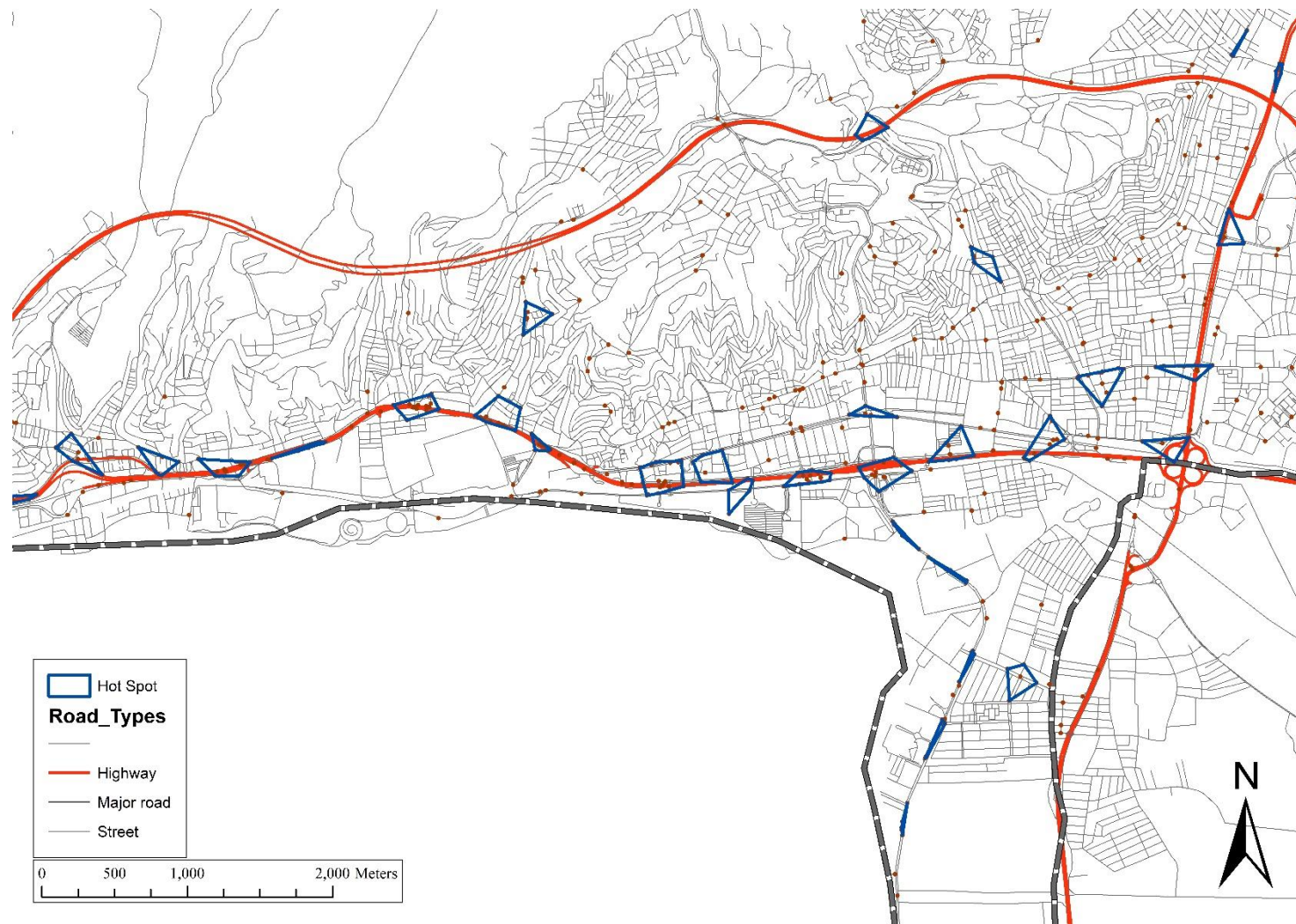
شکل ۴-۷: مقدار $n_{\min} = 10$ و $d = 100$ برای خوشه بندی N_{nh} شهر ایزمیت



شکل ۴-۸: مقدار $n_{\min} = 5$ و $d = 100$ برای خوشه بندی N_{nh} شهر ایزمیت



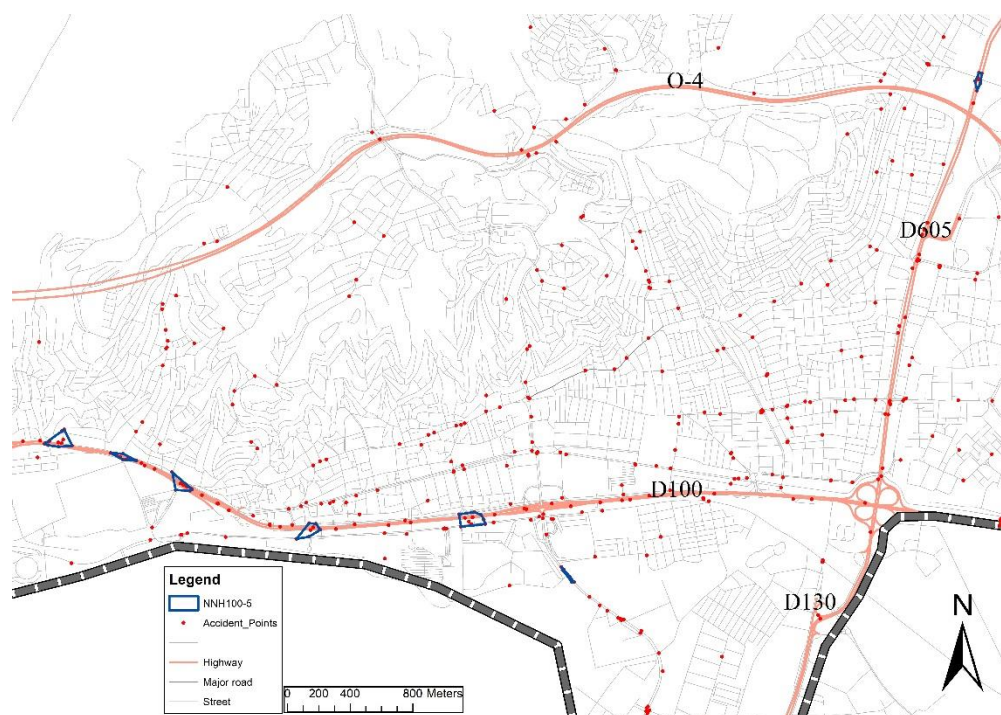
شکل ۴-۹: مقدار $n_{\min} = 10$ و $d = 200$ برای خوشه بندی N_{nh} شهر ایزمیت



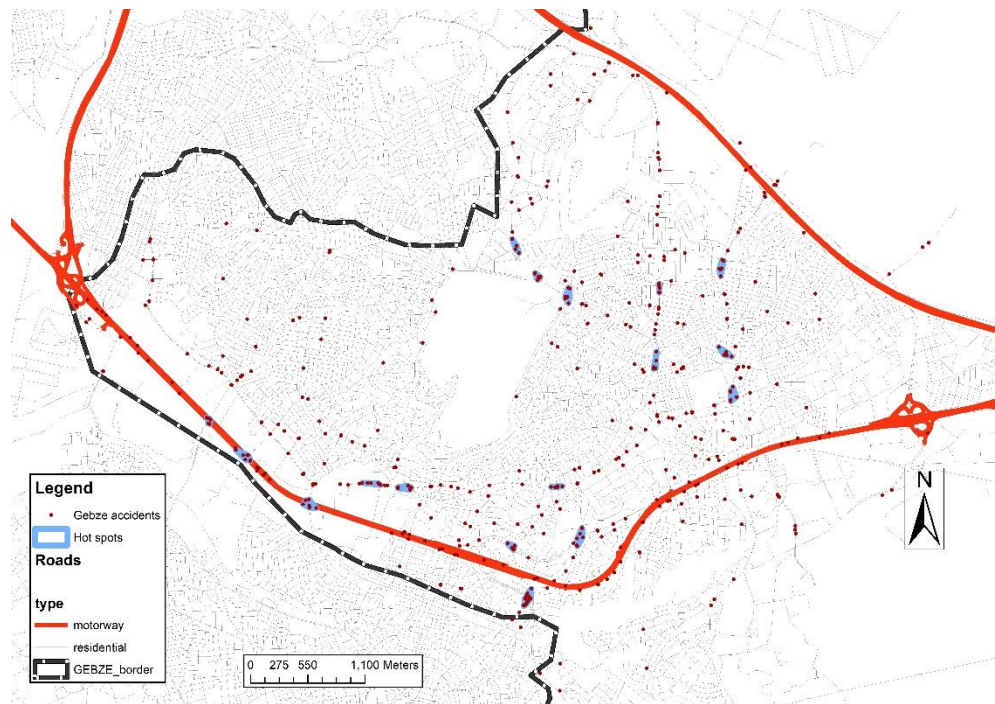
شکل ۴-۱۰: مقدار $n_{\min} = 5$ و $d = 200$ برای خوشه بندی N_{nh} شهر ایزمیت

همانطور که در شکل‌های ۷-۴ تا ۱۰-۷ مشاهده می‌شود، با در نظر گرفتن فاصله آستانه ۱۰۰ متر و n_{min} های ۵ و ۱۰ برای مقایسه مقادیر مختلف n_{min} در هر دو شهر تعداد مختلفی از خوشه‌های حادثه خیز بدست آمد. همانطور که در شکل ۷-۴ و ۸-۴ قابل ملاحظه است به ازای کمتر شدن مقدار n_{min} تعداد خوشه‌ها بیشتر می‌شود. بار دیگر مقدار فاصله آستانه را ۲۰۰ متر در نظر گرفتیم و مقادیر n_{min} های ۵ و ۱۰ قرار دادیم که در شکل ۹-۴ و ۱۰-۴ قابل ملاحظه است. این بار نیز با کمتر شدن n_{min} تعداد خوشه‌ها بیشتر شد و بدین ترتیب هر چه مقدار n_{min} بیشتر باشد خوشه دارای تعداد تصادفات بیشتری خواهد بود و برای بررسی عوامل تصادف از لحاظ آماری قابل قبول‌تر است اما تعداد خوشه‌ها کمتر می‌شود و داده‌های ناخوشه زیاد می‌شود.

پس از مقایسه مقادیر متفاوت d و n_{min} بهترین مقدار برای اینکه بتوان با دیگر روش‌های خوشه‌بندی مقایسه نمود فاصله ۱۰۰ متر و تعداد نقاط ۵ در نظر گرفته شد. چرا که در این مقدار خوشه‌ها شامل نقاط حادثه‌خیز با تراکم بیشتر هستند.



شکل ۴-۱۱: مقدار $n_{min} = 5$ و $d = 100$ برای خوشه بندی N_{nh} شهر ایزمیت

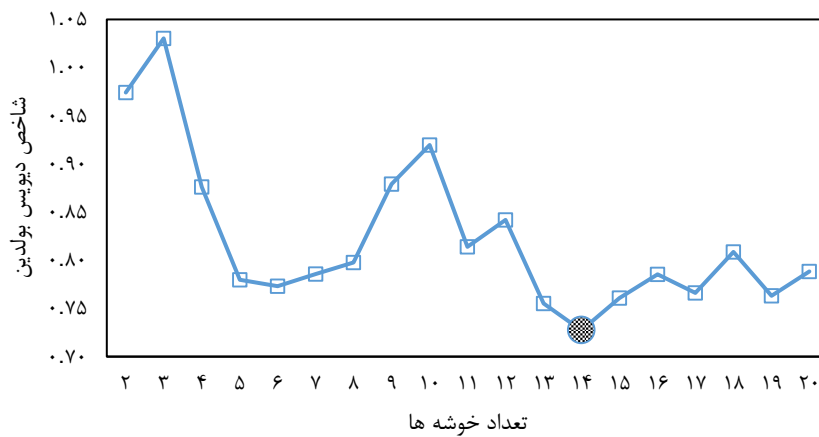


شکل ۴-۱۲: مقدار $n_{min} = 5$ و $d = 100$ برای خوشه بندی Nnh شهر ایزمیت

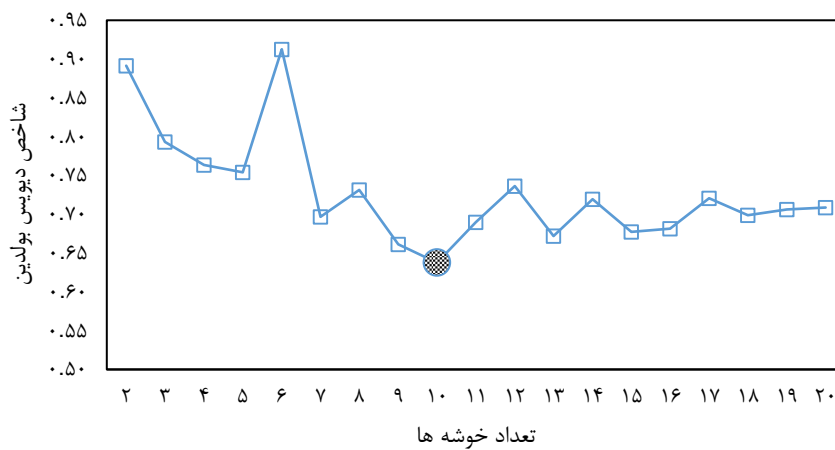
پس از بررسی خوشه‌ها مشاهده می‌شود که اکثر نواحی حادثه‌خیز در نزدیکی یا خود تقاطع‌ها، ورودی و خروجی‌های بزرگراه می‌باشند. تعداد کمی از این نواحی حادثه‌خیز در راه‌های فرعی و دور از مرکز شهر قرار دارند. در شهرهای با عرض معابر کم و توپوگرافی دارای اختلاف ارتفاع مانند گبزه نیز نقاط حادثه‌خیز در تقاطعات درون معابر دسترسی و راه‌های فرعی رخ می‌دهند (شکل ۴-۱۲). چرا که توپوگرافی خاص و بافت شهری گبزه مطمئناً رفتار ترافیکی متفاوتی نسبت به شهر ایزمیت که دارای خیابان‌های عریض است، خواهد داشت. در قسمت پایانی این فصل به تشریح خوشه‌ها مفصلاً خواهیم پرداخت.

۴-۵- روش خوشه بندی K-Means

شاخص اعتبارسنجی دیویس-بولدین^۱ برای بدست آوردن K بهینه در شکل ۴-۱۳ و ۴-۱۴ برای شهر گبزه و ایزمیت نمایش داده شده است. تعداد خوشه مناسب برای گبزه $K=14$ و برای ایزمیت $K=10$ به دست آمد.



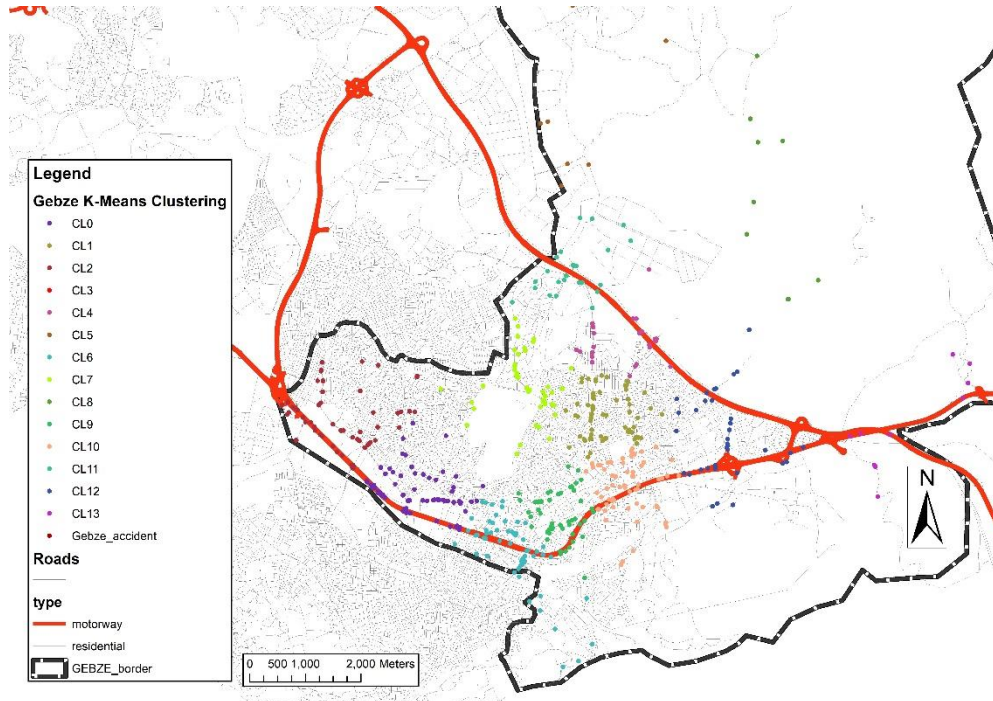
شکل ۴-۱۳: نمودار شاخص اعتبارسنجی دیویس-بولدین برای یافتن تعداد بهینه خوشه‌ها در روش K-Means برای شهر گبزه



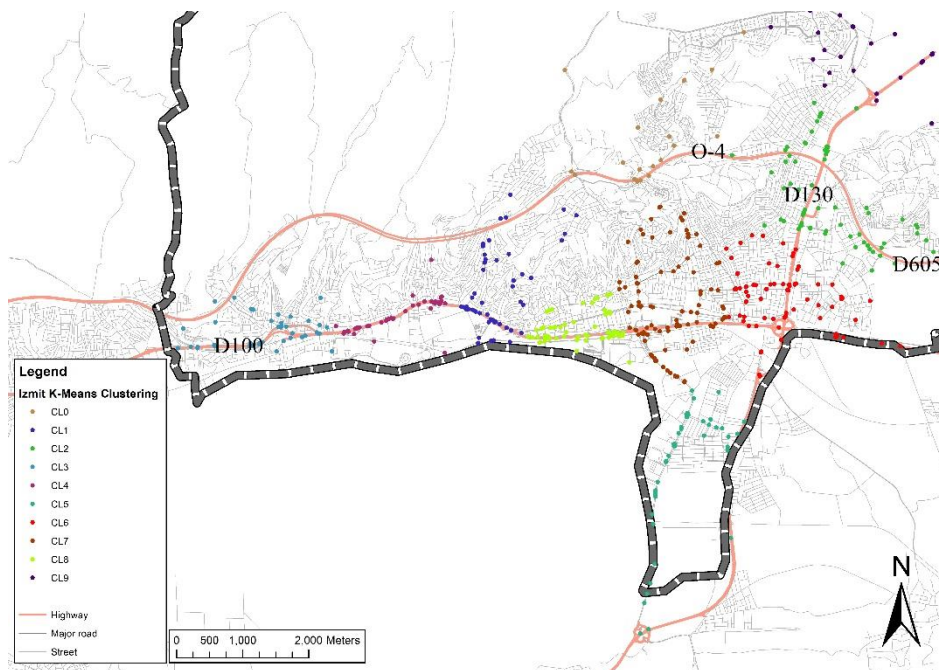
شکل ۴-۱۴: نمودار شاخص اعتبارسنجی دیویس-بولدین برای یافتن تعداد بهینه خوشه‌ها در روش K-Means برای شهر ایزمیت

¹ Davis-Bouldin Index

خروجی خوشه‌بندی به‌دست‌آمده از خوشه‌بندی K-Means به ترتیب برای شهر گبزه و ایزمیت در شکل‌های ۱۵-۴ و ۱۶-۴ قابل مشاهده است.



شکل ۴-۱۵: خوشه‌بندی K-Means شهر گبزه

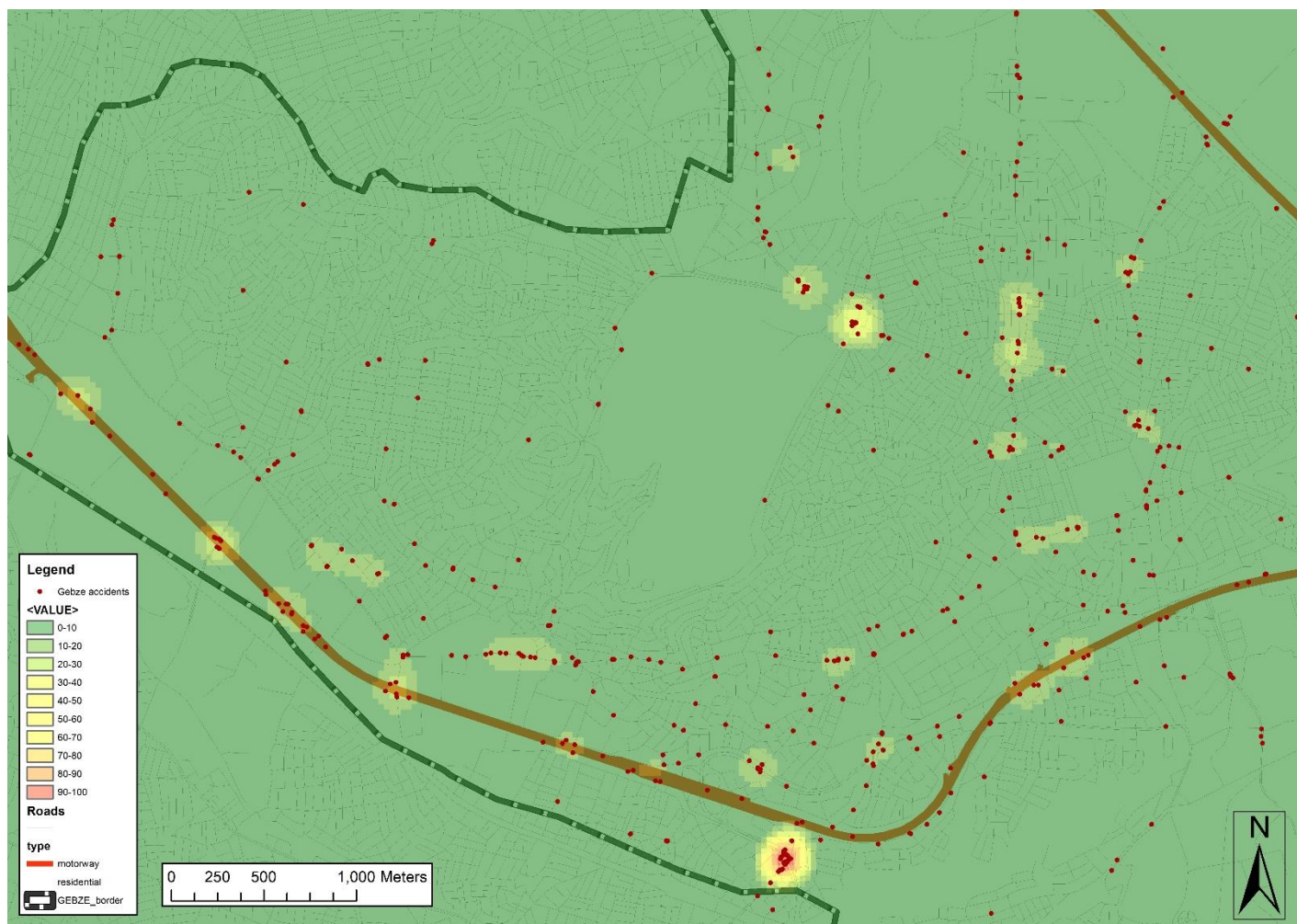


شکل ۴-۱۶: خوشه‌بندی K-Means شهر ایزمیت

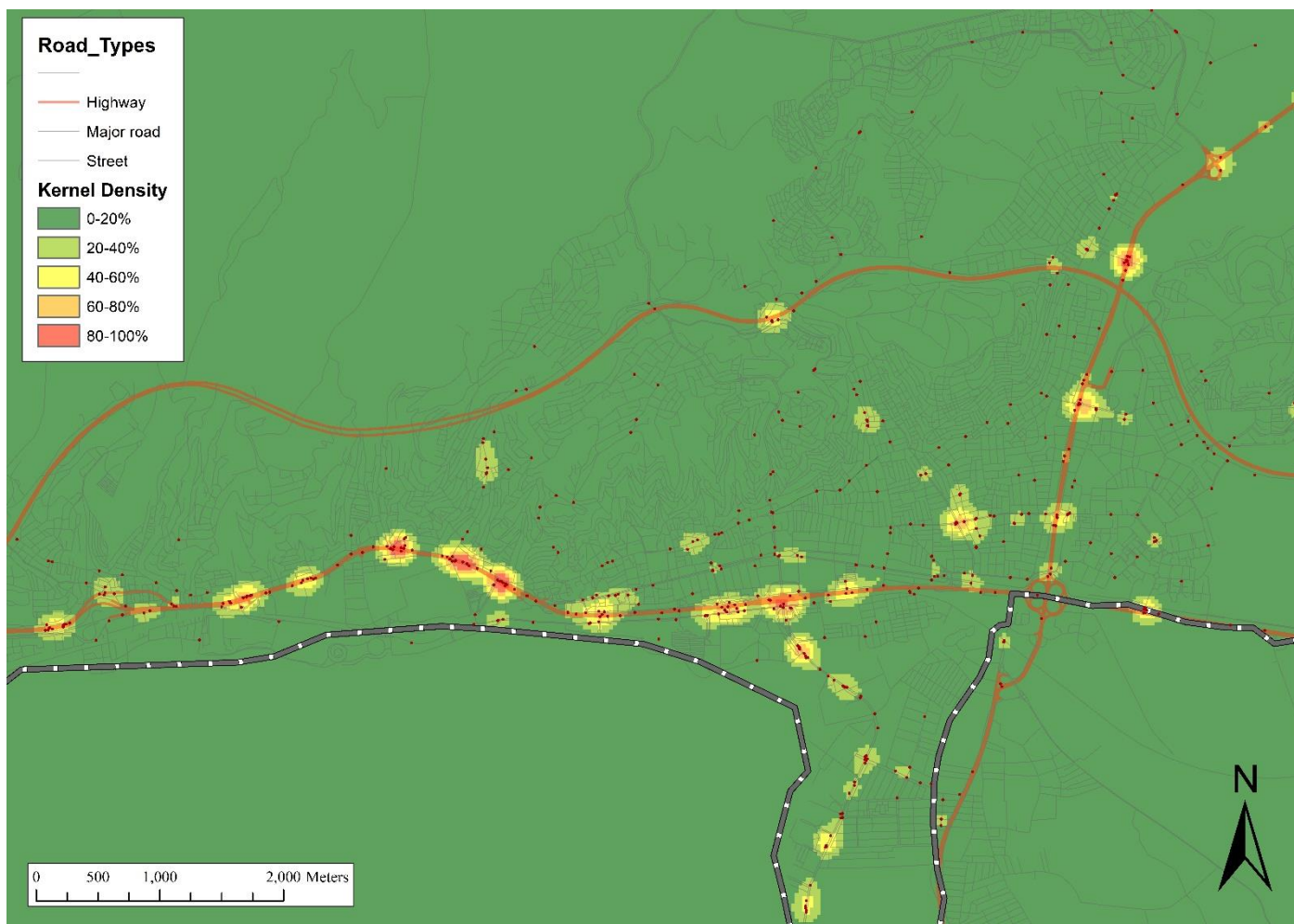
۴-۶- روش تخمین چگالی کرنل

همانطور که در روش تحقیق ذکر شده است، اندازه هر سلول و اندازه پهنای باند (شعاع جستجو) برای محاسبه نقشه KDE از محاسبه نقاط حادثه‌خیز شناسایی شده با روش خوشه بندی Nnh حاصل می شود. بدین ترتیب، کل مساحت convex hullها و حداکثر فاصله بین نقاط در نقاط حادثه‌خیز در نظر گرفته می شود. برای این مطالعه، شعاع جستجو نیز با توجه به نتایج خوشه بندی Nnh انتخاب شده است. برای انجام این کار، حداکثر فاصله تصادفات موجود در خوشه اندازه گیری می شود. حداکثر فاصله بین تصادفات ۲۰۰ متر است. به همین دلیل، برای هر نوع KDE، ۲۰۰ متر به عنوان شعاع جستجو استفاده می شود.

با توجه به این نقشه، مناطق قرمز رنگ که به مقدار تراکم نسبی ۸۰٪-۱۰۰٪ می باشد به عنوان نقاط حادثه‌خیز شناخته می شوند. در شهر اکثر تصادفات در محور D100 بزرگراه ایزمیت به استانبول رخ داده است. این تصادفات هم در خط اصلی بزرگراه و هم در رمپها مشاهده شده است.



شکل ۴-۱۷: نقشه تخمین ترام کرنل با شعاع جستجوی ۲۰۰ متر برای شهر گبزه



شکل ۴-۱۸: نقشه تخمین ترام کرنل با شعاع جستجوی ۲۰۰ متر برای شهر ایزمیت

تعداد نقاط قرار گرفته در مناطق حادثه خیز را محاسبه و مقدار نرخ ضربه و دقت پیش‌بینی را برای روش کرنل نیز محاسبه می‌کنیم. از سوی دیگر، در برخی از شریان‌های کوچک، مناطق متراکم دیده می‌شود که به دلیل ویژگی‌های کاربری زمین مانند مسکونی، تجاری و غیره در آن فعالیت‌های عابر پیاده و مناطق مسکونی قابل مشاهده است برای مثال یکی از نقاط حادثه‌خیز در منطقه مسکونی در شکل ۴-۱۹ قابل مشاهده است.



شکل ۴-۱۹: تقاطعی در منطقه مسکونی که در تمامی روش‌های خوشه‌بندی به عنوان نقطه حادثه‌خیز شناخته شده است.

۴-۷- نتایج روش‌های خوشه‌بندی و بررسی عوامل مؤثر بر تصادفات در نقاط

حادثه‌خیز

با ادغام انواع خوشه‌بندی در محیط GIS نقاط حادثه‌خیز تصادفات را به راحتی می‌توان تشخیص داد. سپس برای درک رابطه بین پارامترهای مؤثر بر تصادفات و تعداد مجروحین در خوشه‌ها در روش

خوشه‌بندی GridBSCAN، با استفاده از آزمون ANOVA در نرم‌افزار SPSS هر خوشه را مورد بررسی قرار دادیم و در هر خوشه تعدادی از پارامترها به عنوان پارامتر معنی‌دار شناخته شد. در این آزمون تعداد افراد جراحی دیده به عنوان داده وابسته و بقیه پارامترهای ذکر شده در فصل سوم به عنوان داده مستقل در نظر گرفته شد.

۴-۷-۱- توصیف خوشه‌های حادثه‌خیز شهر گبزه

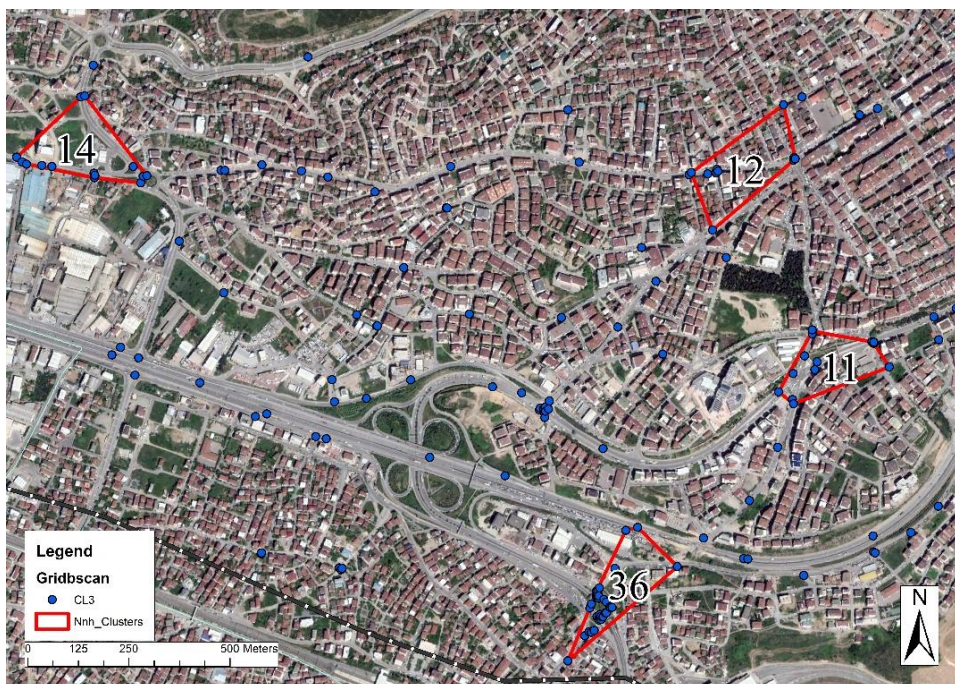
در این قسمت به بررسی و تشریح خوشه‌های حادثه‌خیز شهر گبزه می‌پردازیم و با ادغام خوشه‌ها ویژگی‌های مناطق را می‌کنیم. ابتدا پارامترهای معنادار در خوشه بندی GridBSCAN را از طریق آزمون ANOVA دریافت کرده و پس از ادغام خوشه بندی Nnh و KDE پارامترها را مجدداً بررسی می‌نماییم. همچنین با در اختیار بودن ابزار دیدخیابان در نرم‌افزار گوگل‌ارث به بررسی میدانی شرایط تقاطع می‌پردازیم.

خوشه CL3

این خوشه با بیشترین تعداد تصادف (۱۷۵) در جنوب غربی شهر قرار دارد. با توجه به جدول ۴-۱ دو پارامتر، نوع تقاطع که تصادف در آن محل اتفاق افتاده و وجود یا عدم وجود گاردریل به عنوان داده تأثیر گذار بر تعداد افراد مجروح بدست آمده است. با توجه به شکل ۴-۲۰ نقاط آبی تمام نقاط داخل خوشه CL3 می‌باشد و خطوط قرمز خوشه بندی Nnh و شماره‌ها تعداد تصادفات در داخل آن را نشان می‌دهد. با توجه به اینکه نوع تقاطع رابطه معناداری با تعداد مجروحین دارد و موقعیت تصادفات در نقشه‌ی ماهواره ای نشان می‌دهد که بیشترین تصادفات در تقاطعات رخ داده است با بررسی مجدد تصادفات موجود در خوشه‌های Nnh چند نقطه به عنوان پرحادثه ترین نقاط در داخل این خوشه‌ها بدست آمد.

جدول ۴-۱: نتایج آزمون ANOVA برای خوشه بندی GridBSACN خوشه‌های شهر گبزه

نام خوشه	نقاط درون خوشه	پارامترهای معنی دار با تعداد مجروحین	مجموع مربعات	درجه آزادی	مجدور میانگین	F	Sig.	اتا	مجدور اتا
CL0	۱۶	نوع تصادف روز تصادف (روزهای هفته)	۲/۲۲۳	۴	۰/۵۵۶	۳/۵۶۶	۰/۰۴۲	۰/۷۵۱	۰/۵۶۵
CL1	۱۰	-	-	-	-	-	-	-	-
CL2	۳۶	وجود یا عدم وجود شانه راه طبقه راه (خیابان یا بزرگراه شهری) نوع راه (یکطرفه، دوطرفه جدا شده و جدا نشده)	۷/۲۵۹	۱	۷/۲۵۹	۴/۳۸۴	۰/۰۴۴	۰/۳۳۸	۰/۱۱۴
			۸/۶۱۷	۱	۸/۶۱۷	۵/۳۳۳	۰/۰۲۷	۰/۳۶۸	۰/۱۳۶
			۱۰/۸۸۹	۲	۵/۴۴۴	۳/۴۱۱	۰/۰۴۵	۰/۴۱۴	۰/۱۷۱
CL3	۱۷۵	نوع تقاطع	۹/۰۲۸	۴	۲/۲۵۷	۲/۷۷۹	۰/۰۲۹	۰/۲۴۸	۰/۰۶۱
CL4	۳۱	وجود یا عدم وجود گاردریل	۷/۵۷۸	۱	۷/۵۷۸	۴/۰۴۵	۰/۰۴۷	۰/۱۸۰	۰/۰۳۲
CL5	۱۲۳	نوع تصادف	۲۵/۵۶	۵	۵/۱۱۲	۲/۸۶۵	۰/۰۱۸	۰/۳۳۰	۰/۱۰۹



شکل ۴-۲۰: ترکیب خوشه CL3 از روش GridBSCAN و خوشه‌های حاصل از Nnh

جدول ۴-۲: توزیع نقاط در تقاطعات مختلف در خوشه CL3 شهر گبزه

تقاطع	تعداد	واریانس
چهارراه	۴۲	۰/۲۰۹
مسیر مستقیم	۵۰	۰/۶۸۹
میدان	۳۸	۲/۱۲۹
تقاطع T شکل	۳۸	۰/۴۱۳
تقاطع Y شکل	۷	۰/۲۸۶
مجموع	۱۷۵	۰/۸۴۶

همانطور که در جدول ۴-۲ مشاهده می‌شود تعداد تصادفات در تقاطعات قابل توجه می‌باشد. با توجه به تصاویر دریافتی از طریق نرم‌افزار گوگل ارث^۱ می‌توان به صحت تجمع نقاط در این تقاطعات پی

^۱ Google Earth

برد. همچنین با ابزار دید خیابان^۱ در گوگل ارث می توان از نزدیک این مکان ها را مشاهده نمود. تعدادی از این تقاطعات را در تصاویر زیر می توان مشاهده کرد.



شکل ۴-۲۱: تصاویر هوایی و تصاویر دید خیابان در نرم افزار گوگل ارث

پارامتر گاردریل نیز به تبع به علت وقوع تصادفات در نقاط تقاطع شهری عدم وجود گاردریل به عنوان داده معنی دار به عنوان نتیجه ای بدیهی محسوب می شود.

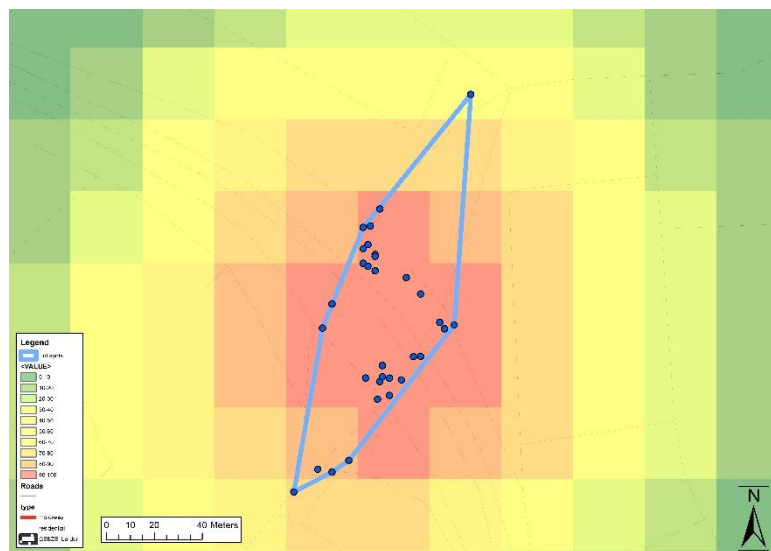
با ادغام خوشه بندی های مختلف در GIS که در شکل ۴-۲۲ قابل مشاهده است می توان نقاط حادثه خیز را دقیق تر مورد بررسی قرار داد. پس از بررسی مجدد داده های درون خوشه بندی حاصل از Nnh برای مثال برای میدانی که به عنوان نقطه حادثه خیز شناخته شده است پارامتر علائم چراغ دار

^۱ Street view

راهنمایی به عنوان داده معنی دار بدست آمد . با توجه به شکل تقاطع و به علت تداخل حرکت راننده‌ها در این میدان تصادفات زیادی در این منطقه اتفاق افتاده است. این مکان به عنوان یکی از پر حادثه‌ترین نقاط شهر گبزه می‌باشد.



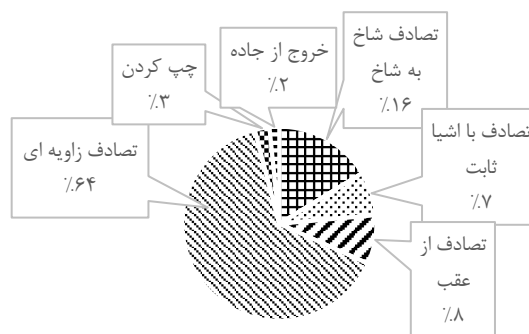
شکل ۴-۲۲: میدان پرتصادف در شهر گبزه



شکل ۴-۲۳: ادغام خوشه بندی KDE, Nnh و GridBSCAN

خوشه CL5

با ۱۲۵ تصادف دارای بیشترین تصادف بعد از خوشه CL3 می‌باشد. در این خوشه پارامتر نوع تصادف با تعداد مجروحین رابطه معناداری دارد. همانطور که در شکل ۴-۲۴ قابل مشاهده است در این خوشه تصادفات بیشتر از نوع تصادف زاویه‌ای یا از بغل می‌باشد.



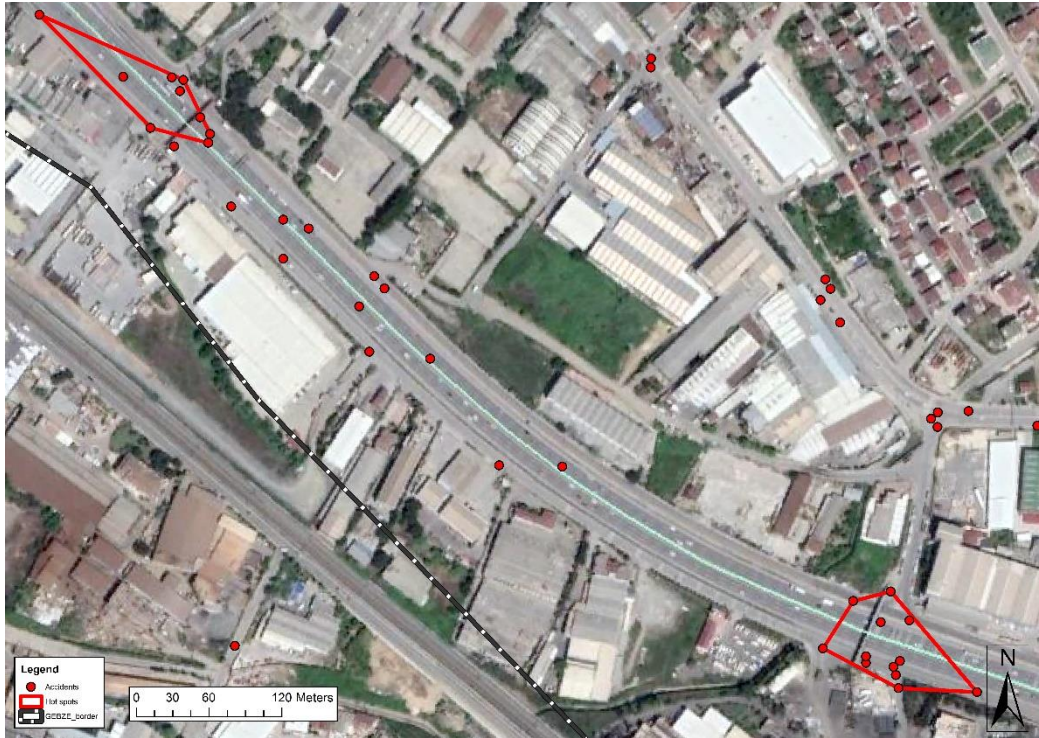
شکل ۴-۲۴: توزیع تصادفات در خوشه CL5 در روش GridBSCAN

پس از بررسی مجدد این خوشه و ادغام دیگر روش‌های خوشه‌بندی پارامتر زمانی شب و روز به عنوان داده معنی دار با مقدار معنی‌داری ۰/۰۱۶ با آلفا ۰/۰۵ بدست آمد. با بررسی تصاویر دید خیابان می‌توان ویژگی‌های محیطی این مناطق را مشاهده کرد. این منطقه از مناطق پرتراکم ساختمانی گبزه است و عرض خیابان ۹ متر و مسیر دو طرفه بدون میانه راه است، و دارای پارک حاشیه‌ای است.

خوشه CL2

تصادفات در این خوشه در محدوده‌ی بزرگراه ایزمیت-استانبول (D100) قرار دارد. برای مثال برای یکی از تقاطعات در این خوشه با توجه به خوشه بندی Nnh که در شکل ۴-۲۴ قابل مشاهده است با بررسی ویژگی‌های تصادف در این منطقه از لحاظ مکانی پارامتر نوع راه و طبقه راه به عنوان داده معناداری به دست آمد و باتوجه به اینکه موقعیت تصادفات این خوشه در محدوده بزرگراه می‌باشد معناداری این پارامترها بدیهی است. اما از لحاظ زمانی بیشتر تصادفات در این منطقه در روز اول کاری

یعنی روز دوشنبه اتفاق افتاده است و این موضوع می‌تواند دلیل مهم‌تری برای حادثه‌خیز بودن این محدوده باشد.



شکل ۴-۲۵: تصادفات خوشه CL2 و خوشه‌های روش Nnh در شهر گبزه

۴-۷-۲- توصیف خوشه‌های حادثه‌خیز شهر ایزمیت

در این بخش به شرح خوشه‌های شهر ایزمیت می‌پردازیم که این شهر نیز به دو بخش شهری و محور بزرگراه استانبول-ایزمیت تقسیم می‌شود. همانطور که در شکل ۴-۲۶ قابل مشاهده است تصادفات بیشتر در قسمت بزرگراه D100 ایزمیت به استانبول قرار دارند و تعداد دیگری از تصادفات نیز در تقاطعات و مناطق پرتراکم مسکونی قرار دارند.



شکل ۴-۲۶: تصویر ماهواره‌ای محور D100 عبوری از شهر ایزمیت و توزیع نقاط تصادف

با توجه به جدول ۴-۳ خوشه‌هایی که دارای بیشترین تعداد تصادفات هستند را مورد بررسی قرار می‌دهیم.

جدول ۴-۳: نتایج آزمون Anova برای خوشه‌های شهر ایزمیت

نام خوشه	نقاط درون خوشه	پارامترهای معنی دار با تعداد مجروحین	مجموع مربعات	درجه آزادی	مجدور میانگین	F	Sig.	اتا	مجدور اتا
CL0	۱۰	-	-	-	-	-	-	-	-
CL1	۴۸	-	-	-	-	-	-	-	-
CL2	۱۰	نوع راه (یکطرفه، دوطرفه جدا شده و جدانشده) طبقه راه (خیابان یا بزرگراه شهری) مسیر مستقیم یا قوس	۰/۷۱۱	۱	۰/۷۱۱	۶/۴۰۰	۰/۰۳۵	۰/۶۶۷	۰/۴۴۴
CL3	۱۰	-	-	-	-	-	-	-	-
CL4	۲۹	-	-	-	-	-	-	-	-
CL5	۲۳	خط کشی‌های راه	۱۱/۴۵۸	۱	۱۱/۴۵۸	۱۶/۲۶۳	۰/۰۰۱	۰/۶۵۲	۰/۴۲۵
CL6	۱۰	شرایط نور مسیر مسیر مستقیم یا قوس	۰/۴	۱	۰/۴	۶/۴۰۰	۰/۰۳۵	۰/۶۶۷	۰/۴۴۴
CL7	۱۶	وجود یا عدم وجود شانه راه	۶/۳۴۹	۱	۶/۳۴۹	۵/۶۸	۰/۰۳۲	۰/۵۳۷	۰/۲۸۹
CL8	۲۴	-	-	-	-	-	-	-	-
CL9	۱۳	تعداد خودرو درگیر تصادف	۱۰/۶۷۴	۲	۵/۳۳۷	۱۳/۰۳۲	۰/۰۰۲	۰/۱۸۵۰	۰/۷۲۳
CL10	۱۱	مسیر صاف یا شیبدار	۲/۶۲۷	۱	۲/۶۲۷	۱۱/۲۶۰	۰/۰۰۸	۰/۷۴۵	۰/۵۵۶
CL11	۵۴	تعداد خودرو درگیر تصادف	۱۵/۸۶۲	۳	۵/۲۸۷	۴/۰۷۴	۰/۰۱۲	۰/۳۳	۰/۱۰۹
CL12	۱۴	طبقه راه (خیابان یا بزرگراه شهری) وجود یا عدم وجود گاردریل شرایط نور مسیر	۳/۸۵۷	۱	۳/۸۵۷	۱۰/۲۸۶	۰/۰۰۸	۰/۶۷۹	۰/۴۶۲
CL13	۱۰	طبقه راه (خیابان یا بزرگراه شهری) وجود یا عدم وجود شانه راه	۸/۸۱۷	۱	۹/۳۴۴	۱۱/۵۹۵	۰/۰۰۹	۰/۷۶۹	۰/۵۹۲

خوشه CL11

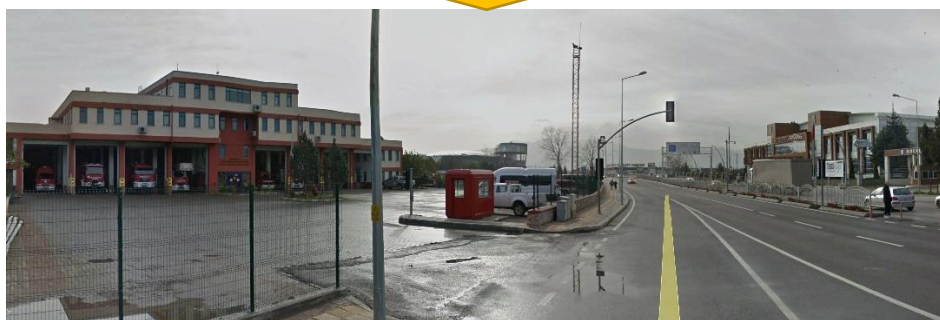
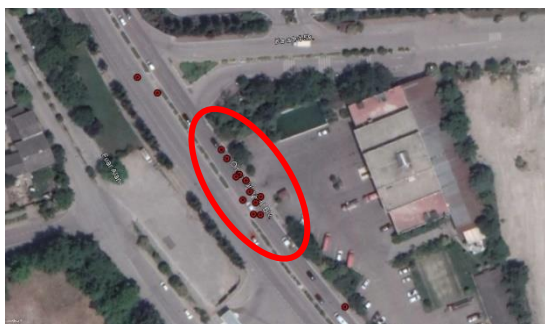
با بیشترین تعداد تصادف (54) در جنوب شهر و محور D100 ایزمیت به استانبول قرار دارد و در اطراف یکی از تقاطعات غیر همسطح اصلی در ورودی شهر می باشد (شکل ۴-۲). در این خوشه تعداد خودروهای درگیر تصادف به عنوان پارامتر مؤثر بر تعداد مجروحین مشخص شده است. همانطور که در جدول زیر قابل مشاهده است تصادفات دو خودرویی آمار قابل توجهی را نشان می دهد.

جدول ۴-۴: آمار توصیفی تعداد خودروهای درگیر تصادف در خوشه CL11

تعداد خودرو درگیر تصادف	N	واریانس
۱	۱۲	۱
۲	۳۲	۱/۱۱۲
۳	۷	۰/۹۰۵
۴	۳	۷
مجموع	۵۴	۱/۵۲۴

یکی از نقاط شناسایی شده درون این خوشه در خیابان تورکچاکال در مقابل آتش نشانی ایزمیت

می باشد. در این نقطه بیشتر تصادفات بین دو خودرو بوده است.



شکل ۴-۲۷: نقطه حادثه خیز در نزدیکی آتش نشانی ایزمیت



شکل ۴-۲۹: تقاطع غضنفر بیلگه شهر ایزمیت خوشه CL6

خوشه CL2

این خوشه در محور اصلی استانبول-ایزمیت در قسمت غربی شهر قرار دارد. در این خوشه پارامتر مسیر قوس به عنوان داده موثر بر تصادفات شناخته شده است. با توجه به طرح هندسی راه این خوشه در مسیر قوس خطرناکی قرار گرفته است که آمار تصادفات در این ناحیه قابل توجه است (شکل ۴-۳۰).



شکل ۴-۳۰: خوشه CL2 شهر ایزمیت

۴-۸- مقایسه روش‌های خوشه‌بندی و شاخص‌های ارزیابی خوشه‌ها

همانطور که در فصل دو اشاره گردید برای بررسی کیفیت خوشه‌بندی باید تعدادی از معیارهای درونی و بیرونی خوشه‌بندی را مورد بررسی قرار داد. در جدول ۴-۵ کیفیت خوشه‌بندی GriDBSCAN که در نرم افزار Elki انجام شد، قابل مشاهده است.

جدول ۴-۵: جدول شاخص‌های خوشه‌بندی GriDBSCAN

ایزیمیت	گبزه	
۰/۲۷۹	۰/۱۱۵	جاکارد
۰/۴۳۷	۰/۲۰۶	اندازه F
۱/۰۰۰	۱/۰۰۰	دقت
۰/۲۷۹	۰/۱۱۵	حساسیت
۰/۲۷۹	۰/۱۱۵	رند
۰/۰۰۰	۰/۰۰۰	رند اصلاحی
۰/۵۲۹	۰/۳۳۹	فولکرز-مالوز
۰/۳۴۰	۰/۲۴۰	NMI
۰/۵۸۰	۰/۴۹۰	مجذور NMI
۰/۴۳۷	۰/۴۲۸	اندازه F
۱/۰۰۰	۱/۰۰۰	حساسیت
۰/۲۷۹	۰/۲۷۲	دقت
۰/۶۷۰	۰/۵۶۷	اندازه F
۰/۵۰۳	۰/۳۹۶	شاخص خلوص
۱/۰۰۰	۱/۰۰۰	شاخص خلوص معکوس
۰/۶۶۲	۰/۵۶۴	اندازه F
۰/۵۰۲	۰/۳۹۴	دقت
۰/۹۷۴	۰/۹۸۹	حساسیت
شاخص جینی ± 0.3641 = میانگین		
میانگین ± 0.3605 =		

همانطور که در جدول ۴-۶ مشاهده می‌شود روش خوشه‌بندی GriDBSCAN در مدت زمان کمتری نسبت به الگوریتم معمولی DBSCAN اجرا می‌شود و برای داده‌هایی همچون داده‌های تصادفات که دارای حجم زیادی هستند، مدت زمان اجرای الگوریتم را به طور قابل توجهی کاهش می‌دهد. تفاوت عمده خوشه‌بندی DBSCAN معمولی با GriDBSCAN در اجرای سریعتر الگوریتم است.

جدول ۴-۶: مقادیر مورد استفاده برای الگوریتم GriDBSCAN

Izmit	Gebze	پارامترهای الگوریتم GriDBSCAN
۲۵۰	۳۵۰	پارامتر (m)Eps
۱۰	۱۵	minPts
۵۰۰	۷۰۰	d شبکه (m)
۱۴	۶	تعداد خوشه‌ها
۱۲	۱۰	GriDBSCAN Runtime (msec)
۱۹	۲۰	DBSCAN Runtime (msec)

در جدول ۴-۶ برای هر شهر مقادیر مورد استفاده برای الگوریتم و نتایج به دست آمده از الگوریتم (تعداد خوشه‌ها، تعداد داده‌های خوشه‌بندی شده و ...) نشان داده شده است. پس از بررسی تمام نتایج مقدار نرخ ضربه و شاخص دقت پیش‌بینی را برای مشخص کردن بهترین روش خوشه‌بندی استفاده شده در تحقیق حاضر در جدول ۴-۷ آمده است.

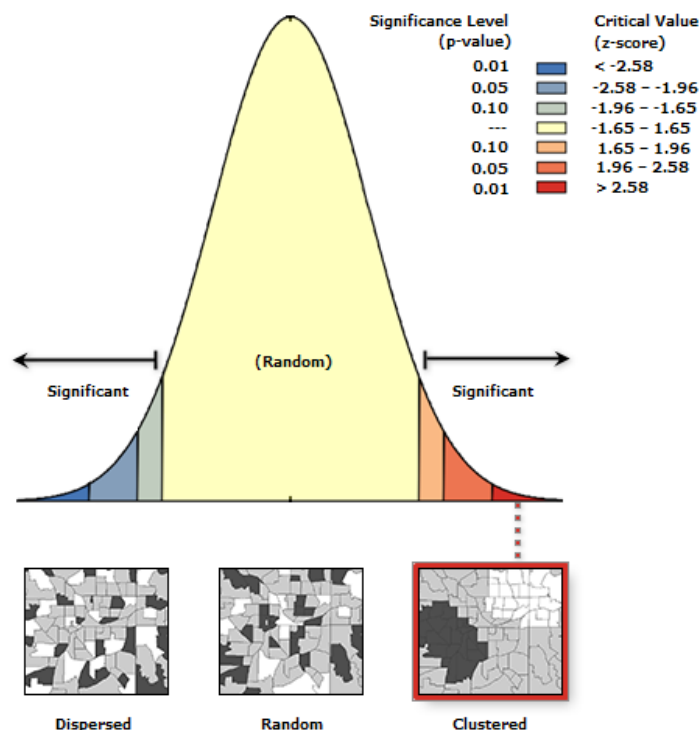
جدول ۴-۷: مقادیر نرخ ضربه و شاخص دقت پیش‌بینی روش‌های مختلف

KDE		Nnh		GriDBSCAN		نتایج حاصل از خوشه‌بندی
ایزمیت	گبزه	ایزمیت	گبزه	ایزمیت	گبزه	
۳۱۴۲۵۲۲۰	۶۷۸۰۲۹۳۸	۳۱۴۲۵۲۲۰	۶۷۸۰۲۹۳۸	۳۱۴۲۵۲۲۰	۶۷۸۰۲۹۳۸	مساحت منطقه مورد مطالعه
۱۷۲۷۹۸۰	۲۵۶۱۵۴۶	۳۱۹۹۲۹	۲۹۱۸۲۸	۱۲۸۰۴۲۲	۶۴۶۵۸۸۳	مساحت خوشه‌بندی شده (متر مربع)
۲۲۵	۲۸۰	۱۷۶	۲۵۴	۲۸۴	۳۹۱	تعداد نقاط خوشه‌بندی شده
۵۷۲	۶۴۷	۵۷۲	۶۴۷	۵۷۲	۶۴۷	تعداد نقاط کل
۳۰/۷۷	۴۳/۲۸	۳۰/۷۷	۳۹/۲۶	۴۹/۶۵	۶۰/۴۳	نرخ ضربه
۲/۹۲	۱۱/۴۵	۳۰/۲۲	۹۱/۲۱	۱/۱۷۸	۶/۳۷۷	شاخص PAI

با توجه به نتایج جدول ۴-۷ با توجه به مقادیر PAI مدل Nnh روش دقیق تری برای شناسایی نقاط حادثه‌خیز بدست آمده است. با استفاده از خوشه‌بندی روش GridBSCAN می‌توان مناطق متراکم تر را پیدا کرده و ویژگی‌های کلی منطقه‌ی تصادف را بدست آورد و سپس با روش Nnh برای شناسایی نقاط حادثه‌خیز اقدام نمود. الگوریتم GridBSCAN برای داده‌های انبوه بسیار مناسب است و برای شناسایی مناطق تراکم داده‌ها کاربردی است و داده‌های پرت را از خوشه‌ها جدا می‌سازد.

۹-۴- نتایج حاصل از توابع موجود در GIS

در این قسمت نتایج توابع خودهمبستگی موران و گتیس‌ارد که برای تعیین الگوی پراکنش داده‌ها است آمده است. با توجه به توضیحات ذکر شده در فصل دو این تابع نشان می‌دهد آیا داده‌های مورد نظر دارای الگوی خوشه‌ای هستند یا نه. در شکل ۴-۳۱ مقادیر شاخص موران، امتیاز استاندارد (Z-Score) و P-Value برای شهر ایزمیت قابل مشاهده است.

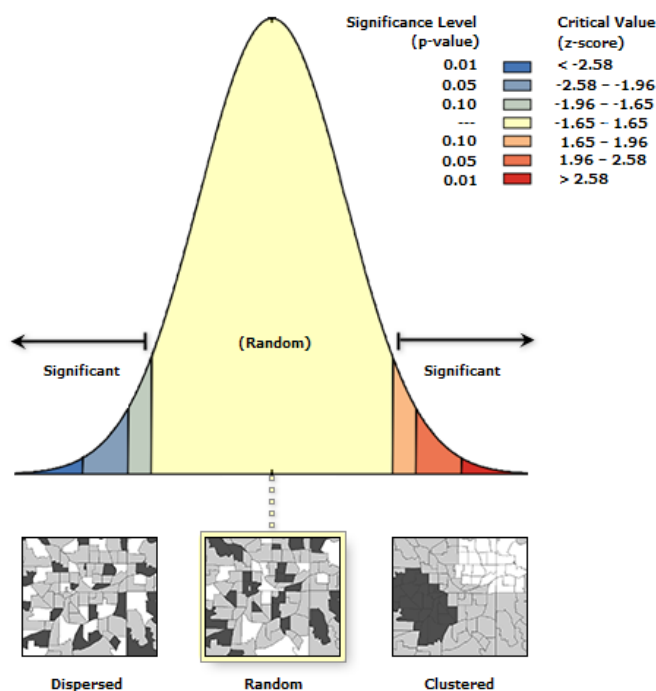


شکل ۴-۳۱: نتیجه تابع خود همبستگی موران برای شهر ایزمیت

جدول ۴-۸: نتایج تابع خودهمبستگی موران برای شهر ایزمیت

Izmit		Gebze	
Moran's Index:	۰/۰۵۷۱۱۷	Moran's Index:	-۰/۰۰۳۳۶۲
z-score:	۳/۱۳۳۸۴۰	z-score:	-۰/۱۲۵۵۸۸
p-value:	۰/۰۰۱۷۲۵	p-value:	۰/۹۰۰۰۵۸

با توجه به مقادیر به دست آمده P-Value که مقدار کمتری دارد و ولی شاخص موران نزدیک به ۱ و امتیاز استاندارد Z مثبت داده‌های مورد استفاده در شهر ایزمیت دارای الگوی خوشه‌ای هستند. تابع موران الگوی خوشه‌ای بودن داده‌های مورد نظر را تایید می‌کند. اما در شهر گبزه که شاخص موران منفی به دست آمده و مقدار P-Value بالاتری دارد، دارای الگوی تصادفی هستند و خوشه‌بندی داده‌ها این شهر الگوی تصادفی دارد.



شکل ۴-۳۲: نتیجه تابع خود همبستگی موران برای شهر ایزمیت

فصل ۵

بحث و نتیجه‌گیری

۵-۱- مقدمه

تاکنون تحقیقات زیادی به منظور ایجاد ارتباط بین مبحث خوشه‌بندی و تحلیل تصادفات انجام گرفته است. همان‌طور که در پیشینه تحقیق بیان شد؛ محققین از روش‌های مختلف داده‌کاوی همچون طبقه‌بندی و خوشه‌بندی برای تحلیل تصادفات استفاده کرده‌اند. در بیشتر مطالعات تصادفات، تصادفات برون‌شهری مدنظر بوده و تصادفات درون‌شهری به علت ساختار پیچیده و به علت شبکه‌ای بودن راه در مناطق شهری کمتر مورد مطالعه قرار گرفته است. اما به علت افزایش تصادفات درون‌شهری در کشورهای در حال توسعه تصمیم بر این شد تا نقاط حادثه‌خیز تصادفات را شناسایی کرده و عوامل مؤثر بر آن‌ها را مورد بررسی قرار دهیم. تحلیل مکانی تصادفات و اهمیت سیستم اطلاعات جغرافیایی در تحلیل‌های مکانی سبب شد که در این تحقیق به تحلیل مکانی تصادفات با الگوریتم‌های خوشه‌بندی مرتبط با بحث‌های مکانی پرداخته شود.

۵-۲- خلاصه نتایج و جمع‌بندی

با توجه به نتایج به دست آمده از روش‌های خوشه‌بندی، روش‌های خوشه‌بندی غیر آماری Nnh، GridBSCAN و روش خوشه‌بندی آماری برآورد تراکم کرنل انتخاب شدند. با توجه به اینکه روش‌های مبتنی بر تراکم برای تحلیل نقاط حادثه‌خیز روشی قابل اطمینان‌تر و مطلوب‌تر هستند و نتایج نسبتاً پایداری را به کاربر ارائه می‌دهند. داده‌های تصادفات با استفاده از این روش‌ها تحلیل شدند. در نهایت مناطق حادثه‌خیز شهری به خوبی شناسایی شدند و این مناطق شناسایی شده با توجه به واقع شدن در مناطق مسکونی یا اداری و یا تقاطعات شهری از لحاظ انطباق با دنیای واقعی نیز در مکان‌های مستعد تصادف قرار گرفته‌اند.

باتوجه به مقدار شاخص دقت پیش‌بینی PAI برای روش‌های خوشه‌بندی انجام شده روش Nnh به عنوان بهترین روش برای شناسایی نقاط حادثه‌خیز بدست آمد. تفسیر خوشه‌ها برای هر کدام از شهرها

نشان داد که هر شهر ویژگی‌های مختص به خود را دارد و نقاط حادثه‌خیز نتیجه‌ی تجمیع عوامل مختلف ترافیکی، هندسی، و دیگر عوامل در یک منطقه است.

در تحقیق حاضر در شهر ایزمیت اغلب نقاط حادثه‌خیز تصادف در محور اصلی D100 ایزمیت- استانبول اتفاق افتاده است و بیشترین تصادفات در این محور در نزدیکی خروجی‌ها و ورودی‌ها بود و همینطور در محور وسط این بزرگراه در نزدیکی ایستگاه‌های اتوبوس و پل عابریپاده تصادفات رخ داده است. در شبکه شهری ایزمیت نیز تصادفات بیشتر در نزدیکی تقاطعات بود. روش‌های مورد استفاده همگی توانستند به خوبی نقاط تمرکز تصادف را خوشه بندی کنند.

در شهر گبزه نیز با توجه به توپوگرافی متنوع و شیب‌دار و همچنین تراکم بالای کاربری‌های مختلف در عرض معابر کم تصادفات به صورت پراکنده در اکثر تقاطعات و خیابان‌های محلی و دسترسی اتفاق افتاده است .

با توجه به توزیع زمانی داده‌ها در هر دو شهر، اکثر تصادفات در زمان اوج ترافیک عصر و سپس در زمان اوج ترافیک صبح اتفاق افتاده است. سازمان‌های ترافیکی و ارگان‌های مرتبط، باید با توجه به عوامل به دست آمده مؤثر بر تصادفات در این دو شهر در جهت رفع این عوامل تصمیم‌گیری کرده و اقدام نمایند.

در انتها با توجه به مقدار شاخص دقت پیش‌بینی بادقت ترین روش، روش Nnh بدست آمد. با توجه به عنوان تحقیق سعی بر آن بود تا روش‌های مختلف خوشه‌بندی تصادفات را مورد بررسی قرار دهیم و روشی کاربردی ارئه دهیم. در این تحقیق الگوریتمی جدید به عنوان GridBSCAN مورد استفاده گردید. این روش نسبت به DBSCAN معمولی از لحاظ زمانی ۵۰٪ درصد سریعتر اجرا می‌شود.

۵-۳- پیشنهادات کارهای آینده

- باتوجه به نبود داده‌هایی همچون حجم ترافیک، عرض دقیق معابر توصیه می‌شود تأثیر این دو عامل بر تصادفات را نیز در تحقیقات آتی در نظر گرفته‌شود.
 - از الگوریتم GridBSCAN برای خوشه‌بندی تصادفات برون‌شهری و برای آزادراه‌ها و بزرگراه نیز استفاده گردد تا مشخص گردد این الگوریتم در کدام مناطق کاربردی‌تر است.
 - با توجه به زمان اجرای بالای این الگوریتم بهتر است برای داده‌های بزرگ^۱ همانند تصادفات یا جرایم در مناطق درون‌شهری که دارای ویژگی‌های مکانی هستند نیز استفاده شود.
 - وضعیت جوی دقیق نیز در ماه‌های مختلف سال و ساعات مختلف روز تغییر می‌کند؛ لذا بررسی نقش عامل فوق با روش‌های حاضر در این تحقیق در تصادفات قابل انکار نیست و باید مورد توجه قرار گیرند.
- ارائه راهکار جهت کاهش تعداد و خسارات تصادفات عابرین پیاده با توجه به اثرات عوامل موثر و انجام مطالعات قبل-بعد با همکاری نهاد‌های اجرایی جهت ارزیابی میزان کارایی این راهکارها.

^۱ Big Data

مراجع:

احدی، م.، زایرزاده، ع.، ۱۳۸۴. پارامترهای مؤثر بر خستگی رانندگان و نقش آن در وقوع تصادفات. اولین کنفرانس بین‌المللی حوادث رانندگی و جاده‌ای.

جمال شهرابی، ۱۳۹۲. داده‌کاوی. جهاد دانشگاهی (دانشگاه صنعتی امیرکبیر).

سجادی، ح.، ۱۳۸۸. تحلیل تصادفات جاده‌ای در سیستم اطلاعات جغرافیایی با تأکید بر خصوصیات جاده و محیط (مطالعه موردی: محور قزوین - رشت). وزارت علوم، تحقیقات و فناوری - دانشگاه شهید بهشتی - دانشکده علوم زمین.

عفتی، م.، رجبی، م.ع.، حکیم پور، ف.، شعبانی، ش.، ۱۳۹۳. تحلیل عوامل مکانی مؤثر بر تمرکز تصادفات در راه‌های برون‌شهری با استفاده از GIS و داده‌کاوی. علوم و فنون نقشه‌برداری سال چهارم ۲، ۸۷-۱۰۲.

محمد سلمانی، مهدی رمضان زاده لسبویی، مسلم دریکوند، فرخ ثابتی، ۱۳۸۷. بررسی عوامل مؤثر بر تصادفات جاده‌ای و ارائه راه‌کارهایی برای کاهش آن؛ مورد مطالعه: منظومه روستایی جنوب خور و بیابانک. پژوهش‌های جغرافیای انسانی سال چهارم ۶۵، ۸۷-۱۰۴.

معمد، م.، ۱۳۸۸. مدل‌سازی تصادفات ترافیکی جهت ارتقاء ایمنی با استفاده از روش‌های آماری پیشرفته. وزارت علوم، تحقیقات و فناوری - دانشگاه صنعتی خواجه‌نصیرالدین طوسی - دانشکده عمران.

Alotaibi, A.S., 2018. Density-based clustering for road accident data analysis. *Int. J. Adv. Appl. Sci.* 5 8 , 113–121.

Anderson, T.K., 2009. Kernel density estimation and K-means clustering to profile road accident hotspots. *Accid. Anal. Prev.* 41 3 , 359–364.

Andersson, A.K., Chapman, L., 2011. The impact of climate change on winter road maintenance and traffic accidents in West Midlands, UK. *Accid. Anal. Prev.* 43 1 , 284–

289.

Association, W.R., 2016. Road Safety Manual: Recommendations from the World Road Association (PIARC).

Bradley, P.S., Mangasarian, O.L., Street, W.N., 1997. Clustering via concave minimization, in: Advances in Neural Information Processing Systems. pp. 368–374.

Chainey, S., Tompson, L., Uhlig, S., 2008. The utility of hotspot mapping for predicting spatial patterns of crime. Secur. J. 21 1–2 , 4–28.

Chang, L.-Y., Chen, W.-C., 2005. Data mining of tree-based models to analyze freeway accident frequency. J. Safety Res. 36 4 , 365–375.

Chen, S., Saeed, T.U., Alinizzi, M., Lavrenz, S., Labi, S., 2019. Safety sensitivity to roadway characteristics: a comparison across highway classes. Accid. Anal. Prev. 123, 39–50.

Chimba, D., Sando, T., 2009. The prediction of highway traffic accident injury severity with neuromorphic techniques. Adv. Transp. Stud. 2009 19 , 17–26.

Delen, D., Sharda, R., Bessonov, M., 2006. Identifying significant predictors of injury severity in traffic accidents using a series of artificial neural networks. Accid. Anal. Prev. 38 3 , 434–444.

Dobra, A., 2002. Classification and regression tree construction. Retrieved Sept. 18, 2011.

Dunham, M.H., 2003. Data Mining: Introductory and Advanced Topics, Pearson Education India. Pearson Education India.

Ester, M., Kriegel, H.-P., Sander, J., Xu, X., 1996. A density-based algorithm for discovering clusters in large spatial databases with noise, in: Kdd. pp. 226–231.

Fowlkes, E.B., Mallows, C.L., 1983. A method for comparing two hierarchical clusterings. J. Am. Stat. Assoc. 78 383 , 553–569.

Gaonkar, M.N., Sawant, K., 2013. AutoEpsDBSCAN: DBSCAN with Eps automatic for

large dataset. *Int. J. Adv. Comput. Theory Eng.* 2 2 , 11–16.

Geurts, K., Wets, G., Brijs, T., Vanhoof, K., 2003. Profiling of high-frequency accident locations by use of association rules. *Transp. Res. Rec. J. Transp. Res. Board* 1840 , 123–130.

Grabusts, P., Borisov, A., 2002. Using grid-clustering methods in data classification, in: *Parallel Computing in Electrical Engineering, 2002. PARELEC'02. Proceedings. International Conference On. IEEE*, pp. 425–426.

Han, J., Pei, J., Kamber, M., 2011. *Data mining: concepts and techniques*. Elsevier.

Ivan, J.N., Wang, C., Bernardo, N.R., 2000. Explaining two-lane highway crash rates using land use and hourly exposure. *Accid. Anal. Prev.* 32 6 , 787–795.

Jaccard, P., 1901. Étude comparative de la distribution florale dans une portion des Alpes et des Jura. *Bull Soc Vaudoise Sci Nat* 37, 547–579.

Kaufman, L., Rousseeuw, P.J., 1990. Partitioning around medoids (program pam). *Find. groups data an Introd. to Clust. Anal.* 68–125.

Kovács, F., Legány, C., Babos, A., 2005. Cluster validity measurement techniques, in: *6th International Symposium of Hungarian Researchers on Computational Intelligence*. Citeseer.

Kumar, S., Toshniwal, D., 2016. Analysis of hourly road accident counts using hierarchical clustering and cophenetic correlation coefficient (CPCC). *J. Big Data* 3 1 , 1–11. doi:10.1186/s40537-016-0046-3

Kunt, M.M., Aghayan, I., Noii, N., 2011. Prediction for traffic accident severity: comparing the artificial neural network, genetic algorithm, combined genetic algorithm and pattern search methods. *Transport* 26 4 , 353–366.

Levine, N., 2010. *CrimeStat III: A Spatial Statistics Program for the Analysis of Crime Incident Locations (version 3.3)*. Ned Levine & Associates: Houston, National Institute of TX.

Loh, W., 2008. *Classification and regression tree methods*. Wiley StatsRef Stat. Ref.

Online.

MacQueen, J., 1967. Some methods for classification and analysis of multivariate observations, in: Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability. Oakland, CA, USA, pp. 281–297.

Mahran, S., Mahar, K., 2008. Using grid for accelerating density-based clustering, in: Computer and Information Technology, 2008. CIT 2008. 8th IEEE International Conference On. IEEE, pp. 35–40.

Moghaddam, F.R., Afandizadeh, S., Ziyadi, M., 2011. Prediction of accident severity using artificial neural networks. *Int. J. Civ. Eng.* 9 1 , 41.

Moreira, A., Santos, M.Y., Carneiro, S., 2005. Density-based clustering algorithms–DBSCAN and SNN. Univ. Minho-Portugal.

Oh, C., Oh, J.-S., Ritchie, S., Chang, M., 2001. Real-time estimation of freeway accident likelihood, in: 80th Annual Meeting of the Transportation Research Board, Washington, DC.

Oh, J.-S., Oh, C., Ritchie, S.G., Chang, M., 2005. Real-time estimation of accident likelihood for safety enhancement. *J. Transp. Eng.* 131 5 , 358–363.

Polat, K., Durduran, S.S., 2011. Subtractive clustering attribute weighting (SCAW) to discriminate the traffic accidents on Konya–Afyonkarahisar highway in Turkey with the help of GIS: A case study. *Adv. Eng. Softw.* 42 7 , 491–500.

Prabha, S.L., Shanavas, A.R.M., 2014. Educational data mining applications. *Oper. Res. Appl. An Int. J.* 1 1 , 1–6.

Prasannakumar, V., Vijith, H., Charutha, R., Geetha, N., 2011. Spatio-temporal clustering of road accidents: GIS based analysis and assessment. *Procedia-Social Behav. Sci.* 21, 317–325.

Rand, W.M., 1971. Objective criteria for the evaluation of clustering methods. *J. Am. Stat. Assoc.* 66 336 , 846–850.

Sabel, C.E., Kingham, S., Nicholson, A., Bartie, P., 2005. Road traffic accident

simulation modelling-a kernel estimation approach, in: The 17th Annual Colloquium of the Spatial Information Research Centre University of Otago, Dunedin, New Zealand. Citeseer, pp. 67–75.

Sagberg, F., 2018. Characteristics of fatal road crashes involving unlicensed drivers or riders: implications for countermeasures. *Accid. Anal. Prev.* 117, 270–275.

Sasaki, Y., 2007. The truth of the F-measure. *Teach Tutor mater* 15, 1–5.

Shafabakhsh, G.A., Famili, A., Bahadori, M.S., 2017. GIS-based spatial analysis of urban traffic accidents: Case study in Mashhad, Iran. *J. traffic Transp. Eng. (English Ed.* 43, 290–299.

Tan, P.-N., Steinbach, M., Kumar, V., 2005. *Introduction to data mining*. 1st.

Tian, R., Yang, Z., Zhang, M., 2010. Method of road traffic accidents causes analysis based on data mining, in: *Computational Intelligence and Software Engineering (CiSE), 2010 International Conference On*. IEEE, pp. 1–4.

Tseng, W.-S., Nguyen, H., Liebowitz, J., Agresti, W., 2005. Distractions and motor vehicle accidents: Data mining application on fatality analysis reporting system (FARS) data files. *Ind. Manag. Data Syst.* 105 9, 1188–1205.

Turkstat, T.S.I., 2018. *Turkish Statistical Institute Road Traffic Accident Statistics 2014*.

World-Health-Organization, 2015. *Global status report on road safety 2015*. World Health Organization.

Yi-ling, Z., Hongbo, X., Shuo, B.A.I., 2008. OPTICS-Plus for Text Clustering. *J. Chinese Inf. Process.* 22 1, 51–60.

Abstract

Based on traffic accident statistics, it is clear that traffic accidents impose the harmful effects and unnecessary costs on society; thus, researchers try to analysis traffic accidents, identify accident locations, and determine significant parameters. Today, data mining techniques, specifically spatial clustering methods, are used to analyze accident data and find their spatial patterns. In this research, a combined method called GridBSCAN algorithm was used for clustering accident data. The DBSCAN algorithm applied for its spatial characteristics and the GRID algorithm used by taking data into a grid mesh increased the accuracy and execution time in big data such as accident data. Other clustering methods such as K-Means, Nnh, KDE and spatial autocorrelation were also used to be compared with GridBSCAN method. The results showed that Nnh algorithm was the most accurate method for spatial clustering of traffic accident points, and GridBSCAN algorithm was also applicable to the separation of high density areas of accidents in high volume data. The clusters obtained from GridBSCAN algorithm properly differentiated the accidents of different urban areas according to the density. By examining the parameters in these clusters, it is possible to reduce the number of accidents and understand the factors affecting them.

Keywords: Urban accidents, data mining algorithms, spatial clustering analysis, GridBSCAN algorithm



Shahrood University of Technology

Faculty of Civil Engineering

M.Sc. Thesis in Transportation Engineering

Presenting a spatial clustering method for crash data

By: Amin Ganjali Khosrowshahi

Supervisor:

Dr. Iman Aghayan

Advisor:

Dr. Abdolahad Choupani

Jan 2019