



دانشگاه صنعتی امیر کبیر
(پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر

گزارش پروژه‌ی پایانی درس یادگیری ماشین

عنوان مقاله:

زیر نمونه برداری برای روش‌های ترکیبی و بدون
نظارت جهت کشف داده‌های پرت به صورت کارآمد
و مؤثر

نام دانشجو: سید احمد نقوی نوزاد

ش-د: ۹۴۱۳۱۰۶۰

استاد درس:

دکتر ناظر فرد

بهار ۱۳۹۵

سید علی

(۱) مقدمه:

در ابتدا تعریف مختصری از داده‌ی پرت ارائه داده و سپس به لزوم کشف داده‌های پرت پرداخته و در ادامه به دو دسته‌بندی کلی از انواع روش‌های کشف داده‌های پرت اشاره خواهیم نمود. در نهایت به بررسی مختصری در مورد رویکرد پیشنهادشده در این مقاله خواهیم پرداخت.

یک داده‌ی پرت به یک داده‌ی ضبط‌شده و یا مجموعه‌ای از داده‌ها اطلاق می‌گردد که بنا به ظاهر نسبت به سایر مجموعه داده ناسازگار بوده و رفتار غیر نرمالی از خود بروز می‌دهند. حال این رفتار می‌تواند نسبت به کلیت مجموعه داده مورد بررسی قرار گیرد (داده‌ی پرت سراسری یا Global Outlier) و یا هم در حیطه‌ی یک همسایگی خاص از داده‌ی مورد نظر (مانند K نزدیک‌ترین همسایه) مورد توجه واقع گردد (داده‌ی پرت محلی یا Local Outlier). در این مقاله ما به طور خاص به دنبال کشف داده‌های پرت محلی با استفاده از روش‌های اصطلاحاً ترکیبی (ensemble) خواهیم بود.

کشف داده‌های پرت از اهمیت ویژه‌ای در بسیاری از کاربردهای عملی، نظیر کشف خطاهای اندازه‌گیری توسط حسگرهای مختلف (که به نوعی با داده‌ی پرت به عنوان نویز رفتار می‌کنند)، کشف سوء استفاده از کارت‌های اعتباری و یا هم یافتن اندازه‌گیری‌های غیرمعارف در داده‌های علمی دارد؛ چرا که رفتار غیرمعارفی که داده‌های پرت نسبت به باقی داده‌ها از خود بروز می‌دهند، می‌تواند موجب بروز مشکلات عدیده گشته و در نتیجه نیازمند توجهات خاصی بوده و البته می‌توانند سبب به وجود آمدن بینش‌های جدیدی نیز نسبت به مجموعه داده‌ی مورد بررسی گردند.

روش‌ها و رویکردهای متعددی جهت کشف داده‌های پرت تاکنون ارائه شده‌اند، که البته هر کدام با توجه به نوع داده‌ی پرتی که بر روی آن متمرکز شده‌اند و یا هم کاربرد و زمینه‌ی خاصی که نیازمند کشف داده‌ی پرت می‌باشد، با دیگران تفاوت دارند. دسته‌ی اول، رویکردهای آماری پارامتریک^۱ می‌باشند، که سعی دارند تا با استفاده از تخمین پارامترهای یک تخمین فرضی، آن تخمین را به مجموعه داده‌ی مورد بررسی نسبت دهند. اما مشکل موجود در مورد این نوع رویکردهای مبتنی بر یک توزیع فرضی خاص، این است که پارامترهای توزیع مدنظر نسبت به حضور داده‌های پرت حساس بوده و ممکن است در محاسبه‌ی آن‌ها دچار خطا شویم. از جمله اثرات مضر که داده‌های پرت بر روی تخمین پارامترهای توزیع مد نظر می‌گذارند، می‌توان به پدیده‌هایی تحت عنوان «پوشش»^۲ و نیز «غرقانیدن»^۳ اشاره نمود. و اما دسته‌ی دوم از رویکردهای مورد استفاده جهت کشف داده‌های پرت، با نام رویکردهای غیرپارامتریک^۴ شناخته می‌شوند، که بر عکس مورد قبلی، توزیع خاصی را به داده‌ها نسبت نداده بلکه سعی در آن دارند تا به صورت آشکارا و یا هم ضمنی، وجهه‌های خاصی از تابع توزیع چگالی احتمالاتی را تخمین بزنند. از جمله روش‌های مبتنی بر این نوع رویکرد، روش‌های «مبتنی بر فاصله» و روش‌های «مبتنی بر چگالی» می‌باشند که البته این روش‌ها نیز می‌توانند مانند مورد قبل، از اثرات مضر مشابه پوشش و غرقانیدن که در نتیجه‌ی عمل تخمین‌زدن فاصله و چگالی صورت می‌گیرد، آسیب ببینند.

و اما در مورد رویکرد مورد استفاده در این مقاله باید گفت که با توجه به این مطلب که تکنیک‌های ترکیبی در مورد کشف داده‌های پرت کمتر مورد مطالعه قرار گرفته‌اند، لذا قصد داریم تا به «زیرنمونه‌برداری» تحت عنوان یک تکنیک خاص بپردازیم، که سبب می‌شود تا تنوع میان انواع روش‌های کشف داده‌ی پرت آشکار گردد. هم‌چنین به صورت آماری و مبتنی بر آزمایشات نشان خواهیم داد که یک روش کشف‌کننده‌ی داده‌ی پرت، بر روی یک

¹ Parametric Statistical Approaches

² Masking

³ Swamping

⁴ Non-parametric Approaches

زیرنمونه به خودی خود، علاوه بر اینکه سبب استنتاج تنوع میان انواع روش‌ها می‌گردد، تحت شرایط خاصی از عملکرد همان روش کشف داده‌ی پرت بر روی کل مجموعه‌داده پیشی می‌گیرد. البته که ساخت یک روش ترکیبی که بر روی زیرنمونه‌های متعدد کار می‌کند، می‌تواند سبب بهبود چشمگیر نتایج گردد؛ و نیز ذکر این نکته ضروری می‌نماید که روش‌های ترکیبی مطرح‌شده تاکنون در مقالات متعدد، روش‌های بانظارت^۵ بوده و تلاش در دسته‌بندی^۶ داده‌ها دارند، ولی در این مقاله تلاش خواهیم داشت تا به صورت آماری، صحت عملکرد روش‌های ترکیبی بی‌نظارت^۷ و مبتنی بر خوشه‌بندی^۸ در کشف داده‌های پرت را به اثبات برسانیم. در پایان باید گفت که استفاده از ترکیبی از روش‌های کشف داده‌های پرت بر روی زیرنمونه‌های متعدد از مجموعه‌داده، با توجه به اندازه‌ی زیرنمونه‌ها^۹ و نیز حجم کلی روش‌های ترکیب‌شده^{۱۰}، به طور کلی نسبت به استفاده از یک روش منفرد کشف داده‌ی پرت بر روی کل مجموعه‌داده عملکرد بهتری خواهد داشت.

۲) شرح روش و پارامترها:

در این مقاله این‌گونه نیست که تنها از مجموعه‌داده نمونه‌برداری کرده و سپس الگوریتم کشف داده‌ی پرت را بر روی زیرنمونه‌های مربوطه و با احتساب حضور سایر داده‌های خارج از زیرنمونه اجرا نمائیم، چرا که با این کار اطلاعات زیادی در مورد ماهیت پرت‌بودن بسیاری از داده‌ها از بین رفته و نیز بسیاری از داده‌ها نیز تنها امتیازی حسب پرت‌بودن از تنها برخی از زیرنمونه‌ها خواهند برد. بلکه در عوض ما در این مقاله، به ازای هر عضو ensemble، یک زیرنمونه از مجموعه داده انتخاب کرده و همسایگی مربوطه را تنها به ازای داده‌های موجود در زیرنمونه برای هر کدام از الگوریتم‌های کشف داده‌ی پرت محاسبه می‌نمائیم، که این روش سبب افزایش سرعت قابل توجهی در مقایسه با سایر روش‌های ensemble گشته و نیز نتایج نهائی بهتری را حاصل می‌نماید.

در این مقاله یک روش پایه را با نام Feature Bagging، به عنوان رقیب روش ensemble معرفی شده در این مقاله مورد استفاده قرار می‌دهیم و برای هر دوی روش‌های ترکیبی (ensemble و Feature Bagging) یک مقدار معین برای تعداد اعضای ensemble یعنی ۲۵ را انتخاب خواهیم نمود. در مورد روش رقیب یعنی Feature Bagging نیز جهت ترکیب نتایج نهائی، از عمل میانگین‌گیری استفاده خواهیم نمود. در مورد روش ensemble نیز از زیرنمونه‌های با اندازه‌های متفاوت جهت ارزیابی نتایج نهائی استفاده نموده و در مورد روش‌های پایه نیز از اندازه‌های مختلف k برای یافتن نزدیک‌ترین همسایه‌ها استفاده می‌کنیم. در نهایت نتایج نهائی را در قالب‌های زیر نشان خواهیم داد: (i) مقدار ثابت k و اندازه‌های مختلف زیرنمونه؛ (ii) مقدار ثابت زیرنمونه و مقادیر مختلف k؛ (iii) مقادیر ثابت برای هر دوی اندازه‌ی زیرنمونه و مقدار k. در ضمن مقداری را برای k انتخاب خواهیم نمود که در مورد روش پایه به نتایج مطلوبی منجر گردد.

هم‌چنین برای ارزیابی نتایج نهائی از ROC AUC^{۱۱} استفاده می‌نمائیم که مقادیر مختلف True Positive Rate را بر حسب مقادیر False Positive Rate رسم نموده و نیز یک معیار عمومی و پرکاربرد جهت ارزیابی روش‌های مختلف کشف داده‌های پرت می‌باشد.

⁵ Supervised

⁶ Classification

⁷ Unsupervised

⁸ Clustering

⁹ Subsample size

¹⁰ Ensemble size

¹¹ Area under the receiver operating characteristic curve (ROC AUC)

(۳) مجموعه داده‌های مورد استفاده

در این جا دو مجموعه داده‌ی مستقل دست‌ساز و مصنوعی^{۱۲} را تولید می‌نمائیم (batch1 و batch2) که هر کدام از متشکل از ۳۰ مجموعه داده‌ی مختلف با ابعاد، تعداد خوشه‌ها و نیز اندازه‌ی متفاوت خوشه‌ها می‌باشند. مجموعه‌ی داده‌های مربوطه را با استفاده از توزیع گاوسیین مخلوط^{۱۳} تولید می‌نمائیم و برای این کار از مقادیر مختلف میانگین و انحراف از معیار که از هر کدام در بازه‌های مشخصی می‌باشند استفاده می‌نمائیم. در نهایت می‌بایست مجموعه داده‌های مربوطه را در فضای چندبعدی چرخش^{۱۴} دهیم که برای این کار نیز ابتدا یک ماتریس A را با داده‌های رندوم نرمال تولید نموده و سپس با ضرب ترانهاده‌ی این ماتریس در خود آن، یک ماتریس وارایانس-کوواریانس به دست می‌آوریم که خاصیت لازم به اصلاح positive-semiDefinite در مورد آن برقرار می‌باشد. بعد از این با استفاده از تابع mvnrnd() یک توزیع گاوسیین مخلوط را در مورد داده‌های مربوط به هر کدام از خوشه‌های یک مجموعه داده‌ی مشخص به دست می‌آوریم. جهت این که داده‌های پرت محلی را در مورد این مجموعه داده‌ها برچسب‌گذاری نمائیم، در ابتدا فاصله‌ی مالهالانوبیس^{۱۵} را در مورد داده‌های موجود درون هر خوشه از مرکز خوشه محاسبه نموده و سپس آن داده‌هایی را که این فاصله برای آن‌ها بیشتر از چارک ۰,۹۷۵ توزیع مربوط به فاصله‌های داده‌های درون خوشه باشد، به عنوان داده‌ی پرت برچسب‌گذاری می‌نمائیم. در مورد مجموعه داده‌های واقعی مورد استفاده نیز، در این جا از مجموعه داده‌های معروف Satimage, Segment (که به سه مجموعه داده‌ی مجزا با داده‌های پرت متفاوت تبدیل خواهد شد)، Wisconsin Breast Cancer (WBC) و نیز مجموعه داده‌ی Waveform Database Generator (waveform) استفاده خواهیم نمود. در مورد هر کدام از این مجموعه داده‌ها در ابتدا آن کلاسی را که تعداد داده‌های کمتری را داراست به عنوان طعمه انتخاب نموده و سپس مانند مجموعه داده‌های مصنوعی که پیش ازین قید گردید، آن داده‌هایی را که فاصله‌ی مالهالانوبیس آن‌ها از مرکز خوشه، بیشتر از چارک ۰,۹۷۵ توزیع مربوط به فاصله‌های داده‌های درون خوشه باشد، به عنوان داده‌ی پرت برچسب‌گذاری می‌نمائیم. در نهایت از الگوریتم‌های کشف داده‌های پرت محلی جهت شناسائی داده‌های پرت در مجموعه داده‌های نامبرده استفاده خواهیم نمود.

(۴) نتایج مربوط به پیاده‌سازی

(۱) مجموعه داده‌های مصنوعی

در این جا مجموعه داده‌ی مصنوعی batch1 متشکل از ۳۰ مجموعه داده‌ی مجزا را به عنوان کاندید برگزیده و دو روش پایه جهت کشف داده‌های پرت محلی با نام‌های LOF^{۱۶} و LoOP^{۱۷} را بر روی آن‌ها پیاده نموده و روش‌های ترکیبی قیدشده پیش ازین، با نام‌های FeatureBagging و ensemble را نیز با پارامترهای مشخص (اندازه‌ی زیرنمونه‌های مختلف و مقدار مشخص $k=3$) بر روی آن‌ها پیاده می‌نمائیم. در این جا مقادیر ROC AUC مختلف به ازای زیرنمونه‌های مختلف را به دست آورده و برای رسم آن‌ها از تابع boxplot() استفاده می‌نمائیم. نتایج بیان شده در مقاله و نیز نتایج حاصله از پیاده‌سازی در ادامه می‌آید:

¹² Synthetic

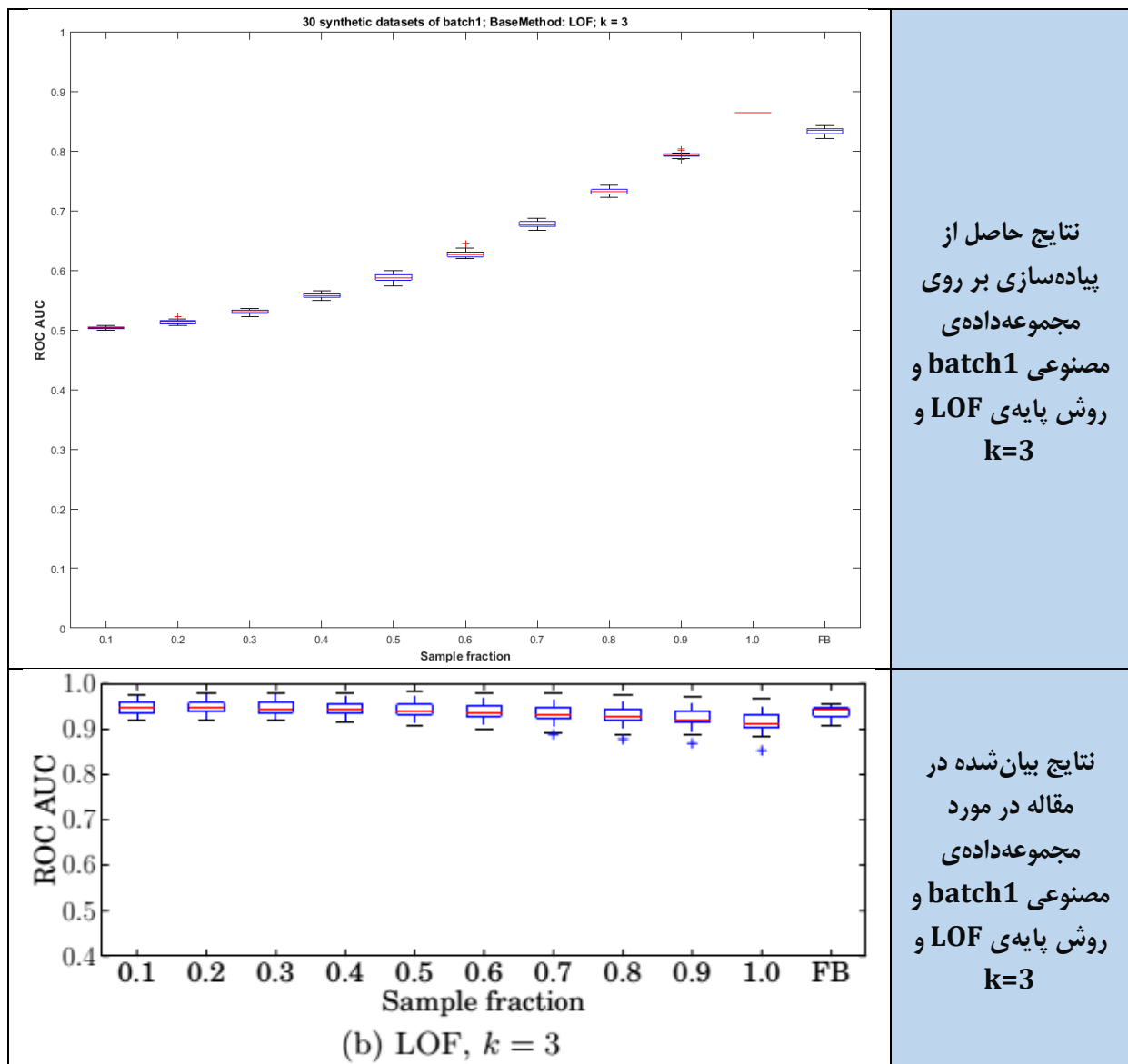
¹³ Gaussian Mixture Distributions

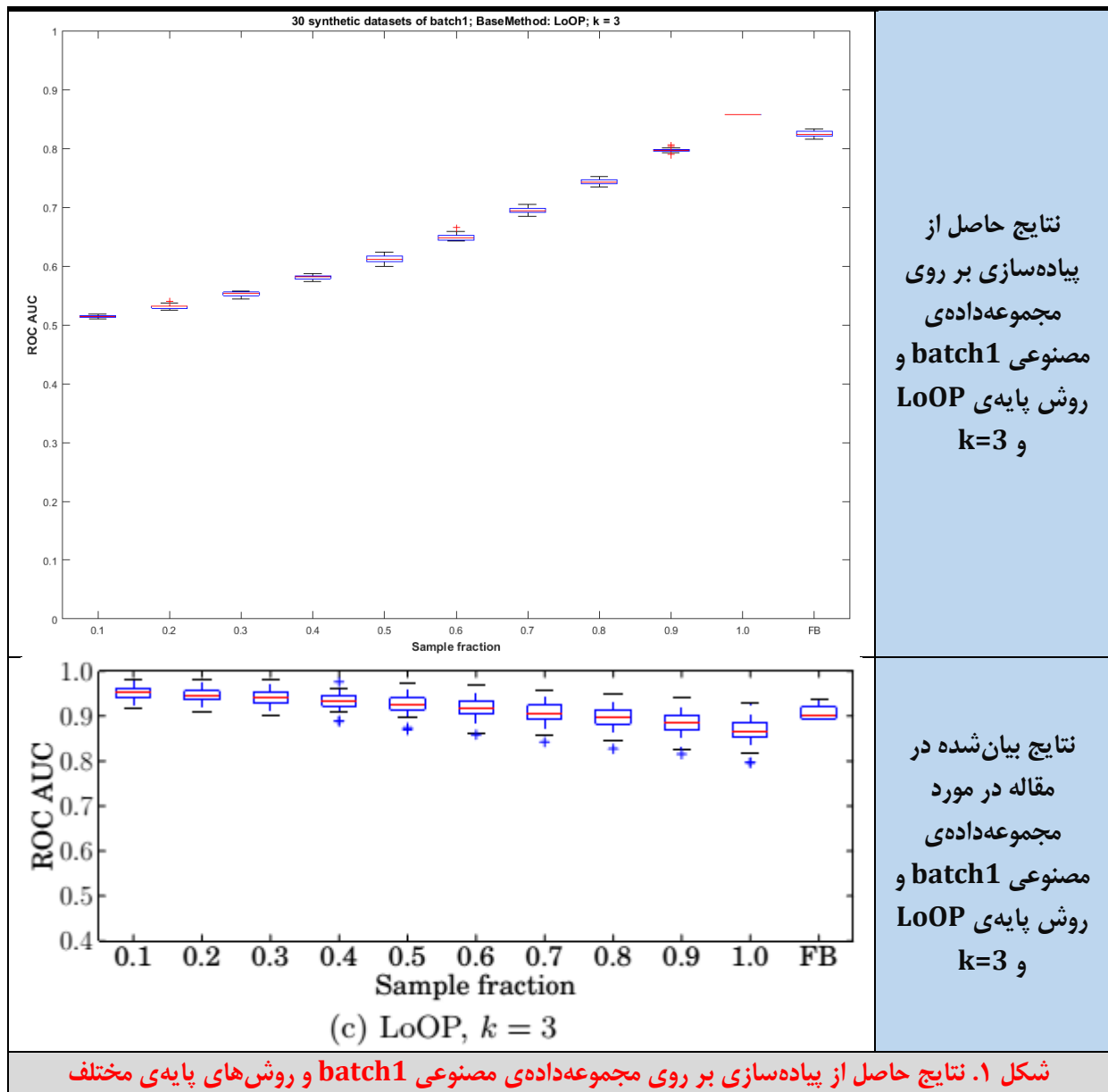
¹⁴ Rotation

¹⁵ Mahalanobis

¹⁶ Local Outlier Factor (LOF)

¹⁷ Local Outlier Probability (LoOP)



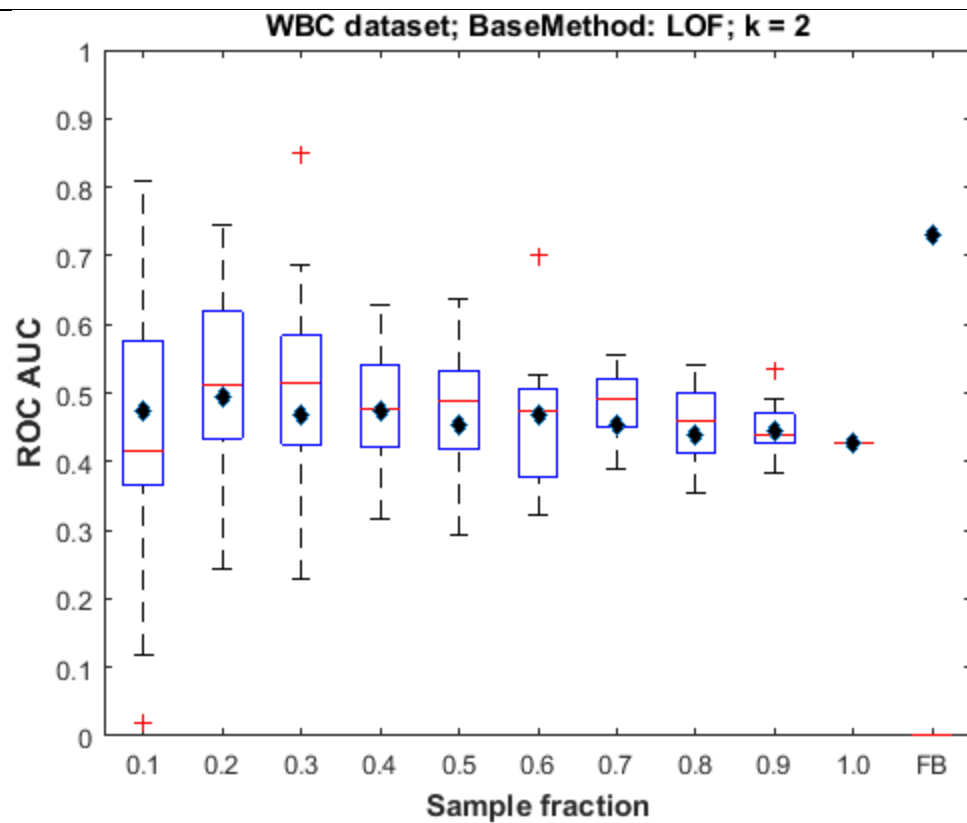


شکل ۱. نتایج حاصل از پیاده‌سازی بر روی مجموعه‌داده‌ی مصنوعی batch1 و روش‌های پایهی مختلف

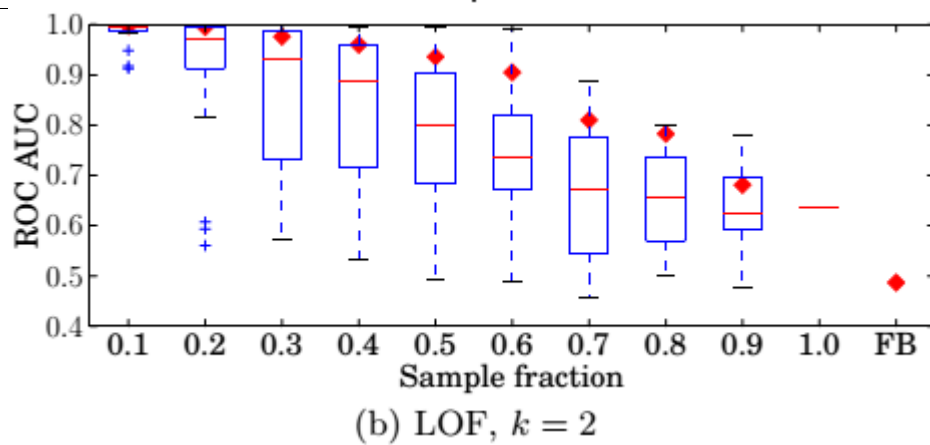
متأسفانه همانطور که قابل مشاهده است، نتایج حاصل از پیاده‌سازی با نتایج قیدشده در مقاله یکسان نمی‌باشد و علت این مسئله را نیز می‌توان در تصادفی‌بودن رویه‌ی تولید داده‌های مصنوعی جستجو نمود. با توجه به نکات قیدشده درمقاله انتظار آن بود که با افزایش اندازه‌ی زیرنمونه، شاهد کاهش محدوده‌ی توزیع مقادیر ROC AUC باشیم که در این‌جا این‌گونه نشد، اما در مورد برخی مجموعه‌داده‌های واقعی که در ادامه خواهد آمد، نتایج خوبی را مطابق مندرجات موجود در مقاله شاهد خواهیم بود.

۲) مجموعه‌داده‌های واقعی

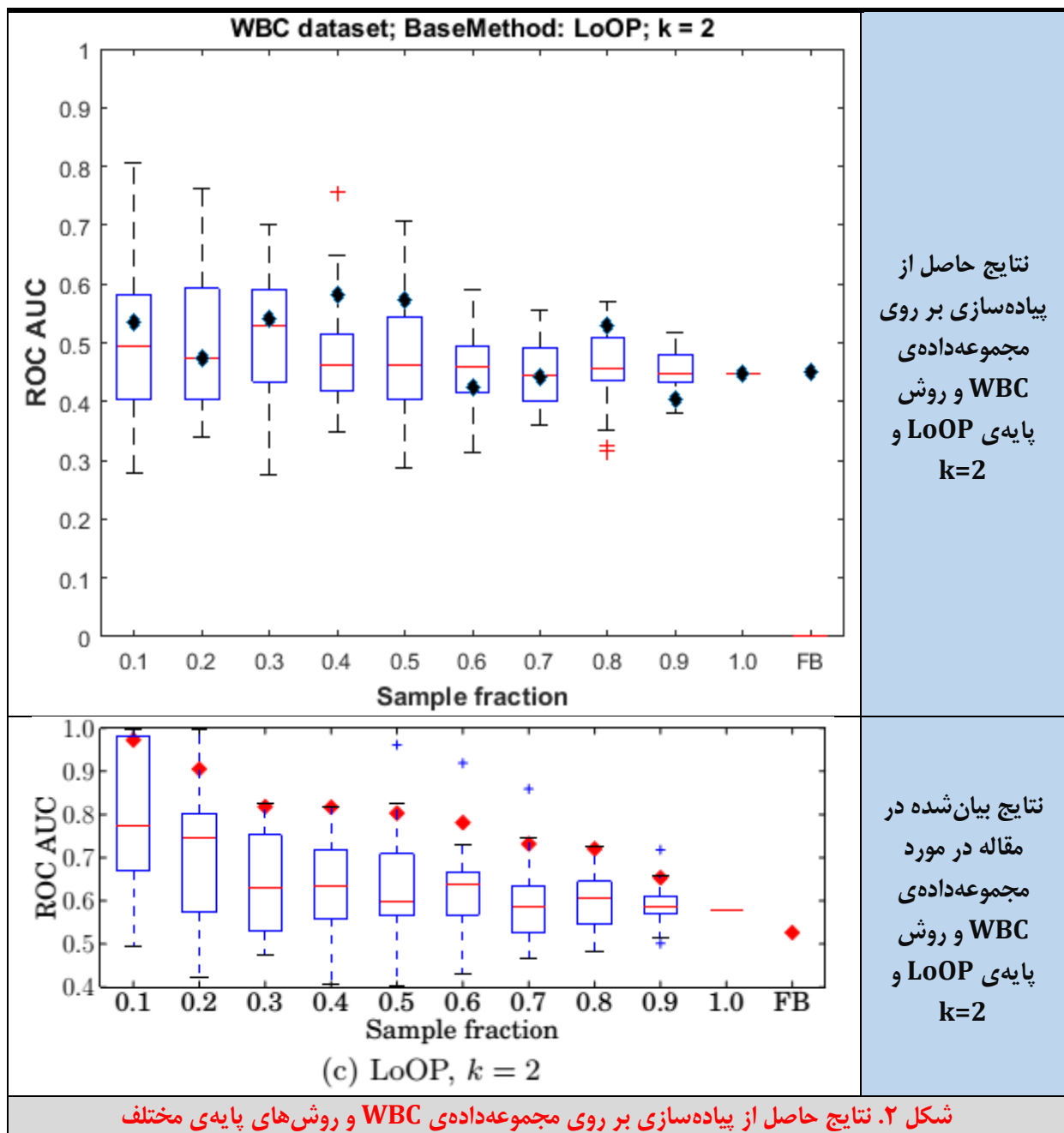
در ابتدا به مجموعه‌داده‌های WBC و Satimage خواهیم پرداخت که تنها در مورد دومی نتایج مطلوبی را شاهد خواهیم بود. نتایج مربوط به مجموعه‌داده‌ی WBC به شرح ذیل می‌باشند:



نتایج حاصل از
پیاده‌سازی بر روی
مجموعه داده‌ی
WBC و روش
پایه‌ی LOF و
 $k=2$

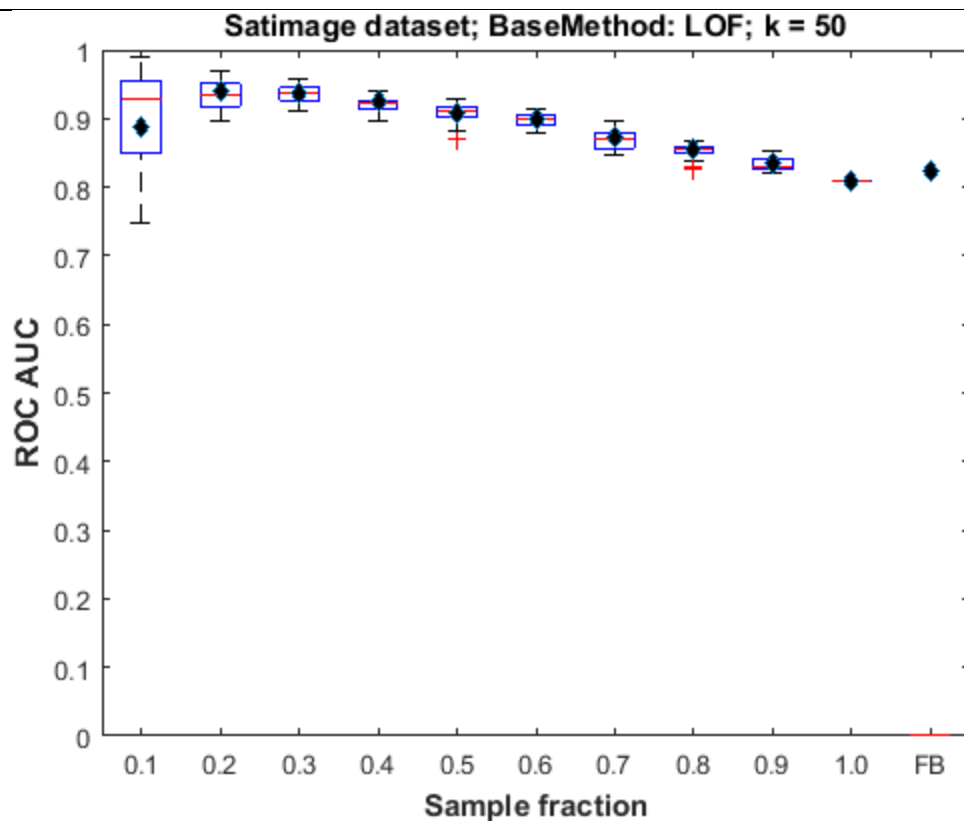


نتایج بیان شده در
مقاله در مورد
مجموعه داده‌ی
WBC و روش
پایه‌ی LOF و
 $k=2$

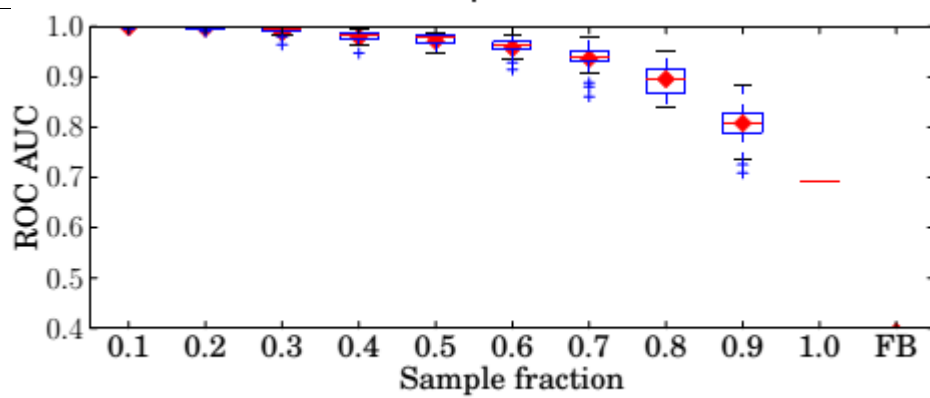


شکل ۲. نتایج حاصل از پیاده‌سازی بر روی مجموعه داده‌ی WBC و روش‌های پایه‌ی مختلف

در ادامه نتایج مربوط به مجموعه داده‌ی Satimage ذکر می‌گردد. لازم به ذکر است که با توجه به بزرگ‌تر بودن مجموعه داده‌ی Satimage، مقدار k را نیز بزرگ‌تر و برابر ۵۰ انتخاب می‌نمائیم. در ادامه داریم:

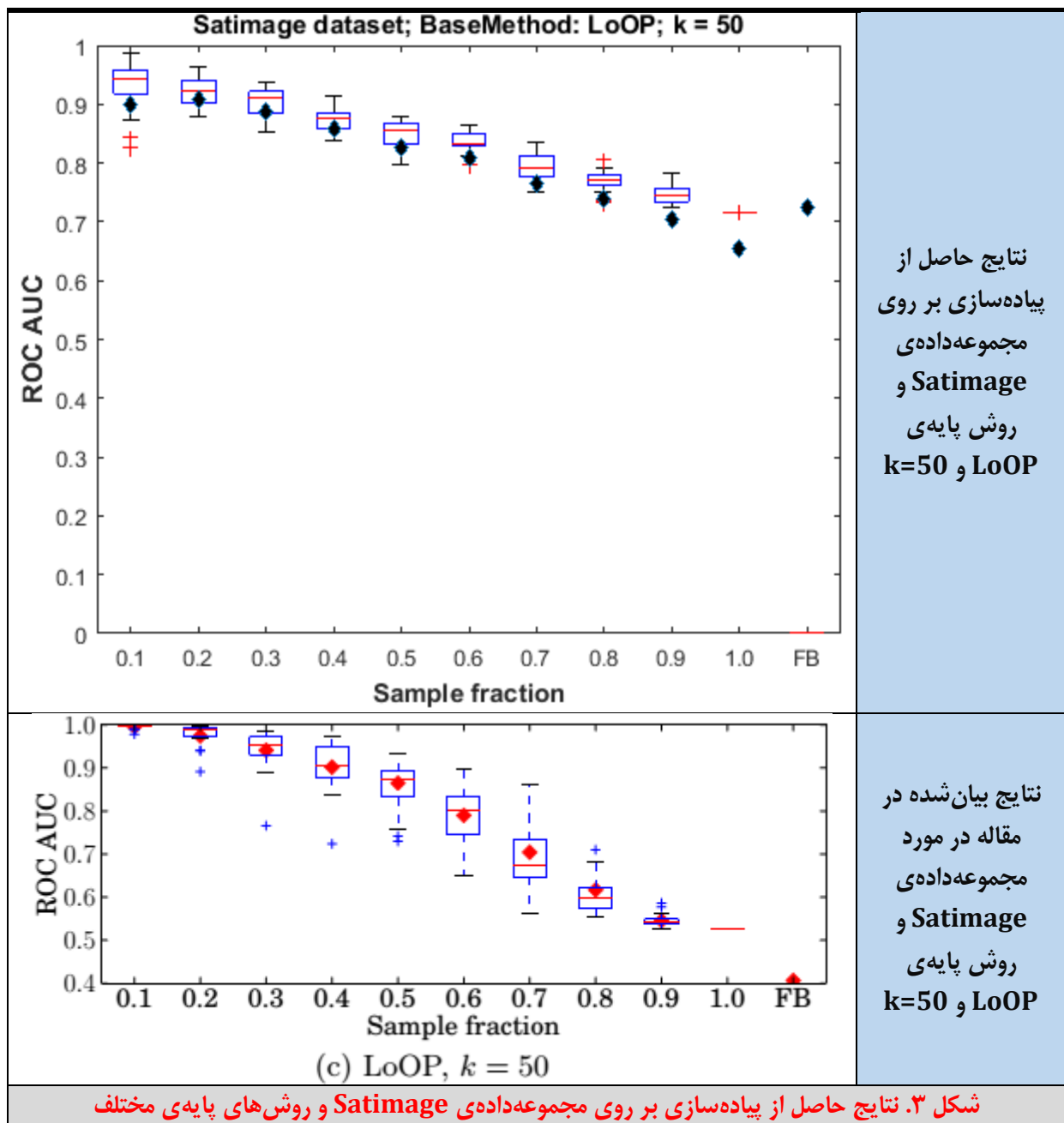


نتایج حاصل از
پیاپی سازی بر روی
مجموعه داده‌ی
Satimage
روش پایه‌ی LOF و
 $k=50$



(b) LOF, $k = 50$

نتایج بیان شده در
مقاله در مورد
مجموعه داده‌ی
Satimage
روش پایه‌ی LOF و
 $k=50$

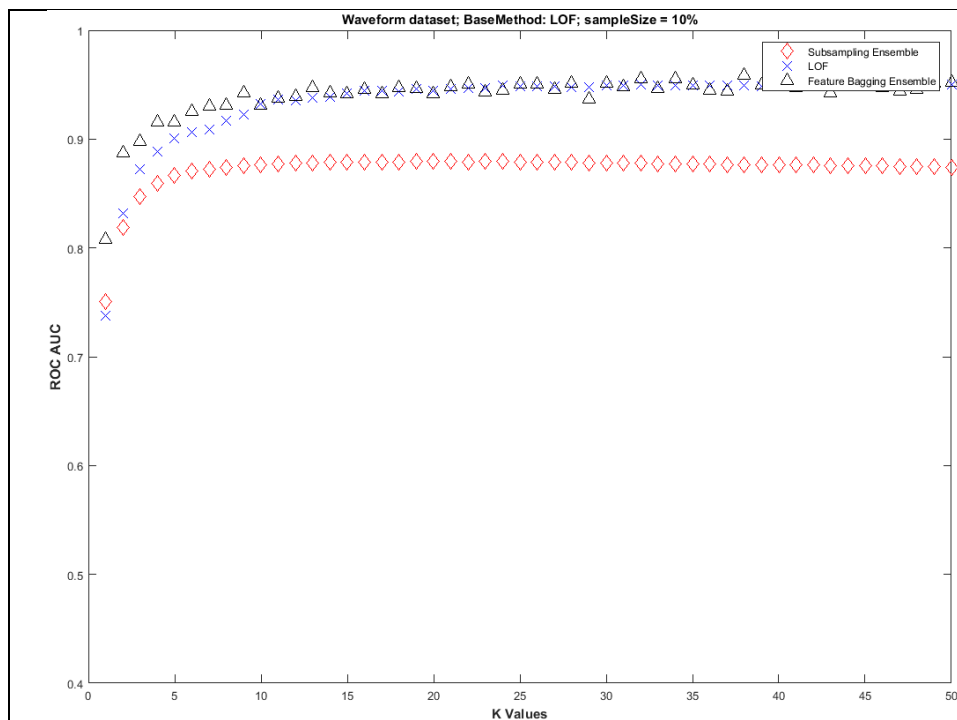


شکل ۳. نتایج حاصل از پیاده‌سازی بر روی مجموعه داده‌ی Satimage و روش‌های پایهی مختلف

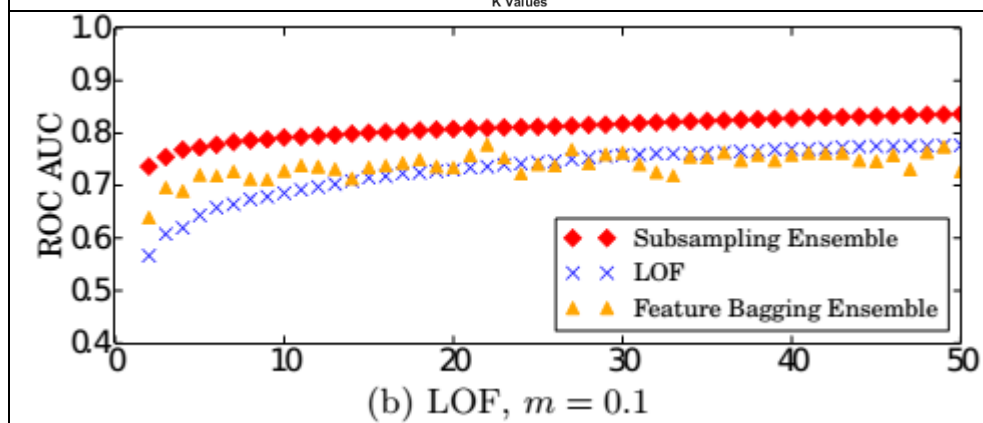
در این قسمت، اندازه‌ی زیرنمونه را ثابت و برابر ۱، ۰ در نظر گرفته و نتایج را بر حسب مقادیر متفاوت k بر روی مجموعه داده‌ی Waveform و برای روش‌های پایه، ترکیب زیرنمونه‌ها^{۱۸} و ترکیب چینش‌های مختلفی از ویژگی‌ها^{۱۹} نمایش می‌دهیم. در ادامه داریم:

¹⁸ Subsampling Ensemble

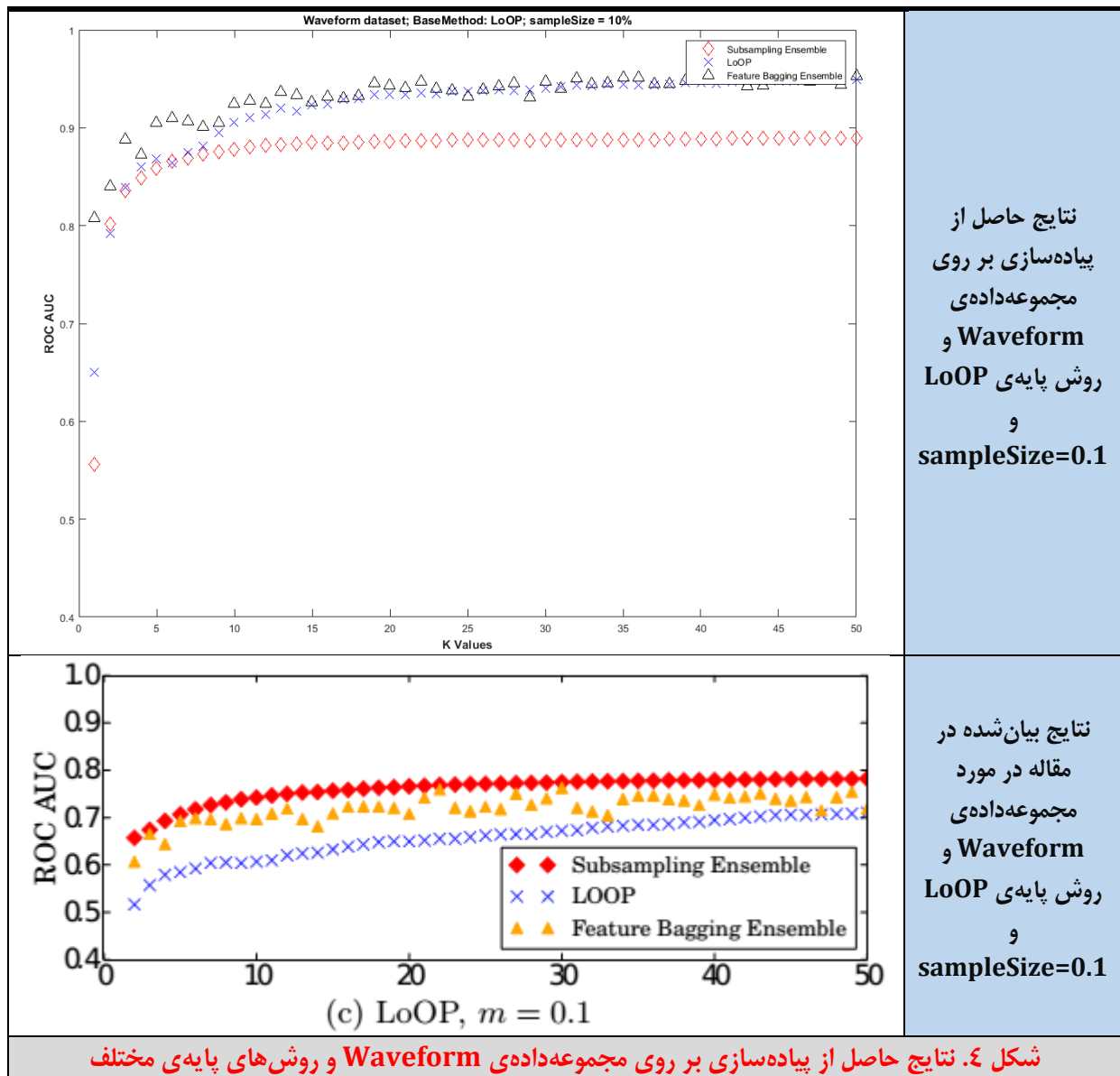
¹⁹ Feature Bagging Ensemble



نتایج حاصل از
پیاده‌سازی بر روی
مجموعه داده‌ی
Waveform
روش پایه‌ی LOF و
sampleSize=0.1



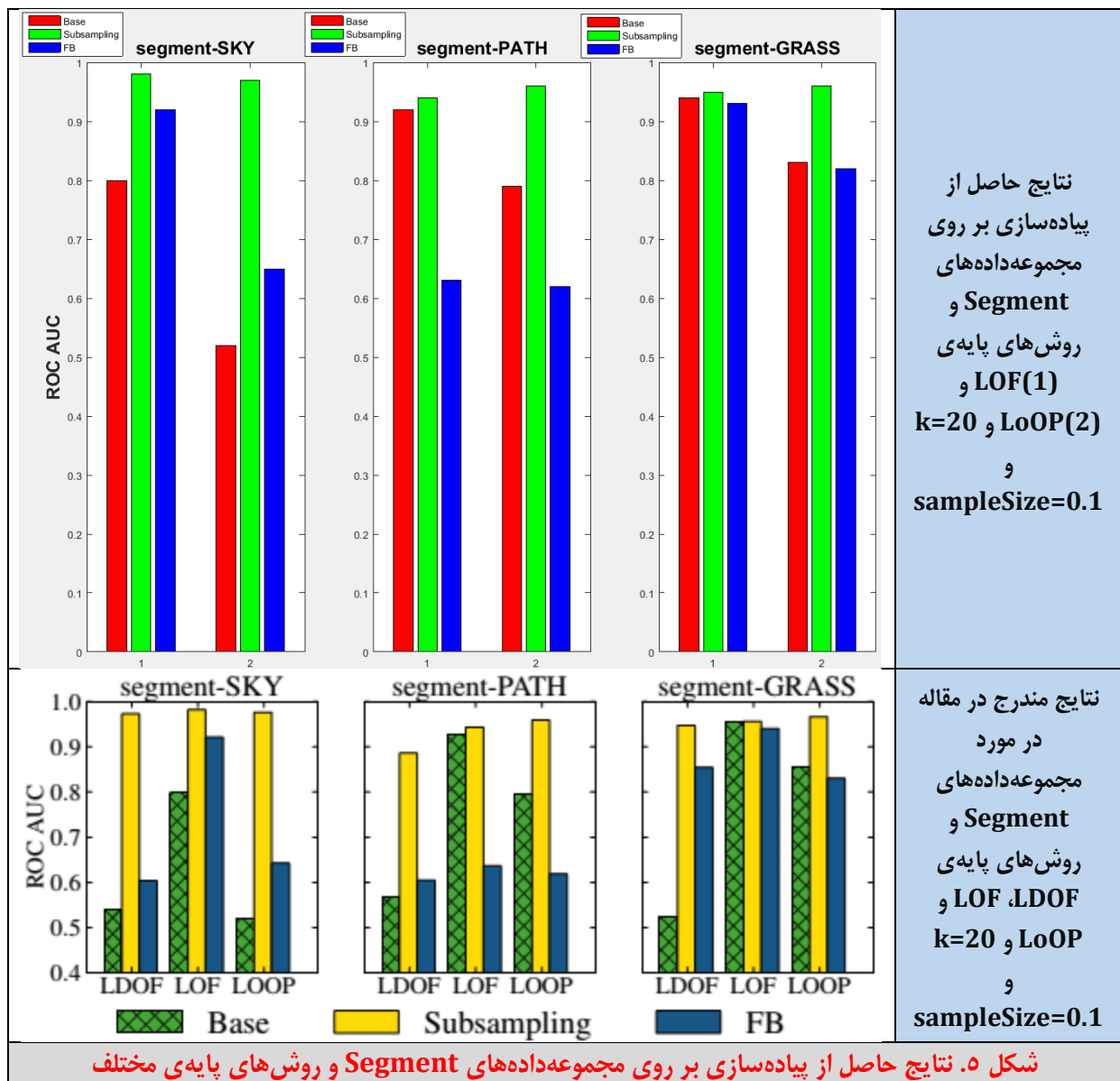
نتایج بیان شده در
مقاله در مورد
مجموعه داده‌ی
Satimage و روش
پایه‌ی LOF و
sampleSize
=0.1



شکل ۴. نتایج حاصل از پیاده‌سازی بر روی مجموعه داده‌ی Waveform و روش‌های پایه‌ی مختلف

علاوه بر اینکه نتایج حاصل از پیاده‌سازی خوشبختانه با نتایج مندرج در مقاله همخوانی بالایی دارند! باید گفت که با افزایش مقدار k شاهد افزایش اندک اما پایدار در مقدار ROC AUC برای روش‌های پایه (LOF و LOOP) و Subsampling Ensemble خواهیم بود و به عبارتی هر دوی این روش‌ها از یک الگوی یکسان و پایدار پیروی می‌نمایند، در حالی که برای روش FeatureBagging شاهد نوسان و واریانس بالاتری در مورد مقادیر ROC AUC هستیم.

در ادامه نتایج حاصله در مورد مجموعه داده‌ی Segment را با مقادیر مشخص برای پارامترها ($k=20$) که در مورد این مجموعه داده و البته روش‌های پایه‌ی مختلف نتایج خوبی را حاصل نموده است؛ و ($subsampleSize=0.1$) به صورتی که در ادامه می‌آید نشان می‌دهیم:



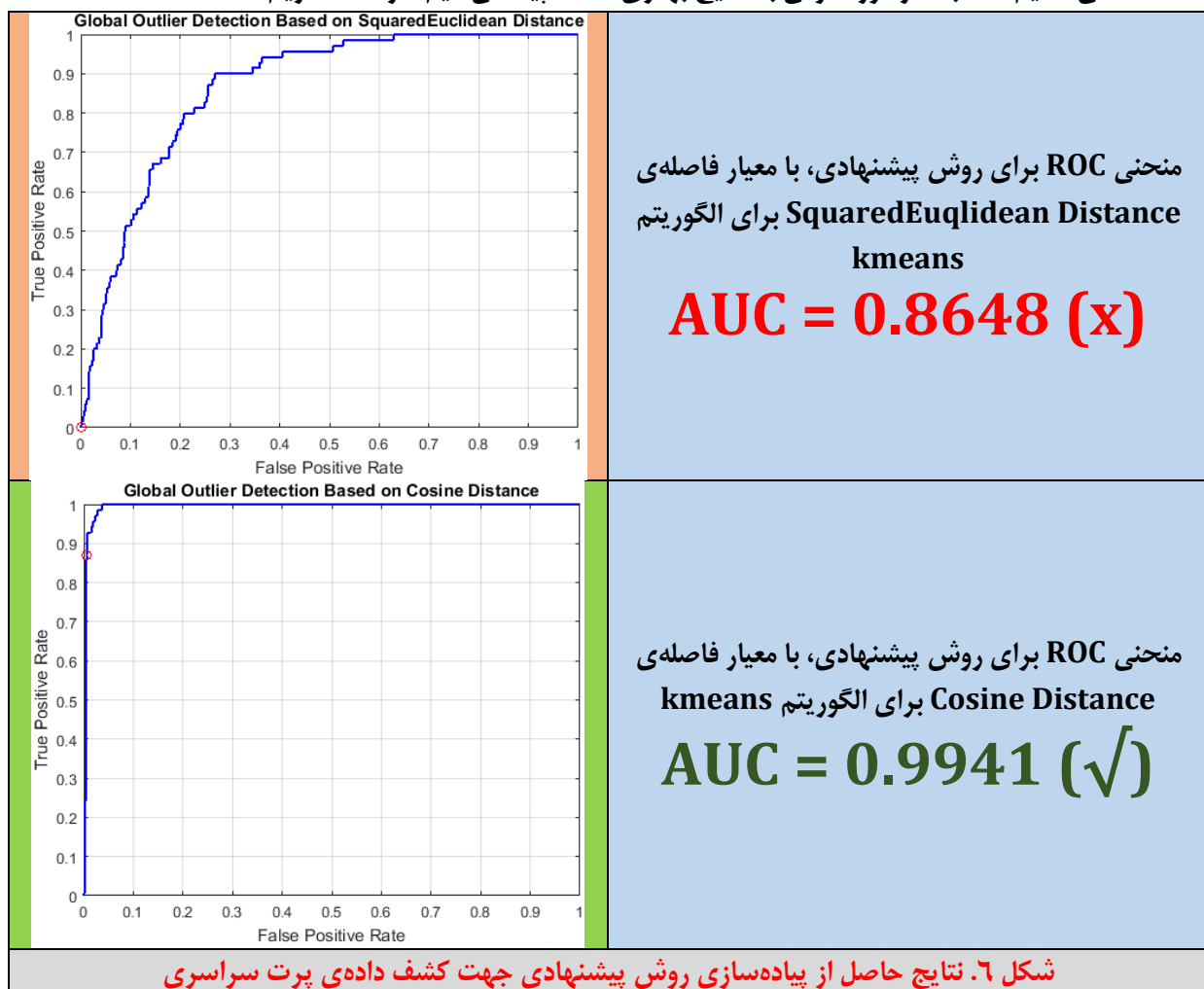
همانطور که قابل مشاهده است، روش زیرنمونه‌برداری معرفی شده در این مقاله به مراتب نتایج بهتری را نسبت به روش‌های پایه و البته روش رقیب FeatureBagging حاصل نموده است، که این خود گواه بر صحت ادعاهای قیدشده در مقاله می‌باشد.

۵) معرفی ایده‌ی پیشنهادی و نتایج پیاده‌سازی

تا این‌جا ما از روش‌های ترکیبی، جهت کشف داده‌های پرت محلی استفاده نمودیم؛ حال قصد داریم تا روشی را جهت کشف داده‌های پرت سراسری معرفی نمائیم. به این ترتیب که ابتدا یک مجموعه داده‌ی مصنوعی را مانند آن‌چه پیش از این قید گردید تهیه می‌نمائیم، تنها با این تفاوت که تعدادی داده‌ی جدید را با مقادیر بسیار دورتر از میانگین هر بعد به مجموعه داده اضافه نموده و آن‌ها را تحت عنوان داده‌ی پرت سراسری برچسب می‌زنیم. حال در عمل پیاده‌سازی ابتدا مجموعه داده را با استفاده از تابع `kmeans()` و مقادیر متفاوت `clustNo` (که همان

تعداد خوشه‌هاست) خوشه‌بندی نموده و سپس برای هر داده یک مقدار احتمالاتی بین ۰ و ۱ را به عنوان ضریب داده‌ی پرت سراسری معرفی می‌نمائیم، که این مقدار از تقسیم فاصله‌ی هر داده از مرکز خوشه‌ای که در آن قرار دارد بر بیشینه‌ی فاصله‌های داده‌های موجود درون خوشه از مرکز به دست می‌آید و البته عددی مابین ۰ و ۱ می‌باشد. در نهایت به ازای خوشه‌بندی‌های مختلف، مقادیر متفاوتی برای این ضرایب خواهیم داشت. سپس جهت به دست آوردن مقدار نهائی برای ضریب مربوطه به ازای هر داده، از روش اول عمق استفاده می‌نمائیم، که در آن ابتدا ماتریس متشکل از ضرایب (سطرها بیانگر داده‌ها و ستون‌ها بیانگر تعداد خوشه‌ها) را به صورت ستونی و هر ستون مجزا از دیگری، در قالب ترتیب نزولی مرتب می‌نمائیم. سپس از سطر اول شروع نموده و سطر به سطر جلو می‌رویم و اولین مقدار ضریب که به ازای هر داده مشاهده می‌کنیم را به عنوان ضریب مربوط به آن داده گزارش می‌نمائیم تا زمانی که به تمامی داده‌ها مقدار ضریب مربوطه نسبت داده شود. به عبارتی در این روش ضریبی برای هر داده از ماتریس مرتب‌شده صورت نزولی انتخاب می‌شود، که در میان خوشه‌بندی‌های مختلف، در بالاترین ردیف ممکن قرار دارد.

در نهایت، بردار ضرایب داده‌ی پرت سراسری را به همراه بردار مربوط به برچسب‌ها به تابع `perfcurve()` می‌دهیم تا منحنی ROC مربوط به روش پیشنهادی را برای ما رسم نموده و مساحت زیر نمودار (AUC) را نیز برای ما فراهم نماید. هر چه این مساحت بیشتر باشد، عملکرد الگوریتم پیشنهادی بهتر بوده است. لازم به ذکر است که در این‌جا ما در عمل خوشه‌بندی با استفاده از `kmeans()`، از دو معیار فاصله‌ی `SquaredEuclidean Distance` و نیز `Cosine Distance` استفاده می‌نمائیم که البته در مورد دومی به نتایج بهتری دست پیدا می‌کنیم. در ادامه داریم:



شکل ۶. نتایج حاصل از پیاده‌سازی روش پیشنهادی جهت کشف داده‌ی پرت سراسری

- [1] **Zimek, Arthur, et al. "Subsampling for efficient and effective unsupervised outlier detection ensembles." *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2013.**
- [2] **Breunig, Markus M., et al. "LOF: identifying density-based local outliers." *ACM sigmod record*. Vol. 29. No. 2. ACM, 2000.**
- [3] **Kriegel, Hans-Peter, et al. "LoOP: local outlier probabilities." *Proceedings of the 18th ACM conference on Information and knowledge management*. ACM, 2009.**