



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد

تشخیص ناهنجاری با استفاده از شبکه های مولد تقابلی

دانشجو

زهرا دهقانیان

اساتید راهنما

دکتر محمد رحمتی

دکتر مریم امیرمزلقانی

بهار ۱۴۰۱







به نام خدا

تاریخ:

تعهدنامه اصالت اثر

## صفحه فرم ارزیابی و تصویب پایان نامه - فرم تأیید اعضاء کمیته دفاع

در این صفحه (هر سه مقطع تحصیلی) باید فرم ارزیابی یا تأیید و تصویب پایان نامه/رساله موسوم به فرم کمیته دفاع برای ارشد و دکترا و فرم تصویب برای کارشناسی، موجود در پرونده آموزشی را قرار دهند.

اینجانب زهرا دهقانین متعهد می‌شوم که مطالب مندرج در این پایان نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آنها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است.

در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است. نقل مطالب با ذکر مآخذ بلامانع است.

در صفحه تعهدنامه اصالت اثر، در قسمت بالا سمت چپ، تاریخ دفاع خود را جایگزین تاریخ نوشته شده کنید.

همچنین در صفحه تعهدنامه اصالت اثر، در خط اول، نام و نام خانوادگی خود را به صورت کامل با نام و نام خانوادگی نمونه، جایگزین کنید. در انتهای متن تعهد، در قسمت امضا نیز باید نام و نام خانوادگی کامل خود را وارد نماید.

## چکیده

یکی از مهم‌ترین فعالیت‌های حوزه تحلیل داده تشخیص ناهنجاری است که در طیف وسیعی از کاربردها همچون تشخیص جعل، کاربردهای پزشکی و سیستم‌های امنیتی به کار گرفته می‌شود. علی‌رغم وجود روش‌های آماری و مبتنی بر یادگیری ماشین، طراحی مدل‌های موثر در تشخیص ناهنجاری در فضای داده پیچیده با ابعاد بالا همچنان به عنوان یک چالش اساسی باقی مانده است. شبکه‌های مولد تقابلی<sup>۱</sup> قادرند تا بر چالش مورد نظر فائق آمده و توزیع داده‌های دنیای واقعی که دارای پیچیدگی و ابعاد بالا هستند را مدل کنند و همین امر سبب می‌شود تا عملکرد امیدوارکننده‌ای در زمینه تشخیص ناهنجاری از خود نشان دهند. در این کار سه شبکه تقابلی CALAD<sup>۲</sup>، RALAD<sup>۳</sup> و RCALAD<sup>۴</sup> با هدف تشخیص ناهنجاری ارائه شده است. اساس کار هر سه مدل پیشنهادی بازسازی داده ورودی با استفاده از شبکه مولد و در ادامه محاسبه میزان اختلاف داده اصلی و بازسازی آن به منظور شناسایی نمونه‌های ناهنجار است. در مدل CALAD با تعریف متغیر جدید  $\hat{Z}_{xx}$  و افزودن تمایزگر ابتکاری  $D_{xxzz}$  چرخه پایداری کامل میان هر دو فضای ورودی و فضای نهان برقرار می‌شود. لازمه شناسایی موثر نمونه‌های ناهنجار بازسازی ضعیف آن‌ها است. هدف از چارچوب پیشنهادی RALAD متمایل کردن تمامی بازسازی‌ها به سمت توزیع داده هنجار است. رویه مورد نظر در این مدل سبب بازسازی ضعیف نمونه‌ها ناهنجار و در نتیجه ایجاد فاصله مناسب میان داده ورودی و بازسازی متناسب با آن می‌شود. در نهایت از ترکیب هر دو ایده مدل جامع RCALAD با هدف حل مسئله تشخیص ناهنجاری در دنیای واقعی ارائه شده است. علاوه بر معماری پیشنهادی، دو امتیاز ناهنجاری جدید نیز در این کار معرفی شده است که در مقایسه با امتیازهای ناهنجاری قبلی قدرت تفکیک‌پذیری بیشتری فراهم می‌آورد. نتایج تجربی بیانگر برتری مدل‌های پیشنهادی در مقایسه با سایر مدل‌های مطرح در زمینه تشخیص ناهنجاری بوده است.

## واژگان کلیدی:

تشخیص ناهنجاری، یادگیری ماشین، شبکه مولد تقابلی، چرخه پایداری کامل، امتیاز ناهنجاری.

<sup>1</sup> Generative Adversarially Networks

<sup>2</sup> Complete Adversarially Learned Anomaly Detection

<sup>3</sup> Regularized Adversarially Learned Anomaly Detection

<sup>4</sup> Regularized Complete Adversarially Learned Anomaly Detection

صفحه	فهرست مطالب
۱	فصل اول: مقدمه.....
۵	۱-۱- ساختار گزارش.....
۶	فصل دوم: مروری بر کارهای پیشین.....
۷	۱-۲- طبقه‌بندی روش‌های تشخیص ناهنجاری از دیدگاه در دسترس بودن برچسب داده.....
۸	۱-۱-۲- تشخیص ناهنجاری با نظارت.....
۸	۲-۱-۲- تشخیص ناهنجاری نیمه‌نظارتی.....
۸	۳-۱-۲- تشخیص ناهنجاری بدون نظارت.....
۹	۲-۲- طبقه‌بندی روش‌های تشخیص ناهنجاری از نظر رویکرد حل مسئله.....
۱۰	۱-۲-۲- روش‌های آماری.....
۱۰	۱-۱-۲-۲- روش‌های پارامتری.....
۱۱	۲-۱-۲-۲- روش‌های غیرپارامتری.....
۱۱	۲-۲-۲- روش‌های یادگیری ماشین.....
۱۲	۱-۲-۲-۲- دسته‌بندی.....
۱۳	۲-۲-۲-۲- نزدیک‌ترین همسایه.....
۱۳	۳-۲-۲-۲- خوشه‌بندی.....
۱۴	۳-۲- دسته‌بندی بر اساس نحوه تشخیص ناهنجاری.....
۱۴	۱-۳-۲- بر اساس فاصله.....
۱۴	۲-۳-۲- دسته‌بندی تک‌کلاسی.....
۱۵	۳-۳-۲- بر اساس بازسازی.....
۱۵	۴-۲- معیارهای ارزیابی روش‌های تشخیص ناهنجاری.....
۱۶	۱-۴-۲- صحت.....
۱۶	۲-۴-۲- بازیابی.....
۱۷	۳-۴-۲- F1-score.....
۱۷	۴-۴-۲- مساحت زیر نمودار مشخصه عملکرد.....
۱۷	۵-۲- شبکه‌های مولد تقابلی و تشخیص ناهنجاری.....
۱۸	۱-۵-۲- شبکه‌های مولد تقابلی.....
۲۲	۱-۱-۵-۲- تحلیل نظری شبکه مولد تقابلی.....
۲۳	۲-۱-۵-۲- مزایا و معایب.....
۲۴	۲-۵-۲- مدل ANOGAN.....
۲۵	۱-۲-۵-۲- یادگیری بدون نظارت متنوع داده‌های طبیعی.....
۲۶	۲-۲-۵-۲- نگاشت تصاویر جدید به فضای نهفته.....
۲۷	۳-۲-۵-۲- تشخیص ناهنجاری.....
۲۸	۴-۲-۵-۲- مزایا و معایب.....

۲۸	..... مدل f-AnoGan
۲۹	..... یادگیری بدون نظارت تصاویر طبیعی
۳۰	..... یادگیری نگاشت سریع از فضای تصویر به فضای نهفته
۳۳	..... شناسایی ناهنجاری
۳۴	..... مزایا و معایب
۳۴	..... مدل ALI
۳۸	..... مقایسه مدل‌های GAN و ALI
۳۹	..... رویکردهای جایگزین برای استنتاج در GAN
۳۹	..... مزایا و معایب
۴۰	..... مدل EGBAD
۴۱	..... مزایا و معایب
۴۲	..... مدل ALICE
۴۲	..... یادگیری تقابلی با اندازه‌گیری اطلاعات
۴۳	..... آنتروپی شرطی
۴۴	..... فرایند یادگیری
۴۴	..... مزایا و معایب
۴۵	..... مدل RCGAN
۴۵	..... منظم‌سازی شبکه مولد و تمایزگر
۴۶	..... پایداری چرخه
۴۷	..... مدل ALAD
۴۸	..... تابع هزینه
۵۰	..... تشخیص ناهنجاری
۵۲	..... جمع‌بندی
۵۳	..... فصل سوم: روش پیشنهادی
۵۷	..... مدل CALAD
۵۸	..... معماری شبکه
۶۱	..... تابع هدف
۶۴	..... مدل RALAD
۶۴	..... معماری شبکه
۶۶	..... تابع هدف
۶۹	..... مدل RCALAD
۷۰	..... معماری شبکه
۷۱	..... تابع هدف
۷۲	..... تشخیص ناهنجاری
۷۴	..... جمع‌بندی



۷۶	فصل چهارم: آزمایش‌ها و نتایج.....
۷۷	۱-۴- دادگان و پیش‌پردازش.....
۷۷	۱-۴-۱- مجموعه داده KDDCup99.....
۷۸	۱-۴-۲- مجموعه داده Arrhythmia.....
۷۸	۱-۴-۳- مجموعه داده Thyroid.....
۷۸	۱-۴-۴- مجموعه داده Musk.....
۷۹	۱-۴-۵- مجموعه داده CIFAR-10.....
۷۹	۱-۴-۶- مجموعه داده SVHN.....
۷۹	۲-۴- تنظیمات مدل.....
۸۰	۳-۴- مدل‌های پایه.....
۸۰	۳-۴-۱- روش OC-SVM.....
۸۰	۳-۴-۲- روش IF.....
۸۱	۳-۴-۳- روش DSEBM.....
۸۱	۳-۴-۴- روش DAGMM.....
۸۱	۳-۴-۵- روش DCAE.....
۸۲	۳-۴-۶- روش DSVDD.....
۸۲	۴-۴- نتایج.....
۸۲	۴-۴-۱- دادگان جدولی.....
۸۳	۴-۴-۱- دادگان تصویری.....
۸۶	۵-۴- بحث.....
۸۶	۵-۴-۱- مطالعه فرسایشی.....
۸۸	۵-۴-۲- انتخاب تابع توزیع جریمه.....
۸۹	۵-۴-۳- ارزیابی کارایی امتیازهای ناهنجاری.....
۹۱	۵-۴-۴- ارزیابی کفایت تمایزگر <i>Dxxzz</i> .....
۹۳	۶-۴- جمع‌بندی.....
۹۴	فصل پنجم: جمع‌بندی، نتیجه‌گیری و کارهای آتی.....
۹۵	۵-۱- جمع‌بندی و نتیجه‌گیری.....
۱۰۲	۵-۲- کارهای آتی.....
۱۰۴	منابع و مراجع.....
۱۰۹	فهرست واژگان انگلیسی به فارسی.....

## صفحه

## فهرست شکل‌ها

شکل ۱-۲: دسته بندی روش‌های تشخیص ناهنجاری.....	۹
شکل ۲-۲: رویکرد کلی شبکه‌های مولد تقابلی.....	۲۰
شکل ۳-۲: شمای کلی روند آموزش کدگذار [31].....	۳۱
شکل ۴-۲: معماری شبکه ALI.....	۳۵
شکل ۵-۲: شمای کلی شبکه ALAD.....	۴۹
شکل ۶-۲: نمونه‌ای از خروجی شبکه ALAD به همراه داده‌های ناهنجار.....	۵۰
شکل ۱-۳: بازسازی نامطلوب نمونه ناهنجار.....	۵۵
شکل ۲-۳: نمایش جریان اطلاعات در مدل CALAD.....	۶۰
شکل ۳-۳: معماری CALAD.....	۶۱
شکل ۴-۳: معماری RALAD.....	۶۵
شکل ۵-۳: تاثیر حضور توزیع $\sigma(x)$ در روند آموزش مدل.....	۶۷
شکل ۶-۳: معماری RCALAD.....	۷۰
شکل ۱-۴: عملکرد مدل RCALAD روی کلاس عدد سه.....	۸۵
شکل ۲-۴: تاثیر توزیع $\sigma(x)$ بر بازسازی نمونه‌های ناهنجار.....	۸۸
شکل ۱-۵: معماری اولیه شبکه مولد تقابلی.....	۹۵
شکل ۲-۵: معماری مدل ALI.....	۹۶
شکل ۳-۵: معماری شبکه ALICE.....	۹۷
شکل ۴-۵: معماری شبکه ALAD.....	۹۸
شکل ۵-۵: معماری شبکه CALAD.....	۹۹
شکل ۶-۵: معماری شبکه RALAD.....	۹۹
شکل ۷-۵: معماری شبکه RCALAD.....	۱۰۰

## صفحه

## فهرست جدول‌ها

جدول ۴-۱: نتایج خروجی مدل پیشنهادی در مقایسه با مدل‌های پایه بر روی مجموعه داده‌های جدولی.....	۸۳
جدول ۴-۲: نتایج خروجی مدل پیشنهادی در مقایسه با مدل‌های پایه بر روی مجموعه داده CIFAR10.....	۸۴
جدول ۴-۳: نتایج خروجی مدل پیشنهادی در مقایسه با مدل‌های پایه بر روی مجموعه داده SVHN.....	۸۵
جدول ۴-۴: تاثیر بخش‌های مختلف پیشنهادی در بهبود نتایج دادگان جدولی.....	۸۶
جدول ۴-۵: تاثیر بخش‌های مختلف پیشنهادی در بهبود نتایج دادگان تصویری.....	۸۷
جدول ۴-۶: تاثیر $\sigma(x)$ های مختلف بر عملکرد مدل RCALAD.....	۸۸
جدول ۴-۷: مقایسه عملکرد امتیازهای ناهنجاری پیشنهادی با سایر امتیازها روی دادگان جدولی.....	۸۹
جدول ۴-۸: مقایسه عملکرد امتیازهای ناهنجاری پیشنهادی با سایر امتیازها روی دادگان تصویری.....	۹۱
جدول ۴-۹: ارزیابی عملکرد مدل در حضور/عدم حضور هر یک از اجزا.....	۹۲
جدول ۵-۱: روند تکامل توابع بهینه‌سازی شبکه‌های مولد تقابلی.....	۹۲

صفحه

فهرست الگوریتم‌ها

الگوریتم ۱-۲: آموزش گرادیان نزولی کوچک دسته‌ای شبکه‌های مولد تقابلی.....	۲۱
الگوریتم ۲-۲: رویه آموزش یادگیری خصمانه استنتاج.....	۳۸
الگوریتم ۱-۳: شبه کد الگوریتم ALAD.....	۵۱
الگوریتم ۲-۳: روند محاسبه‌ی امتیاز ناهنجاری.....	۷۴

## فهرست علائم

داده ورودی	$\pi$	آنتروپی شرطی	$x$
شبکه مولد	CCC	چرخه پایداری کامل	$G$
شبکه تمایزگر	$\hat{x}$	بازسازی $x$	$D$
توزیع شبکه مولد	$\hat{z}_{\hat{x}}$	نگاشت معکوس $\hat{x}$	$p_g$
توزیع داده ورودی	$N$	توزیع نرمال	$p_{data}$
نگاشت معکوس $x$	$\hat{z}$	نگاشت معکوس $x_z$	$z_x$
کدگذار	$\sigma(x)$	توزیع کمکی	$E$
توزیع کدگذار	$A$	امتیاز ناهنجاری	$q$
توزیع یکنواخت	$z$	نمونه از توزیع گاوسی	$U$
خروجی مولد با			$x_z$
ورودی توزیع گاوسی			

## فصل اول:

### مقدمه

هنگام تجزیه و تحلیل دادگان موجود در دنیای واقعی، شناسایی نمونه‌های غیرمشابه با سایر نمونه‌ها امری ضروری به نظر می‌رسد. چنین نمونه‌هایی با عنوان ناهنجاری شناخته می‌شوند و از عملیات شناسایی چنین نمونه‌هایی با عنوان مسئله تشخیص ناهنجاری یاد می‌شود. این مسئله یک بخش حائز اهمیت از زمینه تحقیقاتی داده‌کاوی است چرا که شامل کشف الگوهای جذاب و نادر در داده‌هاست [1]. این مسئله به طور گسترده در آمار و یادگیری ماشین مورد مطالعه قرار گرفته است و با مترادف‌هایی مانند تشخیص داده پرت<sup>۱</sup>، شناسایی نوآوری<sup>۲</sup>، تشخیص انحراف<sup>۳</sup> و استخراج استثنا<sup>۴</sup> نیز یاد می‌شود. تعریف رسمی و مورد قبول این مسئله به صورت زیر است:

"یک ناهنجاری مشاهده‌ای است که به میزانی از سایر مشاهدات منحرف می‌شود که ظن‌هایی را برای این که توسط مکانیسم متفاوتی تولید شده باشد، ایجاد می‌کند [2]."

ناهنجاری‌ها جزو پارامترهای مهم هر مجموعه داده‌ای در نظر گرفته می‌شوند و در دامنه وسیعی از کاربردها تاثیرگذار هستند. به عنوان مثال، الگوی غیر معمول ترافیک در یک شبکه کامپیوتری میتواند به معنای هک شدن رایانه و انتقال داده‌ها به مقصدهای غیرمجاز باشد. رفتار غیر عادی در معاملاتی که توسط کارت‌های اعتباری انجام می‌شوند می‌تواند نشانگر فعالیت‌های اقتصادی با هدف کلاهبرداری باشد [3]، و یا یک ناهنجاری در تصویر MRI ممکن است وجود تومور بدخیم را نشان دهد [4]. تشخیص ناهنجاری به طور گسترده در زمینه‌های کاربردی گوناگونی مانند: پزشکی، بهداشت عمومی، تشخیص کلاهبرداری، سنجش از دور<sup>۵</sup>، تشخیص نفوذ، پردازش تصویر، آسیب‌های صنعتی، شبکه‌های حسگر [5]، رفتار روبات‌ها و داده‌های نجومی به کار گرفته شده است [3].

داده‌های تمامی این مسائل یا از نوع سری زمانی هستند و یا فارغ از زمان می‌باشند. رویکرد مورد استفاده برای حل مسائل مربوط به این دو جنس داده کاملاً متفاوت از یکدیگر هستند. تحقیقات بسیار گسترده‌ای

---

<sup>1</sup> Outlier detection

<sup>2</sup> Novelty detection

<sup>3</sup> Deviation detection

<sup>4</sup> Exception mining

<sup>5</sup> Remote sensing

در حوزه داده‌های سری زمانی صورت گرفته است. به عنوان مثال در مدل رگرسیون<sup>۶</sup> نمونه‌ای که به مقدار زیادی از مدل تعیین شده منحرف شود به عنوان داده ناهنجار شناخته می‌شود [۶]. در روش دیگر با استفاده از مدل ARIMA<sup>۷</sup> مقدار آینده را پیش‌بینی می‌کنند و با محاسبه میزان اختلاف داده پیش‌بینی شده و مقدار واقعی داده به شناسایی نمونه ناهنجار می‌پردازند [۷]. علاوه بر این با ظهور CNN<sup>۸</sup> و RNN<sup>۹</sup> در زمینه پیش‌بینی سری زمانی و ثابت شدن کارایی آن‌ها، توجه‌ها به سمت استفاده از این ساختارها در زمینه تشخیص ناهنجاری جلب شد، نحوه استفاده از این نوع از شبکه‌های عصبی همانند مدل رگرسیون و ARIMA است [۷]. در روش‌های آماری موجود برای مسئله تشخیص ناهنجاری فاصله داده تا توزیع یا مدل طراحی شده محاسبه می‌شود و در صورتی که از یک حد آستانه بیشتر باشد به عنوان ناهنجاری شناخته می‌شود. در دسته دیگر از روش‌های آماری هر نمونه به عنوان یک نقطه در فضای  $n$  بعدی در نظر گرفته می‌شود و حول تمامی نمونه‌ها یک فضای محدب محاسبه می‌کند و فرض می‌کند نمونه‌های ناهنجار در لبه این فضا قرار می‌گیرد [۸]. در روش‌های مبتنی بر فاصله با استفاده از یک معیار فاصله میزان تفاوت نمونه تا سایر داده‌ها سنجیده می‌شود و در صورتی که از یک حد آستانه بیشتر باشد به عنوان ناهنجاری شناخته می‌شود. روش دیگر محاسبه چگالی یک نمونه و همسایه‌های آن است تا بدین وسیله به معیار جدید برای شناسایی نمونه‌های ناهنجار ایجاد شود. این معیار LOF<sup>۱۰</sup> نام دارد و هرچه این معیار بالاتر باشد، احتمال ناهنجار بودن آن نمونه بیشتر است [۹]. به طور خاص داده‌های مورد بررسی در این تحقیق همگی از جنس فارغ از زمان هستند که در فصل دوم به بررسی دسته‌بندی‌های موجود برای حل آن‌ها خواهیم پرداخت.

با هدف تشخیص نمونه‌های ناهنجار موجود در دادگان دنیای واقعی، تاکنون روش‌های متنوعی مورد استفاده قرار گرفته‌اند. به طور کلی تشخیص ناهنجاری بر دو اصل استوار است، شناسایی و مدل کردن رفتار داده هنجار و میزان انحراف از رفتار داده هنجار [10]. به طور کلی می‌توان روش‌های تشخیص

---

<sup>6</sup> Regression

<sup>7</sup> Autoregressive integrated moving average

<sup>8</sup> Convolutional Neural Networks

<sup>9</sup> Recurrent Neural Networks

<sup>10</sup> Local Outlier Factor



ناهنجاری را به دو دسته روش‌های آماری و روش‌های مبتنی بر یادگیری ماشین اشاره کرد. روش‌های آماری اگرچه در برخی از موارد کارایی مناسبی دارند اما عملکرد صحیح و مناسب آن‌ها در گرو صحت پیش‌فرض‌های استفاده شده در همین روش‌هاست و در صورتی که پیش‌فرض‌های اولیه در مورد توزیع داده اشتباه باشد نتایج نهایی ناامید کننده خواهد بود [11]. مزیت اصلی روش‌های مبتنی بر یادگیری ماشین استفاده از تجربه‌های گذشته به منظور انجام پیش‌بینی‌های صحیح در آینده است. این روش‌ها تنها با مشاهده نمونه‌های گذشته به طراحی مدل می‌پردازند و پیش‌فرض خاصی نسبت به توزیع داده ندارند. این دسته از الگوریتم‌ها نیازمند تعداد مناسبی از داده‌ها هستند تا بتوانند مدل پیشنهادی خود را آموزش دهند. از جمله مهم‌ترین مدل‌های موجود در این دسته می‌توان به شبکه‌های عصبی اشاره کرد که سابقه طولانی در زمینه تشخیص ناهنجاری دارد. به عنوان مثال شبکه‌های عصبی کدگذار<sup>۱۱</sup> و خودکدگذار<sup>۱۲</sup> مدلی برای بازسازی داده‌های عادی آموزشی آموزش داده می‌شود و نمونه‌های با خطای بازسازی بالا به عنوان نمونه ناهنجار در نظر گرفته می‌شوند [12].

در سال ۲۰۱۷ از شبکه‌های عصبی مولد تقابلی برای تشخیص ناهنجاری در زمینه تصاویر پزشکی (تصاویر شبکه) و در مقایسه با سایر روش‌ها به موفقیت قابل توجهی دست یافت [13]. نتایج درخشان شبکه‌های مولد تقابلی در عرصه پردازش تصویر و استخراج ویژگی سبب محبوبیت آن در زمینه‌های کاربردی مختلف شده است. به طور خاص این دسته از شبکه‌ها به عنوان یک چهارچوب قدرتمند برای مدل‌سازی مجموعه داده‌های پیچیده با ابعاد بالا شناخته می‌شوند. یک از اصلی‌ترین چالش‌های موجود در استفاده از شبکه‌های مولد تقابلی مقابله با پیچیدگی‌های استنتاج است [14]. در سال‌های اخیر تلاش‌های گسترده‌ای انجام شده است تا با استفاده از شبکه‌های عصبی خود کدگذار در ساختار شبکه‌های عصبی مولد تقابلی از پیچیدگی‌های استنتاج کاسته شده و بر چالش‌های موجود غلبه کنند. علیرغم این تلاش‌ها همچنان ضعف‌هایی در روند یادگیری بلوک‌های موجود در ساختار شبکه‌های مولد تقابلی موجود است و از تمامی ظرفیت موجود به منظور دریافت اطلاعات و آموزش هر چه بهتر مدل استفاده نمی‌شود.

---

<sup>11</sup> Encoder

<sup>12</sup> Autoencoder

## ۱-۱- ساختار گزارش

در فصل بعدی ابتدا به دسته‌بندی روش‌های تشخیص ناهنجاری از دیدگاه‌های مختلف می‌پردازیم، معیارهای ارزیابی مدل‌های تشخیص ناهنجاری را معرفی می‌کنیم و در قسمت انتهایی فصل روش‌های تشخیص ناهنجاری مبتنی بر شبکه‌های مولد تقابلی را بررسی و مرور می‌کنیم. در این قسمت سعی می‌شود تا ضمن دسته‌بندی روش‌های تشخیص ناهنجاری، مروری گذرا بر روش‌های به نسبت قدیمی‌تر نیز انجام شود. در ادامه به طور دقیق‌تر زنجیره‌ای از کارها مورد بحث قرار خواهد گرفت که در طول این زنجیره نقاط ضعف و کمبودهای مدل پیشنهادی بر طرف می‌شود. فصل سوم به معرفی مدل پیشنهادی و روش نوین تشخیص ناهنجاری اختصاص خواهد داشت. در این فصل الگوریتم پیشنهادی، که بر اساس حل یک مسئله بهینه‌سازی با در نظر گرفتن توزیع توام<sup>۱۳</sup> پارامترهای موجود در ساختار شبکه عصبی مولد تقابلی طراحی شده است، به شناسایی نمونه‌های ناهنجار در فضای نهفته<sup>۱۴</sup> و فضای ورودی می‌پردازد. در فصل چهارم عملکرد مدل روی دادگان‌های مختلف آزمایش می‌شود. در این فصل با دو نوع مختلف از دادگان روبرو خواهیم بود، دادگانی جدولی<sup>۱۵</sup> که شامل دادگان KDD و ARRHYTHMIA و دادگان تصاویر که شامل CIFAR-10 و SVHN<sup>۱۶</sup> است. در فصل پنجم و نهایی این گزارش پیشنهادهای موثر به منظور اصلاح و بهبود احتمالی روند آموزش شبکه مولد تقابلی ارائه می‌شود و در گام آخر جمع‌بندی مطالب ارائه شده و نتیجه‌گیری نهایی صورت خواهد پذیرفت.

---

<sup>13</sup> Joint distribution

<sup>14</sup> Latent space

<sup>15</sup> Tabular

<sup>16</sup> Street view house number

فصل دوم:

مروری بر کارهای پیشین

در این فصل ابتدا به دسته‌بندی روش‌های مختلف تشخیص ناهنجاری و مرور روش‌های شاخص هر دسته پرداخته خواهد شد. در گام بعدی معیارهای ارزیابی مدل‌های تشخیص ناهنجاری معرفی می‌شوند. در ادامه بر روی کارهایی که تاکنون در زمینه تشخیص ناهنجاری با استفاده از شبکه‌های مولد تقابلی انجام گرفته‌اند مروری خواهیم داشت. در گام اول این قسمت، مقاله پایه با عنوان شبکه‌های مولد تقابلی مورد بررسی قرار خواهد گرفت در ادامه، کار تشخیص ناهنجاری بدون نظارت با شبکه‌های عصبی تقابلی به منظور راهنمایی عملیات اکتشاف نشانگر<sup>۱</sup> به اختصار AnoGan شرح داده خواهد شد و پس از آن f-AnoGan<sup>۲</sup> که در ادامه کار قبلی و توسط همان نویسندگان انجام شده است بررسی خواهد شد. با توجه به ضعف‌های موجود در ساختار f-AnoGan مقاله مکمل این مدل با نام استنتاج یادگرفته شده به روش تقابلی<sup>۳</sup> به اختصار ALI مرور خواهد شد. در گام بعدی مدل EGBAD<sup>۴</sup> که جزو اولین کارها در زمینه تشخیص ناهنجاری که با الهام از مدل ALI خلق شده است مرور می‌شود. در مرحله بعدی با توجه تضمین نشدن شرط سازگاری حلقه<sup>۵</sup> در ALI مقاله آلیس<sup>۶</sup> مورد اشاره قرار می‌گیرد و در انتها تشخیص ناهنجاری یادگرفته شده به روش تقابلی<sup>۷</sup> که در ادامه کارهای پیشین و همچنین مقاله پایه در این پروژه است به طور دقیق بررسی خواهد شد.

## ۲-۱- طبقه‌بندی روش‌های تشخیص ناهنجاری از دیدگاه در دسترس

### بودن برچسب داده

بیشتر روش‌های تشخیص ناهنجاری در مرحله آموزش خود نیاز دارند تا برچسب نمونه‌های مختلف در دسترس آن‌ها باشد تا بتوانند در خصوص طبیعی یا ناهنجار بودن یک نمونه در مرحله آزمایش

<sup>۱</sup> Unsupervised anomaly detection with generative adversarial networks to guid marker

<sup>۲</sup> Fast unsupervised anomaly detection with generative adversarial networks

<sup>۳</sup> Adversarially Learned Inference

<sup>۴</sup> Efficient GAN-Based Anomaly Detection

<sup>۵</sup> Cycle consistency

<sup>۶</sup> ALICE

<sup>۷</sup> Adversarially Learned Anomaly Detection

تصمیم‌گیری کنند[15]. فرایند تهیه و دستیابی به داده‌های دارای برچسب دقیق شامل طیف گسترده‌ای از عملیات‌های بسیار هزینه‌بر و دشوار است، از این‌رو تکنیک‌های تشخیص ناهنجاری را بر اساس میزان در دسترس بودن برچسب‌ها می‌توان به سه دسته: تشخیص ناهنجاری با نظارت، تشخیص ناهنجاری نیمه نظارتی و تشخیص ناهنجاری بدون نظارت تقسیم کرد[16].

## ۲-۱-۱- تشخیص ناهنجاری با نظارت

در این دسته از روش‌ها هر دو الگوری رفتاری غیرطبیعی و طبیعی مدل می‌شوند. در این مدل‌ها به داده‌های غیر طبیعی برچسب ناهنجاری و به داده‌های طبیعی برچسب عادی می‌زنند. در این رویکرد، برخی از مدل‌ها نمونه‌های ورودی را با نمونه‌های غیرعادی مقایسه می‌کنند و برخی دیگر نمونه‌ها را با نمونه‌های برچسب عادی مقایسه می‌کنند تا بر اساس آن در مورد ماهیت نمونه ورودی تصمیم‌گیری کنند[17].

## ۲-۱-۲ تشخیص ناهنجاری نیمه‌نظارتی

در تشخیص ناهنجاری نیمه‌نظارتی تنها الگوی رفتار طبیعی داده مدل می‌شود و به بیان دیگر تنها به برچسب‌های کلاس عادی نیاز داریم. از نظر کمی این دسته از روش‌ها کاربرد بیشتری نسبت به روش‌های تشخیص ناهنجاری نظارتی دارند[18].

## ۲-۱-۳ تشخیص ناهنجاری بدون نظارت

اساس کار این دسته از روش‌ها همانند روش‌های خوشه‌بندی<sup>۸</sup> است و مدل کلاس داده‌های ناهنجاری را به صورت خودکار از سایر کلاس‌ها تمیز می‌دهد[19]. این روش خوشه‌ای از داده‌ها با رفتار نزدیک به هم پیدا می‌کند و بدین ترتیب عملیات شناسایی ناهنجاری صورت می‌گیرد. اینگونه از مدل‌ها در بسیاری از تشخیص‌ها دچار مشکل می‌شوند چراکه ممکن است نمونه‌های ناهنجار خود باعث ایجاد خوشه‌هایی با

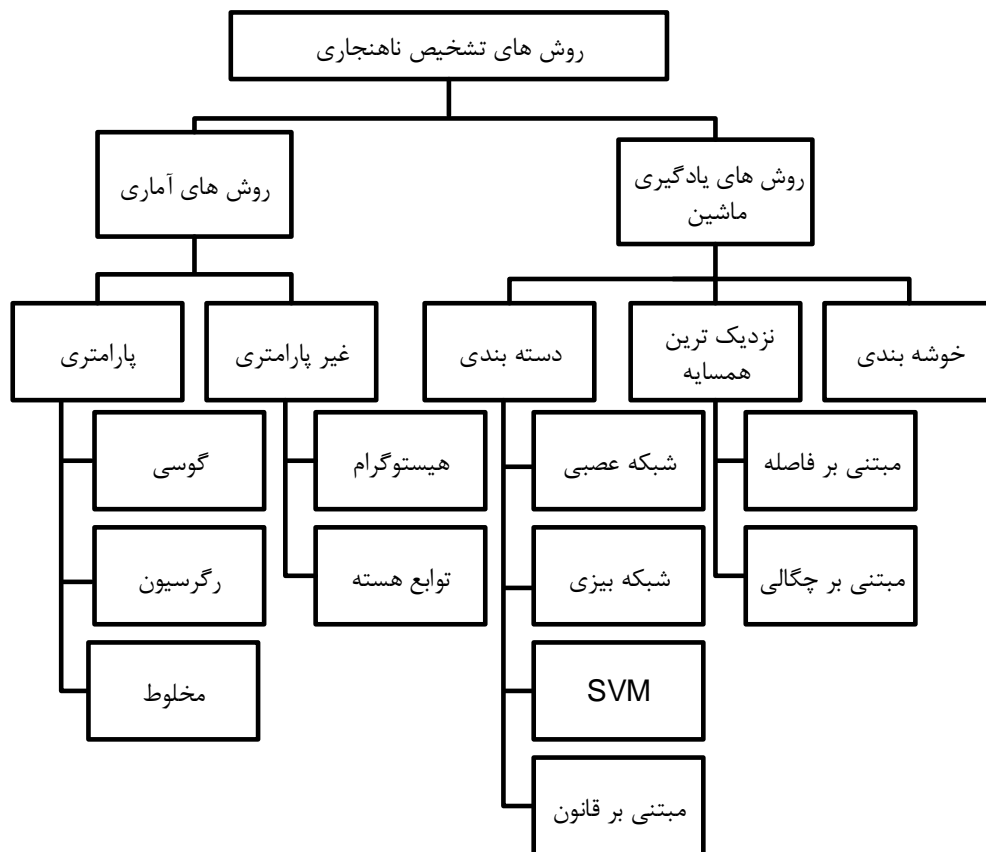
---

<sup>8</sup> Clustering

الگوی مشابه داده‌های عادی شوند، به همین دلیل تکنیک‌های بدون نظارت در تولید نتایج دقیق کارآمد نیستند و اغلب دارای نرخ مثبت کاذب<sup>۹</sup> هستند [20].

## ۲-۲- طبقه‌بندی روش‌های تشخیص ناهنجاری از نظر رویکرد حل مسئله

روش‌های تشخیص ناهنجاری به طور کلی به دو دسته روش‌های آماری و روش‌های مبتنی بر یادگیری ماشین تقسیم می‌شوند. روش‌های آماری خود شامل دو دسته پارامتری و غیرپارامتری هستند و روش‌های یادگیری ماشین شامل خوشه‌بندی، نزدیک‌ترین همسایه و دسته‌بندی است. در شکل ۱-۲ دسته‌بندی این روش‌ها به طور دقیق‌تر به تصویر کشیده شده است. در ادامه هر یک از این روش‌ها بررسی و مرور خواهد شد.



شکل ۱-۲: دسته بندی روش‌های تشخیص ناهنجاری.

<sup>۹</sup> False positive rate

## ۲-۱-۲- روش‌های آماری

تشخیص ناهنجاری با روش‌های آماری به ترتیب بر اساس آمارگان‌های آماری مانند میانگین و انحراف از معیار، توزیع داده‌ها و توابع احتمال (برای ساختن نمایه‌های رفتاری) انجام می‌شود [21]. در اینگونه از روش‌ها بر اساس آزمون‌های آماری هر نوع انحراف از رفتار عادی داده‌ها تشخیص داده می‌شود و داده مورد نظر به عنوان ناهنجاری در نظر گرفته می‌شود. به منظور توسعه مدل‌های آماری در تشخیص ناهنجاری از دو نوع تکنیک پارامتری و غیرپارامتری استفاده می‌شود [22]. تکنیک‌های غیرپارامتری از اطلاعات زمینه‌ای داده‌ها استفاده نمی‌کنند، به بیان دیگر اطلاعاتی از توزیع داده ورودی ندارند در حالی که روش‌های پارامتری با استفاده از همین اطلاعات مدل را طراحی می‌کنند.

### ۲-۱-۲-۲ روش‌های پارامتری

در روش‌های آماری فرض می‌شود داده‌های واقعی بر اساس پارامترهای مشخص از یک توزیع یا تابع خاص تولید می‌شوند، این دسته از روش‌ها خود به سه دسته کلی مدل رگرسیونی، مدل گاوسی و مدل مخلوط تقسیم می‌شوند.

در مدل رگرسیونی داده‌ها بر یک مدل رگرسیونی منطبق می‌شوند و باقی‌مانده<sup>۱۰</sup> هر داده که بر مدل منطبق نیست به عنوان معیار جهت تشخیص ناهنجاری به کار برده می‌شود.

در مدل گاوسی، فرض بر این است که داده‌ها به توزیع گاوسی تعلق دارند و پارامترهای مدل با استفاده از استفاده از تخمین بیشینه درست‌نمایی<sup>۱۱</sup> تعیین می‌شوند. در این مدل‌ها از آزمون‌هایی نظیر آزمون کای-دو<sup>۱۲</sup> جهت شناسایی نمونه ناهنجار استفاده می‌شود [23].

مدل‌های مخلوط خود ترکیبی از سایر مدل‌های پارامتری هستند. چنین مدل‌هایی در برخی از کاربردها عملکرد بسیار موفقی از خود نشان داده‌اند. به عنوان مثال با به کارگیری یک مدل مخلوط از روش‌های

<sup>10</sup> Residual

<sup>11</sup> Maximum likelihood estimation

<sup>12</sup> Chi-square

پارامتری برای تشخیص ناهنجاری‌های شبکه، توانسته‌اند در طی زمان بسیار کوتاهی تمام ناهنجاری‌های موجود در شبکه که توسط سناریوهای مختلف ایجاد شده بودند را شناسایی کنند [11].

## ۲-۱-۲-۲- روش‌های غیرپارامتری

در این روش از نمونه‌های عادی برای تولید مدل استفاده می‌شود و انحراف نمونه از مدل به عنوان امتیاز ناهنجاری در نظر گرفته می‌شود. این روش را می‌توان به دو دسته مدل‌های مبتنی بر هیستوگرام و مدل‌های مبتنی بر هسته تقسیم کرد.

در مدل‌های مبتنی بر هیستوگرام، هیستوگرام بر اساس تقریب از داده‌های عادی تولید می‌شود و برای اگر نمونه ورودی در محدوده‌های خاصی از هیستوگرام قرار گیرد به عنوان ناهنجاری شناخته می‌شود [24].

روش مدل‌سازی مبتنی بر هسته<sup>۱۳</sup> یک تابع تشابه بر اساس نمونه‌های موجود از داده استنباط می‌شود [25]. در اینگونه از مدل‌ها در دسترس بودن نمونه‌های کافی به منظور بازنمایی کامل رفتار مجموعه داده ضروریست چراکه در غیر این صورت دقت مدل کاهش می‌یابد.

## ۲-۲-۲- روش‌های یادگیری ماشین

روش‌های مبتنی بر یادگیری ماشین بر اساس تجربه حاصل از مشاهده نمونه‌های قدیمی و به کارگیری آن، ظرفیت تمایز میان رفتارهای غیرطبیعی و طبیعی داده تا حد مناسبی بهبود می‌بخشند [26]. این طبقه از روش‌ها خود به سه گروه دسته‌بندی، نزدیک‌ترین همسایه و خوشه‌بندی تقسیم می‌شوند [27]، در ادامه به بررسی هر یک از این دسته‌ها می‌پردازیم.

<sup>13</sup> kernel



## ۲-۲-۱- دسته‌بندی

هدف اصلی از روش‌های مبتنی بر دسته‌بندی، اختصاص هر نمونه از داده به یکی از کلاس‌های از پیش تعیین شده بر اساس ویژگی‌های آن نمونه است. از مزیت‌های این دسته از روش‌ها می‌توان به توانایی بالای آن‌ها در تمایز میان کلاس‌های مختلف داده در زمان آزمایش اشاره کرد. از روش‌های متداول که در این دسته جای دارند می‌توان به شبکه‌های بیزی، ماشین بردار پشتیبان<sup>۱۴</sup>، برخی روش‌های مبتنی بر قانون و شبکه‌های عصبی اشاره کرد.

شبکه‌های بیزی در واقع مدل‌های گرافیکی هستند که اتصالات میان نمونه‌های مختلف را بر اساس محاسبه احتمال پیشین<sup>۱۵</sup> یک نمونه از داده به همراه دسته‌ای از پیش‌شروط مورد بررسی و ترجمه قرار می‌دهند. اساس کار این دسته از روش‌ها استفاده از یادگیری با نظارت است.

ماشین‌های بردار پشتیبان از دسته الگوریتم‌های یادگیری با نظارت هستند که در صورت استفاده از هسته نمونه‌ها را فضای با ابعاد بالاتر می‌برند و در فضای جدید نمونه‌ها را به دو کلاس تقسیم می‌کنند. استفاده از هسته زمانی توجه پذیر است که نمونه‌ها در فضای با ابعاد پایین جداپذیر نباشند. این مدل به دلیل استفاده از یک مرز خطی به منظور جداسازی نمونه‌های غیرطبیعی و عادی به عنوان دسته‌بندی خطی شناخته می‌شود [28].

روش‌های مبتنی بر قانون بر اساس یک سری از قواعد رفتار و عملکرد نمونه‌های عادی را می‌آموزد، بنابراین اگر یک نمونه نتواند از این مجموعه قوانین پیروی کند به عنوان نمونه ناهنجار شناخته خواهد شد. از مطرح‌ترین روش‌هایی که در این دسته می‌گنجند می‌توان به درخت تصمیم<sup>۱۶</sup> اشاره کرد [29].

شبکه‌های عصبی رفتار سیستم عصبی انسان را تقلید می‌کنند و شامل مجموعه‌ای از فرایندهای بهم پیوسته هستند که به طور همزمان روی داده عمل می‌کنند. در این دسته از روش‌ها از نمونه‌های عادی برای آموزش شبکه عصبی استفاده می‌شود. از نظر دسترسی به برجسب نمونه‌ها شبکه‌های عصبی را می‌توان مشترک بین دسته‌های یادگیری با نظارت و بدون نظارت در نظر گرفت [30]. یکی از انواع این

<sup>14</sup> Support vector machine

<sup>15</sup> Posterior probability

<sup>16</sup> Decision tree

شبکه‌ها که در سال‌های اخیر به موفقیت چشم‌گیری دست یافته است، شبکه‌های مولد تقابلی<sup>۱۷</sup> هستند. این دسته از شبکه‌ها با موفقیت بر روی داده‌های دنیای واقعی که دارای ابعاد بالا هستند اعمال شده‌اند و چهارچوب یادگیری خصمانه آن‌ها عملکرد مناسبی از خود بر جای گذاشته است. قابلیت این شبکه‌ها نشان‌دهنده ظرفیت آن‌ها برای استفاده در مسئله تشخیص ناهنجاری می‌باشد و به همین جهت به کارگیری شبکه‌های مولد تقابلی در حوزه تشخیص ناهنجاری اخیراً مورد توجه و کاوش قرار گرفته است [31]. در تشخیص ناهنجاری به کمک شبکه‌های مولد تقابلی با استفاده از فرایند آموزش تقابلی رفتار عادی داده مدل می‌شود سپس با اندازه‌گیری امتیاز ناهنجاری روی نمونه‌های مختلف عمل شناسایی نمونه ناهنجار صورت می‌پذیرد. شبکه‌های مولد تقابلی با کمک آموزش و نمونه‌گیری از مدل‌های مولد به نتایج بسیار مناسبی در مقایسه با دیگر روش‌ها دست می‌یابند همچنین این مدل‌ها امکان آموزش داده‌های از دست رفته به کمک الگوریتم‌های یادگیری تقویتی<sup>۱۸</sup> را می‌دهد. بیان جزئیات بیشتر در خصوص این دسته از شبکه‌ها را به به بخش‌هایی که در آینده خواهیم داشت موکول می‌کنیم.

## ۲-۲-۲-۲- نزدیک‌ترین همسایه

روش نزدیک‌ترین همسایه مبتنی بر سنجش فاصله یا تراکم میان داده‌هاست به بیان دیگر امتیاز ناهنجاری مقدار همین فاصله است و بسته به مسئله و میزان در دسترس بودن برچسب‌ها این روش می‌تواند به عنوان روش یادگیری بدون نظارت و یا با نظارت به کار گرفته شود [11].

## ۲-۲-۲-۳- خوشه‌بندی

روش مبتنی بر خوشه‌بندی از دسته روش‌های یادگیری بدون نظارت است که برای شناسایی مجموعه نمونه‌های شبیه به هم به کار برده می‌شود. ناهنجاری‌ها ممکن است تشکیل یک خوشه کوچک بدهند یا در هیچ خوشه‌ای جای نگیرند. این روش در مقایسه با روش‌های مبتنی بر فاصله از پیچیدگی محاسباتی کمتری برخوردار است و در عین حال نقطه ضعف این روش عملکرد نامناسب روی دادگان کوچک است

<sup>17</sup> Generative Adversarial Networks (GAN)

<sup>18</sup> Reinforcement learning

چرا که مدل بینش مناسبی نسبت به داده ندارد و به عنوان مثال برای قسمتی از فضا که برای آن نمونه آموزشی نداریم همواره برچسب ناهنجاری در نظر می‌گیرد در صورتی که ممکن است در حضور تعداد داده کافی برچسب آن نمونه خاص برچسب عادی باشد [11].

## ۳-۲- دسته‌بندی بر اساس نحوه تشخیص ناهنجاری

در این قسمت بر اساس نحوه تشخیص ناهنجاری روش‌های موجود را به دسته بر اساس فاصله، دسته‌بندی تک‌کلاسی و بر اساس بازسازی<sup>۱۹</sup> تقسیم می‌کنیم [32].

### ۳-۲-۱- بر اساس فاصله

از روش‌های مبتنی بر فاصله می‌توان به عنوان یکی از کلاس‌های اصلی روش‌های تشخیص ناهنجاری یاد کرد. در این روش‌ها با استفاده از محاسبه فاصله یک نمونه خاص با نزدیک‌ترین همسایه‌ها و یا نزدیک‌ترین کلاستر، نمونه ناهنجار شناسایی می‌شود. بدیهی‌ست به کار بردن چنین روش‌هایی نیازمند طراحی یا انتخاب معیار فاصله مناسب است [32].

### ۳-۲-۲- دسته‌بندی تک‌کلاسی

در رویکرد دسته‌بندی تک‌کلاسی تنها نمونه‌های عادی به یک دسته‌بند نظیر SVM آموزش داده می‌شوند [33]، در واقع این مدل‌ها یک مرز تصمیم حول نمونه‌های عادی یاد می‌گیرند، در صورتی که نمونه ورودی داخل این مرز قرار گیرد به عنوان نمونه عادی و در غیر این‌صورت به عنوان نمونه ناهنجار شناخته می‌شود.

<sup>19</sup> Reconstruction

## ۲-۳-۳- بر اساس بازسازی

این دسته از الگوریتم‌ها بر اساس میزان خطای بازسازی به شناسایی نمونه‌های ناهنجار می‌پردازند. PCA<sup>۲۰</sup> و الگوریتم‌های مشتق از آن جزو همین دسته روش‌ها هستند. اخیراً اساس کار بیشتر کارهای پژوهشی و کاربردی در زمینه تشخیص ناهنجاری شبکه‌های عصبی هستند و به نظر می‌رسد این شبکه‌ها دارای سابقه طولانی در این زمینه هستند [32]. به عنوان مثال رویکردهای مبتنی بر خودکدگذار و خودکدگذار متغیر<sup>۲۱</sup> روند بازسازی نمونه‌های عادی را فرامی‌گیرند و نمونه‌های با خطاری بازسازی زیاد را به عنوان ناهنجاری در نظر می‌گیرند. مدل‌های مبتنی بر انرژی و مدل‌های ترکیبی گاوسی با خودکدگذار عمیق<sup>۲۲</sup> به طور خاص در زمینه تشخیص ناهنجاری مورد تحقیق و پژوهش قرار گرفته‌اند. چنین روش‌هایی توزیع داده را با استفاده از خودکدگذار یا روش‌های مشابه مدل می‌کنند و سپس بر اساس انرژی و یا ترکیب گاوسی‌ها یک معیار آماری تشخیص ناهنجاری پدید می‌آورند. در سال‌های اخیر از شبکه‌های مولد تقابلی به منظور تشخیص ناهنجاری استفاده شده است. در این مدل‌ها به هنگام آزمایش برای هر نمونه ورودی عمل استنتاج انجام می‌شود و با استفاده از انتشار خطای گرادیان نزولی<sup>۲۳</sup> در شبکه مولد، پارامترهای فضای نهفته بازیابی می‌شوند و با استفاده از این پارامترها می‌توان به شناسایی نمونه‌های ناهنجار پرداخت.

## ۲-۴- معیارهای ارزیابی روش‌های تشخیص ناهنجاری

صرف نظر از رویکرد به کار گرفته شده، تشخیص ناهنجاری از مرحله یادگیری که در آن با استفاده از نمونه‌های آموزشی مدل آموزش داده می‌شود، آغاز می‌شود. پس از اتمام مرحله یادگیری، مدل دسته‌بندی نمونه‌هایی که تاکنون آن‌ها را مشاهده نکرده است را بر اساس معیارهای مورد نظر انجام می‌دهد. نتیجه ارزیابی در میزان تشخیص ناهنجاری‌ها می‌تواند در چهار دسته مختلف گزارش شود که

<sup>20</sup> Principal component analysis

<sup>21</sup> Variational autoencoder

<sup>22</sup> Deep autoencoding gaussian mixture models

<sup>23</sup> Gradient descent

عبارتند از مثبت صحیح<sup>۲۴</sup> به اختصار TP، منفی صحیح<sup>۲۵</sup> به اختصار TN، مثبت کاذب یا FP و منفی کاذب یا FN<sup>۲۶</sup>. در حوزه تشخیص ناهنجاری از معیارهای معمول و استاندارد نظیر صحت<sup>۲۷</sup>، بازیابی<sup>۲۸</sup>، F1-score و مساحت زیر نمودار مشخصه عملکرد<sup>۲۹</sup> و یا به اختصار AUROC که در ادامه به بررسی جزئیات ریاضی هر یک خواهیم پرداخت.

## ۲-۴-۱- صحت

معیار صحت یکی از معیارهای پایه در تمامی مسائل مربوط به دسته‌بندی علی‌الخصوص تشخیص ناهنجاری می‌باشد. این معیار بیانگر این است که چه مقدار از داده‌هایی که به عنوان ناهنجاری تشخیص داده شده‌اند، واقعا ناهنجار هستند. بیان ریاضی این معیار مطابق معادله ۱-۲ می‌باشد.

$$\text{Precision} = \frac{TP}{TP + FP} \quad (1-2)$$

## ۲-۴-۲- بازیابی

این معیار در کنار معیار صحت دید نسبتاً خوبی از وضعیت کلی عملکرد مدل نمایان می‌کند. این معیار بیان می‌کند که چه بخشی از داده‌های ناهنجار کشف شده و مدل توانسته چه نسبتی از این دسته را بازیابی کند. معادله ۲-۲ نحوه محاسبه این معیار را نمایش می‌دهد.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2-2)$$

<sup>24</sup> True positive

<sup>25</sup> True negative

<sup>26</sup> False negative

<sup>27</sup> Precision

<sup>28</sup> Recall

<sup>29</sup> Area Under Curve Receiver Operating Characteristics

## F1-score – ۳-۴-۲

در حالت کلی، یک مدل خوب باید هر دو معیار صحت و بازیابی بالایی داشته باشد و به هنگام مقایسه مدل‌ها، در صورتی که هر دو معیار یادشده یک مدل از دیگری بیشتر باشد، مشخصاً آن مدل عملکرد بهتری داشته است. اما اگر هر دو به صورت همزمان به سمت یک مقدار مناسب همگرا نباشند، انتخاب مدل بهتر با مشکل روبرو می‌شود و نیاز به یک معیار با دید کلی‌تر می‌باشد. معیار F1-score به طور همزمان هر دو جنبه مورد سنجش را ارزیابی می‌کند و معیار جامع‌تری برای مقایسه می‌باشد. بیان ریاضی این معیار مطابق معادله ۳-۲ می‌باشد.

$$F1 - score = \frac{2}{\frac{1}{Recall} + \frac{1}{Precision}} = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall} \quad (3-2)$$

## ۴-۴-۲- مساحت زیر نمودار مشخصه عملکرد

این معیار به منظور تعیین میزان تولنایی مدل در تشخیص ناهنجاری به هنگام افزایش اندازه مجموعه داده تعریف شده است. این معیار با هدف ایجاد اطمینان از توانایی مدل در کنترل تغییرات سریع حجم داده ایجاد شده است و برای محاسبه آن با توجه به جنس مجموعه داده مورد آزمایش روش‌های متفاوتی ارائه شده است [11].

## ۵-۲- شبکه‌های مولد تقابلی و تشخیص ناهنجاری

هدف از مدل‌های یادگیری عمیق، کشف مدل‌های سلسله‌مراتبی قوی است. این مدل‌ها نشان‌دهنده توزیع احتمال انواع داده‌هایی است که در کاربردهای هوش مصنوعی نظیر تصاویر طبیعی، شکل موج صوتی حاوی گفتار به کار می‌رود. برجسته‌ترین موفقیت یادگیری عمیق در طراحی مدل‌های تمایزگر<sup>۳۰</sup> بوده است. این مدل‌ها قادرند تا ورودی با ابعاد بالا را دریافت کنند و کلاس هر یک از نمونه‌ها را با قدرت

<sup>30</sup> Discriminator models

تشخیص خود با دقت بالا مشخص کنند. استفاده از الگوریتم‌های پس‌انتشار<sup>۳۱</sup>، حذف تصادفی<sup>۳۲</sup> و واحدهای خطی تکه‌ای<sup>۳۳</sup> که دارای گرادینان با رفتار مناسب هستند عامل موفقیت چشم‌گیر یادگیری عمیق است.

استفاده کاربردی از الگوریتم‌هایی نظیر تخمین بیشینه درست‌نمایی و الگوریتم‌های مرتبط با آن همراه با چالش‌ها و دشواری‌های زیادی نظیر محاسبات احتمالاتی زیاد و خارج از کنترل است، وجود این چالش و همچنین سختی‌های موجود در استفاده از واحدهای خطی تکه‌ای در حوزه مدل‌های مولد سبب شده است تا مدل‌های مولد عمیق کمتر مورد توجه قرار گیرند. شبکه‌های مولد تقابلی بر این دست از چالش‌ها و دشواری‌ها فائق آمده و نقش مدل‌های مولد عمیق را پررنگ‌تر ساخته است.

## ۲-۵-۱- شبکه‌های مولد تقابلی

شبکه‌های مولد تقابلی اولین بار در سال ۲۰۱۴ توسط آقای گودفلو و همکاران ابداع شد [۲۸]، در این شبکه‌ها یک مدل مولد در برابر یک مدل تمایزگر قرار می‌گیرد، مدل تمایزگر سعی می‌کند میان داده‌های واقعی و داده‌های تولیدی توسط شبکه مولد تمایز ایجاد کند. مدل مولد را می‌توان مانند تیمی فرض کرد که سعی در تولید ارز جعلی با میزان شباهت بسیار بالا به ارز واقعی دارد و در طرف مقابل مدل تمایزگر مشابه پلیس است که سعی در کشف ارز تقلبی دارد. رقابت در این بازی، هر دو تیم را به سمت بهبود روش‌های خود سوق می‌دهد تا اینکه ارز تقلبی از ارز اصلی غیرقابل تشخیص باشد. این چارچوب می‌تواند الگوریتم‌های آموزشی خاصی را برای انواع مختلف از مسائل و مدل‌ها بهینه‌سازی کند. بخش مولد با دریافت نویز تصادفی، از طریق پرسپترون چند لایه نمونه‌هایی با توزیع مشابه داده اصلی تولید می‌کند و مدل تمایزگر با استفاده از مدل پرسپترونی چند لایه تلاش می‌کند تا میان نمونه‌های مختلف تمایز ایجاد کند. در این نوع از تعریف شبکه می‌توان هر دو مدل را با استفاده از الگوریتم‌های پس‌انتشار و

<sup>31</sup> Backpropagation

<sup>32</sup> Dropout

<sup>33</sup> Piecewise linear units

حذف تصادفی ایجاد کرد و برای نمونه‌گیری از مدل مولد تنها از الگوریتم انتشار رو به جلو<sup>۳۴</sup> استفاده کرد و در نتیجه به کارگیری هیچ الگوریتمی نظیر استنتاج تقریبی و یا زنجیره مارکوف ضروری نیست.

در شبکه مولد تقابلی به طور همزمان دو مدل مولد و تمایزگر آموزش داده می‌شود. مدل مولد  $G$  توزیع داده را ضبط می‌کند و مدل تمایزگر  $D$  احتمال این که نمونه از داده‌های تولید شده توسط  $G$  باشد را تخمین می‌زند. تابع هدف برای شبکه مولد  $G$  به حداکثر رساندن احتمال اشتباه شبکه  $D$  است. این بستر منجر به یک بازی دو نفره مانند بازی‌های بیشینه-کمینه<sup>۳۵</sup> می‌شود. در فضای توابع دلخواه  $G$  و  $D$  یک راه حل منحصر به فرد وجود دارد و در این راه حل شبکه مولد  $G$  توزیع داده‌های آموزشی را یاد گرفته است و شبکه تمایزگر  $D$  احتمال را در همه جا برابر مقدار  $1/2$  نشان می‌دهد، به بیان دیگر میان داده‌های واقعی و داده تولید شده توسط شبکه مولد نمی‌تواند تمیز دهد.

مدل‌سازی چارچوب تقابلی به وسیله ایجاد یک مدل چند لایه پرسپترون برای هر دو مدل مولد و تمایزگر انجام می‌شود. برای یادگیری توزیع مولد  $p$  روی داده  $x$ ، یک تابع نویز خالص  $p(z)$  را به عنوان ورودی تعریف می‌کنیم، سپس یک نگاشت از فضای نهفته به فضای داده را به عنوان  $G(z; \theta_g)$  نشان می‌دهیم، در اینجا  $G$  یک تابع مشتق‌پذیر است که توسط یک پرسپترون چند لایه با پارامترهای  $\theta_g$  نمایش داده می‌شود. همچنین برای شبکه تمایزگر  $D$  یک پرسپترون چند لایه لایه  $D(x; \theta_d)$  با یک خروجی اسکالر<sup>۳۶</sup> تعریف می‌کنیم.  $D(x)$  بیانگر احتمال این است که  $x$  از داده‌های اصلی به جای توزیع  $p_g$  ناشی شده باشد. در این میان به شبکه  $D$  آموزش داده می‌شود تا احتمال تخصیص برچسب صحیح به داده‌های واقعی و نمونه‌های تولیدی از  $G$  را به حداکثر برساند. به طور همزمان به شبکه  $G$  آموزش داده می‌شود تا تابع هدف  $\log(1 - D(G(z)))$  را به حداقل برساند. به عبارت دیگر، شبکه‌های  $D$  و  $G$  بازی کمینه-بیشینه دو نفره زیر با تابع  $V(G, D)$  مطابق معادله ۲-۴ انجام می‌دهند.

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim q} [\log D(x)] + \mathbb{E}_{z \sim p} [\log (1 - D(G(z)))] \quad (۲-۴)$$

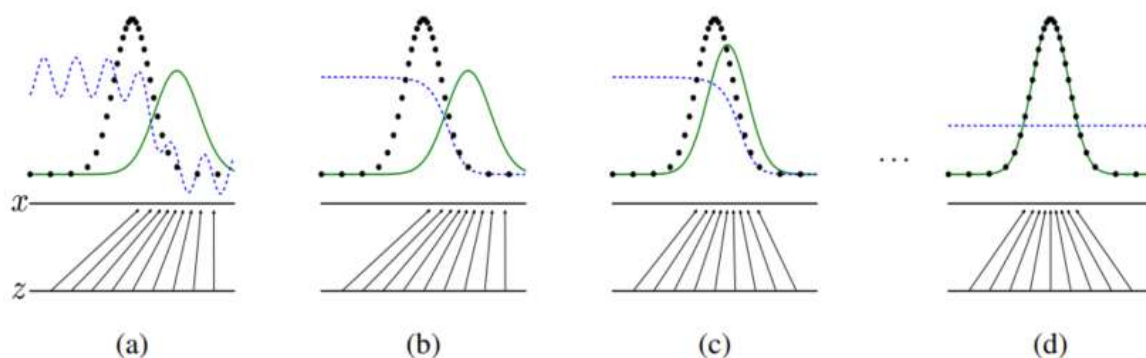
<sup>۳۴</sup> Feed forward

<sup>۳۵</sup> Minimax

<sup>۳۶</sup> Scaler



با تحلیل نظری صورت گرفته بر روی شبکه‌های مولد تقابلی، نشان داده شده است که این شبکه‌ها پتانسیل کافی برای بازیابی توزیع داده‌های اصلی را در قالب شبکه مولد  $G$  دارند. در عمل معادله ۱-۲ ممکن است گرادیان کافی برای آموزش شبکه  $G$  را فراهم نکند. در اوایل یادگیری، هنگامی که شبکه مولد  $G$  ضعیف است، شبکه تمایزگر  $D$  می‌تواند با اطمینان بالا نمونه‌های غیرواقعی را شناسایی کند چرا که نمونه‌های تولیدی با نمونه‌های آموزشی کاملاً متفاوت هستند. در این حالت  $\log(1 - D(G(z)))$  شیب می‌شود. در این حالت می‌توان به جای آموزش  $G$  برای به حداقل رساندن تابع  $\log(1 - D(G(z)))$  می‌توانیم  $G$  را برای به حداکثر رساندن  $D(G(z))$  آموزش دهیم. این تابع در همان نقطه ثابت  $G$  و  $D$  قرار دارد اما گرادیان قوی‌تری در یادگیری فراهم می‌کند. رویکرد کلی شبکه‌های مولد تقابلی در شکل ۲-۲ نشان داده شده است.



شکل ۲-۲: رویکرد کلی شبکه‌های مولد تقابلی [۳۱].

همانطور که در شکل ۲-۲ مشاهده می‌کنید توزیع تمایزگر (خط آبی شکسته) به‌روزرسانی می‌شود تا بتواند نمونه‌های توزیع داده‌های اصلی (خط مشکی نقطه‌چی) از داده‌های تولیدشده توسط توزیع مولد  $p_g$  (خط سبز پیوسته) تمیز دهد. خط افقی پایین بیانگر فضای نهفته است که متغیر  $z$  با توزیع یکنواخت از آن نمونه‌برداری شده است. خط افقی بالا بخشی از فضای داده واقعی  $x$  است. فلش‌های رو به بالا نشان می‌دهد که چگونه تابع  $G$ ،  $z$  را به طور غیریکنواخت به  $x$  نگاشت می‌کند. به مرور زمان  $G$  در مناطق چگال‌تر منقبض می‌شود و در مناطق با چگالی کمتر باز می‌شود. در قسمت (a) شکل ۲-۲ دو بلوک تقابلی نزدیک همگرایی هستند، یعنی  $p$  به  $q$  نزدیک شده است و همچنین  $D$  دسته‌بند تا حدی دقیق می‌باشد. در قسمت (b) در حلقه داخلی الگوریتم،  $D$  آموزش می‌بیند تا بتواند نمونه‌های غیرواقعی را تشخیص دهد و به  $D_G^*(x) = \frac{q(x)}{q(x)+p(x)}$  همگرا شود. در (c) پس از به‌روزرسانی  $G$  گرادیان ناشی از  $D$ ،  $G$  را به گونه‌ای هدایت می‌کند که به سمت مناطقی مایل شود که به توزیع داده واقعی نزدیک‌تر شود. در قسمت آخر (d)

پس از انجام چند مرحله از آموزش اگر  $G$  و  $D$  ظرفیت کافی را داشته باشند به نقطه تعادلی می‌رسند که در آن  $p = q$ . و تمایزگر دیگر قادر نیست میان توزیع داده‌های واقعی و غیرواقعی تفاوتی قائل شود.

به منظور پیاده‌سازی این شبکه‌ها از روش عددی مبتنی بر تکرار استفاده می‌شود. تکمیل بهینه‌سازی شبکه  $D$  در حلقه داخلی مرحله آموزش همراه با چالش‌هایی نظیر هزینه محاسباتی زیاد است و همچنین روی دادگان کوچک منجر به بیش‌برازش<sup>۳۷</sup> است. حال به جای اینکه در هر تکرار هر دو شبکه  $D$  و  $G$  بهینه شوند، به ازای  $k$  مرحله بهینه کردن  $D$  یک مرحله  $G$  بهینه می‌شود. با این کار تا زمانی که  $G$  به اندازه کافی آهسته تغییر کند  $D$  در نزدیکی نقطه بهینه باقی خواهد ماند. شبه کد الگوریتم شبکه‌های مولد تقابلی در الگوریتم ۱-۲ آورده شده است.

$k$  تعداد مراحل اعمال شده تمایزگر ( برای کاهش هزینه محاسبات اینجا عدد یک فرض می‌شود) و  $n$  تعداد تکرار آموزش

for k steps do  
for n steps do

۱. نمونه‌برداری کوچک‌دسته‌ای<sup>۳۸</sup>  $m$  تایی نویز  $\{z^{(1)}, \dots, z^{(m)}\}$  از نمونه‌های نویز  $p(z)$ .
۲. نمونه‌برداری کوچک‌دسته‌ای  $m$  تایی نویز  $\{x^{(1)}, \dots, x^{(m)}\}$  از داده‌های تولید توزیع  $q(x)$ .
۳. بروزرسانی صعودی تمایزگر به وسیله گرادیان تصادفی.

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))]$$

end for

۴. نمونه‌برداری کوچک‌دسته‌ای  $m$  تایی نویز  $\{z^{(1)}, \dots, z^{(m)}\}$  از نمونه‌های نویز  $p(z)$ .
۵. بروزرسانی صعودی مولد به وسیله گرادیان تصادفی.

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)})))$$

end for

الگوریتم ۱-۲: آموزش گرادیان نزولی کوچک دسته‌ای شبکه‌های مولد تقابلی.

<sup>37</sup> Overfitting

<sup>38</sup> Minibatch

## ۲-۵-۱-۱- تحلیل نظری شبکه مولد تقابلی

شبکه مولد  $G$  به طور ضمنی یک تابع توزیع احتمال  $p$  را به عنوان توزیع نمونه‌های  $G(Z)$  تعریف کرده است (توجه داشته باشید  $Z \sim p$ ). در صورتی که ظرفیت مناسب و زمان کافی برای آموزش در اختیار الگوریتم ۱-۲ قرار گیرد در نهایت این الگوریتم توزیع داده ورودی  $q(x)$  را خواهد یافت. نتایج حاصل از شبکه‌های عصبی مولد تقابلی بر اساس تنظیمات غیرپارامتری به دست آمده است. در این بخش نشان می‌دهیم که در بازی کمینه-بیشینه بین دو شبکه  $p = q$  یک بهینه عمومی است. برای این منظور ابتدا تمایزگر بهینه  $D$  را برای هر مولد  $G$  در می‌گیریم.

**قضیه ۱-۲-** برای هر تابع مولد  $G$  ثابت، تابع تمایزگر بهینه  $D$  عبارت است از:

$$D_G^*(x) = \frac{q(x)}{q(x) + p(x)} \quad (۵-۲)$$

**اثبات:** معیار آموزش برای تمایزگر  $D$ ، با توجه به هر مولد  $G$ ، به حداکثر رساندن مقدار  $V(G, D)$  می‌باشد، پس داریم:

$$\begin{aligned} V(G, D) &= \int_x q(x) \log(D(x)) dx + \int_z p(z) \log(1 - D(g(z))) dz \\ &= \int_x q(x) \log(D(x)) dx + p(x) \log(1 - D(x)) dx \end{aligned} \quad (۶-۲)$$

از طرفی می‌دانیم برای هر  $(a, b) \in \mathbb{R}^2$ ، تابع  $y \rightarrow a \log(y) + b \log(1 - y)$  در بازه  $[0, 1]$  در  $\frac{a}{a+b}$  بیشینه است. هم‌چنین می‌دانیم تمایزگر نیاز به تعریف بیرون از مرز  $\text{Supp}(p_{data}) \cup \text{Supp}(p_g)$  ندارد، پس در نتیجه  $D_G^*(x)$  نقطه بهینه برای به حداکثر رساندن  $V(G, D)$  می‌باشد.

می‌توان هدف از آموزش شبکه  $D$  را به حداکثر رساندن لگاریتم درست‌نمایی<sup>۳۹</sup> احتمال  $P(Y = y|x)$  تعبیر کرد، که  $Y$  بیانگر آن است که هر جا  $x$  از توزیع  $q$  باشد  $y=0$  و هر جا از توزیع  $p$  باشد  $y=1$  است. با این تعریف بازی کمینه-بیشینه در معادله ۲-۴ را می‌توان به صورت زیر، بازنویسی کرد:

<sup>39</sup> Log-Likelihood

$$\begin{aligned}
C(G) &= \max_D V(G, D) = \mathbb{E}_{x \sim q} [\log D_G^*(x)] + \mathbb{E}_{z \sim p(z)} [\log (1 - D_G^*(G(z)))] \\
&\quad + \mathbb{E}_{x \sim q} [\log D_G^*(x)] + \mathbb{E}_{x \sim p} [\log (1 - D_G^*(x))] \quad (۷-۲) \\
&= \mathbb{E}_{x \sim q} [\log \frac{q(x)}{q(x) + p(x)}] + \mathbb{E}_{x \sim p} [\log \frac{p(x)}{q(x) + p(x)}]
\end{aligned}$$

**قضیه ۲-۲-** کمینه سراسری  $C(G)$  تنها در حالتی قابل محاسبه است که اگر و تنها اگر  $q = p$  و مقدار  $C(G)$  برابر با  $-\log 4$  باشد.

**اثبات:** طبق قضیه ۱-۲، می‌دانیم هنگامی که  $q = p$  باشد،  $D_G^* = 1/2$  می‌شود. در ادامه نشان دادیم که  $C(G) = \log \frac{1}{2} + \log \frac{1}{2} = -\log 4$  حال برای اثبات این که بهترین مقدار  $C(G)$  این مقدار است، این عبارت را از تعریف  $C(G)$  کم می‌کنیم. طبق همگرایی جنسن-شانون<sup>۴۰</sup> داریم:

$$C(G) = -\log 4 + 2 \cdot JSD(p_{data} || p_g) \quad (۸-۲)$$

می‌دانیم تابع همگرایی جنسن-شانون بین دو توزیع همواره نامنفی است و صفر است اگر و تنها اگر دو توزیع برابر باشند. پس ما نشان دادیم مقدار بهینه برابر با  $-\log 4$  است و نقطه بهینه هنگام برابری دو توزیع رخ می‌دهد. بدین ترتیب اثبات این قضیه نیز به پایان رسید [۲۸].

## ۲-۱-۵-۲- مزایا و معایب

شبکه‌های مولد تقابلی نسبت به مدل‌های قبلی دارای مزایا و معایبی می‌باشد [31]. ایراد این روش آموزش این است که نمایش صریح  $p(x)$  وجود ندارد و  $D$  باید در حین آموزش به خوبی با  $G$  هماهنگ شود و به طور خاص،  $G$  نباید بدون به‌روزرسانی  $D$  خیلی زیاد آموزش داده. اما از طرف دیگر، در هنگام یادگیری این مدل‌ها نیازی به استنباط نیست و می‌توان طیف گسترده‌ای از توابع را در مدل گنجانید. این مزیت در درجه اول محاسباتی است، مدل‌های تقابلی همچنین ممکن است برخی از مزیت‌های آماری

<sup>40</sup> Jensen-Shanon divergence

را از شبکه مولد به دست آورند که مستقیماً با نمونه داده‌ها به روز نمی‌شوند، و فقط با گرادینان‌هایی که از طریق تمایزگر جریان می‌یابند، بروزرسانی می‌شود.

این بدان معنی است که اجزای ورودی مستقیماً در پارامترهای مولد بکارگرفته نمی‌شوند. همچنین یکی دیگر از مزیت‌های شبکه‌های تقابلی این است که آن‌ها می‌توانند توزیع‌های بسیار تیز و لبه‌دار را نشان دهند، در حالی که روش‌های مبتنی بر زنجیره‌های مارکوف نیاز دارند که توزیع تا حدی هموار<sup>۴۱</sup> باشد تا زنجیرها بتوانند میان حالت‌ها جابجا شده و مدل‌سازی را انجام دهند [31].

## ۲-۵-۲- مدل ANOGAN

مدل ANOGAN به منظور مدل‌سازی در زمینه پزشکی و به طور خاص مدل‌سازی وضعیت سلامت موضعی طراحی شده است [13]. این مدل از دسته الگوریتم‌های مولد و بدون نظارت است. با به وجود آمدن این مدل توانایی شبکه‌های مولد تقابلی در ایجاد یک مدل با قدرت بازنمایی بالا در تشریح وضعیت آناتومی ثابت شد. لازم به ذکر است مدل پیشنهادی شبکه تمایزگر و شبکه مولد را به طور همزمان آموزش می‌دهد و با استفاده از هر دو شبکه مشخص می‌کند داده ورودی از جنس داده‌های آموزشی است و یا باید به عنوان ناهنجاری دسته‌بندی شود.

به منظور تشخیص ناهنجاری مدل مورد نظر بازنمایی نمونه‌های متنوع آناتومیکی طبیعی را می‌آموزد. در این کار به جای استفاده از بهینه‌سازی تابع هزینه واحد، از تابع تعادل نش<sup>۴۲</sup> میان هزینه‌ها استفاده شده است که سبب افزایش قدرت بازنمایی و رشد نرخ منفی صحیح<sup>۴۳</sup> به اختصار TNR مدل تولیدی، بهبود روند نگاشت ویژگی<sup>۴۴</sup> و همچنین دستیابی به دقت بالا در طبقه‌بندی داده‌های واقعی از داده‌های غیرواقعی می‌شود. در ادامه چگونگی طراحی مدل و نحوه شناسایی وضعیت و ظواهری که در داده‌های آموزش دیده نشده‌اند شرح داده خواهد شد.

<sup>41</sup> Smooth

<sup>42</sup> Nash cost

<sup>43</sup> True Negative Rate

<sup>44</sup> Feature matching

## ۲-۵-۱- یادگیری بدون نظارت متنوع داده‌های طبیعی

$M$  مجموعه‌ای از تصاویر پزشکی است که هر نمونه آن نمایانگر نمونه‌هایی از تصاویر آناتومی‌های سالم است و با  $I_m$  نمایش داده می‌شود که  $m = 1, 2, \dots, M$ ، که در اینجا  $I_m \in \mathbb{R}^{a \times b}$  است، یعنی اندازه یک تصویر برابر  $a \times b$  است. از هر تصویر  $I_m$ ،  $K$  تکه تصویر دو بعدی  $x_{k,m}$  با ابعاد  $c \times c$  بطور تصادفی از موقعیت‌های مختلف نمونه‌گیری می‌کنیم که منجر به داده‌های  $x_{k,m} \in \mathcal{X}, k = 1, 2, \dots, K$  می‌شود. در طول آموزش، فقط  $I_m$  را در اختیار داریم و برای یادگیری توزیع حاشیه‌ای، که نشان دهنده تنوع تصاویر آموزش است، از یک روش بدون نظارت استفاده می‌شود. برای آزمایش، داریم  $\langle y_n, l_n \rangle$ ، که  $y_n$  تصاویر مشاهد نشده با ابعاد  $c \times c$  استخراج شده از داده  $l_n \in \{0, 1\}$  آرایه‌ای از برچسب‌های حقیقی مبتنی بر تصویر باینری با  $n = 1, 2, \dots, N$  است. این برچسب‌ها فقط در طول آزمایش استفاده می‌شوند، تا کارایی روش تشخیص ناهنجاری ارزیابی شود.

شبکه مولد  $G$  توزیع  $p$  را روی داده  $x$  از طریق نگاشت نمونه‌های  $z$  توسط تابع  $G(z)$  آموزش می‌بیند؛ در واقع بردارهای تک بعدی با توزیع یکنواخت از فضای نهفته  $z$  نمونه‌برداری می‌شوند و به فضای دو بعدی تصویر که در آن تصاویر آناتومی سالم وجود دارند نگاشت می‌شوند. در این تنظیمات، معماری شبکه مولد  $G$  معادل یک کدگذار پیچشی<sup>۴۵</sup> که از پشته‌های پیچشی استفاده می‌کند، در نظر گرفته می‌شود. تمایزگر  $D$  یک CNN استاندارد است که یک تصویر دو بعدی را به یک مقدار  $D(\cdot)$  نگاشت می‌کند. مقدار خروجی تمایزگر  $D(\cdot)$  بیانگر احتمال این است که ورودی تمایزگر، از فضای تصاویر واقعی یعنی فضای نمونه‌های آموزشی نمونه‌برداری شده باشد و یا توسط شبکه مولد تولید شده باشد.  $D$  و  $G$  به طور همزمان از طریق بازی کمینه-بیشینه با تابع  $V(D, G)$  و معادله ۲-۹ بهینه‌سازی می‌شوند [13].

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim q(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z)))] \quad (9-2)$$

در این بازی شبکه تمایزگر آموزش می‌بیند که احتمال اختصاص نمونه‌های واقعی را بیشینه و نمونه‌های تولیدی از  $p$  با برچسب جعلی را کمینه کند. همچنین شبکه مولد  $G$  آموزش می‌بیند همزمان با حداقل

<sup>45</sup> Convolutional

کردن  $V(G) = \log(1 - D(G(z)))$  که معادل با حداکثر کردن  $V(G) = D(G(z))$  است، شبکه تمایزگر  $D$  را فریب دهد. به طور کلی در طول آموزش تقابلی، مولد در تولید تصاویر واقع بینانه و تمایزگر در شناسایی صحیح تصاویر واقعی و تولید شده بهبود می‌یابد.

## ۲-۵-۲-۲- نگاشت تصاویر جدید به فضای نهفته

وقتی آموزش تقابلی به پایان رسید، شبکه مولد یاد می‌گیرد که  $G(z) = x$  را از فضای نهفته با نمایش  $z$  به تصویر واقعی (عادی)  $x$  نگاشت کند. شبکه‌های GAN به‌طور خودکار نگاشت معکوس  $\mu(x) = z$  را انجام نمی‌دهد. فضای نهفته دارای گذار خطی است، بنابراین نمونه‌گیری از دو نقطه نزدیک بهم در فضای نهفته، دو تصویر مشابه بصری نیز ایجاد می‌کند.

با فرض اینکه تصویر  $x$  را برای بررسی داریم، هدف این است که یک نقطه  $z$  را در فضای پنهان پیدا کنیم که مطابق با تصویر  $G(z)$  باشد و از نظر بصری در آن نقطه  $G(z)$  شبیه به تصویر  $x$  باشد و در توزیع حاشیه‌ای  $\mathcal{X}$  قرار داشته باشد. میزان شباهت  $x$  و  $G(z)$  بستگی به این دارد که چه تصویری از توزیع داده  $p_g$  برای آموزش مولد استفاده می‌شود. برای پیدا کردن بهترین  $z$ ، با نمونه‌گیری تصادفی  $z_1$  از توزیع فضای نهفته  $\mathcal{Z}$  شروع می‌کنیم و آن را به شبکه مولد آموزش دیده، برای تولید تصویر  $G(z_1)$  اعمال می‌کنیم. سپس بر اساس تصویر ایجاد شده  $G(z_1)$  یک تابع اتلاف تعریف می‌کنیم، که گرادینان به روزرسانی ضرایب  $z_1$  را فراهم می‌کند و در نتیجه یک موقعیت بروز شده در فضای نهفته  $z_2$  بدست می‌آید. به عبارتی برای پیدا کردن شبیه‌ترین تصویر  $G(z_\Gamma)$ ، نقطه  $z$  در فضای نهفته  $\mathcal{Z}$  در یک فرآیند تکراری از طریق  $\gamma = 1, 2, \dots, \Gamma$  با مراحل پس‌انتشار بهینه می‌شود.

تعریف تابع اتلاف برای نگاشت از تصاویر فضای نهفته شامل دو بخش است [34]، باقی‌مانده خطا<sup>۴۶</sup> و باقی‌مانده تمایز<sup>۴۷</sup>. باقی‌مانده اتلاف شباهت بصری بین تصویر تولید شده  $G(z_\Gamma)$  و تصویر مورد بررسی  $x$  را تقویت می‌کند. باقی‌مانده تمایز، تصویر تولید شده  $G(z_\Gamma)$  را در حاشیه توزیع آموزش دیده قرار می‌دهد. بنابراین، هر دو مؤلفه GAN آموزش دیده، تمایزگر  $D$  و مولد  $G$ ، برای یافتن ضرایب  $z$  از طریق پس‌انتشار مورد استفاده قرار می‌گیرند.

<sup>46</sup> Residual Loss

<sup>47</sup> Discrimination loss

باقی مانده خطا معیار عدم شباهت بصری بین تصویر مورد بررسی  $X$  و تصویر تولید شده  $G(z_\gamma)$  در فضای تصویر اندازه گیری می کند و به صورت معادله ۱۰-۲ تعریف می شود.

$$\mathcal{L}_R(z_\gamma) = \sum |x - G(z_\gamma)| \quad (10-2)$$

با فرض یک مولد کامل  $G$  و یک نگاشت کامل فضای نهفته، برای یک مورد بررسی ایده آل، تصاویر  $X$  و  $G(z_\gamma)$  یکسان هستند. در این حالت باقی مانده خطا برابر با صفر است. باقی مانده تمایز برای شبکه تمایزگر بدین ترتیب طبق معادله ۱۱-۲ تعریف می شود.

$$\mathcal{L}_D(z_\gamma) = \sigma(D(G(z_\gamma)), \alpha) \quad (11-2)$$

در این رابطه  $G(z_\gamma)$  تصویر تولید شده توسط شبکه مولد و  $\sigma$  آنتروپی متقاطع سیگموئید<sup>۴۸</sup> است. هدف نهایی آموزش این است که خروجی شبکه تمایزگر برای تصاویر تولید برابر با ۱ باشد که به معنی این است که به ازای تصاویر تولیدی  $G(z_\gamma)$  هدف نهایی آموزش  $\alpha = 1$  تعریف می شود [13].

## ۲-۵-۳- تشخیص ناهنجاری

در طی شناسایی ناهنجاری ها در داده ی جدید، ابتدا نمونه مورد بررسی جدید  $X$  را به عنوان یک تصویر طبیعی یا غیر عادی ارزیابی می کنیم. تابع اتلافی که برای نگاشت به فضای نهفته مورد استفاده قرار می گیرد، در هر تکرار  $\gamma$  بروزرسانی می شود و سازگاری تصاویر تولید شده  $G(z_\gamma)$  با تصاویر را که در طول آموزش متخاصم مشاهده می شود ارزیابی می شود. بنابراین، این نمره ناهنجاری تناسب تصویر مورد بررسی  $X$  را با مدل تصاویر عادی بیان می کند، این معیار می تواند مستقیماً از تابع اتلاف در معادله ۱۱-۲ بدست آید. بدین ترتیب امتیاز ناهنجاری بر اساس تابع اتلاف مطابق معادله ۱۲-۲ تعریف می شود.

$$A(x) = (1 - \lambda) \cdot \mathcal{L}_R(z_\gamma) + \lambda \cdot \mathcal{L}_D(z_\gamma) \quad (12-2)$$

که در آن به ترتیب  $\mathcal{L}_R(z_\gamma)$  مقدار باقی مانده خطا و  $\mathcal{L}_D(z_\gamma)$  باقی مانده تمایزگر است. این مدل، نمره ناهنجاری بزرگی برای تصاویر غیر عادی بدست می آورد و یک نمره ناهنجاری کوچک بدین معنی است که این تصویر بسیار مشابه تصاویر هنجاریست که قبلاً در طول آموزش دیده شده است. برای تشخیص

<sup>48</sup> Sigmoid cross entropy



ناهنجاری مبتنی بر تصویر، از نمره ناهنجاری  $A(x)$  استفاده می‌شود. علاوه بر این در این جا، برای شناسایی مناطق غیر عادی در یک تصویر، از رابطه  $x_R = |x - G(z_\Gamma)|$  استفاده شده است و بخش‌هایی از تصویر که  $x_R$  بزرگ‌تری دارند به عنوان بخش ناهنجار شناسایی می‌شوند [9].

## ۲-۵-۴- مزایا و معایب

علیرغم نتایج قابل قبول این مدل در تشخیص ناهنجاری در تصاویر پزشکی، همانطور که مشاهده کردیم، برای بررسی تصویر  $x$ ، باید نقطه متناظر با آن در فضای نهفته پیدا شود به طوری که در آن نقطه  $G(z)$  از نظر بصری بیشتر شبیه به تصویر  $x$  باشد. برای پیدا کردن بهترین  $z$ ، باید با نمونه‌گیری تصادفی  $z_1$  از توزیع فضای نهفته  $Z$  شروع کنیم و آن را به شبکه مولد آموزش بدهیم و برای تولید تصویر، به  $G(z_1)$  اعمال می‌کنیم. سپس بر اساس تصویر ایجاد شده  $G(z_1)$  می‌بایست یک تابع اتلاف تعریف کنیم، که گرادینان به روزرسانی ضرایب  $z_1$  را فراهم کند و در نتیجه یک موقعیت بروز شده در فضای نهفته  $z_2$  بدست آورد. به عبارتی برای پیدا کردن شبیه‌ترین تصویر  $G(z_\Gamma)$ ، نقطه  $z$  در فضای نهفته  $Z$  در یک فرآیند تکراری از طریق  $\gamma = 1, 2, \dots, \Gamma$  با مراحل پس‌انتشار بهینه و پیدا می‌شود. این فرایند تکراری مبتنی بر تصادف است و هزینه و پیچیدگی محاسباتی زیادی به مدل تحمیل می‌کند [35]. با توجه به وجود این مشکل می‌توان مکانیزمی طراحی کرد که مدل پیشنهادی فرایند نگاشت معکوس تصاویر در حین آموزش مدل فراگیرد، در ادامه به بررسی همین مدل‌ها می‌پردازیم.

## ۲-۵-۳- مدل f-AnoGan

مدل f-AnoGan در ادامه کار قبلی و توسط همان نویسندگان در سال ۲۰۱۹ ارائه شد [35]. در مدل AnoGan از شبکه‌های عصبی عمیق کانولوشنی مولد تقابلی<sup>۴۹</sup> یا به اختصار DCGAN برای آموزش بدون نظارت شبکه مولد و تمایزگر استفاده شده است. در مدل قبلی برای شناسایی نقطه متناسب با تصویر ورودی بر اساس الگوریتم پس‌انتشار از یک فرایند مبتنی بر تکرار استفاده می‌شد. به هنگام شناسایی ناهنجاری در کاربردهای دنیای واقعی، این فرایند مبتنی بر تکرار از نظر پیچیدگی زمانی

<sup>49</sup> Deep convolutional generative adversarial network

مشکلات قبل توجهی ایجاد می‌کند. f-AnoGan فرایند مبتنی بر تکرار مورد نظر را با یادگیری یک نگاشت معکوس از فضای اصلی به فضای نهفته جایگزین می‌کند. علاوه بر این در ساختار مدل جدید از WGAN<sup>50</sup> به جای DCGAN استفاده شده است.

چهارچوب ارائه شده در این کار شامل دو گام آموزشی روی تصاویر عادی است، در گام اول شبکه مولد تقابلی آموزش می‌بیند و در گام بعدی بر اساس شبکه مولد آموزش دیده شده شبکه کدگذار آموزش می‌بیند. پس از آموزش قسمت‌های مختلف مدل بر اساس استنتاج برای هر تصویر یک امتیاز ناهنجاری محاسبه می‌شود. همانند کار قبلی شبکه مولد روی تصاویر عادی آموزش داده می‌شود و بازنمایی‌های تصاویر سالم در فضای نهفته به دست می‌آید. علاوه بر این کدگذار نیز نگاشت تصاویر به فضای نهفته را آموزش می‌بیند. در ادامه در سه بخش به بررسی شبکه‌های تمایزگر و مولد، کدگذار و نحوه امتیازدهی به داده‌ها به منظور تشخیص داده‌های ناهنجار می‌پردازیم.

## ۲-۵-۳-۱- یادگیری بدون نظارت تصاویر طبیعی

داده‌های آموزشی در این مدل به شکل  $x = x_{k,n} \in \mathcal{X}$  نمایش داده می‌شود. در این نوع از نمایش  $k = 1, 2, \dots, K$  و  $n = 1, 2, \dots, N$  است. به بیان دیگر تعداد  $N$  تصویر پزشکی داریم و از هر تصویر  $I_n$  به صورت تصادفی از مناطق مختلف تصویر تعداد  $K$  نمونه با ابعاد  $S \times S$  انتخاب می‌شود. توجه شود  $I_n \in \mathbb{R}^{u \times v}$  و  $u \gg S$  و  $v \gg S$ . برای ارزیابی مدل از داده تصویر  $y_m$  با ابعاد  $S \times S$  که از میان دادگان آزمایش انتخاب شده است، حاشیه نویسی پیکسل مربوطه به صورت  $a_m \in \{0, 1\}^{S \times S}$  و برچسب مربوطه استفاده می‌شود. این مجموعه داده آزمایش به صورت  $\langle y_m, a_m \rangle$  است و به طور همزمان شامل داده عادی و ناهنجار است.

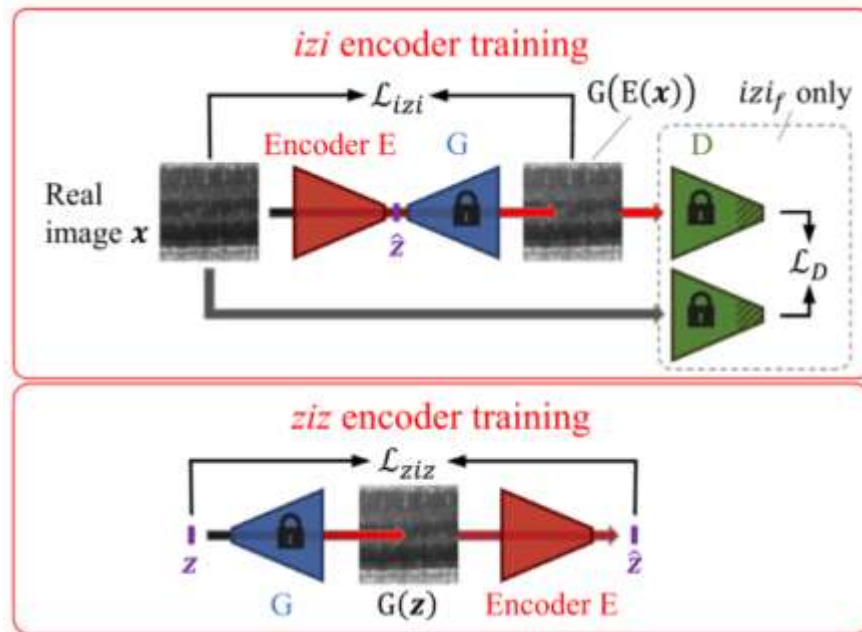
برای یادگیری تنوع موجود در تصاویر از WGAN استفاده می‌شود. این شبکه نگاشت غیرخطی از فضای نهفته  $Z$  به فضای ورودی را یاد می‌گیرد. مانند مدل‌های دیگر در این مدل نیز شبکه مولد و تمایزگر به طور همزمان بهینه می‌شوند. در ابتدا از فضای  $Z \in \mathbb{R}^d$  با ابعاد  $d$  نمونه نوین انتخاب می‌کنیم. در طول فرایند آموزش شبکه مولد تلاش می‌کند تا توزیع شبکه مولد یعنی  $p$  را تا حد امکان به توزیع داده

<sup>50</sup> Wasserstein GAN

ورودی یعنی  $q$  نزدیک کند و در نتیجه بتواند شبکه تمایزگر را به گونه‌ای فریب دهد که نتواند تشخیص دهد از توزیع داده واقعی است و یا توسط شبکه مولد ایجاد شده است. در پایان فرایند آموزش شبکه مولد توزیع داده‌های آموزشی را یاد گرفته است و شبکه تمایزگر می‌تواند تخمین بزند داده تولیدی توسط مولد تا چه اندازه به توزیع داده واقعی نزدیک است.

## ۲-۳-۵-۲- یادگیری نگاشت سریع از فضای تصویر به فضای نهفته

شبکه مولد GAN را می‌توان به صورت  $G(z) = x$  نمایش داد. این نمایش به توانایی نگاشت شبکه مولد از فضای نهفته به فضای داده ورودی اشاره دارد. در روند آموزشی شبکه GAN اولیه هیچ نگاشتی از فضای داده ورودی به فضای نهفته آموزش داده نمی‌شود. در مدل پیشنهادی این مقاله نگاشت معکوس به صورت  $E(x) = z$  نمایش داده می‌شود و این نگاشت توسط یک کدگذار آموزش دیده می‌شود. آموزش این کدگذار با دو معماری مختلف قابل پیاده‌سازی است، روش اول  $z$ -image- $z$  و روش دوم image- $z$ -image نام دارد. روش اول که به صورت خلاصه  $z \rightarrow z$  و روش دوم به طور خلاصه  $z \rightarrow \text{image} \rightarrow z$  نامیده می‌شود. در هر دو معماری از خودکدگذارهای کانولوشنی استفاده می‌شود. از کدگذار  $E$  برای نگاشت معکوس استفاده می‌شود و از شبکه مولد که در واقع یک WGAN با وزن‌های یادگرفته‌شده ثابت است، به عنوان کدگشا استفاده می‌شود. تفاوت دو روش فوق در ترتیب استفاده از کدگذار و کدگشا است. هنگام آموزش کدگذار تنها پارامترهای کدگذار بهینه می‌شوند و پارامترهای شبکه مولد ثابت هستند. معماری‌های مختلف جهت آموزش کدگذار در شکل ۲-۳ قابل مشاهده است.



شکل ۲-۳: شمای کلی روند آموزش کدگذار [35].

در آموزش کدگذار به روش izi خطای  $\mathcal{L}_{izi}$  بر اساس باقی‌مانده<sup>۵۱</sup> از تفاوت تصاویر ورودی واقعی و تصویر بازسازی شده بهینه می‌شوند کدگذار آموزش می‌بیند. در حین آموزش کدگذار به روش  $\mathcal{L}_{izi_f}$  به صورت توأم خطای  $\mathcal{L}_{izi}$  که همان خطای باقی‌مانده میان تصویر ورودی واقعی و تصویر بازسازی شده است به همراه خطای  $\mathcal{L}_D$  که خطای باقی‌مانده روی ویژگی‌های شبکه تمایزگر است، بهینه می‌شود. در آموزش کدگذار به روش  $\mathcal{L}_{ziz}$  خطای باقی‌مانده میان نمونه‌های تصادفی و موقعیت‌های موجود در فضای نهفته بهینه می‌شود.

الگوریتم ziz با معکوس کردن ترتیب کدگذار و کدگشا در ساختار معمول یک خودکدگذار به وجود می‌آید. در هنگام آموزش یک نمونه از فضای نهفته انتخاب می‌شود و با استفاده از شبکه مولد که وزن‌های آن ثابت نگاه داشته شده است به فضای داده واقعی نگاشت می‌شود و کدگذار تلاش می‌کند تا معکوس این نگاشت به فضای نهفته را یاد بگیرد، بنابراین در این روش به هیچ تصویر واقعی یا به بیان دیگر به هیچ نمونه‌ای از فضای داده واقعی نیاز نیست. در واقع در معماری ziz ساختار یک کدگذار از فضای نهفته به فضای نهفته است و در همین حال نگاشت مورد نیاز از فضای نهفته به فضای واقعی داده ورودی ثابت در نظر گرفته

<sup>51</sup> Residual

شده است. تابع هدف این آموزش به صورت خطای  $MSE^{52}$  روی نمونه اولیه  $z$  و مقدار بازسازی شده آن توسط کدگذار تعریف شده است. تابع هدف این معماری را در زیر مشاهده می‌کنید:

$$\mathcal{L}_{ziz}(z) = \frac{1}{d} \| z - E(G(z)) \|^2 \quad (13-2)$$

$d$  در معادله بیانگر ابعاد نمونه‌ها در فضای نهفته است. در این روش کدگذار بر خلاف روش  $izi$  هیچ نمونه‌ای از فضای تصاویر واقعی نمی‌بیند و این مسئله می‌تواند بر آموزش صحیح کدگذار تاثیر منفی بگذارد.

در آموزش کدگشا به روش  $izi$  از ساختار کدگذار استاندارد استفاده می‌شود، بدین صورت که در ادامه کدگذار کدگشا (شبکه مولد) قرار خواهد گرفت. در فرایند آموزش ابتدا نگاشت معکوس از فضای داده واقعی به فضای نهفته توسط کدگذار انجام می‌شود و در ادامه نگاشت از فضای نهفته به فضای داده واقعی توسط کدگشا با ضرایب ثابت صورت می‌پذیرد. ساختار این روش به صورت از فضای واقعی به فضای واقعی است. تابع هدف این روش با استفاده از خطای  $MSE$  بدین شکل پیاده‌سازی می‌شود که خطای باقی‌مانده میان تصویر واقعی و تصویر خروجی مولد کمینه می‌شود، تابع هدف مورد نظر در ادامه آمده است:

$$\mathcal{L}_{izi}(x) = \frac{1}{n} \| x - G(E(x)) \|^2 \quad (14-2)$$

در معادله ۱۴-۲  $\| \cdot \|^2$  بیانگر جمع مربعات خطا در سطح پیکسل میان دو تصویر است. داده‌های آموزشی این روش داده‌های آموزش همان WGAN یعنی داده‌های عادی است. در این روش نمی‌توان به طور مستقیم میزان دقت کدگذار را در فضای نهفته اندازه گرفت و تنها می‌توان به صورت غیرمستقیم نگاشت مربوط به فضای نهفته را به فضای داده واقعی انتقال داده و در این فضا میزان دقت را اندازه گرفت به بیان دیگر میزان دقت به صورت تصویر-تصویر محاسبه می‌شود.

در روش دیگر به نام  $izif$  از تمایزگر نیز استفاده می‌شود. در روش  $izi$  تابع هدف میزان شباهت در فضای تصویر را تحمیل می‌کند. هنگام نگاشت تصاویر جدید ممکن است با نمونه‌هایی روبرو شویم که در مرحله

<sup>52</sup> Mean squared error

آموزش به صورت تنک از فضای نهفته متناظر آن‌ها نمونه گرفته باشیم وقتی نقطه متناظر را به فضای تصویر (فضای داده ورودی) می‌بریم با تصاویر تولیدی دیگر نمی‌توان تمایزگر را متقاعد کرد. در نتیجه این اتفاق تنها کمینه کردن تفاوت تصاویر در سطح پیکسل گاهی اوقات منجر به تولید تصاویر عادی می‌شوند که واقعی به نظر نمی‌رسند ولی هنوز خطای باقی‌مانده کمی حتی برای نمونه‌های ناهنجار دارند و این مورد سبب می‌شود تا دیگر خطای باقی مانده (خطای بازسازی) در فضای داده ورودی دیگر به عنوان معیار مناسب تشخیص ناهنجاری در نظر گرفته نشود.

نویسندگان مقاله دریافتند که باقی‌مانده که خود معیار مورد نظر ما برای تشخیص ناهنجاری است در فضای ویژگی توسط تمایزگر انباشته می‌شود و این عبارت حتما باید در تابع هدف مربوط به آموزش کدگذار گنجانده شود. بنابر این آمارگان تصاویر ورودی و تصاویر خروجی محاسبه می‌شود تا با استفاده از آن‌ها تصاویر خروجی شبکه مولد شبیه تصاویر ورودی بشود و بدین ترتیب روش izif پدید آید. تابع هدف این روش به شکل زیر است:

$$\mathcal{L}_{izif}(x) = \frac{1}{n} \cdot \|x - G(E(x))\|^2 + \frac{k}{n_d} \cdot \|f(x) - f(G(E(x)))\|^2 \quad (2-15)$$

در معادله ۲-۱۵ ویژگی‌های شبکه تمایزگر که در واقع بردار ویژگی لایه‌های میانی این شبکه است با  $f(\cdot)$  و ابعاد این ویژگی با نماد  $n_d$  نمایش داده می‌شود، همچنین  $k$  عامل وزن است. اوزان شبکه تمایزگر همان اوزانی است که در آموزش WGAN یاد گرفته شده‌اند و هنگام آموزش کدگشا ثابت در نظر گرفته شده‌اند. این مدل سبب می‌شود تا هم در فضای تصویر و هم در فضای نهفته کدگذار به جهت مناسبی حرکت کند.

### ۲-۳-۳- شناسایی ناهنجاری

در مرحله آزمایش میزان انحراف تصویر اصلی از تصویر بازسازی شده به منظور تشخیص ناهنجاری اندازه‌گیری می‌شود. تمامی موارد مورد نیاز برای بازسازی تصویر و تشخیص ناهنجاری در هنگام آموزش WGAN و کدگذار یاد گرفته می‌شود. برای محاسبه امتیاز ناهنجاری مستقیماً از تعریف خطای استفاده شده در آموزش کدگذار استفاده می‌شود. امتیاز نهایی که برای تشخیص ناهنجاری در مدل f-AnoGan استفاده می‌شود به صورت زیر است.

$$\mathcal{A}(X) = \mathcal{A}_R(X) + \kappa \cdot \mathcal{A}_D(X) \quad (۱۶-۲)$$

در اینجا  $\mathcal{A}_D(X) = \frac{1}{n_d} \cdot \|f(x) - f(G(E(x)))\|^2$  و  $\mathcal{A}_R(X) = \frac{1}{n} \cdot \|X - G(E(X))\|^2$  ، همچنین  $\kappa$  عامل وزن است. عبارت مورد نظر برای نمونه‌های عادی دارای خطای کمی است و برای نمونه‌های ناهنجار دارای مقدار بزرگی است. از آنجایی که مدل تنها روی نمونه‌های عادی آموزش دیده است نمونه‌های بازسازی شده از نظر بصری شبیه تصویر ورودی هستند. توانایی بازسازی تصویر به طوری که شبیه تصویر ورودی باشد رابطه عکس دارد با درجه یا میزان تمایز ناهنجاری دارد. تصاویر عادی میزان انحراف کمی دارند در حالی که تصاویر ناهنجار که به بازسازی خود نگاشت می‌شوند میزان انحراف زیادی دارند. قدر مطلق خطای باقی‌مانده در سطح پیکسل به صورت زیر تعریف می‌شود.

$$\dot{\mathcal{A}}_R(X) = |X - G(E(X))| \quad (۱۷-۲)$$

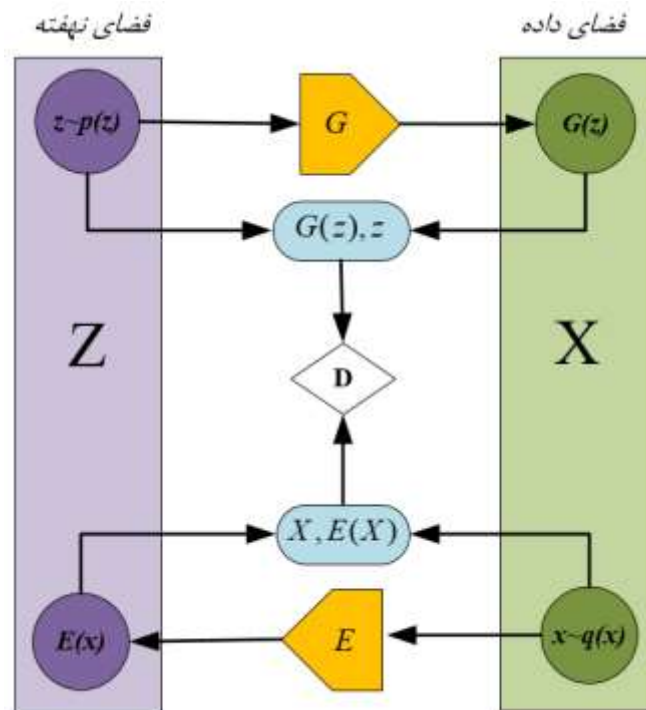
## ۲-۵-۴- مزایا و معایب

f-AnoGan مشکلات مدل قبلی (AnoGan) نظیر پیچیدگی زمانی بالا در هنگام اجرا رفع کرد. همچنین از شبکه مولد تقابلی قوی‌تری در ساختار مدل خود استفاده کرد و به نتایج قابل قبولی نیز دست یافت. در مدل قبلی از یک فرایند تصادفی و مبتنی بر تکرار برای نگاشت معکوس از فضای داده واقعی به فضای نهفته استفاده می‌شد که از نظر زمانی هزینه گزافی را به هنگام اجرا به مدل تحمیل می‌کرد. در مدل جدید f-AnoGan با استفاده از یک کدگشا پارامترهای نگاشت معکوس مورد نظر فراگرفته می‌شود [14]. علی‌رغم این موفقیت این مدل همچنان از مشکلاتی نظیر عدم استفاده توأم از هر دو فضای تصویر و نهفته و نبود چرخه پایداری در حین آموزش رنج می‌برد.

## ۲-۵-۴- مدل ALI

این شبکه در سال ۲۰۱۷ در کنفرانس ICLR معرفی شد [14]. این شبکه‌ها با هدف یادگیری نگاشت معکوس از فضای ورودی  $X$  به فضای نهفته  $Z$  تعریف شد. در این شبکه، علاوه بر شبکه مولد  $G$  که در

معماری اصلی نیز تعریف شده بود، یک کدگذار  $E^{53}$  نیز وجود دارد که از دامنه داده‌های ورودی  $X$  به دامنه ویژگی‌ها  $Z$  می‌برد. بدین ترتیب خروجی بخش مولد یک دوتایی  $^{54}$  است؛ که یکی از دامنه ویژگی‌ها و دیگری از دامنه داده‌های ورودی است. این مدل به طور همزمان شبکه مولد و شبکه استنتاج را با استفاده از یک فرآیند تقابلی به کار می‌برند. شبکه مولد، نمونه‌ها را از یک فضای نهفته آماری به فضای داده‌ها نگاشت می‌کند و شبکه استنتاج نمونه‌های آموزش را از فضای داده به فضای متغیرهای نهفته نگاشت می‌کند. به این صورت یک بازی خصمانه بین دو شبکه انجام می‌شود. در این جا شبکه تمایزگر باید یاد بگیرد تا تفاوت بین جفت ورودی فضای نهفته و فضای داده را تشخیص دهد. شبکه تمایزگر  $D$  در این جا علاوه بر تفکیک در فضای داده، در فضای ویژگی نیز تفکیک می‌کند. به این معنا که تشخیص می‌دهد دوتایی وارد شده، داده واقعی است یا ویژگی تولید شده توسط  $E$  و یا داده جعلی است که توسط  $G$  و به همراه ویژگی‌های  $Z$  درست شده، است. در تصویر زیر چارچوب کلی این الگوریتم، به نمایش درآمده است:



شکل ۲-۴: معماری شبکه ALI.

<sup>53</sup> Encoder

<sup>54</sup> Tuple



دو تابع توزیع احتمال روی  $X$  و  $Z$  در نظر بگیرید:

- $q(x, z) = q(x)q(z|x)$  تابع توزیع تعریف شده برای کدگذار  $E$
- $p(x, z) = p(z)p(x|z)$  تابع توزیع تعریف شده برای کدگشا  $G$

این دو توزیع، توابع توزیع حاشیه‌ای دارند که برای ما آشناست: توزیع حاشیه‌ای کدگذار  $q(x)$  تابع توزیع داده‌های اصلی است و توزیع حاشیه‌ای کدگشا  $p(z)$  معمولاً به عنوان یک تابع توزیع ساده مانند تابع توزیع استاندارد  $p(z) = N(0, I)$  در نظر بگیریم. بدین ترتیب روند تولید  $p(x, z)$  و  $q(x, z)$  معکوس می‌باشد.

هدف اصلی شبکه ALI مطابقت این دو توزیع است. اگر این شرط محقق شود، ما اطمینان حاصل می‌کنیم که تمام توزیع‌های حاشیه‌ای و توزیع‌های شرطی مطابقت دارد. برای دستیابی به این توابع توزیع، یک بازی تقابلی صورت می‌گیرد. جفت  $(x, z)$  از دو توزیع  $q(x, z)$  یا  $p(x, z)$  در نظر گرفته می‌شود و یک شبکه تمایزگر می‌آموزد تا بین این دو خروجی، تمایز قائل شود؛ در حالی که دو شبکه کدگشا و کدگذار می‌آموزند تا این شبکه را فریب دهند. در نهایت تابع هدفی که این بازی بر اساس آن صورت می‌گیرد به صورت زیر است:

$$\min_{G,E} \max_D V(D, G) = \mathbb{E}_{q(x,z)} [\log D(x, E(x))] + \mathbb{E}_{p(x,z)} [\log (1 - D(G(z), z))] \quad (۱۸-۲)$$

ویژگی جالب رویکردهای خصمانه این است که آن‌ها نیازی به محاسبه تابع چگالی شرطی ندارند و تنها نیاز دارند که به نحوی نمونه برداری شوند که این امکان را به وجود آورد که بتواند از پس انتشار گرادیان برای آموزش شبکه استفاده شود. در مورد شبکه ALI، این بدان معنی است که گرادیان‌ها باید از شبکه تمایزگر به شبکه‌های مولد و کدگذار انتشار یابند.

به طور دقیق‌تر شبکه تمایزگر آموزش می‌بیند که بین نمونه‌هایی که از کدگذار  $q(x, z) \sim (\hat{x}, \hat{z})$  و نمونه‌هایی که از کدگشا  $p(x, z) \sim (\hat{x}, \hat{z})$  تولید می‌شود، تمایز بگذارد. شبکه مولد و شبکه کدگذار نیز می‌آموزند که شبکه تمایزگر را فریب دهند؛ یعنی جفت  $x, z$  تولید کنند که  $p(x, z)$  از  $q(x, z)$  غیر قابل تشخیص باشد. شاید در این‌جا سوال مطرح شود که چرا باید شبکه کدگذار تلاش کند تا شبکه تمایزگر را فریب دهد و چرا باید شبکه کدگذار تلاش کند تا در معادله ۱۸-۲ خروجی تمایزگر برای داده‌های هنجار کمینه شود.

برای پاسخ به این سوال باید به این نکته توجه کرد که در حین آموزش شبکه کدگذار، داده آموزشی برای این شبکه وجود ندارد و تنها اطلاعات موجود قابل استفاده، جفت داده‌هایی است که توسط شبکه مولد ساخته شده است و این شبکه بر اساس آموزش و نگاشتی که شبکه مولد استخراج کرده است، نگاشت معکوس را یاد می‌گیرد. در واقع در حین آموزش شبکه کدگذار مفهوم داده آموزشی عوض شده و این شبکه با تلاش به فریب دادن شبکه تمایزگر آموزش می‌یابد. در این جا اساس کار یادگیری، تولید نگاشت معکوس به فضای نهفته  $Z$  برای تصاویر ورودی اصلی شبیه به جفت تولید شده توسط شبکه مولد، می‌باشد. یعنی داده آموزشی برای این شبکه تصاویر تولیدی مولد و متغیر  $Z$  آن است و داده ورودی برای آزمون و تولید خروجی، تصاویر  $X$  اولیه اصلی شبکه هستند.

در الگوریتم ۲-۲ شبکه ALI توصیف شده است. اثبات می‌شود که با فرض یک تمایزگر بهینه، شبکه مولد، واگرایی جنسن-شانون<sup>۵۵</sup> را بین  $p(x, z)$  و  $q(x, z)$  به حداقل می‌رساند.

روش استفاده شده در مدل ALI تنها راه استنتاج در شبکه‌های عصبی مولد تقابلی نیست. راه دیگر برای انجام این کار استفاده از شبکه استنتاج جلورو<sup>۵۶</sup> در ساختار GAN است. در مدل InfoGAN [۳۲] با کمینه کردن اطلاعات متقابل<sup>۵۷</sup> میان مجموعه  $C$  از فضای نهفته و  $X$  به وسیله توزیع کمکی  $Q(c | x)$  نگاشت معکوس را یاد می‌گیرد. InfoGan نیاز دارد تا تابع احتمال پسین<sup>۵۸</sup>  $Q(c | x)$  قابل تخمین و ارزیابی باشد. در مدل ALI تنها نیاز است که بتوان از شبکه استنتاج نمونه گرفت تا بدین وسیله تابع پیچیده توزیع پسین را بازنمایی کرد. عمل انجام شده در این کار مشابه این است که یک کدگذار برای بازسازی  $z$  آموزش دهیم. به عنوان مثال پیدا کردن کدگذار به طوری که  $\mathbb{E}_{z \sim p(z)} [\|z - E(G(z))\|_2^2] \approx 0$  نمونه که در جمله قبل بدان اشاره شد از نظر رویه‌ای شبیه به InfoGAN اما در این روش از یک شبکه مولد با ضرایب ثابت و همچنین تابع توزیع پسین گاوسی با واریانس قطری ثابت استفاده شده است.

<sup>55</sup> Jenesen-Shannon Divergence

<sup>56</sup> Feedforward

<sup>57</sup> Mutual information

<sup>58</sup> Posterior

روند آموزش را می‌توان به دو فاز تقسیم کرد. در فاز اول شبکه مولد تقابلی به صورت معمول آموزش می‌بیند. در فاز دوم کدگشا ثابت در نظر گرفته می‌شود و کدگذار به روش مدل ALI آموزش داده می‌شود. در این روش کدگذار و کدگشا در هنگام آموزش هیچ تعاملی با هم ندارند و در واقع کدگذار بر اساس هر چه کدگشا آموخته است آموزش می‌بیند. مشخص است اگر کدگذار و مولد با هم تعامل داشته باشند روند مدل‌سازی داده بهبود خواهد یافت.

رویه آموزش یادگیری خصمانه استنتاج

مقداردهی اولیه پارامترها  $\theta_g, \theta_d \leftarrow$

Repeat

$x^{(1)}, \dots, x^{(M)} \sim q(x)$  نمونه برداری اولیه از مجموعه داده  $M$

$z^{(1)}, \dots, z^{(M)} \sim p(z)$

$\hat{x}^{(i)} \sim q(z|x=x^{(i)})$ ,  $i = 1, \dots, M$  انتخاب شرطی

$\hat{x}^{(j)} \sim q(z|x=x^{(j)})$ ,  $j = 1, \dots, M$

$\rho_{(q)}^{(i)} \leftarrow D(x^{(i)}, \hat{x}^{(i)})$ ,  $i = 1, \dots, M$  محاسبه پیش‌بینی تمایزگر

$\rho_{(p)}^{(j)} \leftarrow D(\hat{x}^{(j)}, z^{(j)})$ ,  $j = 1, \dots, M$

$\mathcal{L}_d \leftarrow -\frac{1}{M} \sum_{i=1}^M \log(\rho_q^{(i)}) - \frac{1}{M} \sum_{j=1}^M \log(1 - \rho_q^{(j)})$  محاسبه تلفات تمایزگر

$\mathcal{L}_g \leftarrow -\frac{1}{M} \sum_{i=1}^M \log(1 - \rho_q^{(i)}) - \frac{1}{M} \sum_{j=1}^M \log(\rho_q^{(j)})$  محاسبه تلفات مولد

$\theta_d \leftarrow \theta_d - \nabla_{\theta_d} \mathcal{L}_d$  بروزرسانی گرادیان تمایزگر شبکه

$\theta_g \leftarrow \theta_g - \nabla_{\theta_g} \mathcal{L}_g$  بروزرسانی گرادیان مولد شبکه

Until

الگوریتم ۲-۲: رویه آموزش یادگیری خصمانه استنتاج

## ۲-۵-۴-۱- مقایسه مدل‌های ALI و GAN

شبکه ALI شباهت زیادی به شبکه GAN دارد، اما دو تفاوت اساسی با آن دارد:

۱- بخش مولد دارای دو مؤلفه است: بخش کدگذار،  $E(x)$  که نمونه‌های داده  $x$  را به فضای  $z$

نگاشت می‌کند و بخش مولد  $G(z)$  که نمونه‌ها را از  $p(z)$  (منبع منبع نویز) به فضای ورودی نگاشت می‌کند.

۲- بخش تمایزگر به منظور تمایز بین جفت  $(\hat{x} = G(z), z)$  و  $(x, \hat{z} = E(x))$ ، آموزش دیده می‌شود.

## ۲-۵-۴-۲- رویکردهای جایگزین برای استنتاج در GAN

روش استفاده شده در مدل ALI تنها راه استنتاج در شبکه‌های عصبی مولد تقابلی نیست. راه دیگر برای انجام این کار استفاده از شبکه استنتاج جلورو<sup>۵۹</sup> در ساختار GAN است. در مدل InfoGAN<sup>۶۰</sup> [۳۲] با کمینه کردن اطلاعات متقابل<sup>۶۱</sup> میان مجموعه  $C$  از فضای نهفته و  $X$  به وسیله توزیع کمکی  $Q(c | x)$  نگاشت معکوس را یاد می‌گیرد. InfoGan نیاز دارد تا تابع احتمال پسین<sup>۶۲</sup>  $Q(c | x)$  قابل تخمین و ارزیابی باشد. در مدل ALI تنها نیاز است که بتوان از شبکه استنتاج نمونه گرفت تا بدین وسیله تابع پیچیده توزیع پسین را بازنمایی کرد. عمل انجام شده در این کار مشابه این است که یک کدگذار برای بازسازی  $z$  آموزش دهیم. به عنوان مثال پیدا کردن کدگذار به طوری که  $\mathbb{E}_{z \sim p(z)} [\|z - G_z(G_x(z))\|_2^2] \approx 0$  نمونه که در جمله قبل بدان اشاره شد از نظر رویه‌ای شبیه به InfoGAN اما در این روش از یک شبکه مولد با ضرایب ثابت و همچنین تابع توزیع پسین گاوسی با واریانس قطری ثابت استفاده شده است.

روند آموزش را می‌توان به دو فاز تقسیم کرد. در فاز اول شبکه مولد تقابلی به صورت معمول آموزش می‌بیند. در فاز دوم کدگشا ثابت در نظر گرفته می‌شود و کدگذار به روش مدل ALI آموزش داده می‌شود. در این روش کدگذار و کدگشا در هنگام آموزش هیچ تعاملی با هم ندارند و در واقع کدگذار بر اساس هر چه کدگشا آموخته است آموزش می‌بیند. مشخص است اگر کدگذار و کدگشا با هم تعامل داشته باشند روند مدل‌سازی داده بهبود خواهد یافت.

## ۲-۵-۴-۳- مزایا و معایب

اگرچه وجود مکانیزم مناسب جهت نگاشت معکوس از فضای داده واقعی به فضای نهفته سبب موفقیت و عملکرد قابل قبول مدل ALI شده است اما در این مدل هیچ سازوکاری جهت کنترل میزان شباهت تصویر اصلی و تصویر بازسازی شده در آن تعبیه نشده است و به همین جهت در برخی موارد هیچ شباهتی میان

<sup>59</sup> Feedforward

<sup>60</sup> Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets

<sup>61</sup> Mutual information

<sup>62</sup> Posterior

داده اصلی و داده بازسازی شده وجود ندارد. به بیان دیگر پس از یافتن نقطه متناسب با داده مورد نظر در فضای داده واقعی در فضای نهفته، نقطه مورد نظر را به شبکه مولد می‌دهیم و خروجی حاصل را با ورودی ابتدایی مقایسه می‌کنیم. انتظار می‌رود داده ورودی و خروجی در این چرخه شباهت زیادی داشته باشند ولی همانطور که گفته شد در برخی از نمونه‌ها شباهتی میان این دو تصویر وجود ندارد.

## ۲-۵-۵- مدل EGBAD

همانطور که گفته شد شبکه‌های عصبی تقابلی قادرند توزیع‌های پیچیده دنیای واقعی با ابعاد بالا را مدل کنند و همین امر سبب می‌شود تا بتوان از این شبکه‌ها در زمینه تشخیص ناهنجاری نیز کرد. با توجه به تعداد کارهای انگشت شمار در این زمینه، مدل EGBAD را می‌توان از اولین کارها در تشخیص ناهنجاری با استفاده از شبکه‌های عصبی تقابلی به شمار آورد [۳۶].

مدل ارائه شده بر اساس شبکه عصبی تقابلی دوطرفه به اختصار BiGAN<sup>۶۳</sup> بنا نهاده شده است. وظیفه نگاشت معکوس از فضای داده ورودی به فضای نهفته نیز بر عهده کدگذار است. کدگذار، شبکه مولد و تمایزگر در اینجا به طور همزمان آموزش می‌بینند و وجود بلوک کدگذار سبب کاهش هزینه‌های محاسباتی در گام آزمایش می‌شود. برخلاف ساختار استاندارد GAN که در آن تمایزگر تنها تصویر واقعی و تصویر تولیدی شبکه مولد را ورودی می‌گیرد، بازنمایی این تصاویر در فضای نهفته هم به عنوان ورودی به شبکه تمایزگر داده می‌شود.

استراتژی مورد استفاده در گام آموزش مدل مشابه شبکه ALI است که در فصل قبل به تفصیل به توضیح آن پرداختیم. همانطور که در قسمت قبل بررسی شد در این استراتژی آموزشی تاکید بر آن است که شبکه مولد و کدگذار به طور توانمند آموزش داده شوند. تابع هزینه در هنگام آموزش مطابق معادله ۲-۱۹ بهینه می‌شود.

$$V(D, E, G) = \mathbb{E}_{x \sim p_X} [\mathbb{E}_{z \sim p_E(\cdot|x)} [\log D(x, z)]] + \mathbb{E}_{z \sim p_Z} [\mathbb{E}_{x \sim p_G(\cdot|z)} [1 - \log D(x, z)]] \quad (2-19)$$

<sup>63</sup> Bidirectional Generative Adversarial Model

در معادله ۲-۱۹  $p_X(x)$  بیانگر تابع توزیع داده ورودی است،  $p_Z(z)$  بیانگر توزیع نمونه‌ها در فضای نهفته است،  $p_E(z|x)$  و  $p_G(x|z)$  به ترتیب بیانگر تابع توزیع کدگذار و مولد هستند. پس از آموزش مدل نوبت به تعریف معیاری می‌رسد که به وسیله آن بتوان نمونه‌های ناهنجار را تشخیص داد. بیان ریاضی این معیار در ادامه آمده است.

$$A(x) = \alpha L_G(x) + (1 - \alpha) L_D(x) \quad (2-20)$$

در معادله ۲-۲۰  $L_G(x) = \|x - G(E(x))\|_1$  و  $L_D(x)$  را می‌توان به دو روش تعریف کرد. در روش اول تابع خطای آنروپی<sup>۶۴</sup> متقابل استفاده می‌شود. ورودی این تابع در واقع خروجی تمایزگر است و خود شامل خروجی کدگذار و نمونه متناظر با آن است و به صورت  $\sigma(D(x, E(x)), 1)$  تعریف می‌شود. در اینجا خروجی تمایزگر میزان اطمینان آن در قبال واقعی بودن نمونه ورودی است. در روش دوم تعریف  $L_D(x)$  از خطای تطبیق ویژگی<sup>۶۵</sup> استفاده می‌شود. تابع امتیاز این روش به صورت  $L_D(x) = \|f_D(x, E(x)) - f_D(G(E(x)), E(x))\|_1$  تعریف می‌شود.  $f_D$  در واقع لاجیت‌های تمایزگر است و بیان می‌کند ویژگی‌های تصویر بازسازی شده تا چه اندازه شبیه ویژگی‌های تصویر واقعی هستند. توجه شود هر چه مقدار امتیاز محاسبه شده بالاتر باشد نمونه مورد نظر با احتمال بیشتری نمونه ناهنجار است.

## ۲-۵-۵-۱- مزایا و معایب

مدل مورد بحث در این قسمت ثابت کرد می‌توان با تعریف تابع امتیاز مناسب و استفاده از ساختارهای تقابلی به روز ناهنجاری را داده‌های دنیای واقعی شناسایی کرد. در این مدل هیچ نظارتی بر میزان تشابه داده ورودی و داده بازسازی وجود ندارد در حالی که انتظار می‌رود داده ورودی و داده بازسازی شده توسط شبکه مولد برای نمونه‌های هنجار یکسان باشد. در ادامه به بررسی این مشکل و روش ارائه شده برای حل آن پرداخته می‌شود.

<sup>64</sup> Cross entropy

<sup>65</sup> Feature matching loss

## ۲-۵-۶- مدل ALICE

در حالت استاندارد شبکه GAN تنها نگاشت یک طرفه از فضای نهفته به فضای داده بدست می‌آورد، یعنی فاقد مکانیسم معکوس (از فضای داده به فضای نهفته) است و این امر مانع می‌شود که این شبکه‌ها قادر به استنباط باشند. توانایی محاسبه تابع توزیع متغیر نهفته شرطی ممکن است برای تفسیر داده‌ها و برای برنامه‌های پایین دستی (به عنوان مثال، طبقه‌بندی متغیر نهفته) مهم باشد.

تلاش‌های زیادی برای یادگیری همزمان یک مدل دو طرفه کارآمد برای تولید نمونه‌هایی با کیفیت بالا برای هر دو فضای نهفته و داده صورت گرفته است. در میان این طرح‌ها، یکی از طرح‌ها که به موفقیت چشم‌گیری دست یافته است، شبکه یادگیر استنباط خصمانه ALI است [14]. همان‌طور که شرح داده شد، در این مدل در یک چارچوب شبکه مولد متخاصم، شبکه تمایزگر می‌آموزد تا تفاوت بین دو توزیع توأمان را تشخیص دهد.

همان‌طور که ذکر شد، با این که شبکه ALI یک رویکرد جالب و خلاقانه است، اما یک ایراد اساسی دارد؛ این که بازسازی‌های صورت گرفته از داده‌ها در بعضی موارد حتی به داده‌های اصلی شبیه هم نیستند. دلیل این امر این است که شبکه ALI تنها به دنبال مطابقت دو توزیع توأمان است، اما همبستگی بین دو متغیر تصادفی شرطی در هر یک از این توابع مشخص و اعمال نمی‌شود. در نتیجه حاصل، راه‌حل‌هایی می‌شود که هدف ALI را برآورده سازند اما در بازسازی داده‌های مشاهده شده با مشکل روبرو هستند. این شبکه هم‌چنین مشکلاتی در کشف رابطه صحیح جفت‌ها در زمان تغییر دامنه دارد [۳۷].

## ۲-۵-۶-۱- یادگیری تقابلی با اندازه‌گیری اطلاعات

به یاد داریم که تابع هدف در شبکه ALI به صورت معادله ۲-۲۰ بود:

$$\min_{G,E} \max_D V(D,G) = \mathbb{E}_{q(x,z)} [\log D(x, E(x))] + \mathbb{E}_{p(x,z)} [\log (1 - D(G(z), z))] \quad (2-20)$$

نقطه تعادل این معادله هنگامی است که  $q(x, z) = p(x, z)$  باشد. ارتباط بین متغیرهای تصادفی  $x$  و  $z$  توسط ALI محدود و مقید نشده است. در نتیجه، این امکان وجود دارد که توزیع همسان  $p(x, z) = q(x, z)$  برای یک کاربرد خاص نامطلوب باشد. در واقع بسیاری از کاربردها به ثبات چرخه و وجود یک نگاشت معنی‌دار دو طرفه بین دامنه‌ها احتیاج دارند.

جهت مقابله با مشکل توزیع‌های نامطلوب اما برابر، بر روی راه‌حل‌های شبکه ALI باید محدودیتی بر روی توزیع‌های  $q(x, z)$  و  $p(x, z)$  اعمال شود. این کار با کنترل "عدم قطعیت" بین جفت متغیرهای تصادفی، یعنی  $x$  و  $z$  با استفاده از آنتروپی‌های شرطی انجام می‌شود.

## ۲-۵-۶-۲- آنتروپی شرطی<sup>۶۶</sup>

آنتروپی شرطی یک معیار نظریه اطلاعاتی است که عدم قطعیت متغیر تصادفی  $x$  را به شرط متغیر  $z$  با کمک توزیع توامان  $\pi(x, z)$  تعیین می‌کند:

$$H_{\pi}(x|z) \cong -E_{\pi}(x, z)[\log \pi(x|z)] \quad (21-2)$$

$$H_{\pi}(z|x) \cong -E_{\pi}(x, z)[\log \pi(z|x)]$$

عدم قطعیت متغیر  $x$  به شرط متغیر  $z$  با  $H_{\pi}(x|z)$  مرتبط است. در حقیقت، اگر  $H_{\pi}(x|z) = 0$  باشد در این صورت  $x$  به طور قطعی به  $z$  وابسته می‌باشد. به کمک کنترل میزان عدم قطعیت  $q(z|x)$  و  $p(x|z)$  می‌توان راه حل‌های ALI را در توزیع‌های توامانی که نگاشت آن‌ها منجر به نتایج بهتری می‌شود، محدود کرد. در نهایت با افزودن یک عامل تنظیم کننده آنتروپی شرطی، به تابع هدف زیر دست می‌یابیم:

$$V_{ALICE}(D_{xz}, E, G) = V(D_{xz}, E, G) + V_{CE}(E, G) \quad (22-2)$$

$V_{CE}(E, G)$  وابسته به متغیرهای تصادفی توزیع‌های توامان است. در حالت ایده آل، پس از شناسایی تمام نقاط تعادل تابع هدف ALI، می‌توان با محاسبه آنتروپی شرطی آن‌ها، راه‌حل مطلوب را انتخاب کنیم. با این حال، در عمل این راه غیرقابل استفاده است، زیرا ما از قبل به نقاط تعادل دسترسی نداریم. در ادامه یک راه‌حل برای محاسبه آنتروپی شرطی ارائه می‌شود.

<sup>66</sup> Conditional entropy



## ۲-۵-۶-۳- فرایند یادگیری

در نبود تابع توزیع احتمال صریح که برای محاسبه آنتروپی شرطی مورد نیاز است، می توان حدود آنتروپی شرطی را با استفاده از معیار ثبات چرخه<sup>۶۷</sup> محدود کرد. در این جا برای بازسازی  $\hat{x}$  به طریق زیر عمل می شود:

$$\hat{x} \sim p(\hat{x} | z), z \sim q(z | x), x \sim q(x) \quad (23-2)$$

به کمک روال تولید بالا، تلاش می شود تا  $\hat{x}$  با احتمال بالایی شبیه  $x$  اصلی باشد. اثبات می شود که به کمک این روال تولید  $\hat{x}$  ها، حد بالای آنتروپی شرطی  $V_{CE}(E, G)$  می باشد.

نکته حائز اهمیت این است که می توان عامل تنظیم آنتروپی را به تابع هدف شبکه ALI، بدون اعمال تغییرات اضافی دیگری، در روال آموزش این شبکه اضافه کرد. بدین ترتیب تابع بهینه سازی برای شبکه ALICE به صورت ۲-۲۴ خواهد بود.

$$\begin{aligned} \min_{E, G} \max_{D_{xz}, D_{xx}} V_{ALICE} \\ = V_{ALI} + E_{x \sim q(x)} [\log D_{xx}(x, x) + \log 1 \\ - D_{xx}(x, G(E(x)))] \end{aligned} \quad (24-2)$$

ویژگی پایداری چرخش در مقالات پیش از نیز وجود داشته است این ویژگی در این مقالات به کمک نرم درجه<sup>۶۸</sup> یک و دو و داده های واقعی مانند تصاویر محاسبه شده است. وجود تابع اتلاف بر اساس نرم درجه ۲ مبتنی بر پیکسل، سبب می شود که نمونه های خروجی این شبکه ها تصاویر تاری باشند. به همین علت در این شبکه از یک شبکه تمایزگر که اختلاف بین  $x$  ها و  $\hat{x}$  های بازسازی شده را اندازه گیری می کند، استفاده شده است.

## ۲-۵-۶-۴- مزایا و معایب

همانطور که بررسی شد در این کار با استفاده از کدگذار تلاش شد تا شرط پایداری در شبکه های عصبی مولد برقرار شود به بیان دیگر اگر تصویر ورودی را به کدگذار بدهیم و نقطه متناظر در فضای نهفته را

<sup>67</sup> Cycle Consistency<sup>68</sup> L-norm

بدست آوریم و سپس آن را به عنوان ورودی به شبکه مولد بدهیم انتظار داریم نتیجه نهایی شبیه به تصویر اولیه باشد. علی‌رغم عملکرد مناسب این مدل هنوز توزیع‌های توامی وجود دارد که از آن‌ها استفاده نشده است و همین امر سبب می‌شود تا از تمامی اطلاعات موجود استفاده نشود.

## ۲-۵-۷- مدل RCGAN<sup>۶۹</sup>

این مقاله [38] با تمرکز بر ارائه ساختاری مبتنی بر شبکه‌های عصبی تقابلی به منظور هر چه بهتر کردن تشخیص نمونه‌های ناهنجار به وسیله پوشش تمام فضای نهفته و فضای داده ورودی، در سال ۲۰۲۰ در ادامه مقاله ALICE ارائه شد. اساس این کار بر پایه تعریف تابع جریمه، بیان جدیدی از تابع هزینه و همچنین استفاده نوآورانه از تمایزگر در مسئله تشخیص ناهنجاری است. عملکرد مناسب این مدل در خلال بررسی نتایج آن قابل مشاهده است. در بخش‌های بعدی به بررسی جزئیات بیشتر این مدل می‌پردازیم.

## ۲-۵-۷-۱- منظم‌سازی شبکه مولد و تمایزگر

بیشتر مدل‌هایی که اخیراً به تشخیص ناهنجاری به وسیله شبکه‌های عصبی تقابلی پرداخته‌اند مبتنی بر شبکه‌های عصبی دو طرفه هستند، اگرچه این مدل‌ها قادرند تا تصاویر هنجار را با امتیاز پایین ناهنجاری بازسازی کنند اما هیچ ضمانتی در اختصاص نمره بالای ناهنجاری به نمونه‌های ناهنجار وجود ندارد.

برای رفع این محدودیت‌ها و قادر ساختن شبکه‌های عصبی تقابلی برای تشخیص نمونه‌های ناهنجار از هنجار در این کار توزیع جریمه  $t(x)$  به گونه‌ای تعریف می‌شود که  $x \sim t(x)$  باشد، نمونه‌های تولیدی از این توزیع باید به عنوان داده تقلبی توسط تمایزگر شناخته شود. تابع هدف مدل پیشنهادی مطابق معادله ۲-۲۵ به صورت زیر است.

$$\begin{aligned} \min_{E,G} \max_{D_{xz}} V_{\text{ano}}(D_{xz}, G, E) = & \mathbb{E}_{x \sim q(x)} [\log D_{xz}(x, E(x))] \\ & + \mathbb{E}_{z \sim p(z)} [\log(1 - D_{xz}(G(z), z))] \\ & + \mathbb{E}_{x \sim t(x)} [\log(1 - D_{xz}(x, E(x)))] \end{aligned} \quad (2-25)$$

<sup>69</sup> Regularized Cycle onsistent Generative Adversarial Network for anomaly detection

تابع توزیع  $t(x)$  یک تابع چگالی احتمال مانند توزیع گاوسی انتخاب می‌شود. توزیع گاوسی کاربردهای زیادی در زمینه آموزش تقابلی دارد. در انتهای آموزش، شبکه‌های مولد و تمایزگر متمایل به توزیع داده‌های هنجار خواهد بود. در مدل پیشنهادی این مقاله نمونه  $G(E(x))$  که در آن  $x$  نمونه ناهنجار است، نزدیک به توزیع داده هنجار خواهد بود و بنابراین بازسازی  $x$  با خود  $x$  فاصله خواهد داشت. این فاصله سبب می‌شود تا تشخیص نمونه ناهنجار آسان شود. همانطور که در شکل ۳-۳ مشخص است، نمونه‌های ناهنجار به گونه‌ای در مدل هدایت می‌شوند که بازسازی آن‌ها به سمت داده هنجار صفر متمایل شود. چنین رویکردی سبب می‌شود تا میان داده‌های ناهنجار و بازسازی آن‌ها تفاوت زیادی حاصل شود و در نتیجه تشخیص نمونه‌های ناهنجار گارانتی شود. قابل توجه است که الگوریتم ارائه شده سبب می‌شود تا نقاطی که  $q(x)$  در آن‌ها کوچک است پوشش داده شود و این تنظیم مستقل از ارتباط میان  $t(x)$  و توزیع نمونه‌های ناهنجار انجام می‌شود.

## ۲-۵-۷-۲- پایداری چرخه

به منظور ارضای شرط پایداری چرخه تمایزگر  $D_{xx}$  به معماری مورد نظر اضافه شده است. تابع هدف بخش چرخه پایداری همانند معادله ۲-۲۶ است.

$$\min_{E,G} \max_{D_{xx}} V_{\text{cycle}}(D_{xx}, G, E) = \mathbb{E}_{x \sim q(x)} [\log D_{xx}(x, x)] + \mathbb{E}_{x \sim q(x)} [\log(1 - D_{xx}(x, \hat{x}))] \quad (2-26)$$

در معادله ۲-۲۶  $\hat{x} = G(E(x))$  بازسازی داده ورودی  $x$  است. تابع هدف کامل ارائه شده در این کار در نهایت مطابق معادله ۲-۲۷ به شکل زیر است.

$$\min_{E,G} \max_{D_{xz}, D_{xx}} V_{\text{ano}}(D_{xz}, G, E) + V_{\text{cycle}}(D_{xx}, G, E) \quad (2-27)$$

پس از آموزش مدل نوبت به محاسبه امتیاز ناهنجاری می‌رسد، در مدل پیشنهادی این کار تابع امتیاز پیشنهادی مطابق معادله زیر است.

$$A(x) = 1 - D_{xx}(x, G(E(x))) \quad (2-28)$$

امتیاز ناهنجاری  $A(x)$  بیانگر میزان کیفیت بازسازی  $x$  است. الگوریتم ارائه شده در این کار مدل را مجبور به تولید خطای بزرگ برای نمونه‌های ناهنجار می‌کند در حالی که تابع هدف چرخه پلیداری مدل مجبور به تولید بازسازی مناسب برای نمونه‌های هنجار می‌کند. این اختلاف امتیاز میان نمونه‌های هنجار و ناهنجار معیار مناسبی برای تشخیص نمونه‌های ناهنجار است.

## ۲-۵-۸- مدل ALAD

در این بخش، یک روش تشخیص ناهنجاری مبتنی بر شبکه مولد متخاصم را بررسی می‌کنیم که در زمان آزمون بسیار کارآمد است. در این روش به طور هم‌زمان یک شبکه کدگذار را در حین آموزش فرا می‌گیرد و بدین ترتیب استنتاج سریع‌تر و کارآمدتر را در زمان آزمون امکان‌پذیر می‌کند. علاوه بر این در شبکه معرفی‌شده، تکنیک‌هایی که اخیراً برای بهبود بیشتر شبکه کدگذار و تثبیت آموزش شبکه مولد متخاصم ترکیب‌شده و نشان داده شده که این تکنیک‌ها عملکرد و کارایی را در کاربرد تشخیص ناهنجاری بهبود می‌بخشند. آزمایشات روی طیف وسیعی از داده‌های جدولی و تصویری، کارایی و اثربخشی این رویکرد را در عمل نشان می‌دهد [31].

همان‌طور که پیش از این گفته شد، شبکه‌های GAN استاندارد از نمونه‌گیری کارآمد پشتیبانی می‌کنند و روش‌های مختلفی وجود دارد که می‌تواند آن‌ها را برای تشخیص ناهنجاری تطبیق دهد. به عنوان مثال، برای یک نقطه داده  $x$ ، می‌توان از نمونه‌گیری استفاده کرد تا احتمال ناهنجار بودن  $x$  را تخمین زد. تخمین دقیق احتمال به تعداد زیادی نمونه نیاز دارد و در نتیجه محاسبه احتمال، بار محاسباتی سنگینی دارد.

روش دیگر معکوس کردن<sup>۷۰</sup> شبکه مولد برای یافتن متغیرهای نهفته  $z$  است که به معنای به حداقل رساندن خطای بازسازی با تابع هدف گرادیان نزولی تصادفی می‌باشد. این روش هم‌چنین از نظر محاسباتی بسیار پرهزینه است زیرا هر محاسبه گرادیان نیاز به یک پس انتشار از طریق شبکه مولد دارد.

به واسطه بهره‌وری محاسباتی بالا و قابلیت مدل‌سازی داده‌های ابعاد بالا، از شبکه‌های مولد متقابلی به همراه یک شبکه کدگذار  $E$  (که نمونه‌ها را از فضای داده  $x$  به فضای نهفته  $z$  نگاشت می‌کند) استفاده

<sup>70</sup> Invert

می‌شود. نمایش نهفته هر نمونه از فضای داده در چنین مدل‌هایی صرفاً با عبور از شبکه کدگذار انجام می‌شود. همچنین این مدل پیشرفت‌های اخیر که برای بهبود شبکه کدگذار صورت گرفته مانند افزودن یک شبکه تمایزگر برای بهبود سازگاری چرخه  $x \approx G(E(x))$  را شامل می‌شود.

پیش از این توضیح داده شد که شبکه ALI توزیع توامان داده‌ها را به همراه یک شبکه کدگذار مدل می‌کند. این مدل یک شبکه تمایزگر  $D_{xz}$  دارد که  $x$  و  $z$  را به عنوان ورودی می‌گیرد و بررسی می‌کند که این جفت ورودی از کدام منبع – شبکه مولد و یا شبکه کدگذار – تولید شده است.

با این‌که به لحاظ نظری توزیع توامان شبکه مولد و شبکه کدگذار به یک نقطه میل می‌کند، اما در عمل اغلب نتیجه یکسان نیست و لزوماً به یک نقطه همگرا نمی‌شوند و این پدیده سبب نقض پایداری چرخه می‌شود. نبود پایداری چرخه به این معناست که  $x \approx G(E(x))$  باشد. این مشکل برای روش‌های تشخیص ناهنجاری مبتنی بر بازسازی چالش‌های جدی ایجاد می‌کند. برای حل این مشکل، چارچوب ALICE پیشنهاد می‌کند که آنتروپی شرطی را با افزودن یک شبکه تمایزگر بین متغیر  $x$  و بازسازی آن به روش تقابلی برای سازگاری چرخه تقریب بزنیم.

## ۲-۸-۵-۱ تابع هزینه

برای تثبیت آموزش در مدل پایه ALICE، توزیع‌های شرطی را با اضافه کردن یک قید آنتروپی شرطی دیگر تنظیم می‌کنیم و سپس عملیات نرمال‌سازی طیفی را انجام می‌دهیم.

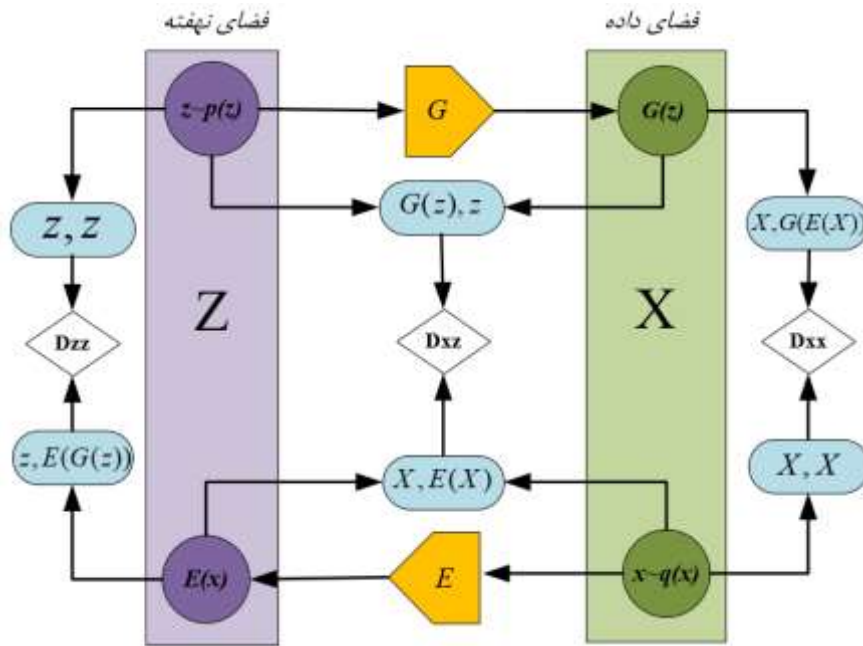
توضیح دقیق‌تر این‌که، در این‌جا فضای نهفته شرطی  $H^\pi(z|x) = -E_{\pi(x,z)}[\log \pi(z|x)]$  را با یک شبکه تمایزگر مخالف دیگر  $D_{zz}$  با نقطه تعادل مشترک تنظیم می‌کنیم که، مطابق افزودن معادله ۲-۲۹ به تابع هزینه در چارچوب ALICE می‌باشد.

$$V(D_{zz}, G, E) = V_{ALICE} + \mathbb{E}_{z \sim p(z)} [\log(D_{zz}(z, z))] + \mathbb{E}_{z \sim p(z)} [\log(1 - D_{zz}(z, G(E(z))))] \quad (2-29)$$

با کنار هم قرار دادن تمامی این اجزا، در نهایت تابع هزینه این شبکه مطابق معادله ۲-۳۰ خواهد بود. شبکه ALAD تلاش می‌کند تا نقطه تعادل این مسئله را آموزش ببیند.

$$\begin{aligned}
& \min_{G,E} \max_{D_{xz}, D_{xx}, D_{zz}} V_{ALAD}(D_{xz}, D_{xx}, D_{zz}, E, G) \\
& = \mathbb{E}_{z \sim p(z), x \sim q(x)} [\log(D_{xz}(x, E(x))) \\
& + \log(1 - D_{xz}(G(z), z)) + \mathbb{E}_{x \sim q(x)} [\log(D_{xx}(x, x)) \\
& + \log(1 - D_{xx}(x, G(E(x)))) + \mathbb{E}_{z \sim p(z)} [\log(D_{zz}(z, z)) \\
& + \log(1 - D_{zz}(z, E(G(z))))]
\end{aligned} \tag{۳۰-۲}$$

در نهایت معماری کلی شبکه ALAD به صورت شکل ۲-۶ خواهد بود.



شکل ۲-۵: شمای کلی شبکه ALAD [۳۲].

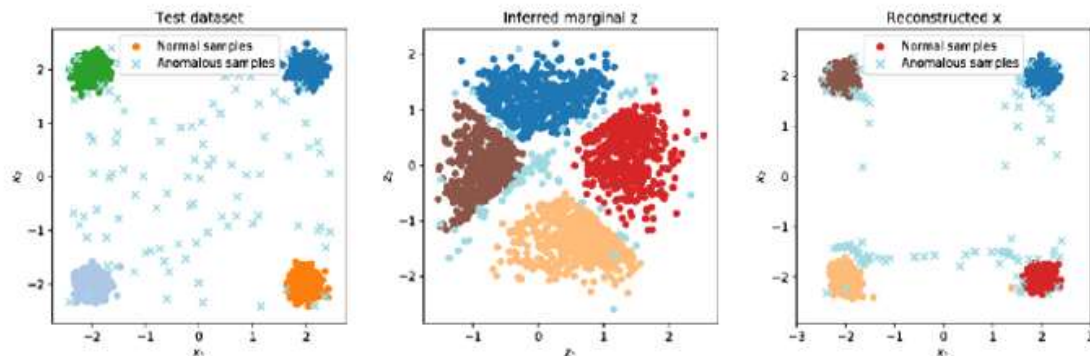
اضافه کردن مرحله نرمال سازی طیفی با انگیزه یادشده مقاله [39] می باشد. در این مقاله نشان داده شده که با افزودن قیود لپسچیتز<sup>۷۱</sup> به تمایزگر شبکه GAN، فاز آموزش تثبیت خواهد شد. در عمل نشان داده شده که با تنظیم مجدد پارامترهای وزن، بهبود بسیار خوبی روی عملکرد شبکه خواهیم داشت. بدین صورت که بزرگترین مقادیر ویژه ماتریس وزن را در شبکه تمایزگر ثابت نگه داریم. این روش از نظر محاسباتی کارآمد است و همچنین آموزش را تثبیت می کند. با آزمایش های صورت گرفته نشان داده شد

<sup>71</sup> Lipschitz Constraints

که افزودن این قیود، نه تنها برای شبکه تمایزگر، بلکه برای شبکه کدگذار نیز سودمند است. قابل توجه است که مدل ALICE شامل این مرحله نمی باشد.

## ۲-۸-۵-۲ تشخیص ناهنجاری

شبکه ALAD یک روش تشخیص ناهنجاری مبتنی بر بازسازی است و بدین صورت عمل می کند که فاصله نمونه از بازسازی را توسط شبکه GAN ارزیابی می کند. نمونه های عادی باید به طور دقیق بازسازی شوند در حالی که نمونه های ناهنجار احتمالاً به طور ضعیف تری بازسازی می شوند. نحوه تشخیص ناهنجاری در شکل ۶-۲ نشان شده است.



شکل ۶-۲: نمونه ای از خروجی شبکه ALAD به همراه داده های ناهنجار [۳۲].

در شکل ۶-۲ ضربدرها نمونه های ناهنجار و دایره های رنگی نمونه های عادی هستند. همان طور که دیده می شود شبکه ALAD تا حدی توانسته برای داده های ناهنجار بازسازی ضعیف داشته باشد

مؤلفه کلیدی دیگر ALAD نمره ناهنجاری است که فاصله بین نمونه های اصلی و بازسازی آن ها را اندازه گیری می کند. انتخاب اولیه ای که به ذهن می رسد، فاصله اقلیدسی بین نمونه های اصلی و بازسازی آن ها در فضای داده است. اما، این معیار ممکن است معیار مطمئنی برای اندازه گیری تشابه نباشد. به عنوان مثال، این معیار در مورد تصاویر می تواند بسیار پرخطا باشد؛ زیرا تصاویر با ویژگی های تصویری مشابه الزاماً از نظر فاصله اقلیدسی نزدیک به یکدیگر نیستند. معیار تعریف شده در این روش از فاصله بین نمونه ها در فضای ویژگی های تمایزگر  $D_{xx}$  محاسبه می شود، که توسط لایه قبل از لاجیت تعریف شده است. از این ویژگی ها هم چنین به عنوان کدهای CNN یاد می شود. به طور دقیق تر می توان گفت با آموزش یک مدل

برای داده‌های عادی و محاسبه  $E$ ،  $G$ ،  $D_{xx}$ ،  $D_{xz}$  و  $D_{zz}$  یک تابع نمره‌دهی را بر اساس خطای بازسازی نرم ۱ مطابق معادله ۲-۳۱ تعریف می‌شود. در این تعریف  $f(.,.)$  تابع فعال‌ساز<sup>۷۲</sup>های لایه قبل از لاجیت و یا همان کد CNN می‌باشد. این نوع تعریف  $A$  به ما این اطمینان را می‌دهد که نمونه به درستی کدگذاری و بازسازی شده و در نتیجه از توزیع داده واقعی می‌باشد.

$$A(x) = ||f_{xx}(x, x) - f_{xx}(x, G(E(x)))||_1 \quad (۲-۳۱)$$

با این تعریف، نمونه‌ها با مقدار بیش‌تر  $A$  با احتمال بالاتری داده ناهنجار خواهند بود. در ادامه در الگوریتم ۲-۳ روال محاسبه  $A(X)$  ارائه می‌شود.

---

الگوریتم محاسبه امتیاز ناهنجاری مدل ALAD

---

$f_{xx}$  لایه ویژگی مربوط به تمایزگر  $D_{xx}$   $x \sim p_{X_{Test}}(x)$ ،  $E$ ،  $G$ ،  $f_{xx}$  ورودی

$A(x)$  خروجی

انجام روال استنتاج

1.  $\tilde{z} \leftarrow E(x)$  کدگذاری نمونه
  2.  $\hat{z} = G(\tilde{z})$  کدگشایی نمونه
  3.  $f_\delta \leftarrow f_{xx}(x, \hat{x})$
  4.  $f_\alpha \leftarrow f_{xx}(x, x)$
  5. بازگرداندن  $||f_\delta - f_\alpha||_1$
  6. اتمام روال محاسبه نمره ناهنجاری
- 

الگوریتم ۲-۳: شبه کد الگوریتم ALAD

---

معیار استفاده شده در این جا از ایده تطابق ویژگی‌های از دست‌رفته الهام گرفته شده است [40]. اما در این جا به جای استفاده از ویژگی‌های محاسبه شده در شبکه تمایزگر GAN استاندارد (که اختلاف را بین نمونه های تولید شده و داده‌های واقعی را محاسبه می‌کند)، از ویژگی‌های محاسبه‌شده در شبکه تمایزگر  $D_{xx}$

---

<sup>72</sup> Activation



استفاده می‌شود. همچنین در این جا به جای استفاده از این معیار در حین آموزش شبکه GAN، از این معیار در هنگام روال استنتاج بهره می‌جوییم.

سوالی که در این جا مطرح می‌شود این است که : چرا نباید از خروجی تمایزگر  $D_{xx}$  به عنوان معیار فاصله استفاده کرد. پاسخ این سوال بدین صورت است که هدف از شبکه تمایزگر  $D_{xx}$  تمایز بین یک جفت نمونه واقعی  $(x, x)$  و بازسازی آن  $(x, G(E(x)))$  می‌باشد و شبکه کدگذار و شبکه مولد داده‌های واقعی و توزیع متغیر نهفته را کاملاً ضبط خواهند کرد. در این حالت  $D_{xx}$  قادر به تفکیک بین نمونه‌های واقعی و نمونه‌های بازسازی شده نخواهد بود و بدین ترتیب یک پیش بینی تصادفی را تولید می‌کند که معیار ناهنجاری مناسبی نخواهد بود.

## ۲-۶- جمع‌بندی

در این فصل ابتداء به دسته‌بندی روش‌های مختلف تشخیص ناهنجاری پرداختیم، در خلال همین دسته‌بندی ها برخی مدل‌های به نسبت قدیمی‌تر و سنتی بررسی شدند. در ادامه اهمیت جایگاه شبکه‌های عصبی تقابلی روشن شد. در گام بعدی معیارهای ارزیابی مدل‌های تشخیص ناهنجاری معرفی شده‌اند. در ادامه پس از بررسی مدل پایه شبکه عصبی تقابلی، زنجیره کارهایی که روی شبکه‌هایی عصبی تقابلی به منظور بهبود صورت گرفته است، بررسی شده است. در ادامه مدل AnoGAN شرح داده شد و با توجه به مشکل آن در نگاشت معکوس از فضای داده ورودی به فضای نهفته، مدل f-AnoGAN که در ادامه کار قبلی است، بررسی شد. با توجه به ضعف‌های موجود در ساختار f-AnoGAN، مقاله مکمل این مدل یعنی شبکه ALI مرور شد. در گام بعدی مدل EGBAD که جزو اولین کارها در زمینه تشخیص ناهنجاری که با الهام از مدل ALI خلق شده است مرور شد. در مرحله بعدی با توجه تضمین نشدن شرط سازگاری حلقه در ALI مقاله ALICE مورد بررسی دقیق‌تر قرار گرفت. در ادامه این روش دو شبکه دیگر ارائه شده است، شبکه RCGAN که ایده اصلی آن متمایل سازی بیشتر توزیع شبکه مولد و کدگذار به سمت توزیع داده‌های هنجار است و شبکه ALAD با هدف تضمین بیشتر چرخه پایداری، یک تمایزگر بین متغیر  $Z$  و بازسازی آن توسط شبکه، به ساختار شبکه ALICE اضافه کرد.

در فصل بعد روش پیشنهادی RCALAD که با تمرکز بر بکارگیری حداکثری جریان اطلاعات موجود در شبکه و همچنین تاکید بر بازسازی ضعیف نمونه‌های ناهنجار، به تفصیل ارائه خواهد شد.

## فصل سوم: روش پیشنهادی

در فصول قبل به بررسی روش‌های تشخیص ناهنجاری با رویکردهای مختلف پرداختیم. همان‌طور که پیش‌تر ذکر شد، با پیشرفت روزافزون زیرساخت‌های محاسباتی و افزایش توان پردازشی، علاقه و توجه محققان حوزه هوش مصنوعی به سمت شبکه‌های عصبی جلب شد. در میان این انواع شبکه‌های عصبی ارائه شده در چند سال اخیر، شبکه‌های مولد تقابلی به نتایج قابل توجه و درخشانی در کاربردهای مختلف دست یافته‌اند. علی‌رغم نتایج قابل دفاع این نوع شبکه‌ها در زمینه‌های پردازش تصویر [13]، پردازش گفتار و پردازش متن [41] در حوزه تشخیص ناهنجاری آن‌طور که شایسته است بدان توجه نشده است. در زمینه شناسایی نمونه‌های ناهنجار در دادگان دنیای واقعی به ندرت می‌توان الگوریتمی یافت که بر مبنای شبکه‌های مولد تقابلی طراحی شده باشد؛ در حقیقت بیشتر ساختارهای تقابلی برای کاربردهای دیگر طراحی شده است و صرفاً همان ساختار در کاربرد تشخیص ناهنجاری مورد استفاده قرار گرفته است، با این حال همین مدل‌های عام‌منظوره، توانسته‌اند به نتایج نسبتاً معقولی دست یابند [32].

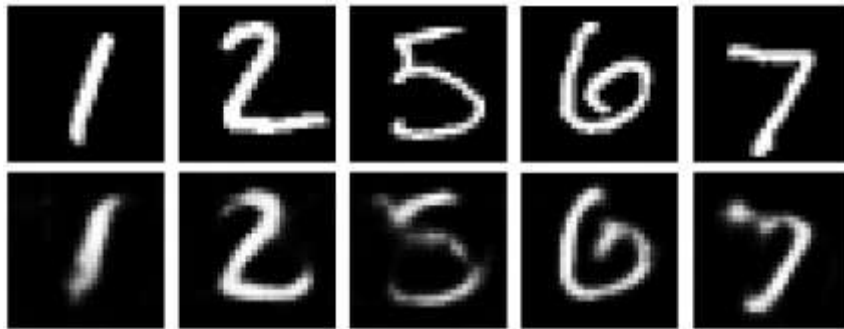
همان‌طور که پیش‌تر از این گفته شد، پس از معرفی شبکه GAN اولیه، مقاله AnoGAN از این ساختار برای کاربرد تشخیص ناهنجاری استفاده کرد. مشکل اصلی این روش، پیچیدگی زمانی بالای روش پیشنهادی به دلیل استفاده از یک ساختار مبتنی بر تکرار<sup>۱</sup> برای یافتن نگاشت معکوس از فضای داده به فضای نهان بود. در سال بعد شبکه ALI دقیقاً با همین هدف، یعنی دستیابی به نگاشت از فضای  $x$  به فضای  $z$  ارائه شد. در این روش یک شبکه کدگذار همزمان با شبکه مولد آموزش داده می‌شود. این شبکه توانست تا حدی مشکل نگاشت معکوس را برطرف کند. در این میان مشکل ساده اما جدی چرخه پایداری ایجاد شد. برای حل این مشکل، شبکه ALICE پیشنهاد افزودن یک تمایزگر بین تصاویر اولیه و بازسازی آن‌ها را داد. شبکه ALAD با افزودن یک تمایزگر دیگر بین نگاشت‌های تصاویر در فضای نهان و  $z$  های بازسازی شده آن‌ها، تلاش کرد تا مفهوم چرخه پایداری را در فضای نهان نیز برقرار کند. اشکال این روش در فرض استقلال بین دو چرخه پایداری فضای داده ورودی و چرخه پایداری فضای نهان بود. تمایزگر افزوده شده در مدل ALAD در واقع در یک چرخه مستقل از چرخه پایداری فضای داده ورودی قرار دارد در حالی که این دو چرخه کاملاً به یکدیگر وابسته می‌باشند. این مشکل در بخش ۳-۳-۱ تحت عنوان مشکل CCC<sup>۲</sup> با جزئیات بیشتر شرح داده خواهد شد. در واقع در این مدل از همه اطلاعات در دسترس

<sup>۱</sup> Iterative

<sup>۲</sup> Complete Cycle Consistency

برای بهبود روند آموزش و تقویت ساختار تقابلی استفاده نشده است، همانطور که می‌دانیم روند تبدیل داده با نگاشت از فضای داده ورودی به فضای نهفته آغاز می‌شود و در ادامه به فضای داده ورودی بازگردانده می‌شود تا با استفاده از میزان تفاوت داده ورودی و داده بازسازی شده امتیاز ناهنجاری محاسبه شود. این ایده در مدل CALAD<sup>۳</sup> ارائه شد و در بخش نتایج عملی نشان داده شد که استفاده از اطلاعات فضای نهفته و فضای داده ورودی به صورت توأم در یک چرخه کامل، سبب بهبود عملکرد و پایداری شبکه‌های عصبی تخصصی می‌شود. قابل توجه است که تاکنون در هیچ یک از کارهای قبلی، از اطلاعات موجود در یک چرخه کامل برای آموزش شبکه استفاده نشده است.

علاوه بر این، در مدل پایه ALAD فرض شده است که با آموزش انجام گرفته بر روی داده هنجار، بازسازی نمونه‌های ناهنجار لزوماً بازسازی ضعیفی خواهد بود در صورتی که هیچ قیدی برای تاکید بر این مهم و متمایل کردن مدل به سمت تولید بازسازی ضعیف در نظر گرفته نشده است. علی‌رغم نتایج قابل قبولی که مدل‌های پیشین ارائه داده‌اند، اما در تمامی آن‌ها با سهل انگاری فرض شده که اگر یک شبکه بر روی داده‌های هنجار آموزش ببیند، لزوماً برای داده‌های هنجار بازسازی خوب و برای ناهنجارها بازسازی ضعیف دارد. اما این فرض لزوماً برقرار نیست و ممکن است شبکه برای تصاویر ناهنجار هم بازسازی نزدیک به تصویر ورودی داشته باشد، همانند تصویر ۱-۳ که در آن کلاس هنجار کلاس صفر می‌باشد و بقیه کلاس‌ها، کلاس ناهنجاری به حساب می‌آیند و همان‌طور که مشاهده می‌کنید، مدل در بازسازی نامناسب نمونه‌های ناهنجار ضعیف عمل کرده و بازسازی بسیار نزدیک و شبیه به کلاس ورودی داشته است.



شکل ۱-۳: بازسازی نامطلوب نمونه ناهنجار.

<sup>3</sup> Complete Adversarially Learned Anomaly Detection

قابل توجه است که این بازسازی نزدیک تصویر ورودی اولیه، عملاً فرض ابتدایی برای روش‌های تشخیص ناهنجاری مبتنی بر بازسازی را نقض کرده و روند تفکیک داده‌ها را مختل خواهد کرد و بدین ترتیب دیگر با این روش داده هنجار و ناهنجار از یکدیگر قابل شناسایی نخواهند بود.

در این بخش همچنین مدل RALAD<sup>4</sup> که مبتنی بر شبکه‌های مولد تقابلی با هدف مقیدسازی مدل برای داشتن بازسازی ضعیف برای داده‌های ناهنجار معرفی می‌شود. در نهایت به کمک ترکیب هر دو ایده و با هدف معرفی یک چارچوب قوی و جامع برای تمامی کاربردهای تشخیص ناهنجاری، مدل RCALAD<sup>5</sup> معرفی شده است. در طراحی این روش تمرکز اصلی بر روی ارائه مدلی است که بتواند در کاربردهای دنیای واقعی نظیر داده‌های پزشکی مورد استفاده قرار بگیرد. مانند بسیاری از الگوریتم‌های مبتنی بر یادگیری، در این جا دو مرحله اصلی آموزش و آزمایش وجود دارد. در قسمت آموزش همانند دیگر چارچوب‌های تقابلی به نوبت بخش مولد و بخش تمایزگر را آموزش می‌دهیم تا هر دو بخش در عین تناسب به نوبت به‌روزرسانی شود. همچنین یک مرحله پیش‌پردازش شامل نرمال‌سازی تصاویر به منظور افزایش دقت مدل نهایی انجام شده است. در مرحله آموزش نمونه‌های ناهنجار را از تصاویر آموزشی حذف می‌کنیم و فقط از آن‌ها در مرحله آزمایش برای ارزیابی مدل پیشنهادی استفاده می‌کنیم. به بیان دیگر مدل تنها توزیع داده‌های عادی را می‌آموزد؛ روند آموزش مدل باید به حدی قدرتمند است که بتواند به خوبی فضای داده‌های ناهنجار را از داده‌های عادی تفکیک کند.

در ادامه و در قسمت تشخیص ناهنجاری، یک ورودی بدون برچسب به ساختار شبکه وارد می‌شود و به کمک اختلاف بازسازی ارائه شده توسط شبکه برای آن ورودی در هر دو فضای داده و فضای نهان، امتیاز ناهنجاری برای هر ورودی محاسبه می‌شود. با توجه به این نکته که شبکه بر روی دادگان نرمال آموزش داده می‌شود انتظار می‌رود که برای داده‌های ناهنجار بازسازی ضعیف‌تری داشته باشد و بدین ترتیب این اختلاف بیشتر شود و در نهایت این داده‌ها امتیاز ناهنجار بیشتری بگیرند. بدین ترتیب با انتخاب یک حد‌آستانه<sup>6</sup> و یا انتخاب یک نسبت معین از میان داده‌های با بیشترین امتیاز، داده‌هایی که ناهنجار هستند شناسایی می‌شوند.

<sup>4</sup> Regularized Adversarially Learned Anomaly Detection

<sup>5</sup> Regularized Complete Adversarially Learned Anomaly Detection

<sup>6</sup> Threshold

در ادامه این بخش به بررسی دقیق‌تر مشکلات پیشنهادی و راه حل ارائه شده و جزئیات هر قسمت از مدل پیشنهادی CALAD، RALAD و RCALAD خواهیم پرداخت. لازم به ذکر است نتایج آزمایش‌ها بر روی هر دو نوع داده تصویر و جدولی بیانگر کارایی و اثربخشی روش‌های پیشنهادی و نمایانگر سازگاری نتایج تئوری و عملی بدست‌آمده برای این مسائل می‌باشد.

### ۳-۱- مدل CALAD<sup>۷</sup>

این بخش به معرفی اولین مدل پیشنهادی اختصاص داده شده است. روش مورد بحث به منظور تشخیص ناهنجاری به دسته روش‌های مبتنی بر بازسازی تعلق دارد و از نظر دسترسی برچسب داده‌ها همانند سایر مدل‌های مبتنی بر شبکه‌های عصبی تقابلی بررسی شده، جزو دسته الگوریتم‌های بدون نظارت به حساب می‌آید. در این روش نیز شبکه مولد موجود در ساختار GAN نداشت از فضای نهان به فضای داده ورودی را فرا می‌گیرد، به بیان دیگر در این قسمت با استفاده از شبکه عصبی تقابلی، توزیع داده هنجار ورودی مدل می‌شود. برای یادگیری نگاشت معکوس از فضای داده ورودی به فضای نهان همانند مدل‌های پیشین از یک شبکه کدگذار استفاده شده است. در ادامه مشابه با مدل ALAD از دو تمایزگر  $D_{xx}$  و  $D_{zz}$  برای تضمین چرخه پایداری در هر دو فضای ورودی و نهان استفاده شده است.

نواوری به کار گرفته شده در این مدل، استفاده از متغیر جدید  $\hat{Z}_x$  و افزودن تمایزگر  $D_{xxzz}$  به ساختار ارائه شده برای تضمین شرط چرخه پایداری کامل در هر دو فضای نهفته و فضای ورودی می‌باشد. در واقع متغیر  $\hat{Z}_x$  بازسازی شبکه از نگاشت تصویر اولیه در فضای نهان است. هدف از تعریف این متغیر جدید دستیابی به یک حلقه کامل نگاشت‌های متوالی در فضای داده ورودی و فضای نهان به صورت وابسته است. قابل توجه است که تمایزگر مورد استفاده در کارهای قبلی، تنها به توزیع‌های مستقل داده‌ها در هر دو فضای نهفته و داده ورودی توجه می‌کردند و عملاً برخی از جریان اطلاعات موجود در شبکه که متعلق به توزیع توأمان متغیرهای فضای داده ورودی و نگاشت متناسب آن در فضای نهفته است، بلا استفاده باقی می‌ماند.

<sup>7</sup> Complete Adversarialy Learned Anomaly Detection

شایان ذکر است که این روال و اثر بکارگیری این نوع تمایزگر توامان یک مرتبه پیش از این ثابت شده است؛ در واقع برای آموزش به هنگام افزودن شبکه کدگذار به ساختار اولیه GAN در هنگام معرفی مدل ALI، دو راه پیش‌رو بود. یک راه افزودن یک شبکه تمایزگر در کنار تمایزگر اولیه موجود، برای تمیز بین متغیرهای فضای پنهان بود و یک راه، تقویت شبکه تمایزگر اولیه و آموزش این شبکه به نحویست که توزیع توامان در هر دو فضای پنهان و واقعی را فرا بگیرند و بتوانند داده‌های آموزشی را از داده‌های تولیدشده توسط این مدل تشخیص دهد.

در مدل ارائه شده، جریان اطلاعات شامل یک فرایند دو مرحله‌ایست، به این ترتیب که ابتدا از روی داده اولیه  $x$  در فضای واقعی یک نگاشت توسط کدگذار به فضای پنهان  $E(x)$  انجام می‌شود و سپس از روی همین داده یک نگاشت معکوس به عنوان بازسازی به فضای داده اولیه توسط شبکه مولد  $G(E(x))$  انجام می‌شود. سپس بار دیگر همین بازسازی به کدگذار فرستاده شده و در واقع بازسازی متغیر در فضای پنهان  $E(G(E(x)))$  بدست می‌آید. در مدل‌های قبلی متغیرها در هر کدام از فضاها با بازسازی آن‌ها به صورت جداگانه مورد بررسی قرار می‌گرفتند و مستقل از هم تمیز داده می‌شدند و بدین ترتیب زنجیره ارتباطات موجود بین این دو فضا نادیده گرفته می‌شد.

بر مبنای استفاده از همین اطلاعات از دست رفته و سابقه بکارگیری توزیع‌های توامان در این زمینه، در این جا مدل CALAD پیشنهاد شده‌است. در این ساختار در کنار چهارچوب اولیه شبکه‌های مولد تقابلی، یک شبکه تمایزگر توامان افزوده شده، تا با به کارگیری بیشترین اطلاعات موجود، مدل به جهت بهتری هدایت شده و سازگاری چرخه‌ها به صورت وابسته به هم بررسی شود؛ یعنی برای آموزش مدل از اطلاعات هر دو فضا به صورت توامان استفاده شود و در نهایت مدل به وزن‌های بهتر و دقت بالاتری دست یابد.

### ۳-۱-۱- معماری شبکه

در این بخش به معرفی تک تک اجزای مدل ارائه‌شده خواهیم پرداخت. در این ساختار همانند کارهای پیشین با هدف کاهش پیچیدگی زمانی، یک کدگذار توام با شبکه مولد در ساختار کلی شبکه عصبی تقابلی آموزش داده می‌شود. نگاشت معکوس از فضای داده ورودی به فضای نهفته به سادگی با تعبیه کدگذار  $E$  در معماری پیشنهادی به دست می‌آید. این کدگذار تنها از فضای داده ورودی نمونه می‌گیرد و

به طور تقریبی بازنمایی متناسب با آن را در فضای نهفته را تولید می‌کند. در اینجا برای آموزش هم‌زمان هر دو شبکه مولد و کدگذار از یک شبکه تمایزگر توامان با نام  $D_{xz}$  استفاده شده است. این تمایزگر بررسی می‌کند که جفت متغیر ورودی متعلق به توزیع داده ورودی  $x$  و نقطه متناظر با آن در فضای نهفته  $E(x)$  است و یا توسط شبکه مولد  $G(z)$  و نمونه‌گیری از فضای نهفته  $z$  تولید شده است.

به منظور ارضای شرط پایداری حلقه در فضای داده ورودی از تمایزگر  $D_{xx}$  افزوده شده است، این تمایزگر به صورت توام نمونه داده ورودی  $x$  و نمونه بازسازی شده متناظر آن را  $\hat{x} = G(E(x))$  به عنوان ورودی دریافت می‌کند. همچنین برای تقویت شرط پایداری حلقه در فضای نهفته، تمایزگر  $D_{zz}$  به این مجموعه اضافه شده است. این تمایزگر شرط پایداری حلقه را در خلال روند تولید نمونه بازسازی شده چک می‌کند. ورودی این تمایزگر نمونه ورودی شبکه مولد از فضای نهفته  $z$  و نمونه بازسازی متناظر با آن در فضای نهان است.

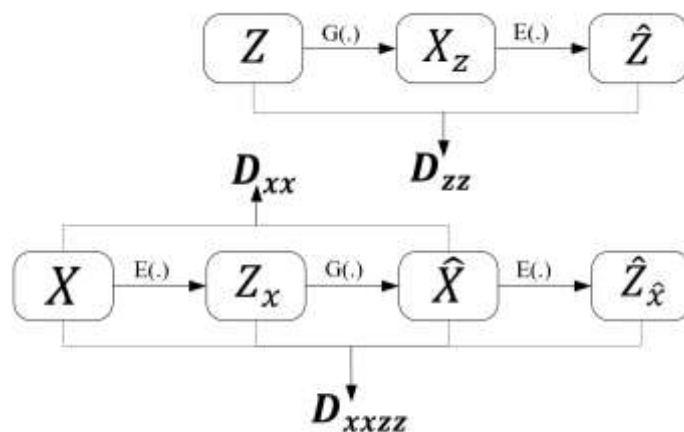
با اضافه شدن تمایزگر  $D_{xxzz}$  به ساختار موجود، تلاش شده است تا از تمامی اطلاعات موجود در یک چرخه کامل به صورت توام استفاده شود یعنی در کنار بررسی هر دو متغیر و بازسازی آن‌ها در همان فضا، توزیع توامان چهارتایی آن‌ها در روند تشخیص نمونه ناهنجار به کار گرفته شود. این شبکه وظیفه تمییز بین نمونه‌های چهارتایی  $(x, x, z_x, z_x)$  و  $(x, G(E(x)), z_x, E(G(z_x)))$  را دارد به بیان دیگر تلاش می‌کند تا  $x$  و بازسازی ارائه شده توسط شبکه و همین‌طور نگاشت تصویر ورودی در فضای نهان  $z_x$  و بازسازی خروجی شبکه مولد توسط کدگذار  $E(G(z_x))$  تا حد امکان به یکدیگر نزدیک کند. هدف از تعبیه این تمایزگر در این ساختار حل مشکل چرخه پایداری کامل<sup>۸</sup> یا به اختصار CCC می‌باشد. تعریف دقیق مسئله CCC در ادامه بررسی می‌شود.

بیان ریاضی مسئله CCC بدین صورت می‌باشد که به ازای هر متغیر  $x$  از فضای ورودی شبکه ابتدا کدگذار نگاشت معکوس به فضای نهفته را تخمین می‌زند که معادل  $E(x) = z_x$  می‌باشد. در مرحله بعد بازنمایی بدست آمده را به شبکه مولد وارد می‌کنیم تا بازسازی شبکه از متغیر ورودی  $G(z_x) = G(E(x)) = \hat{x}$  تولید کند. سپس همین بازسازی را بار دیگر به شبکه کدگذار می‌دهیم تا بازسازی در فضای نهفته را نیز محاسبه شود یعنی  $E(\hat{x}) = E(G(z_x)) = \hat{z}_x$  در این جریان،

<sup>۸</sup> Complete Cycle Consistency



انتظار منطقی از هر شبکه مبتنی بر بازسازی این است که دو متغیر  $\hat{X}$  و  $X$  و همچنین دو متغیر  $\hat{Z}_X$  و  $Z_X$  تا حد امکان کمترین اختلاف را داشته باشند. در مدل ALAD شباهت میان داده ورودی و بازسازی آن و همچنین شباهت  $Z$  و بازسازی آن مستقل از هم و در دو چرخه جداگانه بررسی می‌شد و فرض شده بود که مستقل از هم هستند در حالی که می‌دانیم این دو چرخه کاملاً به یکدیگر وابسته بوده و فرض استقلال این دو مسئله غلط است. در اینجا سعی شده است با بررسی توأم متغیرهای موجود در چرخه CCC در تمایزگر جدید  $D_{xxzz}$  عدم استقلال میان متغیرها مدل شود و جریان اطلاعات موجود در این زنجیره برای بهبود آموزش شبکه و تشخیص هر چه بهتر داده‌های ناهنجار به کار گرفته شود. تفاوت میان ورودی تمایزگر  $D_{xxzz}$  و ورودی تمایزگر  $D_{zz}$  که در مدل ALAD استفاده شده است در شکل ۲-۳ قابل مشاهده است.

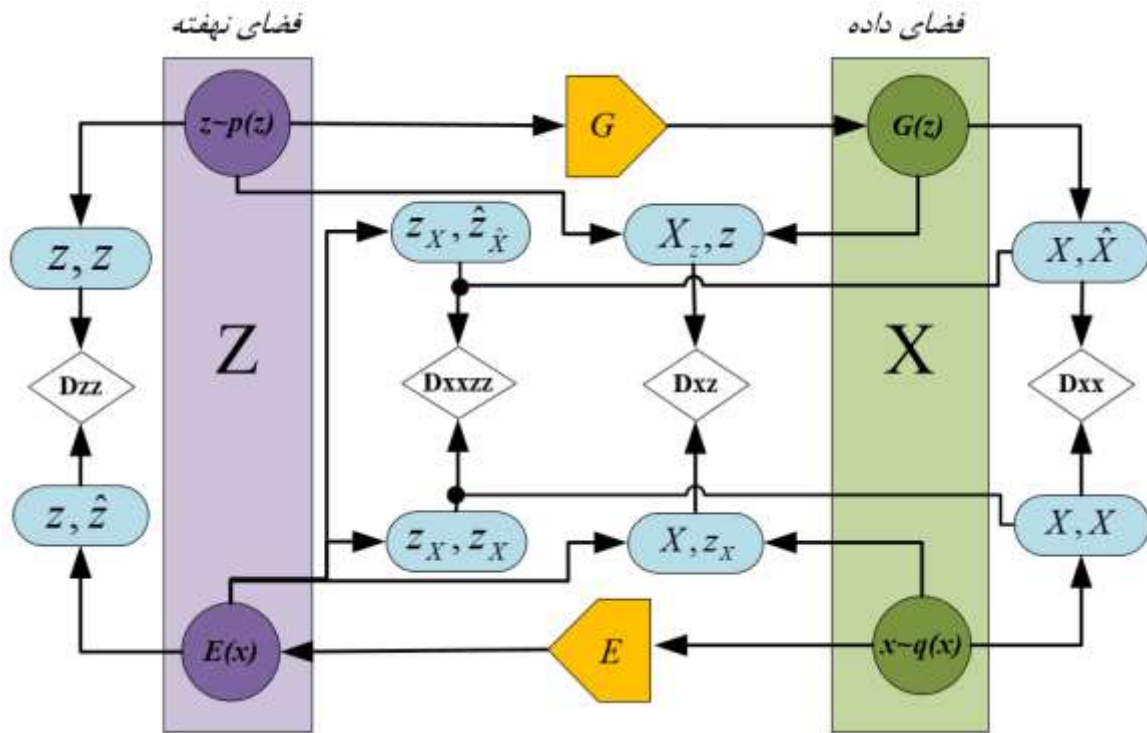


شکل ۲-۳: نمایش جریان اطلاعات در مدل CALAD.

همانطور که در شکل ۲-۳ قابل مشاهده است، چرخه مورد نظر در این مسئله شامل سه گام متوالی است، در کارهای قبلی چرخه کامل CCC وجود نداشت و از اطلاعات به طور کامل در شبکه استفاده نمی‌شد. در حقیقت، مدل‌های پیشین ارائه شده به بررسی مستقل این دو جفت متغیر در فضای جداگانه می‌پرداخت و تمایزگر دید کاملی از جریان اطلاعات و وضعیت داده در هر دو فضا داده ورودی و نهفته به طور همزمان نداشت. قابل توجه است که این متغیر  $\hat{Z}_X$  پیش از این محاسبه نمی‌شده و این چرخه در مدل‌های قبلی تعریف نمی‌شده است.

وجود تمایزگر  $D_{xxzz}$  سبب می‌شود تا ویژگی‌های جدید و عمیق‌تری (به نسبت تمایزگرهای تک گامی  $D_{xx}$  و  $D_{zz}$ ) استخراج شود. این شبکه با در اختیار گرفتن خروجی کل چرخه داده، دید جامعی از

وضعیت تمام قسمت‌های شبکه دارد و با استفاده از تمام این خروجی‌ها به طور همزمان، به ویژگی‌های ترکیبی قوی‌تری برای تمیز نمونه‌های ناهنجار از نمونه داده‌های هنجار دست می‌یابد. جزئیات معماری مدل CALAD در شکل ۳-۳ نمایش داده شده است.



شکل ۳-۳: معماری CALAD.

نامگذاری‌های به کار گرفته شده در شکل ۳-۳ دقیقاً مطابق با توضیحات ابتدای همین بخش می‌باشد. در بخش ۳-۱-۲ به توضیح دقیق‌تر روال آموزش، بررسی جزئیات بلوک‌های موجود و توضیح کامل تابع هدف این مدل خواهیم پرداخت.

### ۳-۱-۲- تابع هدف

همانند سایر مدل‌های پیشین در روند آموزش مدل پیشنهادی از روال‌های آموزش تقابلی استفاده می‌شود، بدیت ترتیب که به صورت متوالی بخش مولد و کدگذار و پس از آن تمایزگرها آموزش داده می‌شوند، یعنی به ترتیب با ثابت نگه داشتن وضعیت بخش مولد و کدگذار، پارامترهای شبکه‌های تمایزگر را به‌روزرسانی می‌کنیم. سپس آموزش این بخش‌ها را متوقف کرده و با توجه به خروجی بهبود یافته

آن‌ها، پارامترهای بخش مولد و کدگذار به روزرسانی می‌شود و این روال بارها و بارها تکرار شده تا مدل به کیفیت خروجی مطلوب دست یابد. هدف از آموزش تمامی مدل‌های مبتنی بر بازسازی از جمله مدل CALAD تولید بازسازی مناسب برای داده‌های هنجار و بازسازی ضعیف برای نمونه داده‌های ناهنجار است.

بیان ریاضی تابع هدف مدل پیشنهادی حاصل از جمع دو بخش کلی  $V_{ano}$  و  $V_{CCC}$  مطابق معادله ۳-۱۴ می‌باشد.

$$\min_{E,G} \max_{D_{xz}, D_{xx}, D_{zz}, D_{xxx}} V_{ano}(D_{xz}, G, E) + V_{CCC}(D_{xxx}, D_{xx}, D_{zz}, G, E) \quad (۱۴-۳)$$

در ادامه به بررسی هدف و جزئیات ریاضی هر یک از این دو بخش خواهیم پرداخت.

اولین بخش از این تابع هدف  $V_{ano}$  می‌باشد. در حالت کلی بخش اول تابع هدف این مسئله یعنی  $V_{ano}$  مطابق معادله ۳-۱۵ تعریف می‌شود.

$$\begin{aligned} \min_{E,G} \max_{D_{xz}} V_{ano}(D_{xz}, G, E) = & \mathbb{E}_{x \sim q(x)} [\log D_{xz}(x, E(x))] \\ & + \mathbb{E}_{z \sim p(z)} [\log(1 - D_{xz}(G(z), z))] \end{aligned} \quad (۱۵-۳)$$

در این معادله، داده‌های هنجار مورد استفاده در مرحله آموزش با تابع توزیع احتمال  $q(x)$  تعریف می‌شود و  $p(z)$  به عنوان تابع توزیع ورودی برای شبکه مولد در نظر گرفته شده است.

بخش دوم تابع هدف مورد استفاده در تابع هزینه این مدل  $V_{CCC}$  است که تلاش می‌کند تا شباهت میان داده تولیدی و بازسازی آن توسط شبکه مولد برای داده‌های هنجار را تضمین کند. در واقع این تابع هدف شرط پایداری چرخه را هم به صورت تک مرحله‌ای و هم به صورت دو مرحله‌ای یعنی چرخه کامل ارضا می‌کند. این بخش مطابق معادله ۳-۱۶ فرمول‌بندی می‌شود.

$$\begin{aligned}
\min_{E, G} \max_{D_{xx}, D_{zz}, D_{xxzz}} V_{ccc}(D_{xxzz}, D_{xx}, D_{zz}, E, G) \\
= \mathbb{E}_{z \sim p(z)} [\log D_{zz}(z, z)] + \mathbb{E}_{z \sim p(z)} [1 - \log D_{zz}(z, E(G(z)))] \\
+ \mathbb{E}_{x \sim q(x)} [\log D_{xx}(x, x)] + \mathbb{E}_{x \sim q(x)} [1 - \log D_{xx}(x, G(E(x)))] \\
+ \mathbb{E}_{x \sim q(x)} [\log D_{xxzz}(x, x, E(x), E(x))] \\
+ \mathbb{E}_{x \sim q(x)} [1 - \log D_{xxzz}(x, G(E(x)), E(x), E(G(E(x))))]
\end{aligned} \quad (۱۶-۳)$$

در معادله ۱۶-۳ که متعلق به بخش چرخه پایداری تابع هدف می باشد از جمع چهار جزء تشکیل شده است، جزء دوم مربوط به تمایزگر  $D_{xx}$  می باشد. این تمایزگر تلاش می کند که شرط چرخه پایداری را به صورت تک گامی با شروع از فضای داده در این مدل ایجاد کند. جزء دوم مربوط به تمایزگر  $D_{zz}$  است که با هدف ایجاد چرخه پایداری تک گامی با شروع از فضای نهفته به مدل اضافه شده است و در نهایت دو عبارت آخر مربوط به تمایزگر جامع  $D_{xxzz}$  است. این تمایزگر در دو گام جریان اطلاعات را دنبال می کند و خروجی تمامی مراحل را بررسی می کند و چرخه پایدار را برای همه مراحل نگاشت داده در مدل ایجاد می نماید. همان طور که در بخش قبل نیز اشاره شد، هدف از تعبیه این تمایزگر در این ساختار، حل مشکل چرخه پایداری کامل یا CCC می باشد.

شاید پس از بررسی دقیق  $V_{ccc}$  این سوال پیش بیاید که آیا با وجود  $D_{xxzz}$  همچنان به دو تمایزگر دیگر نیز احتیاجی هست؟ در پاسخ به این سوال باید گفت که بله احتیاج هست، زیرا همانطور که در بخش ۳-۱ و در شکل ۲-۳ نشان داده شد، این تمایزگر ورودی متفاوتی از تمایزگرهای قبلی دارد و از متغیر  $Z$  که با هدف آموزش شبکه مولد  $G$  و یادگیری نگاشت معکوس از فضای نهان به فضای ورودی است، استفاده نمی کند و اگر دو تمایزگر یادشده حذف شوند، عملاً تضمین چرخه پایداری برای متغیر  $Z$  که اساس بخش  $V_{ano}$  می باشد نیز حذف شده که این اتفاق مطلوب نیست. در بخش ۴-۵-۴ به بررسی تاثیر هر یک از اجزا به تفکیک پرداخته شده است و نشان داده شده که بهترین نتیجه در حضور هر سه تمایزگر حاصل می شود.

## ۳-۲- مدل RALAD<sup>۹</sup>

بخش فعلی به معرفی مدل پیشنهادی دوم اختصاص داده شده است. این مدل بر پایه مدل ALAD بنا نهاده شده است. در مدل‌های مبتنی بر بازسازی فرض بر این است که اگر آموزش و بازسازی برای داده‌های هنجار به خوبی انجام بگیرد، بازسازی داده‌های ناهنجار لزوماً ضعیف و متفاوت از داده اولیه ورودی خواهد بود. در صورتی که این فرض سهل‌انگارانه است و هیچ استلزام یا قید کنترلی برای این مشکل در هیچ یک از مدل‌های پیشین ارائه نشده است. در مدل ارائه شده با افزودن نمونه‌گیری از کل فضا، این استلزام برای بازسازی در فضای هنجار ایجاد و مدل را به سمت فضای بازسازی داده‌های هنجار متمایل<sup>۱۰</sup> کردیم. در بخش‌های آتی به بررسی جزئیات دقیق‌تر این مدل می‌پردازیم.

### ۳-۲-۱- معماری شبکه

روند آموزشی مدل ALAD با همه مزیت‌هایی که نسبت به مدل‌های پیشین دارد، اما از یک مشکل اساسی چشم‌پوشی کرده است؛ مشکل استلزام بازسازی ضعیف. تعریف دقیق این مشکل بدین ترتیب می‌باشد که در تمامی روال‌های آموزش مدل‌های تشخیص ناهنجاری، بازسازی دقیق نمونه‌های هنجار با کمترین خطا به مدل آموزش داده می‌شود و در مرحله آزمایش، نمونه‌های ناهنجار و هنجار به مدل داده می‌شود و همواره فرض می‌شود که برای نمونه‌های هنجار میزان اختلاف تصویر ورودی با تصویر بازسازی شده کم و برای نمونه‌های ناهنجار این اختلاف زیاد خواهد بود.

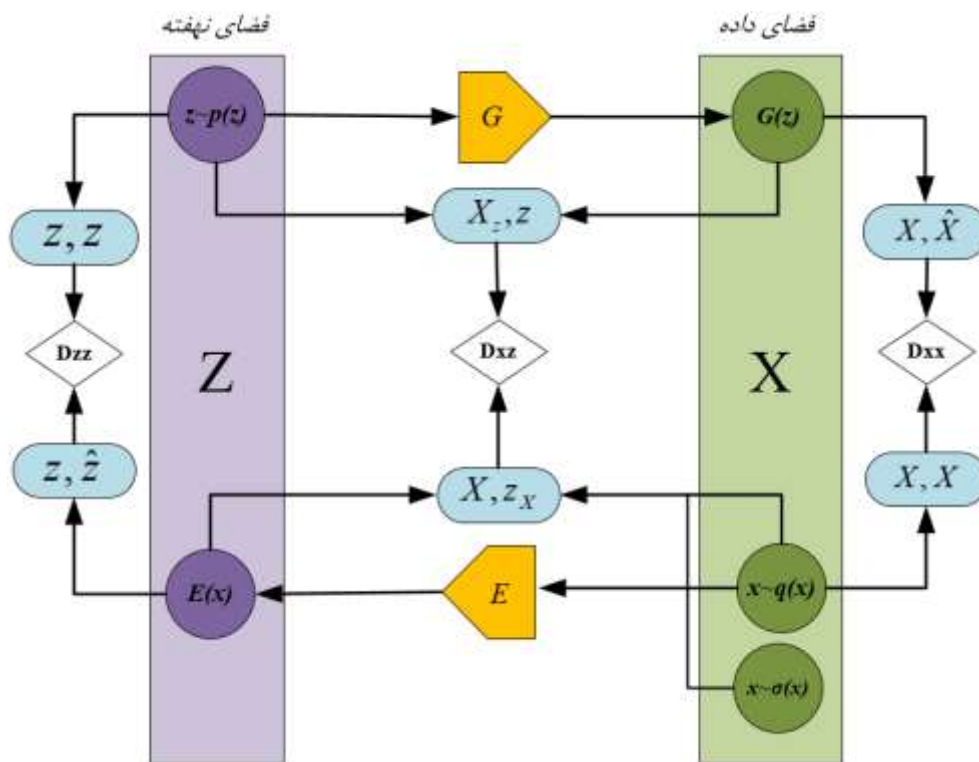
در برخی موارد پیش‌فرض فوق صحیح نیست و نمونه بازسازی شده داده ناهنجار، میزان اختلاف کمی با نمونه ورودی دارد و به همین سبب تشخیص آن به عنوان نمونه ناهنجار دشوار خواهد بود. در واقع در مدل‌های ارائه شده پیشین هیچ استلزامی برای بازسازی ضعیف نمونه ناهنجار وجود ندارد. علت وقوع این امر نگاشت تنک از فضای داده ورودی به فضای نهفته است. در حالت عادی آموزش فضای ورودی تنها به قسمت کوچکی از فضای نهفته نگاشت می‌شود و در نتیجه نمونه‌گیری از فضای نهفته به منظور نگاشت دوباره به فضای ورودی تنک خواهد بود. در زمان مواجهه با نمونه‌های هنجار این امر مشکلی ایجاد نخواهد

<sup>۹</sup> Regularized Complete Adversarial Learned Anomaly Detection

<sup>۱۰</sup> Bias

کرد چرا که فضای متناسب  $Z$  برای نمونه‌های هنجار به خوبی مدل شده است ولی در نمونه‌های ناهنجار با توجه به اینکه مدل تا به حال چنین داده‌های را ندیده است ممکن است نمونه را به نقطه‌ای ناشناخته از فضای نهفته نگاشت کند و در نتیجه بازسازی نمونه ناهنجار نیز ممکن است به نقطه‌ای نامناسب در فضای ورودی نگاشت شود. نگاشت به‌دست آمده از این فرایند، هیچ ضمانتی برای ایجاد بازسازی ضعیف از نمونه ناهنجار ارائه نمی‌دهد.

بلوک  $\sigma(x)$  به منظور پوشش حداکثری فضای نهفته به ساختار مدل ALAD اضافه می‌شود. هدف از تعبیه این بلوک تولید نمونه‌های جدید در فضای داده ورودی و سپس نگاشت آن به فضای نهفته است. انتظار می‌رود در این روند فضای نهفته به شکل مناسب‌تری نسبت به کارهای قبلی پوشش داده شود. نتایج عملی نمایانگر صحت تئوری ارائه شده در این قسمت است. در نهایت شمای کلی مدل پیشنهادی RALAD در شکل ۳-۴ قابل مشاهده است.



شکل ۳-۴: معماری RALAD.

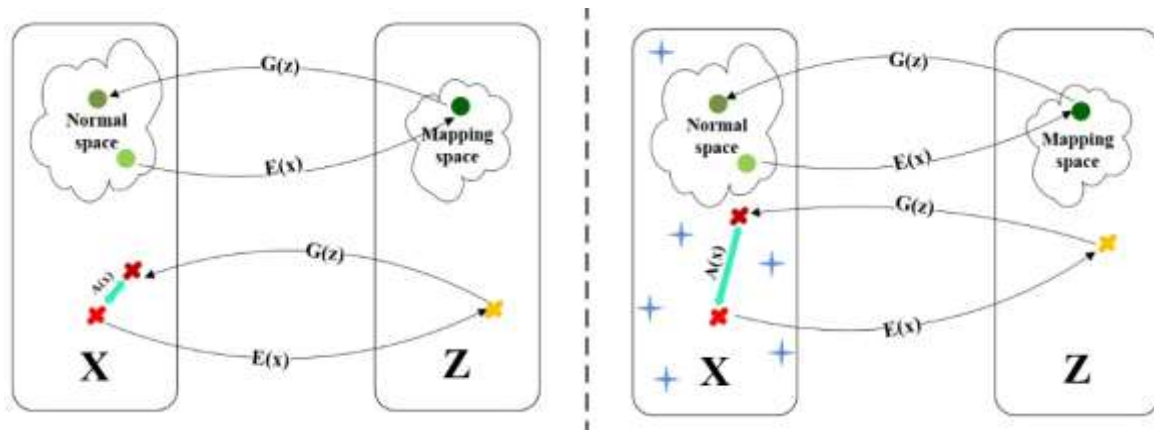
نامگذاری‌های به کار گرفته شده در شکل ۳-۴ مطابق با توضیحات بخش ۳-۱-۱ می‌باشد. در بخش ۳-۲-۲ به توضیح دقیق‌تر روال آموزش و بررسی جزئیات بلوک‌های موجود در مدل خواهیم پرداخت.

### ۳-۲-۲- تابع هدف

تابع هدف و روال آموزش مدل RALAD همانند مدل ALAD است و تنها تفاوت آن وجود توزیع  $\sigma(x)$  است. در حالت کلی مدل‌های مولد تخصصی تلاش می‌کنند تا نزدیک‌ترین بازسازی برای تمامی داده‌ها فارغ از هنجار یا ناهنجار بودن آن‌ها ایجاد شود اما شبکه مورد نیاز برای تشخیص ناهنجاری باید برای داده‌های هنجار بازسازی نزدیک و برای ناهنجارها بازسازی متفاوت از داده ورودی تولید کند و هدف از آموزش تمامی مدل‌های مبتنی بر بازسازی از جمله مدل RALAD تولید بازسازی مناسب برای داده‌های هنجار و بازسازی ضعیف برای نمونه داده‌های ناهنجار است. ایده مدل RALAD برای ایجاد این فاصله بازسازی تمامی داده‌ها در زیرفضای توزیع داده‌های هنجار است.

ممکن است این سوال مطرح شود که در عمل در برخی موارد این رویه نگاشت، سبب نزدیک‌تر شدن فاصله داده ناهنجار ورودی و بازسازی آن شود و این نوع از آموزش خلاف هدف مطلوب عمل کند. در پاسخ به این مسئله باید ذکر کرد که امکان این اتفاق در پاره‌ای از موارد وجود دارد، اما در مقایسه با مدل پایه ALAD که در آن هیچ اطلاعی از وضعیت نگاشت داده‌های ناهنجار وجود نداشت حال با اطمینان بالاتری می‌دانیم که یک فاصله حداقلی میان داده ورودی با بازسازی آن وجود دارد و بنابر نتایج به دست آمده در فصل چهارم ثابت شده است همین فاصله سبب تخصیص امتیاز ناهنجاری مناسب به نمونه‌های ناهنجار می‌شود.

در اینجا برای تاکید بر ایجاد بازسازی ضعیف برای داده ناهنجار، تلاش شده است تا تمامی فضای داده ورودی را به کمک نویز  $\sigma(x)$  پوشش دهیم و مولد شبکه را به سمت بازسازی هر بیشتر نزدیک به توزیع داده هنجار متمایل کنیم. در شکل ۳-۴ چگونگی عملکرد این قسمت از مدل پیشنهادی و نحوه تاثیر آن در فرایند آموزش به تصویر کشیده شده است.



شکل ۳-۵: تاثیر حضور توزیع  $\sigma(x)$  در روند آموزش مدل.

در شکل ۳-۵  $x$  بیانگر فضای داده ورودی و  $z$  بیانگر فضای داده ورودی است. نمونه‌ها توسط مولد  $G$  از فضای داده ورودی به فضای نهفته نگاشت می‌شوند و وظیفه انجام نگاشت معکوس بر عهده کدگذار  $E$  است. دایره‌های سبز رنگ نماد نمونه داده‌های هنجار و ضربدرهای قکد نماد نمونه‌های ناهنجار هستند. علاوه آبی رنگ نشانگر نمونه‌های تولید شده توسط توزیع  $\sigma(x)$  هستند که در تنها مرحله آموزش مورد استفاده قرار گرفته‌اند. فلش فیروزه‌ای مقدار امتیاز ناهنجاری را نشان می‌دهد. همانطور که در شکل ۳-۵ مشاهده می‌شود در صورت عدم حضور  $\sigma(x)$  (در سمت چپ شکل) در روند آموزش، امتیاز ناهنجاری برای نمونه‌های غیرعادی کمتر از زمانی است که از این توزیع استفاده شده است، در واقع در تصویر سمت راست، توزیع  $\sigma(x)$  مدل را به سمت بازسازی همه نمونه‌ها اعم از ناهنجار و هنجار به سمت توزیع داده‌های هنجار متمایل کرده است.

بیان ریاضی تابع هدف مدل پیشنهادی RALAD همانند مدل پیشنهادی قبلی حاصل از جمع دو بخش کلی  $V_{ano}$  و  $V_{CC}$  و مطابق معادله ۳-۱۷ می‌باشد.

$$\min_{E, G} \max_{D_{xz}, D_{xx}, D_{zz}} V_{ano}(D_{xz}, G, E) + V_{CC}(D_{xx}, D_{zz}, G, E) \quad (3-17)$$

تنها تفاوت موجود در این مدل در بخش  $V_{ano}$  نمایان می‌شود. با هدف حل مسئله استلزام بازسازی و توانمندسازی شبکه پیشنهادی برای تمایز هر چه بهتر بین داده‌های هنجار و ناهنجار از توزیع جریمه  $\sigma(x)$  استفاده شده است. ابعاد داده‌های خروجی این توزیع برابر با ابعاد داده ورودی  $x$  می‌باشد و در حالتی که هیچ اطلاعات اضافه‌ای از داده نداشته باشیم، با توجه به این که داده‌ها در مرحله پیش‌پردازش



نرمال شده اند، از توزیع گاوسی نرمال برای پوشش کلی تر فضا استفاده می کنیم. البته اگر اطلاعات اضافه ای به مدل از توزیع داده های ناهنجار داده شود، می توان با بکارگیری به جای توزیع  $\sigma(x)$  اعمال کرد. در حالت کلی بخش اول تابع هدف این مسئله  $V_{ano}$  طبق معادله ۱۸-۳ تعریف می شود.

$$\begin{aligned} \min_{E,G} \max_{D_{xz}} V_{ano}(D_{xz}, G, E) = & \mathbb{E}_{x \sim q(x)} [\log D_{xz}(x, E(x))] \\ & + \mathbb{E}_{z \sim p_g(z)} [\log(1 - D_{xz}(G(z), z))] \\ & + \mathbb{E}_{x \sim \sigma(x)} [\log(1 - D_{xz}(x, E(x)))] \end{aligned} \quad (18-3)$$

در این معادله، داده های هنجار مورد استفاده در مرحله آموزش با تابع توزیع احتمال  $q(x)$  تعریف می شود و  $p_g(z)$  به عنوان تابع توزیع ورودی برای شبکه مولد در نظر گرفته شده است. در این معادله از تابع توزیع جریمه  $\sigma(x)$  برای پوشش بهتر فضای داده ورودی و حل مشکل تنک بودن فضای نهفته استفاده می شود. در این معادله با داده های خروجی از توزیع  $\sigma(x)$  و بازسازی آن ها به عنوان نمونه های تقابلی برخورد می شود و بدین ترتیب با افزودن عبارت سوم یعنی  $\mathbb{E}_{x \sim \sigma(x)} [\log(1 - D_{xz}(x, E(x)))]$  به تابع هزینه کلی، تلاش می شود تا تمامی فضای ورودی به زیرفضا متعلق به توزیع داده های هنجار در فضای نهفته نگاشت شود و با این روش بازسازی های ارائه شده توسط شبکه به سمت توزیع داده هنجار ورودی متمایل شود. توضیح بیشتر این که ابتدا داده تولیدی از توزیع  $\sigma(x)$  به شبکه کدگذار وارد می شود، در طی فرایند آموزش کدگذار می آموزد تمام فضای داده ورودی را به فضای متناسب توزیع داده هنجار در فضای نهفته نگاشت کند، در نتیجه این روند، توزیع داده ورودی شبکه مولد همواره از فضای متناسب با توزیع داده هنجار خواهد بود و خروجی حاصل نیز متعلق به توزیع متناسب با آن در فضای داده ورودی می شود. از آن جایی که شبکه مولد نگاشت از فضای نهفته متعلق به داده هنجار به فضای داده ورودی متناسب را یاد گرفته است، تمامی فضای ورودی فارغ از هنجار یا ناهنجار بودن به فضای داده هنجار نگاشت می شود. پس به ازای داده های ناهنجار اختلاف بازسازی ارائه شده توسط شبکه بیشتر از داده های هنجار خواهد بود. برای اثبات نظری تمایل شبکه های مولد و تمایزگر به سمت داده های هنجار ابتدا تمایزگر و مولد بهینه را آموزش می دهیم. توزیع توام کدگذار روی داده های هنجار به صورت  $q(x, z) = q(x)e(z | x)$  و توزیع توام کدگذار روی داده های ناشی از توزیع جریمه به صورت  $\sigma(x, z) = \sigma(x)e(z | x)$  و

مطابق معادله ۱۹-۳ می‌باشد.  $p(x, z) = p(z)p(x | z)$  است. نقطه بهینه برای تمایزگر  $D_{zz}$  که از معادله ۱۸-۳ بدست می‌آید

$$\begin{aligned} D_{xz}^* &= \frac{q(x, z)}{q(x, z) + \sigma(x, z) + p(x, z)} \\ &= \frac{q(x, z)}{\left(1 + \frac{\sigma(x)}{q(x)}\right) q(x, z) + p(x, z)} \end{aligned} \quad (۱۹ - ۳)$$

در معادله ۱۹-۳ هر دو توزیع داده ورودی و توزیع داده جریمه در نظر گرفته است. از معادله فوق می‌توان نتیجه گرفت بر خلاف GAN استاندارد که تنها روی نمونه‌های عادی آموزش می‌بیند تمایزگر مدل ارائه شده در این کار احتمال بیشتری به نمونه‌های عادی اختصاص می‌دهد و با توجه به این که به نمونه‌های ناهنجار با  $q(x)$  کوچکتر اختصاص می‌یابد، خروجی تمایزگر بهینه برای داده‌های ناهنجار کم می‌شود. نتایج عملی ارائه شده در فصل چهارم گزاره‌های ارائه شده در این قسمت را پشتیبانی می‌کنند.

بخش دوم تابع هدف مورد استفاده در تابع هزینه این مدل  $V_{CC}$  است که دقیقاً همانند  $V_{CC}$  مدل ALAD است. تعریف این بخش در معادله ۲۰-۳ آمده است.

$$\begin{aligned} \min_{E, G} \max_{D_{xx}, D_{zz}} V_{CC}(D_{xx}, D_{zz}, E, G) \\ = \mathbb{E}_{z \sim p(z)} [\log D_{zz}(z, z)] + \mathbb{E}_{z \sim p(z)} [1 - \log D_{zz}(z, E(G(z)))] \\ + \mathbb{E}_{x \sim q(x)} [\log D_{xx}(x, x)] + \mathbb{E}_{x \sim q(x)} [1 - \log D_{xx}(x, G(E(x)))] \end{aligned} \quad (۲۰ - ۳)$$

### ۳-۳-۳ مدل "RCALAD"

در این بخش به بررسی مدل اصلی و کلی پیشنهادی این تحقیق که از ترکیب دو شبکه CALAD و RALAD بدست می‌آید، می‌پردازیم. در این شبکه به هر دو مسئله چرخه پایداری کامل و استلزام بازسازی ضعیف به طور همزمان پرداخته شده است و تلاش شده تا یک چارچوب جامع، کاربردی و سازگار برای تمامی مسائل تشخیص ناهنجاری ارائه شود.

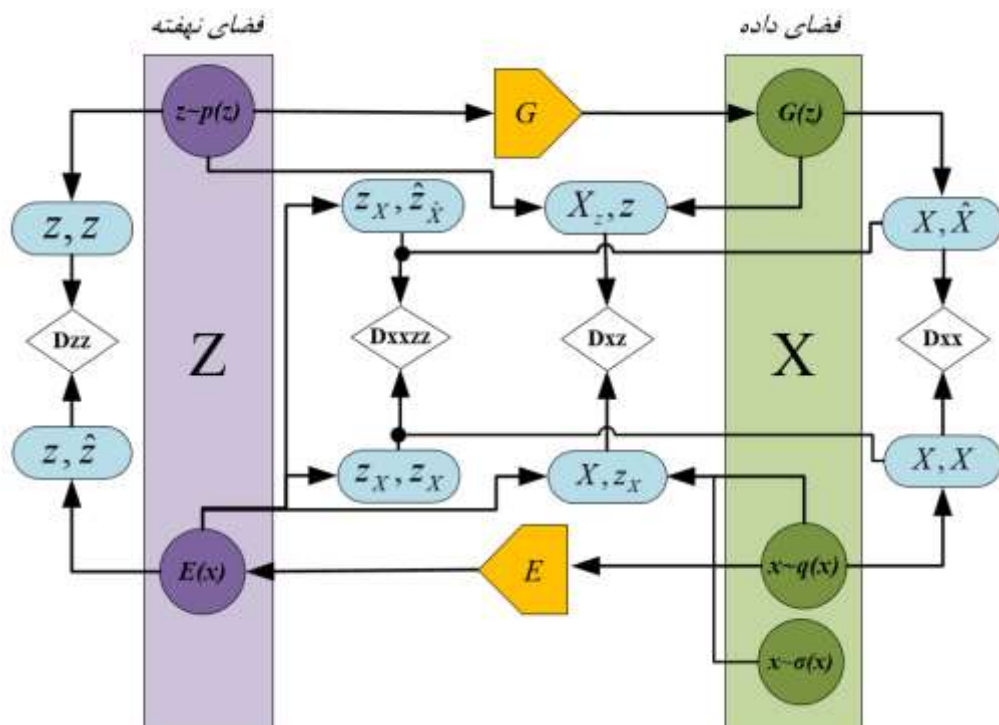
<sup>11</sup> Regularized Complete Adversarial Learned Anomaly Detection

نتایج آزمایش‌ها بر روی هر دو نوع داده تصویر و جدولی بیانگر کارایی و اثربخشی روش پیشنهادی RCALAD می‌باشد و نمایانگر سازگاری نتایج تئوری و عملی بدست‌آمده برای این مسئله است.

### ۳-۳-۱- معماری شبکه

همان‌طور که پیش از این ذکر شد، ایده بکارگرفته شده در این پروژه برای حل مشکل چرخه پایداری کامل، در این‌جا متغیر جدید  $\hat{z}_x$  تعریف و تمایزگر  $D_{xxzz}$  به ساختار کلی افزوده گردیده است که سبب تامین شرط چرخه پایداری کامل میان هر دو فضای ورودی و فضای نهان می‌شود. به جزئیات دقیق تعریف این مسئله، تعریف کامل این متغیر و نحوه آموزش تمایزگر  $D_{xxzz}$  در بخش ۳-۱ به تفصیل پرداخته شده است. همچنین با هدف متمایل سازی بازسازی ارائه شده توسط شبکه به سمت توزیع داده هنجار، از توزیع  $\sigma(X)$  برای نمونه گیری از کل فضا استفاده شده است و با استفاده از این توزیع تلاش شده تا که کل فضای ورودی به شبکه کدگذار نشان داده شود و کدگذار ملزم به ایجاد یک نگاشت نزدیک‌تر به فضای داده هنجار بشود. توضیحات کامل‌تر این موضوع نیز در بخش ۳-۲ شرح داده شده است.

با تجمیع این دو ایده در یک چارچوب در نهایت مدل RCALAD بدست خواهد آمد. شمای کلی مدل پیشنهادی و جزئیات ورودی هر یک از اجزای RCALAD در شکل ۳-۶ قابل مشاهده است.



شکل ۳-۶: معماری RCALAD.

نمادگذاری بکار گرفته شده در شکل ۳-۶ در زیربخش ۳-۱-۱ به طور کامل شرح داده شده است. در بخش بعد به توضیح دقیق تر روال آموزش و بررسی تابع هدف این شبکه خواهیم پرداخت.

### ۳-۲-۳- تابع هدف

تابع هدف مدل پیشنهادی RCALAD همانند مدل های پیشنهادی قبلی حاصل از جمع دو بخش کلی  $V_{ano}$  و  $V_{CCC}$  مطابق معادله ۳-۲۰ می باشد.

$$\min_{E,G} \max_{D_{xz}, D_{xx}, D_{zz}, D_{xxxz}} V_{ano}(D_{xz}, G, E) + V_{CCC}(D_{xxxz}, D_{xx}, D_{zz}, G, E) \quad (۳-۱۷)$$

در بخش  $V_{ano}$  با هدف حل مسئله استلزام بازسازی و توانمندسازی شبکه پیشنهادی برای تمایز هر چه بهتر بین داده های هنجار و ناهنجار از توزیع جریمه  $\sigma(x)$  همانند مدل RALAD استفاده شده است. بخش اول تابع هدف این مسئله یعنی  $V_{ano}$  طبق معادله ۳-۲۱ تعریف می شود.

$$\begin{aligned} \min_{E,G} \max_{D_{xz}} V_{ano}(D_{xz}, G, E) = & \mathbb{E}_{x \sim q(x)} [\log D_{xz}(x, E(x))] \\ & + \mathbb{E}_{z \sim p(z)} [\log(1 - D_{xz}(G(z), z))] \\ & + \mathbb{E}_{x \sim \sigma(x)} [\log(1 - D_{xz}(x, E(x)))] \end{aligned} \quad (۳-۲۱)$$

در این معادله، داده های هنجار مورد استفاده در مرحله آموزش با تابع توزیع احتمال  $q(x)$  تعریف می شود و  $p(z)$  به عنوان تابع توزیع ورودی برای شبکه مولد در نظر گرفته شده است. در این معادله از تابع توزیع جریمه  $\sigma(x)$  برای پوشش بهتر فضای داده ورودی و حل مشکل تنک بودن فضای نهفته استفاده می شود. بخش دوم تابع هدف مورد استفاده در تابع هزینه این مدل  $V_{CCC}$  است که مطابق با  $V_{CCC}$  معرفی شده در مدل CALAD است. در این بخش یک متغیر و تمایزگر جدید برای تضمین چرخه پایداری وابسته هر دو فضا معرفی شده است. تعریف این بخش از تابع هزینه مطابق معادله ۳-۲۲ می باشد.

$$\begin{aligned}
\min_{E,G} \max_{D_{xx}, D_{zz}, D_{xxzz}} V_{ccc}(D_{xxzz}, D_{xx}, D_{zz}, E, G) \\
= \mathbb{E}_{z \sim p(z)} [\log D_{zz}(z, z)] + \mathbb{E}_{z \sim p(z)} [1 - \log D_{zz}(z, E(G(z)))] \\
+ \mathbb{E}_{x \sim q(x)} [\log D_{xx}(x, x)] + \mathbb{E}_{x \sim q(x)} [1 - \log D_{xx}(x, G(E(x)))] \quad (22-3) \\
+ \mathbb{E}_{x \sim q(x)} [\log D_{xxzz}(x, x, E(x), E(x))] \\
+ \mathbb{E}_{x \sim q(x)} [1 - \log D_{xxzz}(x, G(E(x)), E(x), E(G(E(x))))]
\end{aligned}$$

### ۳-۴- تشخیص ناهنجاری

همانطور که گفته شد مدل پیشنهادی در این پروژه به منظور تشخیص ناهنجاری بر اساس بازسازی داده ورودی است. داده هنجار به صورت دقیق و شبیه به داده ورودی بازسازی می شود در حالی که بازسازی نمونه ناهنجار ضعیف خواهد بود. پس از بازسازی نمونه ها، عنصر کلیدی در تشخیص ناهنجاری تعریف امتیاز ناهنجاری با هدف محاسبه فاصله میان نمونه ورودی و خروجی بازسازی شده، توسط شبکه است. اولین انتخاب در زمینه محاسبه میزان فاصله این دو داده فاصله اقلیدسی است ولی در فضای داده این معیار ممکن است به اندازه کافی قابل اتکا نباشد [32]. به عنوان مثال در زمینه تشخیص ناهنجاری در تصاویر ممکن است علی رغم ویژگی های بصری مشابه، فاصله اقلیدسی زیادی داشته باشند.

در این کار به جای محاسبه فاصله میان نمونه ها در فضای داده ورودی، از فضای ویژگی موجود در تمایزگر  $D_{xxzz}$  استفاده می شود. به این منظور خروجی توابع فعالیت<sup>۱۲</sup> لایه یکی مانده به آخر به عنوان ویژگی استفاده می شوند. امتیاز ناهنجاری مورد استفاده به صورت زیر و با استفاده از خطای بازسازی نرم یک و مطابق معادله ۳-۲۳ تعریف می شود.

$$A_{fm}(x) = \|f_{xxzz}(x, x, E(x), E(x)) - f_{xxzz}(x, G(E(x)), E(x), E(G(E(x))))\|_1 \quad (23-3)$$

در معادله ۳-۲۳  $f(\cdot)$  بیانگر تابع فعالیت لایه یکی مانده به آخر در ساختار تمایزگر  $D_{xxzz}$  است.  $A_{fm}(x)$  میزان اطمینان تمایزگر از کیفیت روند کدگذاری و بازسازی داده توسط مولد است که اگر

<sup>12</sup> Activation function

خوب انجام شده باشد در واقع نمونه متعلق به توزیع داده هنجار است که مدل روی آن آموزش دیده است. بنابر مطالب گفته شده هر چه مقدار این معیار بیشتر باشد، اختلاف بازسازی‌ها بیشتر و احتمال ناهنجاری بودن آن ورودی بیشتر است. عملکرد مناسب این معیار در مقایسه با سایر معیارهای تشخیص ناهنجاری در فصل آینده نمایش داده شده است.

معیار پیشنهادی در این قسمت بر اساس روش تطبیق ویژگی<sup>۱۳</sup> و یا به اختصار fm، پایه‌گذاری شده است [40]، در ساختار GAN استاندارد برای تطبیق ویژگی از خروجی شبکه مولد استفاده شده است و در مقاله ALAD [32] از خروجی تمایزگر  $D_{xx}$  برای محاسبه این معیار استفاده شده است.

در مدل ALAD تنها از خروجی تمایزگر  $D_{xx}$  برای شناسایی نمونه‌های ناهنجار استفاده شد و بنابر نتایج به‌دست آمده برای مدل RCALAD این تمایزگر فاقد بخشی از اطلاعات مفید موجود در این مدل است. به منظور بهره‌گیری از تمامی اطلاعات موجود در این مدل برای تشخیص ناهنجاری در این جا یک معیار جدید با نام  $A_{all}$  نیز تعریف شده است. این امتیاز از جمع خروجی هر سه تمایزگر  $D_{xx}$ ،  $D_{zz}$  و  $D_{xxzz}$  تشکیل شده است. تمام تمایزگرهای موجود در مدل پیشنهادی تنها روی نمونه‌های هنجار آموزش دیده‌اند و بازسازی برای تمامی فضای داده ورودی به سمت فضای داده هنجار متمایل شده است، پس انتظار می‌رود تصویر بازسازی شده نمونه ناهنجار و همچنین بازنمایی آن در فضای نهفته که توسط کدگذار تولید می‌شود، ضعیف باشد و تمایزگرهای موجود در مدل این ورودی‌های ناهنجار را شناسایی کنند. بیان ریاضی این معیار در معادله ۳-۲۴ آورده شده است.

$$A_{all}(x) = D_{xxzz}(x, \hat{x}, z_x, \hat{z}_x) + D_{xx}(x, \hat{x}) + D_{zz}(z_x, \hat{z}_x) \quad (۳-۲۴)$$

حال مسئله قابل بررسی این موضوع می‌باشد که آیا معیار  $A_{all}$  حاوی اطلاعات کافی برای تشخیص داده‌های هنجار از ناهنجار می‌باشد یا خیر. پاسخ به این سوال در حالت کلی بله می‌باشد زیرا این تمایزگرها در طی فرایند آموزش یاد می‌گیرند که به اختلاف میان دوتایی  $(x, x)$  و  $(x, \hat{x})$  و همچنین دوتایی  $(z_x, z_x)$  و  $(z_x, \hat{z}_x)$  توجه کنند یعنی هر چه  $\hat{x}$  از  $x$  و یا  $\hat{z}_x$  از  $z_x$  فاصله بگیرد، تشخیص آن برای تمایزگرها ساده‌تر می‌شود. حال با افزودن توزیع  $\sigma(x)$  به این مجموعه و تلاش برای متمایل کردن

<sup>13</sup> Feature Matching



استفاده شد و روند آموزش و چرخه پایداری به بیشینه دقت در بین سایر مدل‌های مبتنی بر شبکه‌های مولد تقابلی رسید.

در گام بعدی به منظور هدایت کردن مدل به سمت تولید بازسازی ضعیف برای نمونه‌های ناهنجار و متمایل سازی تمامی بازسازی‌های ارائه شده توسط شبکه به سمت توزیع هنجار، با الهام از مدل RCGAN از یک توزیع نویز در فضای داده ورودی با نام  $\sigma(x)$  استفاده شد. بدین ترتیب که تلاش شد، با نمونه‌های گذشته توسط شبکه کدگذار همانند نمونه‌های خصمانه تولید شده توسط شبکه مولد برخورد شود و بدین ترتیب نگاشت تمامی فضای ورودی به سمت توزیع داده هنجار متمایل شود.

در گام نهایی دو امتیاز ناهنجاری جدید با نام‌های  $A_{fm}$  و  $A_{all}$  معرفی شدند، همانطور که مشاهده شد معیار اول مبتنی بر خروجی تمایزگرهای موجود بر شبکه تعریف شد بود و اساس کار معیار دوم بر استفاده از تطبیق ویژگی در لایه‌های تمایزگر  $D_{xxxz}$  بنا نهاده شده بود. در فصل آینده کارایی مدل پیشنهادی روی دادگان مختلف بررسی خواهد شد.



## فصل چهارم: آزمایش‌ها و نتایج

به منظور بررسی و ارزیابی کارایی الگوریتم پیشنهادی در تشخیص ناهنجاری در این فصل نتایج عملکرد آن طی چندین روال آزمایشی ارائه می‌شود. در ابتدا به معرفی مجموعه داده‌های آزمایشی و تنظیمات متناسب با هر یک از آن‌ها برای مدل پیشنهادی می‌پردازیم. سپس دیگر الگوریتم‌های تشخیص ناهنجاری مطرح پایه را معرفی کوتاهی خواهیم کرد، در گام بعدی به بررسی نتایج بر روی مجموعه داده‌های جدولی<sup>۱</sup> و تصویری ارائه شده می‌پردازیم و سپس نتایج به‌دست آمده بررسی، مقایسه و تحلیل می‌شود. درنهایت به بحث بیشتر و موشکافی نحوه عملکرد هر یک از امتیازهای معرفی شده در بخش ۳-۳-۳ در مجموعه داده‌های مختلف اشاره می‌شود.

## ۴-۱-۱ دادگان و پیش‌پردازش

برای سنجش عملکرد مدل پیشنهادی و بررسی کارایی آن از جنبه‌های مختلف از مجموعه داده‌هایی با ویژگی متفاوت استفاده می‌شود. در این قسمت مشخصات و ویژگی‌های این مجموعه داده‌ها معرفی خواهد شد.

### ۴-۱-۱-۱ مجموعه داده KDDCup99

مجموعه دادگان KDDCUP یک دادگان جدولی در ارتباط با نفوذ به شبکه‌های کامپیوتری است. این دادگان شامل ۴۹۴۰۲۱ نمونه با ۳۴ ویژگی اسمی<sup>۲</sup> و هفت ویژگی پیوسته<sup>۳</sup> است. در مرحله پیش‌پردازش ویژگی‌های اسمی به روش بازنمایی one-hot کدگذاری می‌شوند و نمونه‌های نهایی هر کدام ۱۲۱ بعد خواهند داشت. داده‌های با برچسب غیرنفوذ با توجه به اینکه حدود بیست درصد مجموعه داده را شامل می‌شود و در اقلیت است به عنوان ناهنجاری در نظر گرفته می‌شود. در مرحله آزمون بیست درصد از داده‌ها با بیشترین امتیاز ناهنجاری به عنوان داده ناهنجار در نظر گرفته می‌شود. معیارهای ارزیابی برای

<sup>1</sup> Tabular

<sup>2</sup> Categorical

<sup>3</sup> Continuous

سنجش این مجموعه داده شامل صحت، بازیابی و F1-score است. توضیحات این معیارها در بخش ۲-۴ به تفصیل مورد بررسی قرار گرفته است.

#### ۴-۱-۲- مجموعه داده Arrhythmia

دادگان جدولی آریتمی قلبی شامل ۴۵۲ نمونه با ۲۷۴ ویژگی است و هر داده می‌تولند به ۱۶ کلاس مختلف دسته‌بندی شود. کوچکترین کلاس‌ها به ترتیب شامل ۳، ۴، ۵، ۷، ۸، ۹، ۱۴ و ۱۵ نمونه هستند و جمعاً ۱۵ درصد از نمونه‌های این دادگان را شامل می‌شوند و همین نمونه‌ها در واقع ناهنجاری هستند. سایر کلاس‌ها به عنوان داده هنجار در نظر گرفته می‌شوند. در اینجا نیز در مرحله آزمایش ۱۵ درصد از دادگان با بیشترین امتیاز به عنوان داده‌های ناهنجار در نظر گرفته می‌شوند. معیارهای ارزیابی برای سنجش این مجموعه داده شامل صحت، بازیابی و F1-score است.

#### ۴-۱-۳- مجموعه داده Thyroid

این دادگان مربوط به بیماری تیروئید و جدولی است و شامل ۳۷۷۲ نمونه در سه کلاس و شامل شش ویژگی پیوسته است. کلاس با برچسب hyperfunction که شامل ۲.۵ درصد از مجموعه داده است به عنوان داده ناهنجار دسته‌بندی می‌شود، بنابراین در مرحله آزمون ۲.۵ درصد از نمونه‌ها با امتیاز ناهنجاری بالا به عنوان داده ناهنجار تشخیص داده می‌شود. در این مجموعه داده ۵۰٪ از نمونه‌های موجود به صورت تصادفی به عنوان داده آموزشی انتخاب شده است. توجه شود که نمونه‌های ناهنجار در تمامی مراحل آموزش مدل از داده‌های آموزشی حذف می‌شود. معیارهای ارزیابی برای سنجش این مجموعه داده همانند دیگر مجموعه داده‌های جدولی شامل صحت، بازیابی و F1-score است.

#### ۴-۱-۴- مجموعه داده Musk

دادگان Musk دادگانی جدولی مربوط به دسته‌بندی شش کلاسی روی مشک مولکولی<sup>۴</sup> شامل ۳۰۶۲ نمونه با ۱۶۶ ویژگی است. کلاس‌های موجود در دسته به نام ۲۱۳ و ۲۱۱ به عنوان داده ناهنجار در نظر

<sup>۴</sup> Musk molecular

گرفته می‌شوند و شامل ۳.۲ درصد دادگان است. در مرحله پیش‌پردازش دو ستون اسمی این مجموعه داده که شامل اسمی مولکول‌ها و ساختار آن‌ها حذف شده است. معیارهای ارزیابی برای سنجش این مجموعه داده مانند بخش‌های گذشته شامل صحت، بازیابی و F1-score است.

#### ۴-۱-۵- مجموعه داده CIFAR-10

این دادگان شامل ۶۰۰۰۰ تصویر  $32 \times 32$  در ده کلاس است. برای آزمایش مدل پیشنهادی روی این دادگان هر بار یک کلاس به عنوان داده هنجار در نظر گرفته می‌شود و نه کلاس دیگر به عنوان داده ناهنجار در نظر گرفته می‌شود. معیار مقایسه مدل پیشنهادی با سایر مدل‌ها در این مجموعه داده AUROC<sup>۵</sup> است. ۸۰ درصد داده به عنوان داده آموزشی در نظر گرفته می‌شود و ۲۰ درصد باقی به عنوان داده آزمون و ۲۵ درصد نمونه‌های آزمون به عنوان داده اعتبارسنجی در نظر گرفته می‌شود. توجه شود نمونه‌های ناهنجار از داده آموزشی و اعتبارسنجی حذف می‌شود.

#### ۴-۱-۶- مجموعه داده SVHN

دادگان SVHN مربوط به دسته‌بندی اعداد بین صفر تا نه روی پلاک خانه‌هاست. این مجموعه داده شامل ۹۹۲۸۹ تصویر  $32 \times 32$  است. رویکرد آموزش، اعتبارسنجی و آزمون در این دادگان همانند دادگان CIFAR10 است.

#### ۴-۲- تنظیمات مدل

در این قسمت به معرفی و بررسی جزئیات به کارگیری و تنظیمات معماری شبکه عصبی مدل پیشنهادی بر روی مجموعه داده‌های مختلف می‌پردازیم. تمامی نتایج گزارش شده در پایان‌نامه با استفاده از چارچوب Tensorflow1 و به زبان پایتون تولید شده است. توان پردازشی مورد نیاز در این کار با استفاده از NVIDIA A100 تامین شده است. تنظیمات آزمون<sup>۶</sup> همانند مدل پایه ALAD می‌باشد.

<sup>۵</sup> Area Under Receiver Operating Curve

<sup>۶</sup> Experimental Setup

جزئیات ساختاری و معماری مدل پیشنهادی در بخش ۳-۲ به تفصیل بررسی شده است. فراپارامترها برای بهینه ساز Adam از کتابخانه Tensorflow1 همانند [32] معادل  $\alpha = 10^{-5}$  و  $\beta_1 = 0.5$  می‌باشد. اندازه هر دسته<sup>۸</sup> در تمامی آزمایش‌ها ۳۲ در نظر گرفته شده است.

### ۴-۳- مدل‌های پایه

در این قسمت چارچوب پیشنهادی با تعداد زیادی از روش‌های تشخیص ناهنجاری مقایسه می‌شود، بخش اعظمی از مدل‌های مورد مقایسه در فصل ۲ مورد بررسی قرار گرفت. در ادامه به طور مختصر مدل‌هایی که پیش از بررسی نشده، شرح داده می‌شود.

#### ۴-۳-۱- روش OC-SVM<sup>۹</sup>

این روش یک مرز حول نمونه‌های هنجار یاد می‌گیرد و نمونه‌هایی که خارج از این مرز قرارگیرند به عنوان ناهنجاری در نظر گرفته می‌شوند. این مرز توسط روش ماشین بردار پشتیبان ایجاد می‌شود و هسته<sup>۱۰</sup> مورد استفاده در این روش RBF<sup>۱۱</sup> است. پارامتر  $\gamma$  برابر نسبت تعداد داده‌های ناهنجار به کل داده‌هاست و  $\gamma$  برابر  $1/m$  قرار داده می‌شود،  $m$  برابر تعداد ویژگی‌های دادگان است [۳۳]، [۴۲].

#### ۴-۳-۲- روش IF<sup>۱۲</sup>

این روش از دسته روش‌های کلاسیک یادگیری ماشین است و به جای مدل کردن توزیع داده هنجار داده ناهنجار را از سایر نمونه‌ها جدا می‌کند. ابتدا در این روش ابتدا تعدادی ویژگی انتخاب می‌شود و یک

<sup>۷</sup> Hyper Parameter

<sup>۸</sup> Batch

<sup>۹</sup> One Class Support Vector Machine

<sup>۱۰</sup> Kernel

<sup>۱۱</sup> Radial Basis Function

<sup>۱۲</sup> Isolated Forest

مقدار تصادفی برای هر ویژگی انتخاب می‌شود تا بتوان داده‌ها را جدا کرد. در ادامه میانگین فاصله هر نمونه تا ریشه به عنوان امتیاز ناهنجاری در نظر گرفته می‌شود [۴۳].

#### ۴-۳-۳- روش DSEBM<sup>۱۳</sup>

اساس کار این روش بر استفاده از انرژی لایه‌هایی که در خودگذارهای حذف نویز به کار برده می‌شوند بنا نهاده شده است. از خطای بازسازی و همچنین خود انرژی به عنوان امتیاز ناهنجاری در این مدل استفاده شده است. DSEBM-r بیانگر تشخیص ناهنجاری با خطای بازسازی و DSEBM-e بیانگر تشخیص ناهنجاری با امتیاز انرژی است [۲۶].

#### ۴-۳-۴- روش DAGMM<sup>۱۴</sup>

این مدل بر اساس خودکدگذارهای مورد استفاده در تشخیص ناهنجاری طراحی شده است. در گام اول مدل یک خودکدگذار را برای تولید فضای نهفته معقول و بازسازی ویژگی‌ها آموزش می‌دهد. شبکه تخمین‌گر دیگری نیز در این مدل آموزش داده می‌شود که به عنوان خروجی پارامترهای GMM که فضای نهفته با ابعاد کوچک را مدل‌سازی می‌کند، تولید می‌کند. در مرحله آزمون میزان درست‌نمایی بازنمایی مدل در فضای نهفته و ویژگی‌های بازسازی شده توسط GMM محاسبه می‌شود و این مقدار به عنوان امتیاز ناهنجاری در نظر گرفته می‌شود [۴۲].

#### ۴-۳-۵- روش DCAE<sup>۱۵</sup>

مدل DCAE یک مدل کلاسیک خودکدگذار است که در آن کدگذار و کدگشا دارای ساختار کانولوشنی هستند. امتیاز ناهنجاری در این مدل نرم  $l_2$  خطای بازسازی است [۴۴].

<sup>13</sup> Deep Structures Energy Based Models

<sup>14</sup> Deep Autoencoding Gaussian Mixture Model

<sup>15</sup> Deep Convolutional Autoencoder

### ۴-۳-۶- روش DSVDD<sup>۱۶</sup>

در این روش یک شبکه عصبی که حجم ابر کره محیط بر داده هنجار را کمینه می‌کند آموزش داده می‌شود. امتیاز ناهنجاری در این مدل فاصله اقلیدسی میان مرکز این ابرکره تا داده ورودی است [۴۵].

### ۴-۴- نتایج

در این بخش به مقایسه مدل پیشنهادی با مدل‌های پایه بخش ۴-۳ روی مجموعه داده‌هایی که در بخش ۴-۱ معرفی شد، می‌پردازیم. در ادامه مقایسه‌ها را در دو بخش دادگان جدولی و تصویری ارائه می‌کنیم.

#### ۴-۴-۱- دادگان جدولی

مدل پیشنهادی بر روی چهار مجموعه داده جدولی شامل KDDCup99، Arrhythmia، Thyroid و MUSK آزمایش شده است. دادگان Thyroid و MUSK به دلیل نرخ پایین ناهنجاریشان انتخاب شده‌اند تا مقاومت<sup>۱۷</sup> مدل در شرایط تفاوت فاحش در میزان نسبت داده‌های ناهنجار به هنجار نیز بررسی شود. در جدول ۴-۱ نتایج حاصل از آزمون مدل‌های پایه و مدل پیشنهادی را با سه معیار استاندارد صحت، بازیابی و F1-score ارزیابی می‌شوند. نتایج زیر مقادیر متوسط به ازای ۳۰ اجرا برای هر یک از مدل‌ها می‌باشد.

<sup>۱۶</sup> Deep Support Vector Data Description

<sup>۱۷</sup> Robustness

جدول ۴-۱: نتایج خروجی مدل پیشنهادی در مقایسه با مدل‌های پایه بر روی مجموعه داده‌های جدولی.

Model	KDDCUP			Arrhythmia			Thyroid			Musk		
	Prec.	Recall	F <sub>1</sub>	Prec.	Recall	F <sub>1</sub>	Prec.	Recall	F <sub>1</sub>	Prec.	Recall	F <sub>1</sub>
IF	92.16	93.73	92.94	51.47	54.69	53.03	<b>70.13</b>	<b>71.43</b>	<b>70.27</b>	47.96	47.72	47.51
OC-SVM	74.57	85.23	79.54	53.97	40.82	45.18	36.39	42.39	38.87	—	—	—
DSEBMr	85.12	64.72	73.28	15.15	15.13	15.10	4.04	4.03	4.03	—	—	—
DSEBMe	86.19	64.46	73.99	46.67	45.65	46.01	13.19	13.19	13.19	—	—	—
AnoGAN	87.86	82.97	88.65	41.18	43.75	42.42	44.12	46.87	45.45	3.06	3.10	3.10
DAGMM	92.97	94.22	93.69	49.09	50.78	49.83	47.66	48.34	47.82	—	—	—
ALAD	94.27	<b>95.77</b>	95.01	50.00	53.13	51.52	22.92	21.57	22.22	58.16	59.03	58.37
DSVDD	89.81	94.97	92.13	35.32	34.35	34.79	22.22	23.61	23.29	—	—	—
RCALAD	<b>95.36</b>	95.62	<b>95.49</b>	<b>58.82</b>	<b>62.50</b>	<b>60.60</b>	53.76	51.53	52.62	<b>62.96</b>	<b>63.33</b>	<b>63.14</b>
error bar	0.28	0.29	0.28	6.6	6.8	5.8	4.3	2.7	2.8	5.06	2.53	2.62

نتایج جدول ۴-۱ بیانگر این مسئله است که مدل RCALAD در سه مورد از چهار مجموعه داده جدولی آزمایشی با اختلاف بهتر عمل کرده است و تنها در مجموعه داده Thyroid از مدل IF دقت کمتری به دست آورده است. با نگاه دقیق‌تر در این مجموعه داده می‌توان از علت این پدیده آگاه شد، ویژگی‌های این مجموعه داده شامل نتایج آزمایشگاهی برای بررسی بیماری تیروئید می‌باشد و شامل ۳۵ ویژگی است. اما ثابت شده است که وجود/عدم وجود این بیماری تنها تحت تاثیر دو ویژگی T3 و T4 است و به همین دلیل مدل IF که مبتنی بر انتخاب ویژگی مهم‌تر و با ارزش‌تر است، در این مسئله بهتر عمل می‌کند. یک ایده کلی برای بهبود نتایج مدل پیشنهادی روی این دادگان می‌تواند به کارگیری مدل IF در مرحله پیش‌پردازش و با هدف استخراج ویژگی‌های قوی‌تر باشد و پس از این مرحله از این ویژگی‌ها برای آموزش مدل RCALAD استفاده شود.

#### ۴-۴-۱- دادگان تصویری

در این بخش عملکرد مدل پیشنهادی روی مجموعه داده تصویری شامل SVHN و CIFAR10 آزمایش شده است. نتایج آزمایش به تفکیک برای هر کلاس از داده‌ها مطابق جدول ۴-۲ و ۴-۳ گزارش می‌شود. این نتایج حاصل از میانگین ۳ مرتبه اجرای هر یک از مدل‌ها می‌باشد. همچنین میانگین معیار AUROC برای هر مجموعه داده در جدول ۴-۴ آورده شده است.



جدول ۴-۲: نتایج خروجی مدل پیشنهادی در مقایسه با مدل‌های پایه بر روی مجموعه داده CIFAR10.

Normal	DCAE	DSEBM	DAGMM	IF	AnoGAN	ALAD	RCALAD
Airplane	59.1 ± 5.1	41.4 ± 2.3	56.0 ± 6.9	60.1 ± 0.7	67.1 ± 2.5	64.7 ± 2.6	<b>72.8 ± 0.8</b>
auto.	<b>57.4 ± 2.9</b>	57.1 ± 2.0	56.0 ± 6.9	50.8 ± 0.6	54.7 ± 3.4	45.7 ± 0.8	50.2 ± 0.3
Bird	48.9 ± 2.4	61.9 ± 0.1	53.8 ± 4.0	49.2 ± 0.4	52.9 ± 3.0	67.0 ± 0.7	<b>72.6 ± 0.2</b>
Cat	58.4 ± 1.2	50.1 ± 0.4	51.2 ± 0.8	55.1 ± 0.4	54.5 ± 1.9	59.2 ± 1.1	<b>64.2 ± 0.9</b>
Deer	54.0 ± 1.3	73.2 ± 0.2	52.2 ± 7.3	49.8 ± 0.4	65.1 ± 3.2	72.7 ± 0.6	<b>74.9 ± 0.5</b>
Dog	<b>62.2 ± 1.8</b>	60.5 ± 0.3	49.3 ± 3.6	58.5 ± 0.4	60.3 ± 2.6	52.8 ± 1.2	60.1 ± 1.1
Frog	51.2 ± 5.2	68.4 ± 0.3	64.9 ± 1.7	42.9 ± 0.6	58.5 ± 1.4	69.5 ± 1.1	<b>75.3 ± 0.4</b>
Horse	58.6 ± 2.9	53.3 ± 0.7	55.3 ± 0.8	55.1 ± 0.7	<b>62.5 ± 0.8</b>	44.8 ± 0.4	56.6 ± 0.2
Ship	<b>76.8 ± 1.4</b>	73.9 ± 0.3	51.9 ± 2.4	74.2 ± 0.6	75.8 ± 4.1	73.4 ± 0.4	<b>77.5 ± 0.3</b>
Truck	<b>67.3 ± 3.0</b>	63.6 ± 3.1	54.2 ± 5.8	58.9 ± 0.7	66.5 ± 2.8	43.2 ± 1.3	52.6 ± 0.6
Mean	59.4	60.3	54.4	55.5	61.8	59.3	<b>65.7</b>

همان‌طور که در جدول ۴-۲ مشاهده می‌شود، مدل پیشنهادی RCALAD می‌تواند در نگاه کلی (۷) کلاس از ۱۰ کلاس داده را) دیگر مدل‌های پایه از جمله ALAD و RCGAN را در مجموعه داده CIFAR10 مغلوب کند و حتی در کلاس‌هایی که به بهترین نتیجه دست نیافته است، به نتایج قابل قبولی دست یافته است. در ادامه عملکرد مدل پیشنهادی بر روی مجموعه داده SVHN مورد ارزیابی قرار گرفته است.

همان‌طور که از جدول ۴-۳ مشخص است اگرچه مدل پیشنهادی RCALAD در بیشتر کلاس‌ها عملکرد بهتری نسبت به سایر مدل‌های پایه دارد برای برخی از کلاس‌ها قادر به تشخیص نمونه‌های ناهنجار نیست. مقایسه وضعیت کلی هر یک از این مدل‌ها بر روی هر دو مجموعه داده در یک نگاه در سطر آخر جدول ۴-۳ انجام شده است. همان‌طور که می‌بینیم عملکرد کلی مدل پیشنهادی RCALAD از دیگر مدل‌ها بهتر است.

جدول ۴-۳: نتایج خروجی مدل پیشنهادی در مقایسه با مدل‌های پایه بر روی مجموعه داده SVHN.

Normal	OCSVM	DSEBMr	DSEBMe	IF	ANOGAN	ALAD	RCALAD
0	$52.0 \pm 1.6$	$56.1 \pm 0.2$	$53.4 \pm 1.8$	$53.0 \pm 0.6$	$57.3 \pm 0.4$	$58.7 \pm 0.9$	<b><math>60.4 \pm 0.1</math></b>
1	$48.6 \pm 5.3$	$52.3 \pm 0.9$	$52.1 \pm 0.3$	$51.2 \pm 0.9$	$57.0 \pm 0.8$	<b><math>62.8 \pm 1.7</math></b>	$59.2 \pm 0.3$
2	$49.7 \pm 7.7$	$51.9 \pm 0.8$	$51.8 \pm 0.4$	$52.3 \pm 0.1$	$53.1 \pm 0.4$	<b><math>55.2 \pm 2.3</math></b>	$54.9 \pm 0.1$
3	$50.9 \pm 1.4$	$51.8 \pm 0.4$	$51.7 \pm 0.5$	$52.2 \pm 0.3$	$52.6 \pm 0.4$	$53.8 \pm 3.3$	<b><math>55.8 \pm 1.9</math></b>
4	$48.4 \pm 5.2$	$52.5 \pm 0.1$	$52.4 \pm 0.2$	$49.1 \pm 0.6$	$53.9 \pm 0.5$	$58.0 \pm 0.1$	<b><math>58.5 \pm 0.2</math></b>
5	$51.1 \pm 2.6$	$52.4 \pm 2.3$	$52.3 \pm 2.6$	$52.4 \pm 0.8$	$52.8 \pm 0.1$	$56.1 \pm 0.9$	<b><math>56.2 \pm 0.4</math></b>
6	$50.1 \pm 3.9$	$52.1 \pm 1.8$	$52.2 \pm 1.8$	$51.8 \pm 0.2$	$53.2 \pm 0.0$	$57.4 \pm 0.6$	<b><math>59.4 \pm 0.5</math></b>
7	$49.6 \pm 1.3$	$53.4 \pm 0.9$	$55.3 \pm 1.1$	$52.0 \pm 0.4$	$55.0 \pm 0.0$	<b><math>58.8 \pm 0.3</math></b>	$58.0 \pm 0.4$
8	$45.0 \pm 4.2$	$51.9 \pm 0.3$	$52.5 \pm 0.6$	$52.3 \pm 0.8$	$52.2 \pm 0.7$	<b><math>55.2 \pm 0.4</math></b>	<b><math>56.1 \pm 0.5</math></b>
9	$52.5 \pm 3.9$	$55.8 \pm 1.7$	$52.7 \pm 1.4$	$53.7 \pm 0.6$	$53.1 \pm 0.1$	$57.3 \pm 0.6$	<b><math>58.3 \pm 0.2</math></b>
Mean	50.2	52.9	52.4	51.6	54.0	57.3	<b>57.7</b>

عملکرد ضعیف مدل پیشنهادی روی عدد سه می‌تواند به دلیل شباهت ظاهری میان عدد سه به عدد پنج و دو در زبان انگلیسی باشد. برای مشاهده نحوه عملکرد مدل روی کلاس عدد سه در شکل ۴-۴ را مشاهده کنید.



شکل ۴-۱: عملکرد مدل RCALAD روی کلاس عدد سه.

در شکل ۴-۱ ردیف اول داده هنجار، ردیف دوم بازسازی داده هنجار، ردیف سوم داده ناهنجار و ردیف چهارم بازسازی داده ناهنجار است.

## ۴-۵- بحث

در این بخش به بررسی کارایی و تاثیر جز به جز هر یک از عناصر موجود در مدل می‌پردازیم. سپس توابع توزیع جریمه مختلف و تاثیر انتخاب هر کدام بحث قرار می‌گیرد و نشان داده می‌شود که نتایج بدست آمده وابستگی چندانی به یک تابع توزیع جریمه خاص ندارند. در نهایت دو امتیاز ناهنجاری ارائه شده مقایسه و حوزه عملکرد متناسب با هر یک بررسی می‌شود.

۴-۵-۱- مطالعه فرسایشی<sup>۱۸</sup>

در این قسمت تاثیر جزء به جزء قسمت‌های مختلف مدل را بر دقت نهایی به دست آمده بررسی می‌کنیم. آزمایش‌ها در این‌جا در شرایط حضور و عدم حضور هر جزء تکرار می‌شوند و نتایج حاصل از آن‌ها با هم مقایسه می‌شود تا میزان تاثیر هر قسمت به طور جداگانه مشخص شود. نماد  $D_{xxzz}$  نشان‌دهنده افزودن همین تمایزگر به مدل پایه ALAD می‌باشد. و همچنین نماد  $\sigma(x)$  به معنای اضافه شدن توزیع کمکی برای پوشش فضای  $X$  می‌باشد.

جدول ۴-۴: تاثیر بخش‌های مختلف پیشنهادی در بهبود نتایج دادگان جدولی.

Model	Precision	Recall	F1 score
KDD99			
Baseline (ALAD)	$0.942 \pm 0.008$	<b><math>0.957 \pm 0.006</math></b>	$0.950 \pm 0.007$
Baseline + $D_{xxzz}$ (CALAD)	<b><math>0.959 \pm 0.004</math></b>	$0.957 \pm 0.007$	<b><math>0.958 \pm 0.005</math></b>
Baseline + $\sigma(x)$ (RALAD)	$0.943 \pm 0.005$	$0.955 \pm 0.004$	$0.949 \pm 0.004$
Baseline + $D_{xxzz} + \sigma(x)$ (RCALAD)	$0.953 \pm 0.007$	$0.956 \pm 0.005$	$0.954 \pm 0.006$
Arrhythmia			
Baseline (ALAD)	$0.500 \pm 0.049$	$0.531 \pm 0.047$	$0.515 \pm 0.048$
Baseline + $D_{xxzz}$ (CALAD)	$0.574 \pm 0.021$	$0.605 \pm 0.022$	$0.575 \pm 0.021$
Baseline + $\sigma(x)$ (RALAD)	$0.546 \pm 0.035$	$0.565 \pm 0.039$	$0.555 \pm 0.037$
Baseline + $D_{xxzz} + \sigma(x)$ (RCALAD)	<b><math>0.588 \pm 0.42</math></b>	<b><math>0.625 \pm 0.41</math></b>	<b><math>0.606 \pm 0.41</math></b>
Thyroid			
Baseline (ALAD)	$0.229 \pm 0.067$	$0.215 \pm 0.067$	$0.222 \pm 0.067$
Baseline + $D_{xxzz}$ (CALAD)	$0.529 \pm 0.071$	<b><math>0.518 \pm 0.075</math></b>	$0.523 \pm 0.073$

<sup>18</sup> Ablation studies

Baseline + $\sigma(x)$ (RALAD)	$0.431 \pm 0.039$	$0.457 \pm 0.043$	$0.443 \pm 0.041$
Baseline + $D_{xxxz} + \sigma(x)$ (RCALAD)	<b><math>0.537 \pm 0.054</math></b>	$0.515 \pm 0.057$	<b><math>0.526 \pm 0.055</math></b>
Musk			
Baseline (ALAD)	$0.500 \pm 0.068$	$0.531 \pm 0.070$	$0.515 \pm 0.069$
Baseline + $D_{xxxz}$ (CALAD)	$0.574 \pm 0.026$	$0.605 \pm 0.027$	$0.575 \pm 0.026$
Baseline + $\sigma(x)$ (RALAD)	$0.546 \pm 0.051$	$0.565 \pm 0.051$	$0.555 \pm 0.051$
Baseline + $D_{xxxz} + \sigma(x)$ (RCALAD)	<b><math>0.629 \pm 0.011</math></b>	<b><math>0.633 \pm 0.016</math></b>	<b><math>0.631 \pm 0.013</math></b>

نتایج جدول ۴-۴ نشان می‌دهد اضافه کردن تمایزگر  $D_{xxxz}$  به طور قابل ملاحظه‌ای کارایی مدل را روی دادگان جدولی افزایش داده است. همچنین این تمایزگر نتایج روی مجموعه داده‌های تصویری CIFAR-10 و SVHN را نیز به شکل مناسبی بهبود داده است. همان‌طور که در جدول ۴-۵ مشاهده می‌شود، نتایج مجموعه داده جدولی درخشان‌تر از مجموعه داده‌های تصویری می‌باشد. شایان ذکر است که اگرچه بهبود حاصل‌شده در دادگان تصویری کمتر می‌باشد، اما در این جنس مجموعه داده‌ها حتی بهبودهای جزئی‌تر از این هم بارزش بوده و قابل توجه می‌باشد.

جدول ۴-۵: تاثیر بخش‌های مختلف پیشنهادی در بهبود نتایج دادگان تصویری.

Model	AUROC
SVHN	
Baseline (ALAD)	$0.573 \pm 0.016$
Baseline + $D_{xxxz}$ (CALAD)	$0.576 \pm 0.014$
Baseline + $\sigma(x)$ (RALAD)	$0.568 \pm 0.018$
Baseline + $D_{xxxz} + \sigma(x)$ (RCALAD)	<b><math>0.577 \pm 0.019</math></b>
CIFAR-10	
Baseline (ALAD)	$0.593 \pm 0.017$
Baseline + $D_{xxxz}$ (CALAD)	$0.634 \pm 0.018$
Baseline + $\sigma(x)$ (RALAD)	$0.642 \pm 0.012$
Baseline + $D_{xxxz} + \sigma(x)$ (RCALAD)	<b><math>0.657 \pm 0.016</math></b>

برای مشاهده هر چه بهتر تاثیر توزیع  $\sigma(x)$  بر بازسازی نمونه‌های ناهنجار شکل ۲-۴ آورده شده است. همانطور که در شکل مشخص است این توزیع مدل را متمایل به تولید کلاس هنجار صفر به ازای همه ورودی‌های ناهنجار کرده است.



شکل ۲-۴: تاثیر توزیع  $\sigma(x)$  بر بازسازی نمونه‌های ناهنجار.

#### ۲-۵-۴- انتخاب تابع توزیع جریمه

در این زیربخش به آزمایش تابع توزیع‌های مختلف  $\sigma(x)$  و میزان تاثیر آن‌ها بر نتایج نهایی مدل می‌پردازیم. توزیع‌های مورد استفاده در این جا دو توزیع گاوسی به صورت  $N(0, I)$ ،  $N(0, 2I)$  و یک توزیع یکنواخت به صورت  $U(-1, +1)$  است. با استفاده از هر کدام از توابع توزیع یاد شده در مدل ارائه شده همانطور که در جدول ۴-۶ مشخص است، بهبودهای پایداری حاصل شده است. تابع توزیع  $\sigma(x)$  به منظور تقلید و یا تخمین توزیع داده ناهنجار طراحی نشده است و می‌تواند مستقل از توزیع داده ناهنجار حتی در شرایطی که متفاوت با آن است تاثیر مورد نظر خود را بر مدل بگذارد.

جدول ۴-۶: تاثیر  $\sigma(x)$  های مختلف بر عملکرد مدل RCALAD.

$t(x)$	KDDCUP			Arrhythmia		
	Prec.	Recall	$F_1$	Prec.	Recall	$F_1$
$\mathcal{N}(0, I)$	<b>0.629</b>	<b>0.633</b>	<b>0.631</b>	<b>0.588</b>	0.625	0.606
$\mathcal{N}(0, 2I)$	0.626	<b>0.633</b>	0.629	0.580	0.629	0.603
$\mathcal{U}(-1, 1)$	0.608	0.604	0.606	0.584	<b>0.633</b>	<b>0.607</b>

## ۴-۵-۳- ارزیابی کارایی امتیازهای ناهنجاری

در این قسمت امتیازهای ناهنجاری پیشنهادی در این پروژه با دیگر معیارهای مبتنی بر بازسازی مقایسه می‌شود. خروجی خام تمایزگرها را در این جا لاجیت<sup>۱۹</sup> می‌نامیم و غرض از ویژگی<sup>۲۰</sup>، ویژگی‌های تولیدی تمایزگر در لایه‌های پنهان است.  $x$  نمونه ورودی،  $z_x$  نگاشت این نمونه در فضای نهان،  $\hat{x} = G(E(x))$  بازسازی تولیدی همان نمونه و  $\hat{z}_x$  نگاشت تصویر بازسازی شده در فضای نهان به وسیله مدل RCALAD است. امتیازهای مورد استفاده در این بخش به شرح زیر است.

$$L_1 : A(x) = \|x - \hat{x}\|_1$$

$$L_2 : A(x) = \|x - \hat{x}\|_2$$

$$\text{Logits} : A(x) = \log(D_{xx}(x, \hat{x}))$$

$$\text{Features} : A(x) = \|f_{xx}(x, x) - f_{xx}(x, \hat{x})\|_1 \quad (۴-۱)$$

$$\text{FM} : A_{fm}(x) = \|f_{xxzz}(x, x, z_x, z_x) - f_{xxzz}(x, \hat{x}, z_x, \hat{z}_x)\|_1$$

$$\text{ALL} : A_{ALL}(x) = D_{xxzz}(x, \hat{x}, z_x, \hat{z}_x) + D_{zz}(z_x, \hat{z}_x) + D_{xx}(x, \hat{x})$$

در ادامه ارزیابی نتایج هر یک از امتیازها روی دادگان جدولی مطابق جدول ۴-۷ قابل مشاهده است.

جدول ۴-۷: مقایسه عملکرد امتیازهای ناهنجاری پیشنهادی با سایر امتیازها روی دادگان جدولی.

Model	Precision	Recall	F1 score
KDD99			
$L_1$	$0.9081 \pm 0.0638$	$0.9108 \pm 0.0638$	$0.9094 \pm 0.0638$
$L_2$	$0.9011 \pm 0.0155$	$0.9004 \pm 0.0157$	$0.9007 \pm 0.0156$
Logits	$0.9169 \pm 0.0162$	$0.9168 \pm 0.0164$	$0.9168 \pm 0.0163$
Features	$0.9127 \pm 0.0029$	$0.9177 \pm 0.0039$	$0.9151 \pm 0.0034$
Features_xxzz	<b><math>0.9327 \pm 0.0017</math></b>	<b><math>0.9377 \pm 0.0017</math></b>	<b><math>0.9301 \pm 0.0017</math></b>
Logits_all	$0.9231 \pm 0.0018$	$0.9207 \pm 0.0018$	$0.9218 \pm 0.0018$
Arrhythmia			
$L_1$	$0.3529 \pm 0.0148$	$0.3750 \pm 0.0164$	$0.3636 \pm 0.0256$
$L_2$	$0.3529 \pm 0.0107$	$0.3750 \pm 0.0108$	$0.3636 \pm 0.0107$

<sup>19</sup> Logit

<sup>20</sup> Feature

Logits	$0.5588 \pm 0.0334$	$0.5937 \pm 0.0386$	$0.5757 \pm 0.0359$
Features	$0.2325 \pm 0.0029$	$0.2500 \pm 0.0029$	$0.2424 \pm 0.0029$
Features_xxzz	$0.4411 \pm 0.0013$	$0.4687 \pm 0.0013$	$0.4545 \pm 0.0013$
Logits_all	<b><math>0.6176 \pm 0.0208</math></b>	<b><math>0.6562 \pm 0.0221</math></b>	<b><math>0.6363 \pm 0.0214</math></b>

Thyroid			
$L_1$	$0.4981 \pm 0.0028$	$0.4908 \pm 0.0024$	$0.4994 \pm 0.0024$
$L_2$	$0.5011 \pm 0.0330$	$0.5004 \pm 0.0318$	$0.5007 \pm 0.0324$
Logits	$0.4969 \pm 0.0142$	$0.4968 \pm 0.0144$	$0.4968 \pm 0.0143$
Features	$0.5127 \pm 0.0119$	$0.5177 \pm 0.0119$	$0.5151 \pm 0.0119$
Features_xxzz	$0.5227 \pm 0.0083$	$0.5123 \pm 0.0083$	$0.5174 \pm 0.0083$
Logits_all	<b><math>53.76 \pm 0.0029</math></b>	<b><math>51.53 \pm 0.0029</math></b>	<b><math>52.62 \pm 0.0029</math></b>

Musk			
$L_1$	$0.5979 \pm 0.0103$	$0.5931 \pm 0.0109$	$0.5954 \pm 0.0106$
$L_2$	$0.6008 \pm 0.0021$	$0.6018 \pm 0.0028$	$0.6013 \pm 0.0024$
Logits	$0.5868 \pm 0.0124$	$0.5897 \pm 0.0127$	$0.5882 \pm 0.0125$
Features	$0.5824 \pm 0.0011$	$0.5883 \pm 0.0019$	$0.5883 \pm 0.0015$
Features_xxzz	$0.6111 \pm 0.0481$	$0.6187 \pm 0.0468$	$0.6148 \pm 0.0474$
Logits_all	<b><math>62.96 \pm 0.0013</math></b>	<b><math>63.33 \pm 0.0013</math></b>	<b><math>63.14 \pm 0.0013</math></b>

همانطور که در جدول ۴-۷ مشاهده می‌شود روی دادگان جدولی خروجی تمایزگر  $D_{xxzz}$  دارای بهترین نتایج به نسبت سایر امتیازهای ناهنجاری است. با توجه به اینکه تعداد ویژگی‌ها روی دادگان جدولی به نسبت دادگان تصویری کمتر است تمایزگر  $D_{xxzz}$  قادر به تشخیص مناسب نمونه‌های ناهنجار است. نتایج روی دادگان تصویری در جدول ۴-۸ قابل مشاهده است.

جدول ۴-۸: مقایسه عملکرد امتیازهای ناهنجاری پیشنهادی با سایر امتیازها روی دادگان تصویری.

Anomaly Score	AUROC
SVHN	
$L_1$	$0.5778 \pm 0.0141$
$L_2$	$0.5636 \pm 0.0251$
Logits	$0.5369 \pm 0.0785$
Features	$0.5763 \pm 0.0367$
Logits_all	$0.5768 \pm 0.0251$
Features_xxzz	$0.5778 \pm 0.0161$
CIFAR-10	
$L_1$	$63.41 \pm 0.0321$
$L_2$	$63.27 \pm 0.0782$
Logits	$62.97 \pm 0.0643$
Features	$63.12 \pm 0.0368$
Logits_all	$64.77 \pm 0.0227$
Features_xxzz	$65.73 \pm 0.0194$

همانطور که در جدول ۴-۸ مشخص است عملکرد امتیاز مبتنی بر ویژگی‌ها روی دادگان تصویری بسیار مناسب است، این می‌تواند به این دلیل باشد بردار ویژگی‌ها برای هر عکس نسبت به بردار ویژگی‌های موجود در دادگان جدولی بزرگتر است و استفاده از پارامترهای بیشتر به منظور تشخیص ناهنجاری سبب بهبود عملکرد امتیاز مورد نظر شده است.

#### ۴-۵-۴- ارزیابی کفایت تمایزگر $D_{xxzz}$

در این بخش به بررسی سوالی که در بخش ۳-۱-۲ مطرح شده می‌پردازیم، توضیح دقیق مسئله این بخش بدین ترتیب است: آیا با وجود  $D_{xxzz}$  همچنان به دو تمایزگر دیگر نیز احتیاجی هست؟ در واقع



$D_{xxzz}$  برای تامین شرط پایداری چرخه کامل کافی نیست؟ آیا دیگر تمایزگرها اطلاعاتی در مدل استخراج کرده و یا اضافه بوده و می‌توان حذف شوند؟

برای پاسخ‌گویی به این سوالات مطابق جدول ۴-۹ به ترتیب هر یک از تمایزگرها را از مدل کنار گذاشته و عملکرد مدل را گزارش می‌کنیم.

جدول ۴-۹: ارزیابی عملکرد مدل در حضور/عدم حضور هر یک از اجزا.

Model	$D_{zz}$	$D_{xx}$	$D_{xxzz}$	Prec.	Recall	$F_1$
<b>KDD99</b>						
ALAD	✓	✓	×	$0.942 \pm 0.008$	<b><math>0.957 \pm 0.006</math></b>	$0.950 \pm 0.007$
ALI + $D_{xxzz}$	×	×	✓	$0.938 \pm 0.007$	$0.951 \pm 0.010$	$0.944 \pm 0.009$
ALI + $D_{zz}$ + $D_{xxzz}$	✓	×	✓	$0.946 \pm 0.005$	$0.955 \pm 0.004$	$0.950 \pm 0.004$
ALICE + $D_{xxzz}$	×	✓	✓	$0.941 \pm 0.005$	$0.954 \pm 0.008$	$0.947 \pm 0.006$
CALAD	✓	✓	✓	<b><math>0.959 \pm 0.004</math></b>	$0.957 \pm 0.007$	<b><math>0.958 \pm 0.005</math></b>
RCALAD	✓	✓	✓	$0.953 \pm 0.007$	$0.956 \pm 0.005$	$0.954 \pm 0.006$
<b>Arrhythmia</b>						
ALAD	✓	✓	×	$0.500 \pm 0.049$	$0.531 \pm 0.047$	$0.515 \pm 0.048$
ALI + $D_{xxzz}$	×	×	✓	$0.522 \pm 0.054$	$0.529 \pm 0.049$	$0.525 \pm 0.052$
ALI + $D_{zz}$ + $D_{xxzz}$	✓	×	✓	$0.571 \pm 0.033$	$0.582 \pm 0.028$	$0.576 \pm 0.031$
ALICE + $D_{xxzz}$	×	✓	✓	$0.543 \pm 0.052$	$0.561 \pm 0.044$	$0.551 \pm 0.048$
CALAD	✓	✓	✓	$0.574 \pm 0.021$	$0.605 \pm 0.022$	$0.575 \pm 0.021$
RCALAD	✓	✓	✓	<b><math>0.588 \pm 0.42</math></b>	<b><math>0.625 \pm 0.41</math></b>	<b><math>0.606 \pm 0.41</math></b>

مطابق نتایج تئوری، افزودن تمایزگر  $D_{xxzz}$  به چارچوب کلی و در کنار دیگر تمایزگرها بالاترین کارایی را داشته است. پس از آن حذف  $D_{xx}$  ضربه کمتری به مدل می‌زند زیرا بخشی از اطلاعاتی که استخراج می‌کند، توسط تمایزگر  $D_{xxzz}$  پوشش داده می‌شود. همان‌طور که نتایج نشان می‌دهد، نتیجه این بخش این است که تمایزگر  $D_{xxzz}$  به تنهایی کافی نیست و این سه تمایزگر در کنار هم بیشترین کارایی را دارند و تمایزگر  $D_{xxzz}$  به تنهایی نمی‌تواند تمامی جنبه‌ها را دیده و اطلاعات مورد نیاز را استخراج کند.

## ۴-۶- جمع‌بندی

در این فصل ابتدا مجموعه داده‌های استفاده شده در این پروژه معرفی شد. همچنین تنظیمات مورد نیاز برای هر مجموعه داده در تقسیم نمونه‌ها به مجموعه آموزشی و آزمون و چگونگی انتخاب تعداد نمونه‌های ناهنجار در هر دادگان مورد بررسی قرار گرفت. در گام بعدی مدل‌های پایه‌ای که به منظور مقایسه با مدل پیشنهادی به کار برده شده‌اند به طور مختصر بررسی شدند. در قدم بعدی مدل پیشنهادی با سایر مدل‌ها مقایسه و برتری کلی آن ثابت شد سپس میزان مشارکت هر قسمت از مدل پیشنهادی بر بهبود نهایی بررسی شد. در مرحله بعدی توابع مختلف به عنوان توزیع  $\sigma(x)$  امتحان شدند و ثابت شد بهبود حاصل از این توزیع مستقل از توزیع داده‌های ناهنجار است. سپس امتیازهای ناهنجاری پیشنهادی در مقایسه با سایر امتیازهای ناهنجاری مبتنی بر بازسازی آزموده شدند و کارایی آن‌ها تایید شد. در گام آخر به این سوال پاسخ داده شد که با وجود تمایزگر  $D_{xxzz}$  به دیگر تمایزگرها نیاز است یا خیر، که نتایج آزمایش‌ها نشان داد که بهترین نتیجه در حضور هر سه تمایزگر بدست می‌آید.

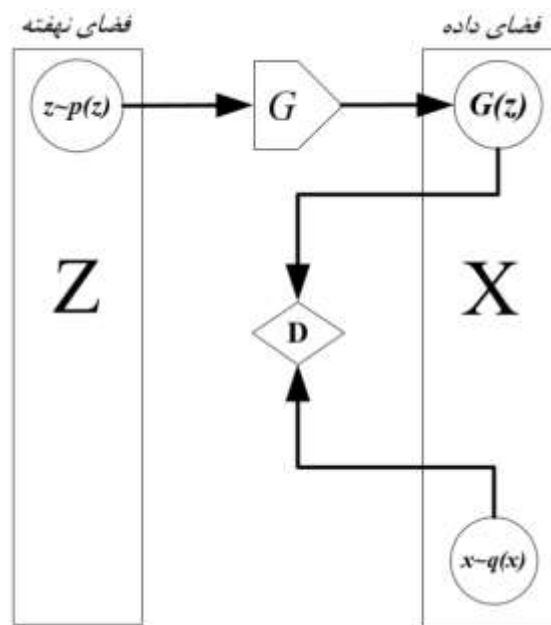
فصل پنجم:

جمع‌بندی، نتیجه‌گیری و کارهای آتی

در بخش پایانی به بر مرور و جمع‌بندی هر آنچه در این پروژه گفته شد می‌پردازیم و خط مش کلی کارهای آتی را مشخص می‌کنیم.

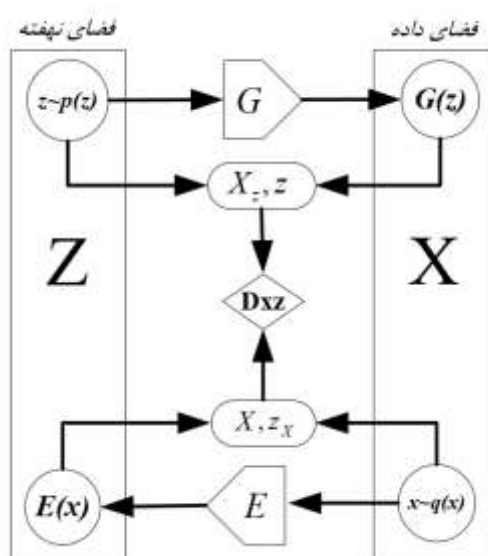
## ۵-۱- جمع‌بندی و نتیجه‌گیری

در بخش اول به اهمیت مسئله تشخیص ناهنجاری و کاربردهای آن در دنیای واقعی پرداختیم و سپس در همین فصل جایگاه شبکه‌های مولد تقابلی در زمینه تشخیص ناهنجاری مشخص شد. در فصل دوم به دسته‌بندی روش‌های که تا کنون در تشخیص ناهنجاری به کار برده شده‌اند از دیدگاه‌های مختلف و همچنین معرفی معیارهای پرکاربرد در ارزیابی مدل‌های تشخیص ناهنجاری پرداختیم. در خلال بررسی‌ها لزوم وجود روش‌هایی برای مدل‌سازی داده‌هایی پیچیده و با ابعاد بالا احساس و در نتیجه توجه‌ها به سمت شبکه‌های مولد تقابلی که قادر به انجام این مهم هستند معطوف شد. در ادامه همین فصل به بررسی تعریف و اصول این نوع از شبکه‌ها پرداخته شد و روند تکاملی و محلی‌سازی آن‌ها با هدف تشخیص ناهنجاری مورد اشاره قرار گرفت. همانطور که مشاهده شد شبکه مولد تقابلی که معماری آن در شکل ۵-۱ مشخص است در سال ۲۰۱۴ معرفی شد [۳۱].



شکل ۵-۱: معماری اولیه شبکه مولد تقابلی.

مطابق با آنچه در فصل دو بیان شد با هدف تشخیص ناهنجاری علاوه شبکه مولد که وظیفه نگاشت از فضای نهفته به فضایی داده ورودی را بر عهده دارد نیازمند فرایند استنتاجی برای نگاشت معکوس از فضای داده ورودی به فضای نهفته هستیم. مدل AnoGan با استفاده از یک فرایند مبتنی بر تکرار نقطه متناظر با داده ورودی را در فضای نهفته محاسبه می‌کرد. مشکل مدل مورد بحث پیچیدگی محاسباتی بالا و همچنین تا حد زیادی تصادفی بودن فرایند آن است. مدل‌های ALI و f-anogan و EGBAD با استفاده از یک کدگذار نگاشت معکوس به فضای نهفته آموخته می‌شود. معماری مدل ALI در شکل ۲-۵ قابل مشاهده است [۱۴]، [۳۵]، [۳۶].

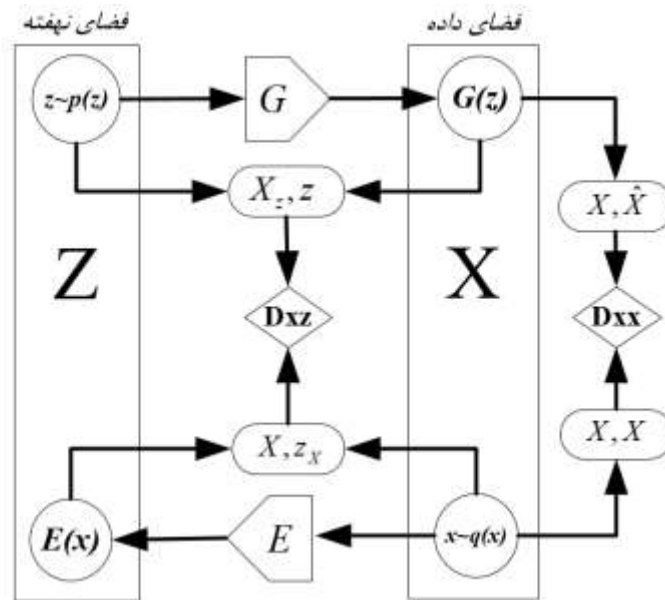


شکل ۲-۵: معماری مدل ALI.

این کارها اگرچه در جزئیات متفاوتند اما چالش اصلی آن‌ها یکسان است. در ادامه با توجه به اینکه هیچ ساختاری برای کنترل مشابهت میان تصویر ورودی و تصویر بازسازی شده توسط شبکه مولد تا به حال وجود نداشته است مدل ALICE به رفع این نقصان پرداخته است. مدل ALICE با اضافه کردن یک تمایزگر که داده ورودی و بازسازی آن را به عنوان ورودی دریافت می‌کند مشکل مورد نظر را که شرط پایداری حلقه نام داشت برطرف نمود. معماری ALICE در شکل ۳-۵ آورده شده است [۳۷].

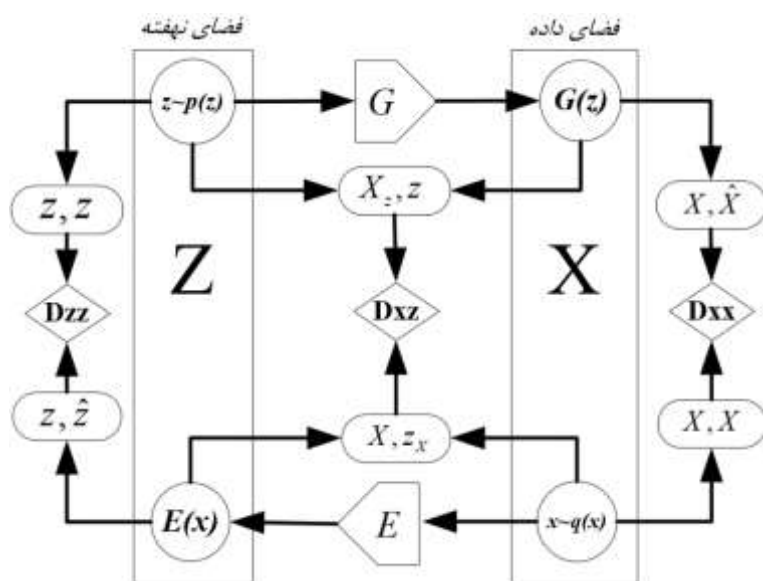
در فصل سه ابتدا دو مدل پایه مورد استفاده در این شبکه به طور دقیق مورد بررسی قرار گرفت. مدل ALAD بر پایه مدل قبلی بنا نهاده شده است و در معماری خود برای تضمین بیشتر پایداری حلقه تمایزگر دیگری برای فضای نهفته اضافه کرده است و همچنین ورودی‌های تمایزگرهای موجود در شبکه

را به صورت توأم در نظر گرفته است. معماری پیشنهادی ALAD سبب افزایش بازدهی در زمان آزمایش شده است و همچنین روند آزمایش را تثبیت کرده است. معماری این شبکه نیز در ۴-۵ نمایش داده شده است [۳۲].



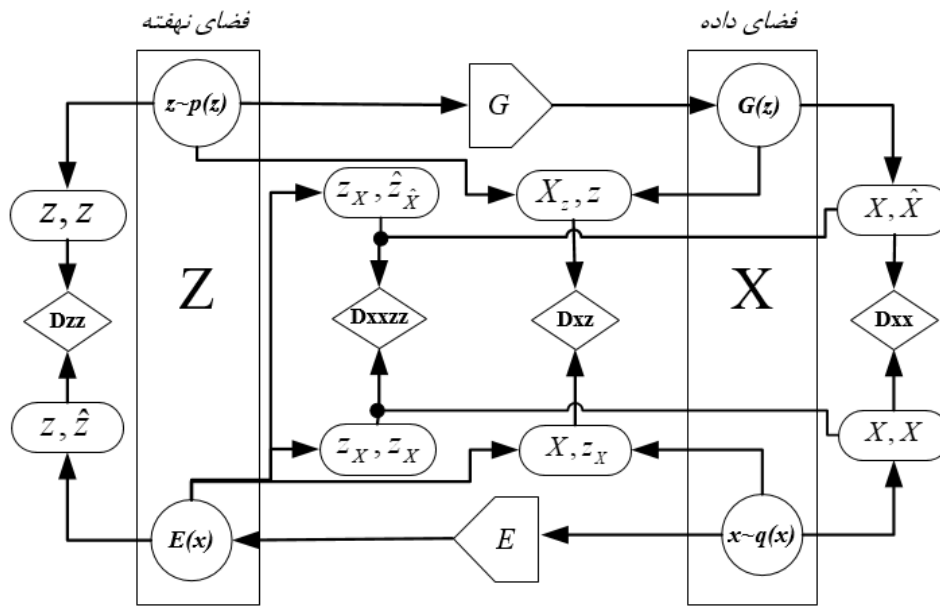
شکل ۵-۳: معماری شبکه ALICE.

مدل پایه دیگر که در این فصل بررسی شده است، مدل RCGAN است [۳۸]. این شبکه با هدف ضمانت بازسازی ضعیف برای نمونه‌های ناهنجار تابع توزیع  $t(x)$  را به ساختار تقابلی ALICE اضافه کرد تا مدل را به گونه‌ای متمایل به سمت بازسازی تمامی نمونه‌های ورودی به فضای داده هنجار در فضای داده ورودی کند. با انجام این کار فاصله میان داده ناهنجار و بازسازی آن زیاد خواهد شد و در نتیجه شناسایی نمونه ناهنجار ساده‌تر خواهد بود.



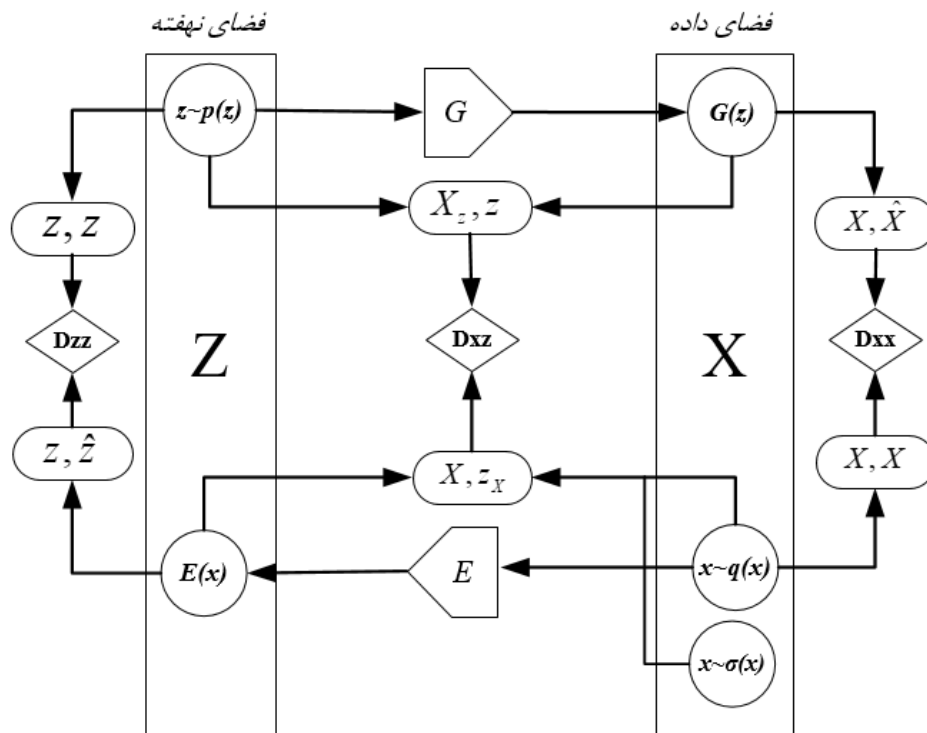
شکل ۵-۴: معماری شبکه ALAD.

پس از این قسمت‌ها نوبت به معرفی مدل‌های پیشنهادی این پروژه می‌رسد. همانطور که در فصل سه گفته شد به منظور تقویت قدرت تشخیص تمایزگر اطلاعات حاصل از روند دگرديسی داده ورودی در تمامی مراحل چرخه، باید توسط تمایزگر قابل دسترس باشد. چرخه مورد نظر در این مسئله شامل سه گام متوالی است، در کارهای قبلی از خروجی‌های چرخه اطلاعات به طور کامل در شبکه استفاده نمی‌شد. به منظور پوشش این نقص تمایزگر  $D_{xxzz}$  به ساختار تقابلی قبلی اضافه شد. نتیجه افزودن این تمایزگر مدل CALAD شد. معماری این مدل در شکل ۵-۵ نمایش داده شده است.



شکل ۵-۵: معماری شبکه CALAD.

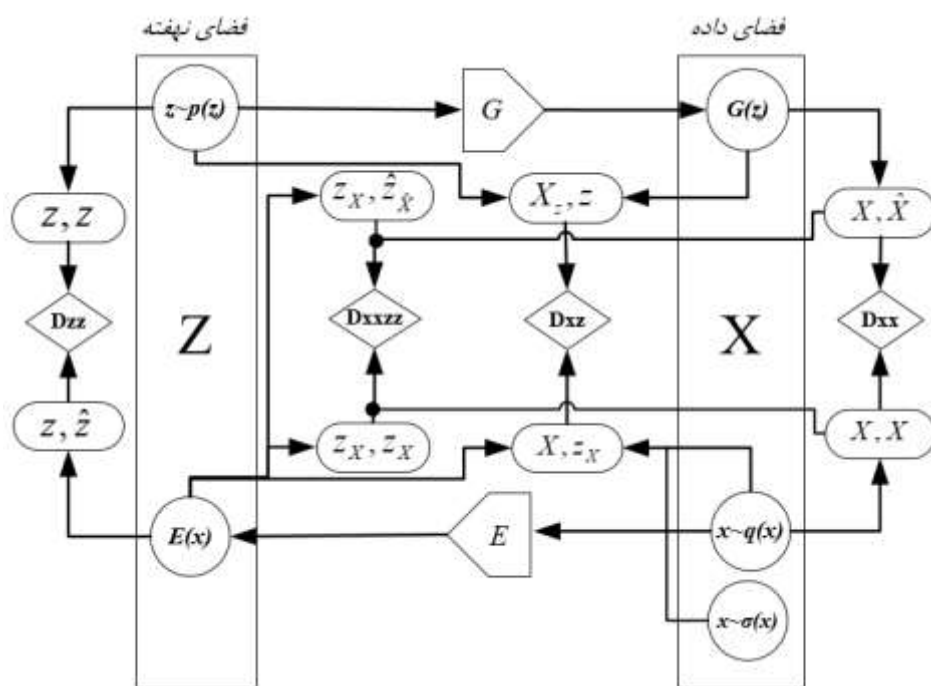
دیگر مدل معرفی شده در این تحقیق با استفاده از توزیع  $\sigma(x)$  برای ضمانت بازسازی ضعیف داده ناهنجار استفاده می‌کند. ساختار مدل RALAD در شکل ۵-۶ نشان داده شده است.



شکل ۵-۶: معماری شبکه RALAD.



در نهایت مدل جامع RCALAD که از تجمیع هر دو ایده بکارگیری متغیر  $\hat{Z}$  و همچنین توزیع  $\sigma(X)$  بدست می‌آید و همانطور که در بخش چهارم نشان داده شد، بیشترین کارایی د میان مدل‌ها را بدست می‌آورد. معماری این مدل در شکل ۷-۵ نمایش داده شده است.



شکل ۷-۵: معماری شبکه RCALAD.

توابع بهینه‌سازی هر یک از شبکه‌های مورد بحث در جدول ۵-۱ آمده است. همانطور که مشخص است توابع بهینه‌سازی این شبکه‌ها کاملاً در امتداد هم و در راستای رفع نقاط ضعف کارهای قبلی هستند.

جدول ۵-۱: روند تکامل توابع بهینه‌سازی شبکه‌های مولد تقابلی.

نام شبکه	تابع بهینه‌سازی
GAN	$\min_G \max_D V_{GAN}(D, G)$ $= \mathbb{E}_{x \sim q(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log (1 - D(G(z)))]$
ALI	$\min_G \max_D V_{ALI}(D, G)$ $= \mathbb{E}_{x \sim q(x)} [\log D(x, G_z(x))] + \mathbb{E}_{z \sim p(z)} [\log (1 - D(G_z(z), z))]$
ALICE	$\min_{E, G} \max_{D_{xz}, D_{xx}} V_{ALICE}$ $= V_{ALI} + \mathbb{E}_{x \sim q(x)} [\log D_{xx}(x, x) + \log 1 - D_{xx}(x, G(E(x)))]$
ALAD	$\min_{G, E} \max_{D_{xz}, D_{xx}, D_{zz}} V_{ALAD}(D_{xz}, D_{xx}, D_{zz}, E, G)$ $= V_{ALICE} + \mathbb{E}_{z \sim p(z)} [\log(D_{zz}(z, z)) + \log(1 - D_{zz}(z, E(G(z))))]$
CALAD	$\min_{G, E} \max_{D_{xxzz}, D_{xz}, D_{xx}, D_{zz}} V_{CALAD}(D_{xxzz}, D_{xz}, D_{xx}, D_{zz}, E, G)$ $= V_{ALAD} + \mathbb{E}_{x \sim q(x)} [\log D_{xxzz}(x, x, E(x), E(x))]$ $+ \mathbb{E}_{x \sim q(x)} [1 - \log D_{xxzz}(x, G(E(x)), E(x), E(G(E(x))))]$
RALAD	$\min_{G, E} \max_{D_{xz}, D_{xx}, D_{zz}} V_{RALAD}(D_{xz}, D_{xx}, D_{zz}, E, G)$ $= V_{ALAD} + \mathbb{E}_{x \sim \sigma(x)} [\log(1 - D_{xz}(x, E(x)))]$
RCALAD	$\min_{G, E} \max_{D_{xxzz}, D_{xz}, D_{xx}, D_{zz}} V_{RCALAD}(D_{xxzz}, D_{xz}, D_{xx}, D_{zz}, E, G)$ $= V_{ALAD} + \mathbb{E}_{x \sim \sigma(x)} [\log(1 - D_{xz}(x, E(x)))]$ $+ \mathbb{E}_{x \sim q(x)} [\log D_{xxzz}(x, x, E(x), E(x))]$ $+ \mathbb{E}_{x \sim q(x)} [1 - \log D_{xxzz}(x, G(E(x)), E(x), E(G(E(x))))]$

نتایج عملی که در فصل چهارم مشاهده کردیم بیانگر کارایی مدل RCALAD در زمینه تشخیص ناهنجاری است. از مجموع مطالب گفته شده تا اینجا می‌توان چنین برداشت کرد استفاده توام از اطلاعات موجود در ساختارهای تقابلی سبب بهبود عملکرد آن‌ها می‌شود. علاوه بر این وجود توزیع  $\sigma(x)$  که

مستقل از توزیع داده ناهنجار است سبب می‌شود که مدل به سمت تولید بازسازی‌ها در فضای داده ناهنجار سوق داده شود.

## ۵-۲- کارهای آتی

در کارهای آتی تلاش خواهد شد با روش‌های معرفی شده در مقاله [40] آموزش شبکه به شکل هر چه بهتر صورت پذیرد. در این کار با استفاده از تطبیق ویژگی به جای آن که شبکه مولد روی تمایزگر آموزش بیش از حد ببیند، تلاش می‌شود تا آمارگان توزیع داده ورودی نیز به شبکه مولد آموزش داده شود. علاوه بر این با استفاده از روش تمایز کوچک دسته‌ای<sup>۱</sup> مولد را مجبور به تولید خروجی‌های متفاوت خواهیم کرد تا کار شبکه تمایزگر سخت‌تر و روند آموزش بهبود یابد. میانگین‌گیری تاریخی<sup>۲</sup> از دور باطل حول یک نقطه بهینه جلوگیری می‌کند و انتظار می‌رود در صورت استفاده از این روش به بهینه محلی مناسب‌تری دست یابیم. همچنین روش نرمال‌سازی مجازی دسته<sup>۳</sup> سبب می‌شود تا نمونه‌های موجود در یک دسته مستقل از هم شوند و در نتیجه روند بهینه‌سازی شبکه عصبی بهبود یابد. کارایی این روش روی DCGAN ثابت شده است. انتظار می‌رود با پیاده‌سازی روش‌های نام برده شده در این بخش، روند تشخیص ناهنجاری به حالت بهینه نزدیک‌تر شود.

روش دیگر که برای بهبود روش آموزش در مدل RCALAD پیشنهاد می‌شود، یادگیری ضریب اهمیت برای هر یک از تمایزگرها با توجه به جنس مسئله است. در واقع با توجه به نتایجی که در بخش ۴-۵-۳ بدست آمد و دیده شد که میزان اهمیت تمایزگرها در مسائل از جنس مختلف (داده جدولی و تصویری) متفاوت است، بنظر می‌رسد که با یادگیری ضریب اهمیت برای هر یک از تمایزگرها در تابع هزینه و استفاده از آن‌ها برای محاسبه امتیاز ناهنجاری، می‌توان به دقت‌های بالاتری دست یافت. دیگر ایده به کار گرفته شده در حوزه تشخیص ناهنجاری در سال‌های اخیر، استفاده از داده‌های کمکی یعنی نمونه‌های ناهنجار شناخته شده (هرچند تعداد آن‌ها بسیار کم باشد) با استفاده از روش‌های یادگیری روی داده‌های

<sup>۱</sup> Minibatch discrimination

<sup>۲</sup> Historical averaging

<sup>۳</sup> Virtual batch normalization

نامتوازن<sup>۴</sup> می‌باشد. در [۴۶] دو تابع هزینه جدید با هدف یادگیری نامتوازن در کاربرد تشخیص ناهنجاری در شبکه‌های مولد تقابلی با نام‌های Patch loss و Anomaly adversarial loss معرفی شده که هر دوی این تابع‌ها قابلیت به کارگیری در چارچوب معرفی شده RCALAD را دارند.

---

<sup>۴</sup> Imbalanced

## منابع و مراجع

- [1] X. Shu, L. Cheng, and S. J. Stolfo, "Anomaly Detection as a Service".
- [2] D. M. Hawkins, *Identification of Outliers*. Netherlands: Springer, 1980. doi: 10.1007/978-94-015-3994-4.
- [3] C. Jiang, J. Song, G. Liu, L. Zheng, and W. Luan, "Credit Card Fraud Detection: A Novel Approach Using Aggregation Strategy and Feedback Mechanism," *IEEE Internet of Things Journal*, vol. 5, no. 5, pp. 3637–3647, 2018, doi: 10.1109/JIOT.2018.2816007.
- [4] X. Dai and M. Bikdash, "Distance-based outliers method for detecting disease outbreaks using social media," *Conference Proceedings - IEEE SOUTHEASTCON*, vol. 2016-July, 2016, doi: 10.1109/SECON.2016.7506752.
- [5] S. A. Haque, M. Rahman, and S. M. Aziz, "Sensor anomaly detection in wireless sensor networks for healthcare," *Sensors (Switzerland)*, vol. 15, no. 4, pp. 8764–8786, 2015, doi: 10.3390/s150408764.
- [6] H. S. Wu, "A survey of research on anomaly detection for time series," *2016 13th International Computer Conference on Wavelet Active Media Technology and Information Processing, ICCWAMTIP 2017*, no. 1, pp. 426–431, 2017, doi: 10.1109/ICCWAMTIP.2016.8079887.
- [7] K. Choi, J. Yi, C. Park, and S. Yoon, "Deep Learning for Anomaly Detection in Time-Series Data: Review, Analysis, and Guidelines," *IEEE Access*, vol. 9, pp. 120043–120065, 2021, doi: 10.1109/ACCESS.2021.3107975.
- [8] I. Ruts and P. J. Rousseeuw, "Computing depth contours of bivariate point clouds," *Computational Statistics and Data Analysis*, vol. 23, no. 1, pp. 153–168, 1996, doi: 10.1016/S0167-9473(96)00027-8.
- [9] S. Papadimitriou, H. Kitagawa, P. B. Gibbons, and C. Faloutsos, "LOCI: Fast outlier detection using the local correlation integral," *Proceedings - International Conference on Data Engineering*, pp. 315–326, 2003, doi: 10.1109/ICDE.2003.1260802.
- [10] F. Sönmez, M. Zontul, O. Kaynar, and H. Tutar, "Anomaly Detection Using Data Mining Methods in IT Systems: A Decision Support Application,"

- Sakarya University Journal of Science*, vol. 22, no. 4, pp. 1–1, 2018, doi: 10.16984/sofenbilder.365931.
- [11] G. Muruti, F. A. Rahim, and Z. A. Bin Ibrahim, “A survey on anomalies detection techniques and measurement methods,” *2018 IEEE Conference on Application, Information and Network Security, AINS 2018*, no. 1, pp. 81–86, 2019, doi: 10.1109/IISA.2018.8631436.
- [12] M. Ahmed, A. Naser Mahmood, and J. Hu, “A survey of network anomaly detection techniques,” *Journal of Network and Computer Applications*, vol. 60, pp. 19–31, 2016, doi: 10.1016/j.jnca.2015.11.016.
- [13] T. Schlegl, P. Seeb, S. Waldstein, U. Schmidt-Erfurth, and G. Langs, “Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery,” *International Conference on Information Processing in Medical Imaging*, vol. 2, pp. 146–157, 2017, doi: 10.1007/978-3-319-59050-9.
- [14] V. Dumoulin *et al.*, “Adversarially learned inference,” *5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings*, pp. 1–18, 2017.
- [15] H. Issa and M. A. Vasarhelyi, “Application of Anomaly Detection Techniques to Identify Fraudulent Refunds,” *SSRN Electronic Journal*, 2012, doi: 10.2139/ssrn.1910468.
- [16] R. Kaur and S. Singh, “A survey of data mining and social network analysis based anomaly detection techniques,” *Egyptian Informatics Journal*, vol. 17, no. 2, pp. 199–216, 2016, doi: 10.1016/j.eij.2015.11.004.
- [17] N. Görnitz, M. Kloft, K. Rieck, and U. Brefeld, “Toward supervised anomaly detection,” *Journal of Artificial Intelligence Research*, vol. 46, pp. 235–262, 2013, doi: 10.1613/jair.3623.
- [18] R. N. Reza Hassanzadeh, “A SemiSupervised GraphBased Algorithm for Detecting Outliers in OnlineSocialNetworks,” pp. 577–582, 2013.
- [19] D. Kwon, H. Kim, J. Kim, S. C. Suh, I. Kim, and K. J. Kim, “A survey of deep learning-based network anomaly detection,” *Cluster Computing*, vol. 22, pp. 949–961, 2019, doi: 10.1007/s10586-017-1117-8.

- [20] Z. Zhao, C. K. Mohan, and K. G. Mehrotra, "Adaptive sampling and learning for unsupervised outlier detection," *Proceedings of the 29th International Florida Artificial Intelligence Research Society Conference, FLAIRS 2016*, pp. 460–465, 2016.
- [21] D. Digitalcommons@uri and Y. Chae, "Representing Statistical Network-Based Anomaly Detection by Representing Statistical Network-Based Anomaly Detection by Using Trust Using Trust," 2017.
- [22] M. A. Rassam, A. Zainal, and M. A. Maarof, "Advancements of data anomaly detection research in Wireless Sensor Networks: A survey and open issues," *Sensors (Switzerland)*, vol. 13, no. 8, pp. 10087–10122, 2013, doi: 10.3390/s130810087.
- [23] G. Thatte, U. Mitra, and J. Heidemann, "Parametric methods for anomaly detection in aggregate traffic," *IEEE/ACM Transactions on Networking*, vol. 19, no. 2, pp. 512–525, 2011, doi: 10.1109/TNET.2010.2070845.
- [24] J. Wu, W. Zeng, and F. Yan, "Hierarchical Temporal Memory method for time-series-based anomaly detection," *Neurocomputing*, vol. 273, pp. 535–546, 2018, doi: 10.1016/j.neucom.2017.08.026.
- [25] S. Zou, Y. Liang, H. V. Poor, and X. Shi, "Unsupervised nonparametric anomaly detection: A kernel method," *2014 52nd Annual Allerton Conference on Communication, Control, and Computing, Allerton 2014*, pp. 836–841, 2014, doi: 10.1109/ALLERTON.2014.7028541.
- [26] D. B. Araya, K. Grolinger, H. F. ElYamany, M. A. M. Capretz, and G. Bitsuamlak, "An ensemble learning framework for anomaly detection in building energy consumption," *Energy and Buildings*, vol. 144, pp. 191–206, 2017, doi: 10.1016/j.enbuild.2017.02.058.
- [27] M. H. Bhuyan, D. K. Bhattacharyya, and J. K. Kalita, "Network anomaly detection: Methods, systems and tools," *IEEE Communications Surveys and Tutorials*, vol. 16, no. 1, pp. 303–336, 2014, doi: 10.1109/SURV.2013.052213.00046.
- [28] M. S. Mohd Pozi, M. N. Sulaiman, N. Mustapha, and T. Perumal, "Improving Anomalous Rare Attack Detection Rate for Intrusion Detection System Using Support Vector Machine and Genetic Programming," *Neural Processing Letters*, vol. 44, no. 2, pp. 279–290, 2016, doi: 10.1007/s11063-015-9457-y.

- 
- [29] R. Ul Islam, M. S. Hossain, and K. Andersson, "A novel anomaly detection algorithm for sensor data under uncertainty," *Soft Computing*, vol. 22, no. 5, pp. 1623–1639, 2018, doi: 10.1007/s00500-016-2425-2.
  - [30] A. H. Moghaddam, M. H. Moghaddam, and M. Esfandyari, "Stock market index prediction using artificial neural network," *Journal of Economics, Finance and Administrative Science*, vol. 21, no. 41, pp. 89–93, 2016, doi: 10.1016/j.jefas.2016.07.002.
  - [31] I. Goodfellow *et al.*, "Generative Adversarial Nets," in *Advances in Neural Information Processing Systems*, Oct. 2014, pp. 2672–2680. doi: 10.1109/ICCVW.2019.00369.
  - [32] H. Zenati, M. Romain, C. S. Foo, B. Lecouat, and V. Chandrasekhar, "Adversarially Learned Anomaly Detection," *Proceedings - IEEE International Conference on Data Mining, ICDM*, vol. 2018-Novem, pp. 727–736, 2018, doi: 10.1109/ICDM.2018.00088.
  - [33] B. Schölkopf, R. Williamson, A. Smola, J. Shawe-Taylor, and J. Piatt, "Support vector method for novelty detection," *Advances in Neural Information Processing Systems*, pp. 582–588, 2000.
  - [34] R. A. Yeh, C. Chen, T. Yian Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with deep generative models," *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, vol. 2017-Janua, pp. 6882–6890, 2017, doi: 10.1109/CVPR.2017.728.
  - [35] T. Schlegl, P. Seeböck, S. M. Waldstein, G. Langs, and U. Schmidt-Erfurth, "f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks," *Medical Image Analysis*, vol. 54, pp. 30–44, 2019, doi: 10.1016/j.media.2019.01.010.
  - [36] H. Zenati, C. S. Foo, B. Lecouat, G. Manek, and V. R. Chandrasekhar, "Efficient GAN-Based Anomaly Detection," 2018, [Online]. Available: <http://arxiv.org/abs/1802.06222>
  - [37] C. Li *et al.*, "ALICE : Towards Understanding Adversarial Learning for Joint Distribution Matching arXiv : 1709 . 01215v2 [ stat . ML ] 5 Nov 2017," no. Nips, pp. 1–22, 2017.



- [38] Z. Yang, I. S. Bozchalooi, and E. Darve, "Regularized Cycle Consistent Generative Adversarial Network for Anomaly Detection".
- [39] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," 2018.
- [40] T. Salimans, I. Goodfellow, V. Cheung, A. Radford, and X. Chen, "Improved Techniques for Training GANs," pp. 1–10.
- [41] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath, "Generative Adversarial Networks: An Overview," *IEEE Signal Processing Magazine*, vol. 35, no. 1, pp. 53–65, 2018, doi: 10.1109/MSP.2017.2765202.
- [42] S. Zhai, Y. Cheng, W. Lu, and Z. Zhang, "Deep structured energy based models for anomaly detection," *33rd International Conference on Machine Learning, ICML 2016*, vol. 3, pp. 1742–1751, 2016.
- [43] F. Tony Liu, K. Ming Ting, and Z.-H. Zhou, "Isolation Forest ICDM08," *Icdm*, 2008, [Online]. Available: <https://cs.nju.edu.cn/zhoush/zhoush.files/publication/icdm08b.pdf%0Ahttps://cs.nju.edu.cn/zhoush/zhoush.files/publication/icdm08b.pdf?q=isolation-forest>
- [44] A. Makhzani and B. Frey, "Winner-take-all autoencoders," *Advances in Neural Information Processing Systems*, vol. 2015-Janua, pp. 2791–2799, 2015.
- [45] L. Ruff *et al.*, "Deep one-class classification," *35th International Conference on Machine Learning, ICML 2018*, vol. 10, pp. 6981–6996, 2018.
- [46] J. Kim, K. Jeong, H. Choi, and K. Seo, "GAN-Based Anomaly Detection In Imbalance Problems." *Lecture Notes in Computer Science*, 12540 LNCS, 128–145. [https://doi.org/10.1007/978-3-030-65414-6\\_118](https://doi.org/10.1007/978-3-030-65414-6_118), 2020.

## فهرست واژگان انگلیسی به فارسی

استخراج استثنا.....Exception mining	A
تنظیمات آزمون.....Experimental Setup	Ablation studies.....مطالعه فرسایشی
<b>F</b>	Activation function.....تابع فعال ساز
منفی کاذب.....False negative	Area Under Curve Receiver Operating
نرخ مثبت کاذب.....False positive rate	Characteristics.....مساحت زیر نمودار مشخصه
تطبیق ویژگی.....Feature matching	عملکرد
خطای تطبیق ویژگی.....Feature matching loss	Autoencoder.....خودکدگذار
انتشار رو به جلو.....Feed forward	<b>B</b>
<b>G</b>	Backpropagation.....پس انتشار
Generative Adversarial Networks.....شبکه مولد	Batch.....دسته
تقابلی	Bias.....متمایل
Gradient descent.....گرادیان نزولی	Bidirectional.....دوطرفه
<b>H</b>	<b>C</b>
Historical averaging.....میان گیری تاریخی	Categorical.....اسمی
Hyper Parameter.....فراپارامتر	Clustering.....خوشه بندی
<b>I</b>	Complete Cycle Consistency.....چرخه پایداری کامل
Invert.....معکوس کردن	Conditional entropy.....آنترپی شرطی
Iterative.....مبتنی بر تکرار	Continuous.....پیوسته
<b>J</b>	Cycle consistency.....چرخه پایداری
distribution Joint.....توزیع توام	<b>D</b>
<b>K</b>	Decision tree.....درخت تصمیم
kernel.....هسته	Deviation detection.....تشخیص انحراف
<b>L</b>	Discrimination loss.....خطای تمایزگر
Latent space.....فضای نهفته	Discriminator models.....مدل های تمایزگر
Log-Likelihood.....لگاریتم درست نمایی	Dropout.....حذف تصادفی
	<b>E</b>
	Encoder.....کدگذار

Reinforcement learning..... یادگیری تقویتی	<b>M</b>
Remote sensing..... سنجش از راه دور	Maximum likelihood estimation..... بیشینه تخمین
Residual..... باقی مانده	درست‌نمایی
Residual Loss..... خطای باقی مانده	Mean squared error..... میانگین مجموع مربعات خطا
Robustness..... مقاومت	Minibatch..... کوچک دسته‌ای
<b>S</b>	Minibatch discrimination..... تمایز کوچک دسته‌ای
Scaler..... اسکالر	Minimax..... بیشینه-کمینه
Sigmoid cross entropy..... آنترپی متقاطع سیگموئید	Mutual information..... اطلاعات متقابل
Smooth..... هموار	<b>N</b>
Support vector machine..... ماشین بردار پشتیبان	Novelty detection..... شناسایی نوآوری
<b>T</b>	<b>O</b>
Tabular..... جدولی	Outlier detection..... تشخیص داده پرت
Threshold..... حد آستانه	Overfitting..... بیش‌برازش
True negative..... منفی صحیح	<b>P</b>
True Negative Rate..... نرخ منفی صحیح	Piecewise linear units..... واحدهای خطی تکه‌ای
True positive..... مثبت صحیح	probability Posterior..... احتمال پسین
Tuple..... دوتایی	Precision..... صحت
<b>V</b>	Principal component analysis..... تحلیل مولفه اصلی
Variational autoencoder..... خودکدگذار متغیر	<b>R</b>
Virtual batch normalization..... نرمال کردن مجازی	Recall..... بازیابی
دسته	Reconstruction..... بازسازی

## Abstract

Anomaly detection is a significant and hence well studied problem in field of data analysis which is used in a wide range of applications such as fraud detection, medical application and cyber security systems. Despite the existence of statistical and machine learning-based methods, designing effective models for anomaly detection in complex high-dimensional data space remains a major challenge. As generative adversarial networks are able to handle this challenge and model the complex high-dimensional distribution of real-world data. as a result it can operate promisingly in field of anomaly detection. In this work we propose CALAD<sup>1</sup>, RALAD<sup>2</sup> and RCALAD<sup>3</sup> models to detect anomalies. Our reconstruction based method reconstruct the input data through generative network and compute reconstruction error to find anomalous example. In the CALAD model defining new variable  $\hat{\mathbf{z}}_{\hat{\mathbf{x}}}$  and using an innovative discriminator  $\mathbf{D}_{xxzz}$ , complete cycle consistency between input space and hidden space is established. Poor reconstruction for anomalous data is a prerequisites in reconstruction based models. RALAD aims to bias the model towards normal data distribution. This bias leads to poor reconstruction of anomalous data and as a result the distance between anomalous input data and its reconstruction will increase. With combining these two ideas, comprehensive RCALAD model is proposed. In addition, two new anomaly score are proposed which provide high resolution power in contrast to other anomaly scores. Finally, experimental results demonstrate the effectiveness of our approach by showing the results if outperforming the current state of the art approaches in terms of the average area under the ROC<sup>4</sup> and F1-score.

**Key Words:** anomaly detection, machine learning, generative adversarial networks, reconstruction error, anomaly score

---

<sup>1</sup> Complete Adversarially Learned Anomaly Detection

<sup>2</sup> Regularized Adversarially Learned Anomaly Detection

<sup>3</sup> Regularized Complete Adversarially Learned Anomaly Detection

<sup>4</sup> Receiver Operating Characteristic



**Amirkabir University of Technology**  
**(Tehran Polytechnic)**

**Department of Computer Engineering and Information Technology**

**Master Thesis**

# **Anomaly Detection with Generative Adversarial Network**

**By**  
**Zahra Dehghanian**

**supervisor**

**Dr. Mohammad Rahmati**  
**Dr. Maryam Amirmazlaghani**

**May 2022**