

رویکرد مبتنی بر یادگیری عمیق برای شناسایی حمله تزریق داده جعلی در شبکه‌های هوشمند

علیرضا سلطانی

کارشناسی ارشد برق، دانشکده مهندسی برق و کامپیوتر، دانشگاه تربیت مدرس، تهران، ایران

alireza_soltani@modares.ac.ir

دکتر حمیدرضا بقایی کاشی

استادیار گروه برق، دانشکده مهندسی برق و کامپیوتر، دانشگاه تربیت مدرس، تهران، ایران (نویسنده مسئول)

hrbaghaee@modares.ac.ir

دکتر محمودرضا حقی فام

استاد گروه برق، دانشکده مهندسی برق و کامپیوتر، دانشگاه تربیت مدرس، تهران، ایران

haghifam@modares.ac.ir

چکیده

با گسترش شبکه‌های توزیع هوشمند و افزایش تهدیدات سایبری، ایمن‌سازی سیستم‌های قدرت به ضرورتی حیاتی تبدیل شده است. حملات تزریق داده جعلی از جمله مخرب‌ترین تهدیدات هستند که با دور زدن مکانیزم‌های تشخیص سنتی، ثبات سیستم را به خطر می‌اندازند. در این پژوهش، یک مدل تشخیصی مبتنی بر **Graph Autoencoder** ارائه شده است که با رویکردی پیشگیرانه به شناسایی این حملات در شبکه توزیع ۳۳ IEEE باسه می‌پردازد. مدل پیشنهادی در چارچوب یک معماری یادگیری عمیق طراحی گردیده و با بهره‌گیری از قابلیت‌های شبکه عصبی گرافی، الگوهای عادی عملکرد شبکه را آموخته و از طریق محاسبه خطای بازسازی، حملات را شناسایی می‌نماید. نتایج حاکی از آن است که مدل ارائه‌شده بالای ۹۸ درصد قادر به تشخیص مؤثر ناهنجاری‌ها بوده و منجر به بهبود عملکرد سیستم حفاظتی می‌گردد.

تاریخ دریافت:

۱۴۰۴/۰۹/۱۲

تاریخ پذیرش:

۱۴۰۴/۱۱/۰۶

کلمات کلیدی:

شبکه هوشمند
تشخیص ناهنجاری
حمله سایبری
یادگیری عمیق

۱. مقدمه

شبکه قدرت یکی از زیرساخت های بنیادین اقتصادی - اجتماعی است که سلامت، پایداری و دسترسی پیوسته آن برای عملکرد جوامع مدرن حیاتی است. در سال های اخیر، تحول فناوری اینترنت صنعتی اشیاء (IIoT) و ظهور شبکه های هوشمند^۲ موجب ادغام گسترده فناوری اطلاعات و ارتباطات با شبکه قدرت شده است؛ این ادغام به طور چشمگیری کارآمدی و قابلیت اطمینان شبکه را افزایش داده است ولی همزمان چالش های امنیتی و سایبری نوینی نیز ایجاد کرده است [۱].

یکی از ستون های عملیاتی مرکز کنترل سیستم قدرت، فرآیند تخمین حالت^۳ است که از داده های خام گردآوری شده توسط سیستم های SCADA^۴ وضعیت سیستم را محاسبه می کند و نتایج آن ورودی بسیاری از تابع های اجرایی مانند پخش بار اقتصادی و تحلیل پایداری است. از این رو، هرگونه دست کاری یا اختلال در نتایج برآورد حالت می تواند به خطاهای تصمیم گیری، تزریق بار نامناسب، یا در نهایت وقوع رویدادهای بحرانی (مانند تلفات و قطعی) منجر شود [۲].

حملات FDI^۵

حملات تزریق داده جعلی (FDI) نوعی حمله ای است که به طور خاص بر سامانه برآورد حالت هدف گیری می کند؛ مهاجم با تغییر مقادیر اندازه گیری ارسالی از سوی شمارنده ها یا PMU^۶ ها می تواند مقدارهای ورودی تخمین گر حالت را به گونه ای دست کاری کند که خطای ساختاری در برآورد حالت ایجاد شود، در حالی که مکانیزم های متداول تشخیص داده های بد (BDD^۷) را دور می زند. کارهای ابتدایی نشان دادند که با داشتن دانشی کافی از مدل شبکه، مهاجم می تواند حمله ای کاملاً پنهانی^۸ بسازد که با ساختار ریاضی برآورد حالت سازگار باشد و به همین دلیل به سختی قابل کشف است. بعداً نشان داده شد حتی با دسترسی ناقص یا دانش جزئی درباره پیکربندی شبکه، حملات قابل تأثیر و تا حدی غیرقابل کشف قابل طراحی هستند و بنابراین دفاع صرفاً مبتنی بر BDD ناکافی است [۳].

در این مقاله با استفاده روش های مبتنی بر هوش مصنوعی برای پیشگیری و تشخیص حملات تزریق داده جعلی در شبکه نمونه (IEEE ۳۳ bus) می پردازیم.

مقایسه با روش های ارائه شده در مراجع:

در [۴] یک حمله ی FDI جعبه سیاه را ارائه می دهد که در آن بردار حمله با استفاده از یک شبکه ی مولد تخاصمی (GAN) تولید می شود تا اختلال هایی مشابه داده های عادی سیستم ایجاد کند. سپس، این بردار حمله از طریق دستکاری ماژول های اندازه گیری شبکه به سیستم تزریق می شود. این حمله بلافاصله پس از وقوع یک خطا در سیستم اجرا می شود تا هم پنهان ماندن از مکانیزم های تشخیص داده بد افزایش یابد و هم

^۱ Industrial Internet of Things

^۲ Smart grids

^۳ State estimation

^۴ Supervisory Control And Data Acquisition

^۵ False Data Injection Attack

^۶ Phasor Measurement Unit

^۷ Bad Data Detection

^۸ Stealthy

اثر تخریبی آن بیشینه شود، در [۵]، [۶]، [۷] و [۸] با استفاده از دانش ساختار شبکه (ماتریس ژاکوبین) حمله پس از دور زدن BDD ساخته شده است. تخمین گر حالت حمله را تشخیص نمی‌دهد. در جدول ۱ این منابع از نظر دقت تشخیص حملات از داده های سالم بررسی می‌شوند.

مدل	شبکه	دقت	فراخوانی	امتیاز F ₁
AlexNet ^[۴]	IEEE ۵۷-bus	۹۸٫۴۰٪	-	-
AlexNet ^[۴]	IEEE ۱۴۵-bus	۹۷٫۱۰٪	-	-
SFTCN ^۱ [۵]	IEEE ۱۴-bus	۹۹٫۶۸٪	۹۸٫۸۳٪	۹۹٫۳۶٪
SFTCN ^[۵]	IEEE ۱۱۸-bus	۹۷٫۲۰٪	۹۱٫۸۵٪	۹۴٫۲۶٪
A-BiTG ^۱ [۶]	IEEE ۱۴-bus	۹۵٫۴۷٪	۹۷٫۲۷٪	۹۶٫۳۷٪
GMM ^۱ [۷]	IEEE ۱۴-bus	۹۵٪<	-	-
LSTM ^۱ [۸]	IEEE ۳۰-bus	۹۴٫۲۴٪	-	-
LSTM ^[۸]	IEEE ۱۴-bus	۹۴٫۵۲٪	-	-

جدول (۱) مقایسه دقت تشخیص روش های مختلف

در روش پیشنهادی در این پروژه بدون دانستن مشخصات شبکه (ماتریس ژاکوبین) با استفاده از روش زیرفضا حمله طراحی شده است و حمله قادر به عبور از BDD مرسوم شبکه است، از مدل تست های متعدد گرفته شده است و دقت تشخیص حمله تقریباً بالاتر از ۹۹ درصد است و حتی اگر یکی از اندازه گیری ها مورد حمله باشد، مدل قادر به تشخیص آن است.

۲. مبانی نظری

در اغلب سیستم های قدرت، وضعیت سیستم با یک بردار حالت x (مثلاً بردار ولتاژها و زاویهها) مدل می‌شود و سنجشها (اندازه گیری ها) Z توسط رابطه ۱ نمایش داده می‌شود.

$$z = H(x) + e \quad \text{رابطه ۱:}$$

که H و e به ترتیب نشان دهنده ماتریس ژاکوبین شبکه و بردار خطای اندازه‌گیری تصادفی است که معمولاً به عنوان توزیع گاوسی مستقل با میانگین صفر (طبق رابطه ۲) فرض می‌شود [۱].

$$\mathcal{N}(\cdot, \sigma^2 I) \sim e \quad \text{رابطه ۲:}$$

حمله تزریق داده جعلی را می‌توان به صورت افزودن یک بردار حمله a به اندازه‌گیری‌ها مدل کرد:

$$H(x) + e + a = a + z = \hat{z} \quad \text{رابطه ۳:}$$

که \hat{z} همان بردار حاصل از حمله است.

^۹ Spatial Feature-based Temporal Convolutional Network

^{۱۰} Attention-based Bidirectional Temporal Gating

^{۱۱} Gaussian Mixture Model

^{۱۲} Long-Short Term Memory

هدف حمله‌ی تزریق داده‌ی جعلی گمراه کردن اپراتور سیستم است تا وضعیت تخمینی و مخدوش شده $x = x + c$ را به‌عنوان یک تخمین معتبر در نظر بگیرد، که در آن $c \neq 0$. انحراف حالت سیستم قدرت است. برای دستیابی به این هدف، مهاجم مقادیر اندازه‌گیری دریافتی در مرکز کنترل را به $\hat{z} = z + a$ تغییر می‌دهد [۱].

روش‌های متعارف تشخیص ناهنجاری مبتنی بر باقیمانده، باقیمانده اندازه‌گیری (R) را با یک آستانه از پیش تعیین شده (τ) مقایسه می‌کنند تا بررسی کنند که آیا اندازه‌گیری‌های بی‌کیفیت (یعنی خراب یا مخدوش شده) در سیستم وجود دارند یا خیر [۹]. در این حالت، آشکارساز^{۱۳} به محض اینکه شرط زیر (رابطه ۴) برقرار شود، وجود یک حمله را اعلام می‌کند:

$$R = z - Hx > \tau \quad \text{رابطه ۴:}$$

برای دور زدن مکانیسم BDD، بردار حمله باید ساختاریافته باشد، به‌گونه‌ای که $a = Hc$. در چنین مواردی (R باقیمانده^{۱۴}) در رابطه ۴ بدون تغییر می‌ماند:

$$\|z + a - H(x + c)\| = \|z - Hx\| \quad \text{رابطه ۵:}$$

و در نتیجه (طبق رابطه ۵) حمله می‌تواند از BDD عبور کند. بر این اساس، اپراتور سیستم قدرت به اشتباه $x + c$ را یک تخمین معتبر در نظر می‌گیرد و بنابراین یک بردار خطای c معرفی می‌شود.

حمله زیر فضا^{۱۵}:

حملات زیرفضا مبتنی بر این ایده هستند که مهاجم می‌تواند بدون دانش کامل از پارامترهای سیستم، تنها با تحلیل داده‌های تاریخی، ساختار شبکه را یاد گرفته و حملات غیرقابل تشخیص طراحی کند [۱۰].

فرض شود که U یک ماتریس پایه برای فضای ستونی $R(H)$ باشد و U_1 زیر ماتریس از U است که با حذف سطرهای مربوط به سنسورهای مورد تهاجم (S_A) به‌دست می‌آید.

شرط امکان‌پذیری حمله غیرقابل مشاهده:

یک حمله غیرقابل مشاهده امکان‌پذیر است اگر و تنها اگر U_1 رتبه ستونی کامل نداشته باشد.

اگر U_1 رتبه کامل داشته باشد، آنگاه تنها بردار v که در $U_1 v = 0$ صدق کند، بردار صفر است.

اگر U_1 رتبه کامل نداشته باشد، آنگاه فضای پوچ $N(U_1)$ دارای بعد مثبت است [۱۰].

$N(U_1)$ همان فضای پوچ ماتریس U_1 است؛ یعنی مجموعه‌ای از تمام بردارهایی که اگر در U_1 ضرب شوند، حاصل بردار صفر می‌شود، اگر این فضا بعد مثبت داشته باشد، امکان ساخت حمله غیرقابل تشخیص وجود دارد.

^{۱۳} Detector

^{۱۴} Residual

^{۱۵} Subspace Attack

الگوریتم طراحی حمله:

$\{Z_1, Z_2, \dots, Z_k\}$ مقادیر بردار اندازه‌گیری‌های گذشته هستند که مهاجم از آن‌ها استفاده می‌کند و مجموعه سنسور‌هایی که مهاجم برای حمله انتخاب می‌کند، S_A هستند.

برای مرکز‌سازی داده‌ها:

$$\bar{Z} = \frac{1}{k} \sum_{i=1}^k Z_i \quad \text{رابطه ۶:}$$

$$Z_c^{(i)} = Z_i - \bar{Z} \quad \text{رابطه ۷:}$$

که $Z_c^{(i)}$ بردار مرکز‌سازی شده اندازه‌گیری‌ها و \bar{Z} بردار میانگین داد‌ها هستند، برای ساده‌سازی فرض می‌کنیم بردارهای نویز e_1, \dots, e_k مستقل و دارای توزیع مشابه^{۱۶} هستند، بردارهای حالت X_1, \dots, X_k نیز وضعیت مشابهی مانند بردارهای نویز دارند و دارای ماتریس کوواریانس معین مثبت Σ_x هستند، بردارهای نویز و بردارهای حالت نیز ناهمبسته هستند. در این صورت، ماتریس کوواریانس Z به صورت زیر محاسبه می‌شود:

$$\Sigma_Z \triangleq E[(Z_i - E[Z_i])(Z_i - E[Z_i])^T] = H\Sigma_x H^T + \sigma^2 I \quad \text{رابطه ۸:}$$

با توجه به اینکه $\hat{\Sigma}_Z$ ماتریس مثبت معین است در نتیجه معکوس پذیر است و چون H^T نیز رتبه کامل سطری دارد، ضرب آن در سمت راست رتبه ستونی را تغییر نمی‌دهد پس:

$$\mathcal{R}(H\hat{\Sigma}_x H^T) = \mathcal{R}(H) \quad \text{رابطه ۹:}$$

طبق رابطه ۹ زیر فضای پوششی (\mathcal{R}) در H معادل زیر فضای پوششی $H\hat{\Sigma}_x H^T$ است.

بنابراین، در عمل می‌توانیم یک ماتریس پایه برای $\mathcal{R}(H)$ را با اعمال تجزیه مقدار منفرد روی ماتریس کوواریانس نمونه‌ای $\hat{\Sigma}_Z$ تخمین بزنیم:

$$\hat{\Sigma}_Z = \frac{1}{k-1} \sum_{i=1}^k Z_c^{(i)} (Z_c^{(i)})^T \quad \text{رابطه ۱۰:}$$

که k تعداد نمونه‌ها و T نشان‌دهنده ترانهاده ماتریس هستند.

محاسبه مقادیر منفرد (SVD)^{۱۷}:

$$\hat{\Sigma}_Z = U\Lambda V^T \quad \text{رابطه ۱۱:}$$

که U ماتریس بردارهای ویژه چپ، Λ ماتریس مقادیر ویژه و V^T ترانهاده ماتریس بردارهای ویژه راست هستند.

^{۱۶} i.i.d: Independent And Identically Distributed

^{۱۷} Singular Value Decomposition

ماتریس \hat{U}_1 حاصل از حذف سطرهای مربوط به S_A در U است، مقادیر منفرد این ماتریس به صورت زیر است:

$$\hat{U}_1 = \tilde{U}\tilde{A}\tilde{V}^T \quad \text{رابطه ۱۲:}$$

اگر حمله امکان پذیر باشد، برای هر بردار غیرصفر $v \in N(\hat{U}_1)$ بردار $a = Uv$ یک بردار حمله غیر قابل مشاهده است [۱۰].

این بردار حمله در هر ردیف از $a \neq 0$ به داده های واقعی شبکه (در اینجا داده های ولتاژ) اضافه می شود و اندازه گیری آن سنسور آلوده (S_A) می شود.

این حمله به گونه ای طراحی شده که فقط داده سنسورهای حمله شده را تغییر دهد ولی در محاسبات باقیمانده سیستم تشخیص سنتی غیر قابل مشاهده باقی بماند، مدل شبکه عصبی ارائه شده با یادگیری وابستگی های پیچیده تر شبکه، می تواند این حمله هوشمند را از طریق خطای بازسازی بالا شناسایی کند.

۳. روش تحقیق

هدف روش، طراحی و ارزیابی یک چارچوب داده محور برای شناسایی و مکان یابی حمله سایبری در سامانه اندازه گیری شبکه توزیع، حملات تزریق داده جعلی است. رویکرد مقاله مبتنی بر یک مدل یادگیرنده گرافی است که ورودی های گره ای را بازسازی می کند و از خطاهای بازسازی برای برچسب گذاری نمونه ها و شناخت ناهنجاری ها استفاده می نماید. این بخش روند کار آزمایشی و روش شناختی را تشریح می کند. [۱۱]

برای مدل سازی و شناسایی ناهنجاری های ناشی از حملات تزریق داده جعلی از چارچوب GAE^{18} مبتنی بر GCN^۹ استفاده شد. در این طراحی، یک رمزگذار جهت یادگیری نمایش نهفته^{۲۰} هر گره و یک رمزگشا^{۲۱} مبتنی بر GCN جهت بازسازی ویژگی های گره ها به کار رفته است. ایده پایه این است که شبکه پس از آموزش روی داده های سالم، الگوهای همبستگی فضایی بین باس ها (گره ها) و ویژگی های آن ها را یاد می گیرد؛ در صورت وقوع حمله، خطای بازسازی گره (ها) افزایش می یابد که می تواند به عنوان شاخص تشخیص و مکان یابی مورد استفاده قرار گیرد [۱۲]. [۱۳]

ماتریس مجاورت A : در شبکه اتصال بین باس ها اگر بین گره (باس) نوز اتصال فیزیکی (خط) وجود داشته باشد؛ $A_{i,j} = 1$ و در غیر این صورت $A_{i,j} = 0$ خواهد شد.

ماتریس خود همراه \hat{A} : طبق رابطه ۱۳ افزودن ماتریس همانی (I) یعنی هر گره خودش نیز در میانگین گیری همسایه ها لحاظ می شود، این عمل اجازه می دهد اطلاعات خود گره به محاسبه ویژگی جدید اضافه شود.

$$\hat{A} = A + I \quad \text{رابطه ۱۳:}$$

^{۱۸} Graph Autoencoder

^{۱۹} Graph Convolutional Network

^{۲۰} embedding

^{۲۱} decoder

$$\widehat{D}_{ii} = \sum_j \widehat{A}_{ij} \quad \text{رابطه ۱۴:}$$

که \widehat{D} جمع سطرهای \widehat{A} ؛ نشان‌دهنده تعداد همسایگان هر گره است.

ماتریس نرمال شده \widehat{A} : طبق رابطه ۱۵ نرمال سازی \widehat{A} باعث می‌شود که از ناپایداری مقیاس جلوگیری شود (از تسلط گره‌های با درجه بالا در انتشار پیام جلوگیری می‌کند) و در بسیاری از پیاده‌سازی‌های GCN این فرم را مبنای لایه کانولوشن قرار می‌دهند.

$$\widehat{A} = \widehat{D}^{-1/2} \widehat{A} \widehat{D}^{-1/2} \quad \text{رابطه ۱۵:}$$

لایه پایه GCN: هر لایه GCN ویژگی‌های گره‌ها را با میانگین‌گیری وزنی از همسایگان و تبدیل خطی $W^{(l)}$ ترکیب می‌کند، سپس تابع غیرخطی σ اعمال می‌شود.

$$H^{(l+1)} = \sigma(\widehat{A} H^{(l)} W^{(l)} + b^{(l)}), \quad H^{(0)} = X \quad \text{رابطه ۱۶:}$$

که در رابطه ۱۶، $H^{(l)}$ ؛ نشان‌دهنده خروجی لایه l ، $W^{(l)}$ ؛ ماتریس وزن لایه l ، $b^{(l)}$ ؛ ماتریس بایاس^{۲۲}، σ ؛ تابع فعال سازی^{۲۳} X ؛ ماتریس ویژگی ورودی، هستند [۱۴].

در این پژوهش از تابع فعال ساز ReLU برای لایه‌ها استفاده شده است.

پیاده‌سازی ارائه‌شده شامل سه کلاس کلیدی است:

GCNEncoder

این مؤلفه مسئول یادگیری نمایش فشرده‌شده گره‌ها در فضای نهفته است. معماری آن شامل دو لایه گراف کانولوشنال (GCN) متوالی است که ورودی اولیه ماتریس ویژگی‌های گره‌های X و ماتریس مجاورت A را دریافت کرده و بردارهای نهفته Z را برای هر گره تولید می‌کند.

GCNDecoder

این بخش مسئول بازسازی ویژگی‌های اولیه از بردارهای نهفته است. در این پیاده‌سازی، رمزگشا از دو لایه گراف کانولوشنال معکوس^{۲۴} تشکیل شده که بردارهای نهفته Z را به مقادیر بازسازی‌شده ولتاژ \widehat{X} نگاشت می‌کند.

GraphAutoEncoder

این کلاس با ترکیب رمزگذار و رمزگشا، مدل کامل را تشکیل می‌دهد، خروجی‌های این مدل شامل بردارهای نهفته Z و مقادیر بازسازی‌شده ولتاژ \widehat{X} است [۱۳]، [۱۵].

^{۲۲} Bias

^{۲۳} Activation Function

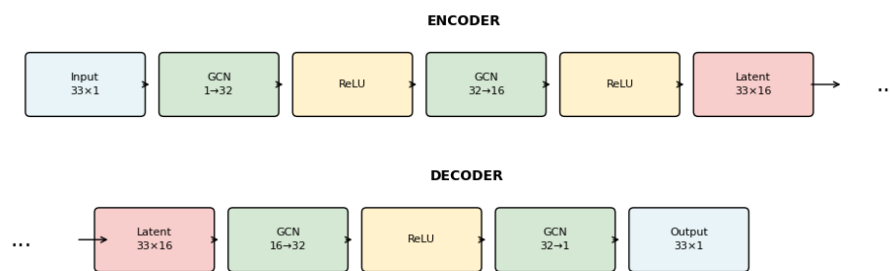
^{۲۴} Deconvolution

تابع هزینه^{۲۵}:

با فرض ثابت بودن توپولوژی شبکه در شرایط عادی، در پیاده‌سازی فعلی تنها بازسازی ویژگی‌ها (ولتاژها) انجام می‌شود؛ بنابراین تابع هزینه اصلی میانگین مربعات خطا^{۲۶} بین X و \hat{X} است:

$$\mathcal{L}_X = \|X - \hat{X}\|_F^2 = \sum_{i=1}^N |x_i - \hat{x}_i|^2 \quad \text{رابطه ۱۷:}$$

که در رابطه ۱۷ نشان‌دهنده مقادیر واقعی ولتاژها، \hat{x}_i مقادیر ولتاژهای ساخته شده توسط مدل و \mathcal{L}_X تابع هزینه هستند.



شکل (۱)، نمودار شبکه عصبی GAE پروژه

شکل ۱، نمودار شبکه عصبی به کار رفته در این پروژه را نشان می‌دهد که، Encoder از دو لایه GCN با توابع فعال‌ساز ReLU تشکیل شده که ابعاد ویژگی‌ها را از ۱ به ۳۲ و سپس به ۱۶ کاهش می‌دهد، Decoder نیز شامل دو لایه GCN است که فرآیند معکوس را برای بازسازی ویژگی‌های اولیه انجام می‌دهد.

(نمایش) نهفته^{۲۷} (z):

بردار نهفته‌ی هر گره است که توسط Encoder تولید می‌شود خروجی z با ابعاد، ابعاد بردار نهفته در تعداد نودها (در پیاده‌سازی ارائه شده ۳۳×۱۶) این بردار بیانگر فشردگی اطلاعات محلی هر باس بوده و از طریق تجمیع ویژگی‌های همسایگان در لایه‌های GCN، مؤلفه‌های ساختاری و ویژگی‌محور شبکه را کدگذاری می‌کند، Decoder همین نمایش نهفته را برای بازسازی ولتاژ گره‌ها به کار می‌گیرد، فرآیند آموزش با کمینه‌سازی خطای بازسازی، به مدل امکان می‌دهد تا یک نمایش فشردگی و معنادار از وضعیت عادی شبکه را فرا گیرد.

^{۲۵} Loss Function

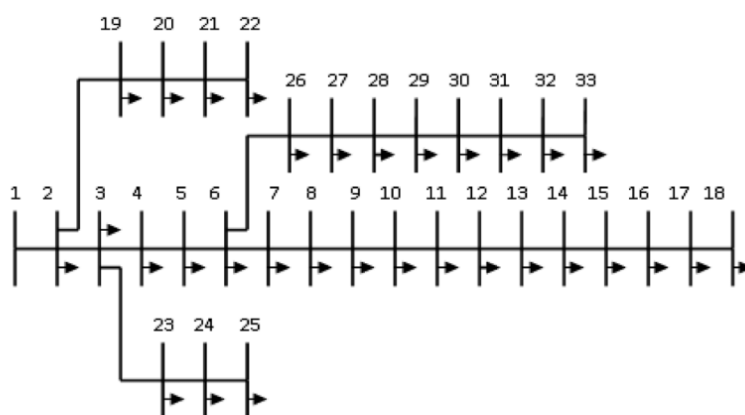
^{۲۶} MSE: Mean Square Error

^{۲۷} Latent

در این پروژه برای تشخیص حملات، این شبکه ابتدا الگوهای عادی عملکرد شبکه را از داده‌های آموزشی یاد می‌گیرد. هنگامی که داده‌های حمله به مدل ارائه می‌شوند، به دلیل انحراف از الگوهای عادی آموخته شده، خطای بازسازی مدل به طور قابل توجهی افزایش می‌یابد. این خطای بالا به عنوان شاخصی برای تشخیص حملات مورد استفاده قرار می‌گیرد [۱۶].

۴. توصیف داده‌ها

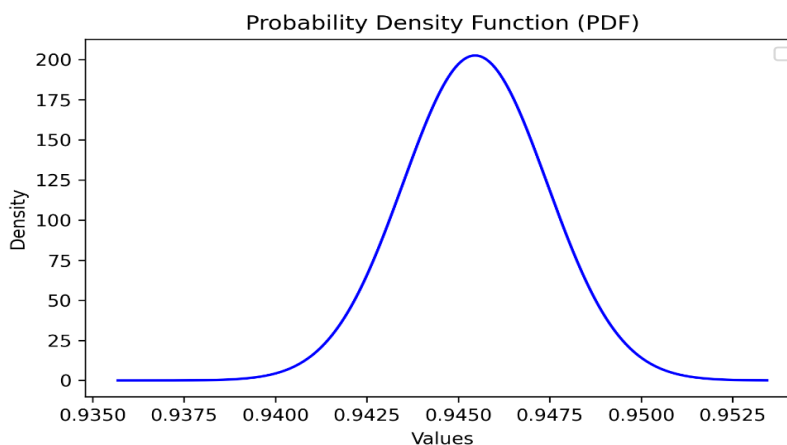
مجموعه داده مورد استفاده در این پژوهش شامل مجموعه‌ای از سناریوهای شبیه‌سازی شده عملیاتی از یک شبکه توزیع نمونه (شبکه IEEE ۳۳ باسه) است.



شکل (۲). شبکه ۳۳ باسه IEEE

شکل ۲ شبکه ۳۳ باسه IEEE و نحوه اتصال باس‌ها به یکدیگر و بارهای شبکه را نشان می‌دهد. تمامی پارامترهای این شبکه (بار، امپدانس خطوط، طول خطوط و...) طبق استاندارد IEEE در نرم افزار دیگ سایلنت پیاده‌سازی شده‌اند.

هر سناریو نمایانگر وضعیت شبکه در یک زمان ثابت است و شامل مشاهدات گره‌ای و یالی استاندارد می‌باشد. نحوه بدست آمدن سناریو‌ها از تغییر بارهای شبکه به صورت توزیع نرمال در بار مرجع شبکه IEEE است (بارها با ضریب ۰.۶ تا ۱.۱ برابر بار اصلی شبکه هستند و توزیع نرمال با میانگین ۰.۸۵ و انحراف از معیار ۰.۱۲۵ دارند). داده‌ها به صورت مجموعه‌ای از فایل‌های جداسازی شده نگهداری می‌شوند و برای ارزیابی مدل به مجموعه‌های آموزش، اعتبارسنجی و تست تقسیم شده‌اند. در این پژوهش از ولتاژهای شبکه و در ۶۰۰۰ سناریو نمونه برداری شده است، چون شبکه ۳۳ باسه است، بردار ورودی شبکه به صورت 33×6000 است، و داده‌ها بدون تغییر از پخش بار بدست می‌آیند. (به طوری که هر سطر نمایانگر یک سناریو و هر ستون مربوط به ولتاژ یک باس خاص می‌باشد). داده‌ها در نرم افزار اکسل و به صورت per-unit شده ذخیره شده‌اند. در شکل ۳ تابع چگالی احتمال ولتاژهای بدست آمده از سناریو‌ها مشاهده می‌شود، این تابع چگالی احتمال مربوط به باس ۹ شبکه تست می‌باشد.



شکل (۳). تابع چگالی احتمال ولتاژ باس ۹

نمونه‌های آلوده به صورت کنترل شده تولید شدند: برای هر نمونه چند سنسور (گره) به طور تصادفی انتخاب و اختلال جمع‌شونده‌ای (طبق روش زیر فضا) با دامنه‌ای منطقی به اندازه‌گیری‌ها اضافه شد تا مقادیر در چارچوب فیزیکی باقی بمانند. (توسط سیستم‌های سستی تشخیص خطا یا داده آلوده تشخیص داده نشوند). الگوهای متنوعی و با شدت‌های مختلف ساخته شد و داده‌های تولید شده به عنوان داده‌های آلوده به حمله تزریق داده جعلی در تست مدل استفاده شدند، حداکثر دامنه حمله اعمال شده ۰٫۱ پریونیت به داده‌های ولتاژ است، زمان شروع حمله به شرایط خاصی بستگی ندارد (مثلاً در حین خطا، قبل و بعد از خطا و...) و در هر زمانی قابل اعمال است.

۵. نتایج اجرای مدل

برای تشخیص ناهنجاری‌های ناشی از حملات تزریق داده جعلی، مدل اصلی مورد استفاده یک Graph Autoencoder است که ورودی‌های گره‌ای و یالی را به فضای پنهان نگاشت و سپس با یک decoder بازسازی مقادیر اصلی را انجام می‌دهد. خطای بازسازی در سطح باس‌ها به عنوان معیار تشخیص حمله استفاده شده است. در این پژوهش ساختار گراف و اتصالات شبکه تست برای شبکه عصبی تعریف و معین شده است. شبکه عصبی با شناخت ساختار و رابطه داده‌های سناریو‌های آموزش، آموزش می‌بیند که در مواجهه با سناریو‌های جدید کدام باس‌ها سالم و کدام باس‌ها مورد حمله تزریق داده جعلی قرار گرفته‌اند. Graph Autoencoder به دلیل توانایی در یادگیری بازنمایی‌های موثر از ساختار و داده‌های شبکه قدرت، و مکانیسم تشخیص ناهنجاری بر اساس خطای بازسازی، گزینه مناسبی برای تشخیص حملات تزریق داده جعلی است. در این شبکه از مدل بدون نظارت^{۲۸} شبکه عصبی استفاده شده است که برای آموزش این شبکه فقط از داده‌های سالم و نرمال شبکه که با تغییر بارهای شبکه بدست آمده‌اند استفاده شده است، برای آموزش مدل در حدود ۸۵ درصد داده‌ها برای آموزش و ۱۵ درصد آن‌ها به اضافه قسمت حمله برای تست مدل استفاده شدند، البته برای تست بیشتر داده‌های بیشتری علاوه بر داده‌های اصلی در نظر گرفته شده‌اند که در ادامه به آن می‌پردازیم.

^{۲۸} Unsupervise

داده‌های سالم در یک زیرمجموعه محدود از فضای ورودی یا زیرفضا قرار می‌گیرند و GAE پس از آموزش، توانایی بازسازی دقیق نمونه‌هایی را دارد که در نزدیکی این زیرفضا قرار دارند. حملات تزریق داده جعلی، به طور کلی، مقادیر ورودی را به بیرون از آن زیرفضا می‌برند (یا آن‌گونه تغییر می‌دهند که نمونه دیگر در توزیع واقعی داده‌های سالم نیست). بنابراین زمانی که یک نمونه آلوده به مدل اعمال می‌شود، decoder قادر به بازسازی دقیق آن نخواهد بود و خطای بازسازی (به ویژه در گره‌های تحت تاثیر حمله) به صورت معناداری بزرگ می‌شود (از آستانه بیشتر می‌شود). از این اختلاف^{۲۹} می‌توان برای تشخیص نمونه آلوده استفاده نمود.

تعیین آستانه:

برای تشخیص اینکه یک نمونه (سناریو) ناهنجار است یا خیر، از خطای بازسازی نمونه‌ای (طبق روابط ۱۸ و ۱۹) استفاده می‌کنیم:

$$MSE^{(s)} = \frac{1}{N} \sum_{i=1}^N |x_i^{(s)} - \widehat{x}_i^{(s)}|^2 \quad \text{رابطه ۱۸:}$$

$$\tau = Q_p \left(\{MSE^{(v)}\}_{v \in \mathcal{V}} \right) \quad \text{رابطه ۱۹:}$$

آستانه تصمیم τ بر مبنای صدک امپیریک p ^{۳۰} از توزیع $MSE^{(s)}$ روی مجموعه ارزیابی بدون حمله (سالم) \mathcal{V} تعیین می‌شود، قاعده تشخیص به این صورت است که، نمونه s را ناهنجار فرض می‌کنیم اگر رابطه ۲۰ برقرار باشد [۱۷].

$$MSE^{(s)} > \tau \quad \text{رابطه ۲۰:}$$

در این پیاده سازی مقدار $p = 0.95$ در نظر گرفته شد.

شبکه عصبی با `batch_size` برابر ۱۲۸ و تعداد `epochs` برابر ۱۰۰ اجرا شد و بهترین نتیجه مدل با `train_loss` برابر ۰.۰۰۰۰۰۶۷۷۴۸ ذخیره شد.

برای اینکه داده های بدست آمده قابل اعتماد باشند، مدل در سه `seed[0,1,42]` امتحان شده است.

این شبکه عصبی بارها برای انواع سناریو های حملات تزریق داده جعلی تست شد، روابط و نتایج پارامتر های سنجش عملکرد شبکه به صورت زیر ارائه می‌گردد:

$$Precision^{۳۱} = TP^{۳۲} / (TP + FP^{۳۳}) \quad \text{رابطه ۲۱:}$$

^{۲۹} Reconstruction Residual

^{۳۰} Empirical Percentile

^{۳۱} دقت

^{۳۲} True Positive

^{۳۳} False Positive

رابطه ۲۲: $Recall^{۳۴} = TP / (TP + FN^{۳۵})$

رابطه ۲۳: $F1 - Score = ۲ \times (Precision \times Recall) / (Precision + Recall)$

رابطه ۲۱، دقت: نسبت نمونه‌های مثبت صحیح به تمام نمونه‌های پیش‌بینی شده به عنوان مثبت.

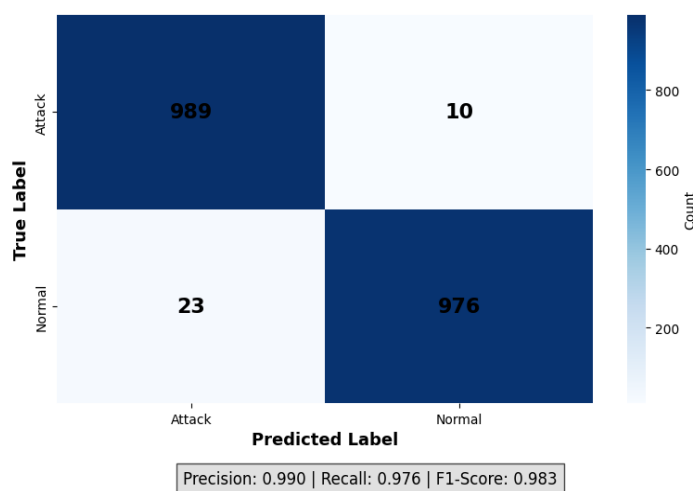
رابطه ۲۲، حساسیت (فراخوانی): نسبت نمونه‌های مثبت صحیح به تمام نمونه‌های مثبت واقعی.

جدول (۲). معیار های بررسی مدل

seed	دقت	فراخوانی	امتیاز F1
۰	۹۹,۳٪	۹۷,۵٪	۹۸,۴٪
۱	۹۸,۶٪	۹۷,۳٪	۹۷,۹٪
۴۲	۹۹,۰٪	۹۸,۱٪	۹۸,۵٪
میانگین	۹۸,۹٪	۹۷,۶٪	۹۸,۳٪

در شکل ۴ ماتریس درهم‌ریختگی^{۳۶} (میانگین) نمایش داده شده است، که نمایان گر تعداد حملات و توانایی مدل در تشخیص حملات و داده های نرمال از یکدیگر است.

CONFUSION MATRIX



شکل (۴). ماتریس درهم‌ریختگی مدل

^{۳۴} فراخوانی

^{۳۵} False Negative

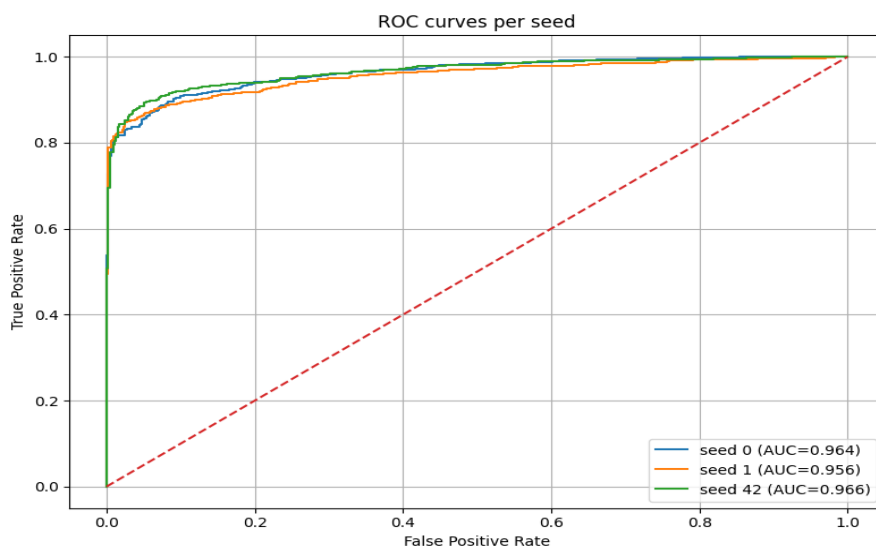
^{۳۶} Confusion matrix

معیار دیگر برای سنجش توانایی مدل در تشخیص حملات، منحنی عملکرد سیستم^{۳۷} است.

$$\text{FPR} = \text{FP}/(\text{FP} + \text{TN}) \quad \text{رابطه ۲۴:}$$

$$\text{AUC-ROC}^{38} = \int_0^1 \text{TPR}(f) d(\text{FPR}(f)) \quad \text{رابطه ۲۵:}$$

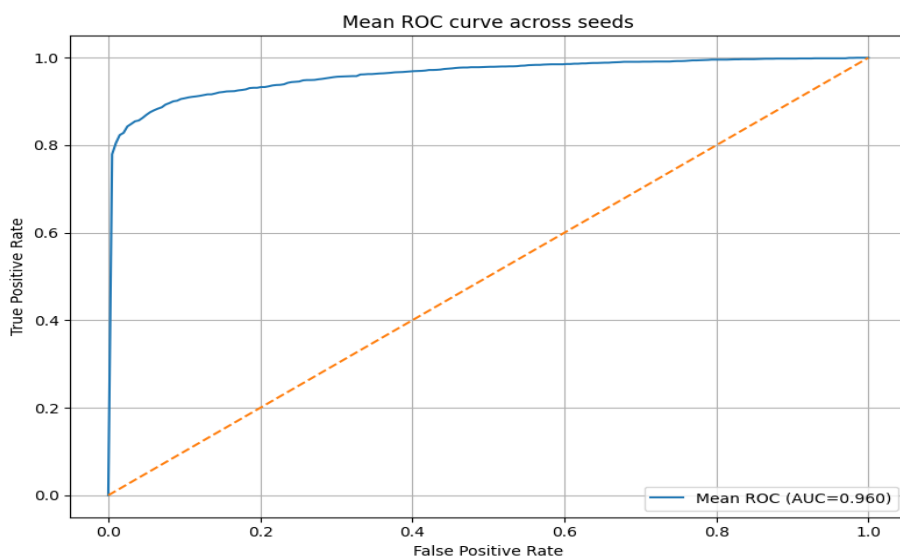
طبق روابط ۲۴ و ۲۵، ROC یک منحنی دوبعدی است که عملکرد تشخیص یک مدل طبقه‌بندی دودویی را با نمایش نرخ مثبت واقعی (TPR) در مقابل نرخ مثبت کاذب (FPR) در سطوح مختلف آستانه تصمیم‌گیری نشان می‌دهد. معیار AUC-ROC، که مقداری بین ۰ و ۱ دارد، احتمال این را اندازه‌گیری می‌کند که مدل بتواند به طور تصادفی یک نمونه مثبت را نسبت به یک نمونه منفی رتبه‌بندی کند. مقدار $\text{AUC}=1$ نشان‌دهنده تشخیص کامل و $\text{AUC}=0.5$ معادل یک طبقه‌بندی کننده تصادفی است. در شکل ۵ و ۶ منحنی ROC و مقدار AUC برای سیدها و مقدار میانگین نمایش داده شده است.



شکل (۵)، منحنی ROC به ازای هر سید

^{۳۷} ROC: Receiver Operating Characteristic

^{۳۸} Area under the Receiver Operating Characteristic Curve



شکل (۶)، منحنی ROC میانگین سیدها

۶. نتیجه گیری و پیشنهادات

در این پژوهش، یک چارچوب نوین تشخیص ناهنجاری مبتنی بر شبکه خود-رمزگذار گرافی (GAE) برای مقابله با تهدیدات پیچیده‌ای چون حملات تزریق داده جعلی در شبکه‌های توزیع هوشمند ارائه شد. هسته مرکزی این روش، استفاده از قابلیت‌های یادگیری عمیق گرافی برای مدلسازی و یادگیری الگوهای رفتاری شبکه قدرت در شرایط عادی است. مدل پیشنهادی با دریافت داده‌های ولتاژ باس‌ها و ساختار توپولوژیک شبکه، اقدام به یادگیری یک نمایش فشرده و بهینه از حالت سالم سیستم می‌نماید. معیار کلیدی تشخیص در این طرح، محاسبه خطای بازسازی است؛ به گونه‌ای که هرگونه انحراف در داده‌های عملیاتی از الگوی عادی آموخته‌شده، منجر به افزایش محسوس این خطا شده و به عنوان نشانه‌ای از یک حمله یا ناهنجاری شناسایی می‌گردد.

ارزیابی‌های تجربی بر روی شبکه استاندارد IEEE ۳۳-bus و با استفاده از یک مجموعه داده جامع متشکل از ۶۰۰۰۰ سناریو، کارایی این چارچوب را به اثبات رساند. مدل پیشنهادی با دستیابی به دقت بالای ۹۸ درصد و فراخوانی بالای ۹۷ در قالب F۱-Score بالای ۹۸ درصد تبلور یافته است، حاکی از توانایی آن در شناسایی کم‌خطای حملات است. همچنین، نتایج حاصل از ماتریس درهم‌ریختگی به وضوح نشان می‌دهد که مدل به خوبی قادر به تمایز بین داده‌های عادی و داده‌های تحت حمله است و از نرخ هشدار کاذب (False Positive) پایینی برخوردار می‌باشد. این مدل نیازی به پیش‌دانش دقیق از الگوهای حمله ندارد و صرفاً با یادگیری وضعیت سالم سیستم، قادر به تشخیص هرگونه انحراف مخرب می‌باشد. این ویژگی، چارچوب ارائه‌شده را به یک راهکار امنیتی پیشگیرانه و مقیاس‌پذیر برای پیاده‌سازی در محیط‌های عملیاتی شبکه‌های هوشمند تبدیل می‌کند.

پیشنهادات برای پژوهش‌های آینده:

- افزایش مقاومت در برابر حملات پیشرفته: آموزش مدل با داده‌های آلوده به حملات غیرخطی و چندمرحله‌ای برای افزایش تاب‌آوری در شرایط واقعی.
- استقرار عملیاتی: پیاده‌سازی مدل در محیط‌های شبیه‌سازی بلادرنگ^{۳۹} و ارزیابی عملکرد آن در مقیاس بزرگ.
- ترکیب داده‌ها با نوین: استفاده از نوین جمعی با داده‌های اصلی برای آموزش مدل برای مقاومت مدل در برابر نوین شبکه.
- استفاده از منابع تولید پراکنده: جایگذاری منابع تولید پراکنده در شبکه اصلی و آزمون مدل در حضور این منابع

منابع

- [۱] S. Wang, S. Bi, and Y. J. A. Zhang, "Locational Detection of the False Data Injection Attack in a Smart Grid: A Multilabel Classification Approach," *IEEE Internet Things J*, vol. ۷, no. ۹, pp. ۸۲۱۸-۸۲۲۷, Sep. ۲۰۲۰, doi: ۱۰.۱۱۰۹/JIOT.۲۰۲۰.۲۹۸۳۹۱۱.
- [۲] G. Liang, J. Zhao, F. Luo, S. R. Weller, and Z. Y. Dong, "A Review of False Data Injection Attacks Against Modern Power Systems," Jul. ۰۱, ۲۰۱۷, *Institute of Electrical and Electronics Engineers Inc.* doi: ۱۰.۱۱۰۹/TSG.۲۰۱۵.۲۴۹۵۱۳۳.
- [۳] S. Bi and Y. J. Zhang, "Using covert topological information for defense against malicious attacks on DC state estimation," *IEEE Journal on Selected Areas in Communications*, vol. ۳۲, no. ۷, pp. ۱۴۷۱-۱۴۸۵, ۲۰۱۴, doi: ۱۰.۱۱۰۹/JSAC.۲۰۱۴.۲۳۳۲۰۵۱.
- [۴] Z. Liu, M. Liu, Q. Wang, and Y. Tang, "False Data Injection Attacks on Data-Driven Algorithms in Smart Grids Utilizing Distributed Power Supplies," *Engineering*, vol. ۵۱, pp. ۶۲-۷۴, Aug. ۲۰۲۵, doi: ۱۰.۱۰۱۶/j.eng.۲۰۲۴.۱۱.۰۲۵.
- [۵] X. Wang, M. Hu, X. Luo, and X. Guan, "A detection model for false data injection attacks in smart grids based on graph spatial features using temporal convolutional neural networks," *Electric Power Systems Research*, vol. ۲۳۸, Jan. ۲۰۲۵, doi: ۱۰.۱۰۱۶/j.epsr.۲۰۲۴.۱۱۱۱۲۶.
- [۶] W. He, W. Liu, C. Wen, and Q. Yang, "Detection of False Data Injection Attacks on Smart Grids Based on A-BiTG Approach," *Electronics (Switzerland)*, vol. ۱۳, no. ۱۰, May ۲۰۲۴, doi: ۱۰.۳۳۹۰/electronics۱۳۱۰۱۹۳۸.

^{۳۹} Real-Time

- [۷] P. Hu, W. Gao, Y. Li, M. Wu, F. Hua, and L. Qiao, "Detection of False Data Injection Attacks in Smart Grids Based on Expectation Maximization," *Sensors*, vol. ۲۳, no. ۳, Feb. ۲۰۲۳, doi: ۱۰,۳۳۹۰/s۲۳۰۳۱۶۸۳.
- [۸] F. Zhang and Q. Yang, "False data injection attack detection in dynamic power grid: A recurrent neural network-based method," *Front Energy Res*, vol. ۱۰, Sep. ۲۰۲۲, doi: ۱۰,۳۳۸۹/fenrg.۲۰۲۲,۱۰۰۵۶۶۰.
- [۹] Y. Liu, P. Ning, and M. K. Reiter, "False data injection attacks against state estimation in electric power grids," in *ACM Transactions on Information and System Security*, May ۲۰۱۱. doi: ۱۰,۱۱۴۵/۱۹۵۲۹۸۲,۱۹۵۲۹۹۵.
- [۱۰] J. Kim, L. Tong, and R. J. Thomas, "Subspace methods for data attack on state estimation: A data driven approach," *IEEE Transactions on Signal Processing*, vol. ۶۳, no. ۵, pp. ۱۱۰۲-۱۱۱۴, Mar. ۲۰۱۵, doi: ۱۰,۱۱۰۹/TSP.۲۰۱۴,۲۳۸۵۶۷۰.
- [۱۱] H. Haider, A.-M. Zaid Ali, and -Yemen Ayeda Al-Hmadi, "False Data Injection Attacks (FDIA) detection by Deep Learning Techniques in Smart Grids: survey."
- [۱۲] A. Takiddin, R. Atat, M. Ismail, O. Boyaci, K. R. Davis, and E. Serpedin, "Generalized Graph Neural Network-Based Detection of False Data Injection Attacks in Smart Grids," *IEEE Trans Emerg Top Comput Intell*, vol. ۷, no. ۳, pp. ۶۱۸-۶۳۰, Jun. ۲۰۲۳, doi: ۱۰,۱۱۰۹/TETCI.۲۰۲۲,۳۲۳۲۸۲۱.
- [۱۳] T. N. Kipf and M. Welling, "Variational Graph Auto-Encoders," Nov. ۲۰۱۶, [Online]. Available: <http://arxiv.org/abs/۱۶۱۱,۰۷۳۰۸>
- [۱۴] W. Liao, B. Bak-Jensen, J. Radhakrishna Pillai, Y. Wang, and Y. Wang, "A Review of Graph Neural Networks and Their Applications in Power Systems," *Journal of Modern Power Systems and Clean Energy*, vol. ۱۰, no. ۲, pp. ۳۴۵-۳۶۰, ۲۰۲۲, doi: ۱۰,۳۵۸۳۳/MPCE.۲۰۲۱,۰۰۰۰۵۸.
- [۱۵] A. Howe, D. Peasley, and M. Papa, "Graph Autoencoders for Detecting Anomalous Intrusions in OT Networks Through Dynamic Link Detection," in *2024 IEEE 21st Consumer Communications & Networking Conference (CCNC)*, IEEE, Jan. ۲۰۲۴, pp. ۱-۶. doi: ۱۰,۱۱۰۹/CCNC۵۱۶۶۴,۲۰۲۴,۱۰۴۵۴۸۴۱.
- [۱۶] C. Zhang and J. W. Jung, "Enhanced Graph Autoencoder for Graph Anomaly Detection Using Subgraph Information," *Applied Sciences (Switzerland)*, vol. ۱۵, no. ۱۵, Aug. ۲۰۲۵, doi: ۱۰,۳۳۹۰/app۱۵۱۵۸۶۹۱.
- [۱۷] R. Chalapathy and S. Chawla, "Deep Learning for Anomaly Detection: A Survey," Jan. ۲۰۱۹, [Online]. Available: <http://arxiv.org/abs/۱۹۰۱,۰۳۴۰۷>

A Deep Learning-Based Approach for Detecting False Data Injection Attacks in Smart Grids

Authors:

Alireza Soltani¹, Hamidreza Baghaee-Kashi², Mahmoudreza Haghifam³

¹ [M.Sc.](#) Student, Electrical Engineering Department, Faculty of Electrical and Computer Engineering, Tarbiat Modares University, Tehran, Iran.

Email: alireza_soltani@modares.ac.ir

² Assistant Professor, Electrical Engineering Department, Faculty of Electrical and Computer Engineering, Tarbiat Modares University, Tehran, Iran.

³ Professor, Electrical Engineering Department, Faculty of Electrical and Computer Engineering, Tarbiat Modares University, Tehran, Iran.

Abstract

The proliferation of smart distribution networks and the escalation of cyber threats have rendered the security of power systems critically imperative. Among the most disruptive threats, False Data Injection (FDI) attacks compromise grid stability by evading conventional detection mechanisms. This research proposes a **Graph Autoencoder-based detection model** that proactively identifies such attacks in the IEEE ۳۳-bus distribution system. Designed within a deep learning architecture, the model harnesses the representational capabilities of Graph Neural Networks (GNNs) to learn normal operational patterns of the grid and subsequently detect anomalies through reconstruction error computation. The results demonstrate that the proposed model effectively identifies intrusions with **over ۹۸% accuracy**, thereby substantially enhancing the performance of the protection system.

Submission Date:

۲۰۲۵/۱۱/۰۳

Accepted Date:

۲۰۲۶/۰۲/۲۵

Keywords:

Smart grid

Anomaly Detection

Cyber-Attack

Deep Learning