

کاهش داده‌های مورد نیاز برای آموزش مدل‌های یادگیری عمیق بر اساس خوشه‌بندی داده‌ها و کاربرد آن در وارون‌سازی یک‌بعدی مگنتوتلوریک

مهدی رحمانی جویباری^۱ و بنفشه حبیبیان‌دهکردی^{۲*}

^۱ دانشجوی دکترا، موسسه ژئوفیزیک دانشگاه تهران، تهران، ایران

^۲ استادیار، موسسه ژئوفیزیک دانشگاه تهران، تهران، ایران

(دریافت: ۱۴۰۳/۰۱/۰۸، پذیرش: ۱۴۰۳/۰۷/۰۴)

چکیده

رویکردهای یادگیری عمیق داده‌محور با چالش تولید داده‌هایی به تعداد زیاد و با کیفیت بالا و بار محاسباتی سنگین و زمان آموزش طولانی تحمیل شده توسط آن روبرو هستند. علاوه بر این، در صورتی که بعد از جداسازی تصادفی داده‌ها به سه مجموعه آموزش، اعتبارسنجی و آزمایش، توزیع آماری یکسانی برای آنها به دست نیاید، به دلیل رفتار نامنظم منحنی خطای آموزش و اعتبارسنجی، تعمیم‌پذیری خوبی حاصل نمی‌شود. در این پژوهش با استفاده از رویکرد مبتنی بر خوشه‌بندی اولیه داده‌ها و اختصاص درصد مشخصی از هر خوشه به سه مجموعه، و با پیمایش نتایج پیش‌بینی، کمینه داده مورد نیاز برای وارون‌سازی با رویکرد یادگیری عمیق ارائه می‌گردد. با اعمال آزمون‌های آماری نشان داده می‌شود که داده‌هایی که با این رویکرد جداسازی شده‌اند، دارای توزیع یکسان در سه مجموعه هستند. یک مدل یادگیری عمیق مبتنی بر معماری U-Net برای وارون‌سازی یک‌بعدی داده‌های مگنتوتلوریک آموزش داده می‌شود. به این منظور از یک مدل ژئوالکتریکی پنج لایه که شرایط یک میدان زمین‌گرایی را شبیه‌سازی می‌کند، استفاده شده‌است. آموزش شبکه با تعداد متفاوت داده‌هایی که با روش گفته شده جداسازی شده‌اند، تکرار و عملکرد آن با معیارهای کمی و کیفی متفاوتی سنجیده می‌شود. با پیمایش نتایج وارون‌سازی با داده‌های آزمایشی یکسان بر مدل‌های آموزش دیده با درصد داده‌های مختلف می‌توان بدون اینکه از دقت شبکه کاسته شود، به میزان ۵۰ درصد تعداد داده‌های مورد نیاز برای آموزش مدل یادگیری عمیق و بنابراین زمان آموزش را کاهش داد. در مواجهه با داده‌های پیچیده‌تر، واقعی‌تر و نویزی قطعا جداسازی تصادفی رهیافت مناسبی برای تشکیل سه مجموعه نیست. هرچه شرایط پیچیده‌تر و تعداد ویژگی‌ها بیشتر باشد، جداسازی تصادفی راهکار نامناسب‌تری است؛ چراکه تفاوت توزیع‌های آماری سه مجموعه بیشتر می‌شود؛ و در نتیجه تعمیم‌پذیری کاهش و تعداد داده‌های مورد نیاز افزایش می‌یابد. در این صورت استفاده از خوشه‌بندی راهکار مناسبی برای یکسان‌سازی توزیع آماری سه مجموعه و کاهش تعداد داده‌هاست.

کلمه‌های کلیدی: خوشه‌بندی، مگنتوتلوریک، یادگیری عمیق، وارون‌سازی

۱ مقدمه

در سال‌های اخیر، نقش روش‌های مبتنی بر هوش مصنوعی و به‌ویژه یادگیری عمیق در جنبه‌های مختلف کار با داده‌های ژئوفیزیکی بسیار پررنگ شده است. روش مگنتوتلوریک، یک روش ژئوفیزیکی الکترومغناطیسی با کاربردهای متنوع در زمینه‌هایی همچون مطالعات تکتونیکی (برای مثال کامپو و همکاران، ۲۰۲۲) و اکتشاف منابع زمین‌گرمایی و هیدروکربنی (برای مثال سگویا و همکاران، ۲۰۲۱؛ میری و همکاران، ۲۰۲۱) است. برای وارون‌سازی و تفسیر سنتی سونداژهای MT، متداول‌ترین روش‌های مورد استفاده توسط کانستبل و همکاران (۱۹۸۷)، پارکر و بوکر (۱۹۹۶)، اسمیت و بوکر (۱۹۸۸) و فیشر و همکاران (۱۹۸۱) ارائه شده‌اند. در پژوهش‌های بعدی از تانسور فاز (کالدول و همکاران، ۲۰۰۴) برای وارون‌سازی یک‌بعدی مگنتوتلوریک استفاده شده است (یونگه، ۲۰۱۱). العلی و همکاران (۲۰۲۰) رهیافت جدیدی برای تعیین تعداد بهینه لایه‌ها ارائه دادند. پارامترهای مدل یک‌بعدی شامل مقاومت ویژه الکتریکی و ضخامت لایه‌ها در قالب یک مسئله وارون بدووضع از روی داده‌ها تعیین می‌شوند. روش‌های یادگیری عمیق به‌عنوان رهیافت‌های آماری مبتنی بر داده قابل کاربرد برای حل مسائل وارون بدووضع هستند. این وارون‌سازی یک مسئله رگرسیون شدیداً غیرخطی است و شبکه‌های عمیق تر توابع غیرخطی را بهتر تقریب می‌زنند. مزیت بهره گرفتن از لایه‌های متعدد در شبکه‌های عمیق، تولید توابع نگاشت غیرخطی با پیچیدگی بالاست. نمونه‌های متعددی از کاربرد روش‌های یادگیری عمیق در مدل‌سازی‌های الکترومغناطیسی (اوه و همکاران، ۲۰۲۰؛ پازیرف، ۲۰۱۹؛ شهریاری و همکاران، ۲۰۲۰) و مگنتوتلوریک (لیو و همکاران، ۲۰۲۱ و ۲۰۲۲؛ لیائو و همکاران، ۲۰۲۱ الف و ب)، گرایش فزاینده به استفاده از این روش‌ها در مدل‌سازی‌های ژئوفیزیکی را نشان می‌دهد.

روش‌های داده‌محور یادگیری عمیق برای یادگیری ویژگی‌های یک سیستم، به داده‌های آموزشی متکی هستند؛ به طوری که فرایند یادگیری بر اساس تجربه و با شناسایی الگوی پیچیده موجود در داده‌ها محقق می‌شود. هرچه مدل پیچیده‌تر باشد، داده‌های آموزشی بیشتری مورد نیاز است. اگر کمیت و کیفیت نمونه‌ها مناسب نباشد و یا مسئله مورد نظر را به خوبی نمایندگی نکنند، تعمیم‌پذیری و کارایی مدل تحت شرایط متغیر با مشکل مواجه می‌شود. بار محاسباتی لازم برای تولید این حجم از داده شبیه‌سازی شده برای آموزش بسیار بالاست؛ گرچه مدل آموزش‌دیده در تولید خروجی به سرعت عمل می‌کند. بنابراین چالش اول، تولید داده‌های کافی و با کیفیت بالاست که معرف سناریوهای واقعی باشند. چالش بعدی، آموزش شبکه با این حجم از داده است که مستلزم انجام محاسبات سنگین و صرف زمان طولانی است. چالش پایانی هم تعمیم‌پذیری مدل آموزش داده شده است.

اولین گام پس از تولید داده در یادگیری عمیق، جداسازی داده‌ها در قالب سه مجموعه آموزش، اعتبارسنجی و آزمایش است. اگر این جداسازی به شکل تصادفی انجام شود، تضمینی وجود ندارد که سه مجموعه داده از توزیع آماری یکسانی برخوردار باشند. در این پژوهش رویکرد خوشه‌بندی داده‌ها و اختصاص درصد مشخصی از هر خوشه به سه مجموعه برای جداسازی به کار گرفته می‌شود. بطور خاص از معماری U-Net برای وارون‌سازی یک‌بعدی داده‌های مگنتوتلوریک استفاده شده است. شبکه U-Net بر اساس مطالعه‌ای که در آن سه معماری U-Net، Residual Network (Res-Net) و Variational Autoencoder (VAE) به‌عنوان سه مدل یادگیری عمیق از لحاظ دقت و تعمیم‌پذیری در وارون‌سازی داده‌های MT مقایسه شده‌اند (رحمانی جوینانی و همکاران، ۲۰۲۴)، انتخاب شده است. به این منظور در مطالعه مورد اشاره، معیارهای متنوعی به کار گرفته

در محیط یک‌بعدی، امپدانس الکتریکی یک کمیت عددی است و به سیستم مختصات بستگی ندارد؛ غالباً به شکل مقاومت ویژه ظاهری (ρ_a) و فاز (φ) بیان می‌شود:

$$\rho_a = \frac{1}{\omega\mu_0} |Z|^2, \quad \varphi = \tan^{-1} \left(\frac{\text{Im}Z}{\text{Re}Z} \right) \quad (2)$$

که μ_0 نفوذپذیری فضای آزاد و ω بسامد زاویه‌ای است.

برای تولید مثال‌های آموزشی از یک مدل ژئوالکتریکی که با حذف تغییرات جانبی از مدل به کاررفته توسط چن و همکاران (۲۰۱۲) حاصل شده‌است و شرایط یک میدان زمین‌گرمایی را شبیه‌سازی می‌کند، استفاده شده‌است. این مدل از پنج لایه شامل روباره، کلاهک رسی، مخزن، ناحیه گذار و سنگ بستر تشکیل شده‌است. آنها کران‌های اولیه برای پارامترهای مدل را بر اساس تجارب خود از میدان‌های زمین‌گرمایی انتخاب کرده‌اند و با استفاده از این کران‌ها می‌توان بازه‌هایی را برای تغییرات مقادیر مقاومت ویژه الکتریکی و عمق لایه‌ها تعیین کرد. مقادیر مقاومت ویژه الکتریکی به ترتیب از سطحی‌ترین به عمیق‌ترین لایه به‌طور تصادفی در بازه‌های $[60, 140]$ ، $[1/6, 2/4]$ ، $[20, 460]$ ، $[64, 96]$ و $[30]$ اهم-متر توزیع شده‌اند. مقادیر ضخامت هم به‌طور تصادفی در بازه‌های $[450, 550]$ ، $[170, 430]$ ، $[1770, 2630]$ و $[1070, 2930]$ متر در نظر گرفته شده‌اند (جدول ۱). با استفاده از الگوریتم پیشرو MT پاسخ‌های این مدل ژئوالکتریکی یک‌بعدی به‌طور تحلیلی در بازه فرکانسی ۱۰۰-۰/۰۱ هرتز و در ۱۳ فرکانس که به‌طور یکنواخت در مقیاس لگاریتمی توزیع شده‌اند، محاسبه و در مجموع ۵۰۰۰۰۰ نمونه داده تولید شد. آماده‌سازی یا پیش‌پردازش داده‌ها شامل مراحل همچون نرمال‌سازی و استانداردسازی به منظور مقیاس کردن مجدد متغیرهای ورودی و خروجی قبل از آموزش یک مدل شبکه است. متغیرهای ورودی با مقادیر خیلی پراکنده یا با واحدهای مختلف، در نهایت باعث می‌شوند مقادیر وزن تغییرات شدید داشته باشند و فرایند یادگیری ناپایدار شود. برای

شدند و عملکرد معماری‌های متفاوت به‌طور عددی و تحلیلی مورد ارزیابی قرار گرفتند. نتایج حاصل، نشان‌دهنده برتری نسبی معماری U/Net در بازیابی مدل‌های لایه‌ای مقاومت ویژه الکتریکی بود. این مطالعه مقایسه‌ای را می‌توان به‌طور خاص مبنایی برای مدل‌سازی MT با رهیافت یادگیری عمیق و بررسی قابلیت‌های آن تلقی کرد. ضمناً ضخامت لایه‌ها متغیر و بخشی از مجموعه پارامترهای خروجی در نظر گرفته شده‌اند. ولی در پژوهش حاضر، با استفاده از رویکرد خوشه‌بندی داده‌ها قبل از ورود به مرحله آموزش، به برخی از چالش‌های بسیار مهم مرتبط با یادگیری عمیق داده‌محور پرداخته شده‌است. بعد از خوشه‌بندی مناسب، با اعمال آزمون آماری KS (Kolmogorov-Smirnov: KS) نشان داده می‌شود که در داده‌هایی که با این رویکرد جداسازی شده‌اند، سه مجموعه دارای توزیع آماری یکسان هستند. با پیمایش کمی و کیفی نتایج وارون‌سازی داده‌های مگنتوتلوریک با درصد‌های مختلف داده‌ای، نشان داده می‌شود که تعداد نمونه‌های مورد نیاز برای آموزش شبکه به شکل موثری، با حفظ دقت پیش‌بینی و یا حتی افزایش آن، کاهش می‌یابد.

۲ روش پژوهش

۲-۱ تولید داده‌های آموزشی

در روش مگنتوتلوریک از نوسان‌های میدان‌های الکترومغناطیسی طبیعی که بطور همزمان ثبت می‌شوند، برای کاوش ساختار ژئوالکتریکی زیرسطحی استفاده می‌شود. در حالت کلی تانسور امپدانس (Z) به عنوان مهمترین تابع تبدیل مگنتوتلوریک، رابطه بین مولفه‌های افقی میدان‌های الکتریکی (E) و مغناطیسی (H) را بیان می‌کند:

$$Z = \frac{E}{H} \quad (1)$$

(۴) پسپردازش می‌شوند تا به بازه اولیه بازگردانده شوند و مقادیر خروجی مطلوب را ارائه دهند.

$$\hat{m}_j = \text{anti_log}(\log(m_{max}) \times m_j) \quad (۴)$$

در مطالعات قبلی، به دلیل مشاهده عملکرد بهتر برای شبکه مورد نظر، فقط از مقاومت ویژه ظاهری به‌عنوان داده ورودی استفاده شده‌است (لیائو و همکاران، ۲۰۲۲؛ لیو و همکاران، ۲۰۲۱)؛ اما در این تحقیق مقادیر فاز هم به داده‌های ورودی اضافه شده‌اند و عملکرد شبکه طراحی شده را در انجام وارون‌سازی ارتقا داده‌اند. در این مرحله از یادگیری باناظر و طرحواره پس انتشار خطا استفاده می‌شود و بنابراین ورودی‌ها یا مشاهدات به همراه خروجی‌های مربوطه در قالب نمونه‌های آموزشی به شبکه داده می‌شوند.

سرعت بخشیدن به همگرایی شبکه، داده‌های ورودی (d_i) و داده‌های خروجی متناظر با آنها (m_j) به‌صورت زیر نرمال‌سازی می‌شوند تا در بازه $[0, 1]$ واقع شوند:

$$\begin{cases} \hat{d}_i = \frac{\log(d_i)}{\log(d_{max})} \\ \hat{m}_j = \frac{\log(m_j)}{\log(m_{max})} \end{cases} \quad (۳)$$

که d_{max} و m_{max} به ترتیب معرف بیشینه مقادیر عددی داده ورودی و خروجی متناظر با آن هستند. گرچه روش‌های مختلفی برای نرمال‌سازی وجود دارد، این روش برای داده‌های MT مناسب‌تر به نظر می‌رسد. پس از اتمام مدل‌سازی شبکه، نتایج پیش‌بینی شده (m_j) توسط رابطه

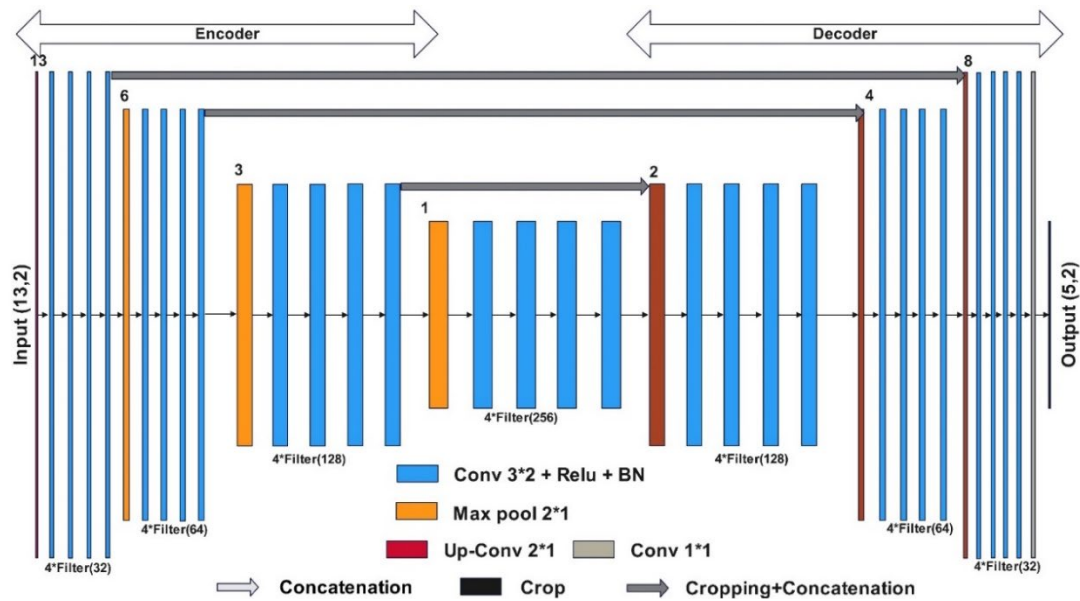
جدول ۱. پارامترهای مورد استفاده برای تولید مدل ژئوالکتریکی یک‌بعدی از مخزن زمین‌گرمایی (برگرفته از مطالعه انجام شده توسط چن و همکاران، ۲۰۱۲).

لایه‌ها	بازه مقادیر مقاومت ویژه الکتریکی (اهم-متر)	بازه عمقی (متر)
لایه اول	[۶۰,۱۴۰]	[۴۵۰,۵۵۰]
لایه دوم	[۱/۶,۲/۴]	[۱۷۰,۴۳۰]
لایه سوم	[۲۰,۴۶۰]	[۱۷۷۰,۲۶۳۰]
لایه چهارم	[۶۴,۹۶]	[۱۰۷۰,۲۹۳۰]
لایه پنجم	[۳۰]	

۲-۲ معماری شبکه

شبکه‌های عصبی پیچشی (Convolutional Neural Network: CNN) در حال حاضر به‌طور گسترده مورد استفاده قرار می‌گیرند. تفاوت شبکه‌های CNN نسبت به شبکه‌های سنتی، استفاده از پیچش به جای ضرب ماتریسی است. دو اثر مشخص این کار، ایجاد اتصالات تنک و به اشتراک‌گذاری پارامتر است. یکی از پرکاربردترین شبکه‌های پیچشی، شبکه U-Net است. این معماری در ابتدا توسط رونبرگر و همکاران (۲۰۱۵) برای تقسیم‌بندی تصاویر پزشکی توسعه داده شد. اساس آن بخش‌های رمزنگار و رمزگشایی هستند که بین آنها پل ارتباطی برقرار شده و جریان اطلاعات را تسهیل می‌کند. در معماری U-

Net به‌کاررفته در این پژوهش، از ۴ فیلتر با اندازه ۳۲، ۴ فیلتر با اندازه ۶۴، ۴ فیلتر با اندازه ۱۲۸ و ۴ فیلتر با اندازه ۲۵۶ استفاده شده‌است (شکل ۱). برای مسئله وارون-1 DMT، داده‌های ورودی و خروجی ماتریس دوبعدی هستند؛ ماتریس خروجی به‌صورت داده‌های مقاومت ویژه الکتریکی و ضخامت و ماتریس ورودی، به‌صورت داده‌های مقاومت ویژه ظاهری و فاز تعریف شده‌اند. با وجود اینکه این معماری در مقایسه با سایر مدل‌های یادگیری عمیق با مجموعه داده‌های محدود نیز موثر است، همچنان چالش مربوط به جمع‌آوری حجم زیادی از داده‌ها و منابع سخت‌افزاری و زمان مورد نیاز برای انجام آن پابرجاست.



شکل ۱. معماری U-Net مورد استفاده.

نرخ یادگیری (Learning Rate)، الگوریتم بهینه‌سازی، تابع فعال‌ساز و تابع اتلاف از جمله فرآیندهایی هستند که قبل از فرآیند یادگیری انتخاب و بعد از اتمام چرخه آموزش در صورت نیاز بر اساس منحنی خطای آموزش و اعتبارسنجی تغییر می‌کنند. در آموزش این شبکه با معماری U-Net، تعداد تکرار چرخه آموزش برابر با ۱۵، اندازه دسته‌های نمونه‌گیری برابر با ۳۲، تعداد و اندازه فیلتر برابر با ۴ فیلتر به اندازه‌های ۳۲، ۶۴، ۱۲۸ و ۲۵۶، نرخ یادگیری برابر با ۰/۰۰۰۱، نوع تابع بهینه‌ساز از نوع Adam، تابع فعال‌ساز از نوع ReLU و تابع اتلاف از نوع میانگین مربعات خطا (MSE) در نظر گرفته شده است.

از تابع اتلاف برای محاسبه میزان خطا، از بهینه‌ساز Adam برای به روزرسانی وزن‌ها و از تابع فعال‌ساز ReLU برای انتقال وزن‌ها به لایه بعدی استفاده می‌شود. مجموعه پارامترهایی که در حین فرآیند یادگیری بهینه می‌شوند، تابع نگاشت یا عملگر شبه‌وارونی را تشکیل می‌دهند که خروجی مدل یادگیری عمیق را با اعمال بر داده‌های ورودی تولید می‌کند (کیم و ناکاتا، ۲۰۱۸). در تابع بهینه‌ساز Adam برخلاف طرحواره کلاسیک کاهش

۲-۳ آموزش شبکه یادگیری عمیق

توانایی نگاشت غیرخطی رویکردهای عمیق باعث می‌شود شبکه با تنظیم پارامترهای خود، قادر به ایجاد نگاشت بین داده‌های مقاومت ویژه الکتریکی ظاهری و فاز از یک طرف و مدل لایه‌ای مقاومت ویژه الکتریکی از طرف دیگر باشد. این پارامترها با به‌روزرسانی موثر، ویژگی‌های اصلی را از داده‌های ورودی فرا می‌گیرند. مقادیر اولیه پارامترها به شکل تصادفی انتخاب و در چرخه آموزش به‌روزرسانی می‌شوند. میزان کارایی روش‌های یادگیری ماشین به‌عنوان رویکردهای داده‌محور همچنین به فرآیندهایی (hyperparameters) بستگی دارد که قبل از فرآیند یادگیری انتخاب می‌شوند. آموزش شبکه نوعی فرآیند بهینه‌سازی است و فرآیندها بر اساس خطای اعتبارسنجی انتخاب می‌شوند. فرآیندها ویژگی الگوریتم یادگیری هستند؛ در مقابل پارامترها ویژگی مدل هستند و بطور محض در حین آموزش، زمانی که الگوریتم سعی در یادگیری نگاشت بین ورودی و خروجی را دارد، تخمین زده می‌شوند. تعداد تکرار چرخه آموزش (Epoch)، اندازه دسته‌های نمونه‌گیری (Batch Size)، تعداد و اندازه فیلتر،

حاصل کرد. همچنین می‌توان از این منحنی‌های خطا در تشخیص توزیع آماری یکسان داده‌های آموزش و اعتبارسنجی استفاده کرد (گودفلو و همکاران، ۲۰۱۶).

۲-۴ منحنی‌های خطا در چرخه آموزش یادگیری عمیق

به‌طور کلی، منحنی خطا نموداری است که میزان خطا را برحسب تکرار در چرخه آموزش نشان می‌دهد. ترسیم منحنی خطا برای مجموعه داده آموزشی، ایده‌ای از چگونگی یادگیری مدل و ترسیم منحنی خطا برای مجموعه داده اعتبارسنجی، ایده‌ای از چگونگی تعمیم‌پذیری مدل ارائه می‌دهد. در ارزیابی چرخه آموزش، ترسیم همزمان این دو منحنی معمول است و از شکل و پویایی منحنی خطا می‌توان برای تشخیص رفتار یک مدل یادگیری عمیق استفاده کرد.

۱-۴-۲ منحنی‌های خطای کم‌برازش

کم‌برازش هنگامی اتفاق می‌افتد که مدل یادگیری عمیق به اندازه کافی پیچیده نباشد که بتواند رابطه‌ی میان ورودی و خروجی را یاد بگیرد و الگوی غالب در داده‌ها را تشخیص دهد؛ در نتیجه، در آموزش داده‌های آموزشی و در نتیجه اعتبارسنجی، عملکرد ضعیفی دارد. اگر مدلی نتواند به خوبی به داده‌های جدید تعمیم داده شود، نمی‌توان از آن برای طبقه‌بندی یا پیش‌بینی استفاده کرد. شکل منحنی خطا در مدل کم‌برازش معمولاً به دو صورت است:

۱- منحنی خطای آموزش بدون توجه به فرایند آموزش ثابت می‌ماند.

۲- منحنی خطای آموزش تا پایان چرخه آموزش همچنان کاهش می‌یابد.

چندین روش برای جلوگیری از کم‌برازش رخ داده پیشنهاد می‌شود که باعث افزایش پیچیدگی و تنوع مدل می‌شود و امکان آموزش موفقیت‌آمیز مدل فراهم می‌گردد.

گرادیان تصادفی (Stochastic Gradient Descent: SGD)، نرخ یادگیری ثابت نیست و به‌صورت تدریجی کم می‌شود. این رویه، واپاشی نرخ یادگیری نام دارد. یکی از مهم‌ترین مزیت‌های بهینه‌ساز Adam این است که نیاز کمتری به تنظیم دستی نرخ یادگیری دارد و در بسیاری از موارد، تنظیمات پیش‌فرض ($\alpha = 0.001$, $\beta_1 = 0.9$ و $\beta_2 = 0.999$) بهترین عملکرد را ارائه می‌دهند. بهینه‌ساز Adam به خوبی با داده‌های متفاوت و شرایط مختلف آموزشی سازگار است و توانایی آن در تطبیق نرخ یادگیری برای پارامترهای مختلف، به‌ویژه در مواجهه با داده‌های نامتعادل یا پراکنده بسیار ارزشمند است (کینگما و بار، ۲۰۱۵). مقدار اولیه نرخ یادگیری بر اساس سعی و خطا برابر با 0.001 انتخاب شد.

هنگام آموزش یک شبکه عمیق با یادگیری مبتنی بر گرادیان و طرحواره پس‌انتشار خطا (Backpropagation Error)، انباشت مشتقات کوچک منجر به مشکل محوشدگی گرادیان (Vanishing Gradient) می‌شود که به مدل ناپایدار و ناتوان در یادگیری موثر می‌انجامد. تابع فعال‌ساز ReLU، مشکل محوشدگی و انفجار گرادیان را که برای توابعی همچون سیگموئید و تانژانت هذلولی رخ می‌دهد، ندارد. همچنین این تابع محدودیتی برای مقدار بیشینه ورودی ایجاد نمی‌کند. از آنجاییکه این تابع فقط تعدادی از گره‌ها را فعال می‌کند، از لحاظ محاسباتی از کارایی بهتری نسبت به سایر توابع فعال‌ساز در شبکه‌های عمیق برخوردار است (نیر و همکاران، ۲۰۱۰).

یک ابزار تشخیصی و پرکاربرد در ارزیابی چرخه یادگیری در مسائل رگرسیون، پایش روند منحنی‌های خطای آموزش و اعتبارسنجی بعد از اتمام چرخه آموزش است. با محاسبه خطا پس از هر چرخه آموزش می‌توان نمودارهایی از عملکرد یادگیری بر مجموعه داده آموزشی و اعتبارسنجی ترسیم و از سه حالت بیش‌برازش (Overfit)، کم‌برازش (Underfit) و برازش مناسب (Goodfit) اطلاع

و در واقع برازش مناسب، هدف الگوریتم یادگیری است. بهترین انتخاب در تعداد تکرار چرخه آموزش، زمانی است که روند کاهش منحنی‌های خطای آموزش و اعتبارسنجی ناملموس باشد و فاصله بین آنها به کمترین مقدار خود برسد.

۲-۴-۴ تشخیص توزیع یکسان داده‌های آموزش و اعتبارسنجی

از منحنی‌های خطا می‌توان در تشخیص توزیع یکسان داده‌های آموزش و اعتبارسنجی استفاده کرد. در صورتی که ویژگی‌های آماری دو مجموعه داده آموزش و اعتبارسنجی یکسان نباشند برازش مناسب در چرخه آموزش اتفاق نمی‌افتد. اگر تعداد نمونه‌ها در یک مجموعه داده نسبت به مجموعه داده دیگر بسیار کم باشد ویژگی‌های آماری دو مجموعه یکسان نخواهد بود. معمولاً در جداسازی داده‌ها در یادگیری عمیق، سهم بیشتر به داده‌های آموزش اختصاص دارد و تعداد نمونه‌ها در مجموعه داده اعتبارسنجی کم است. در این حالت اگر ویژگی‌های آماری مجموعه داده اعتبارسنجی در مقایسه با مجموعه داده آموزشی یکسان نباشد، مجموعه داده اعتبارسنجی اطلاعات کافی به منظور ارزیابی توانایی مدل برای تعمیم‌پذیری را ارائه نمی‌دهد. در این حالت روند منحنی خطای آموزش، نزولی، اما روند منحنی خطای اعتبارسنجی نامنظم خواهد بود (شکل ۲).

۲-۴-۲ منحنی‌های خطای بیش‌برازش

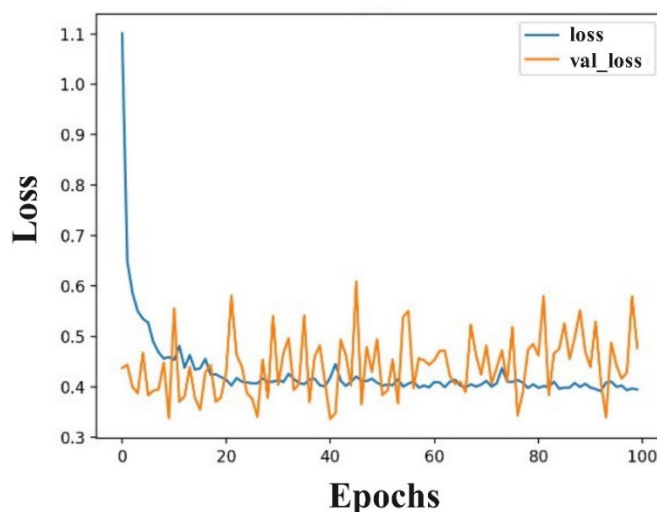
برازش بیش از حد هنگامی اتفاق می‌افتد که مدل، ویژگی‌های داده‌های آموزشی را به جای یادگیری، حفظ کرده باشد؛ یعنی بیش از حد روی آن آموزش دیده باشد؛ در نتیجه، این مدل فقط در مجموعه‌ی داده‌های آموزشی مفید خواهد بود و نمی‌توان آموزش را به مجموعه‌ی داده‌های دیگر که هنوز آن‌ها را ندیده است تعمیم داد. بنابراین خطای تعمیم‌پذیری افزایش می‌یابد. این افزایش خطای تعمیم‌پذیری را می‌توان با عملکرد مدل در مجموعه داده اعتبارسنجی اندازه‌گیری کرد. شکل منحنی‌های خطای اعتبارسنجی در مدل بیش‌برازش تا یک نقطه کاهش و دوباره شروع به افزایش می‌یابد. نقطه عطف (Inflection point)، نقطه‌ای است که می‌توان آموزش را در آن متوقف کرد. زیرا تکرار آموزش پس از آن نقطه، برازش بیش از حد را نشان می‌دهد. روش‌های مختلفی برای جلوگیری از آن وجود دارد از جمله: تقویت داده‌ها (Data Augmentation)، منظم سازی (L2 Regularization)، حذف تصادفی (Drop Out)، توقف زودهنگام (Early Stopping) و ذخیره بهترین وزن‌ها. در اینجا از منظم سازی L2 با مقدار $0/0001$ و ذخیره بهترین وزن‌ها استفاده شده است.

۲-۴-۳ منحنی‌های یادگیری برازش مناسب

در برازش مناسب، فاصله بین خطای یادگیری آموزش و اعتبارسنجی کمینه بوده و هر دو منحنی روند نزولی دارند

جدول ۲. تعداد داده‌ها در سه مجموعه آموزش، اعتبارسنجی و آزمایش با رهیافت تصادفی.

نسبت جداسازی تصادفی	تعداد داده آموزش	تعداد داده اعتبارسنجی	تعداد داده آزمایش
۹۸/۱/۱	۴۹۰۰۰۰	۵۰۰۰	۵۰۰۰
۹۰/۵/۵	۴۵۰۰۰۰	۲۵۰۰۰	۲۵۰۰۰
۷۰/۱۵/۱۵	۳۵۰۰۰۰	۷۵۰۰۰	۷۵۰۰۰



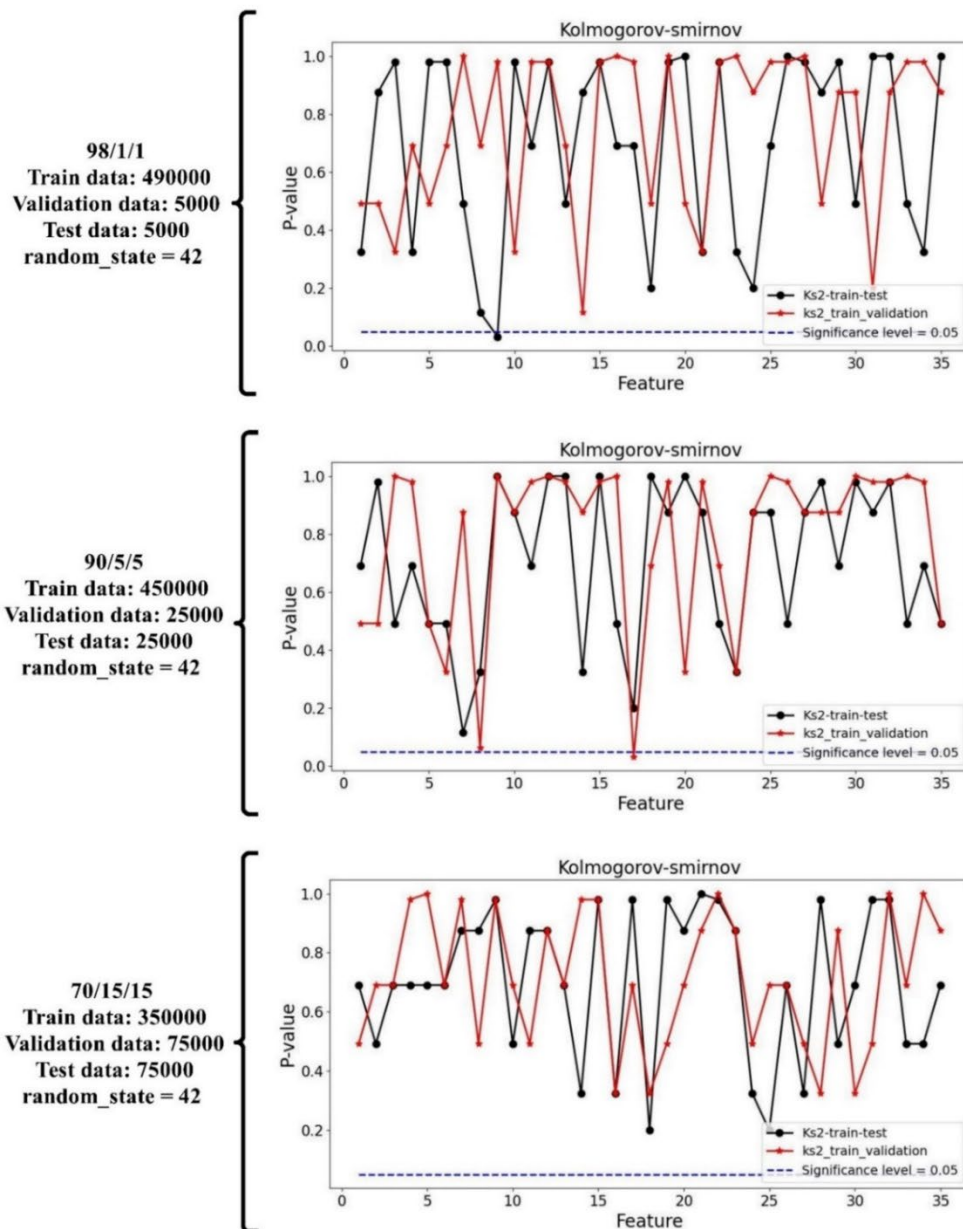
شکل ۲. رفتار نامنظم در منحنی خطای اعتبارسنجی به دلیل یکسان نبودن توزیع آماری داده‌های اعتبارسنجی و آموزش.

۳ خوشه‌بندی داده‌ها

یکی از مسائلی که در آموزش شبکه با آن مواجه هستیم، انتخاب نسبت جداسازی داده‌های آموزشی، اعتبارسنجی و آزمایش است. بر اساس مطالعات مرتبط با یادگیری عمیق، اگر تعداد داده‌ها زیاد باشد معمولاً نسبت جداسازی به صورت ۹۸/۱/۱ و اگر تعداد داده‌ها کم باشد، به شکل ۷۰/۱۵/۱۵ در نظر گرفته می‌شود (اولانی مورینا، ۲۰۲۲)؛ ولی هیچ معیار مشخصی برای نحوه توزیع داده‌ها در سه مجموعه وجود ندارد و عموماً به شکل تصادفی انجام می‌شود. حالت ایده‌آل این است که توزیع آماری سه مجموعه یکسان باشد؛ در غیر اینصورت منحنی خطای اعتبارسنجی رفتار نامنظمی را نشان می‌دهد.

دو راهکار برای یکسان‌سازی توزیع آماری سه مجموعه داده وجود دارد: ۱- تعداد داده‌ها زیاد و نسبت جداسازی کم شود. ۲- با استفاده از روش‌های خوشه‌بندی توزیع آماری آنها یکسان شود. ابتدا این موضوع بررسی می‌شود که آیا با نسبت‌های مختلف امکان دستیابی به توزیع آماری یکسان برای سه مجموعه داده وجود دارد. آزمون KS یکی از آماره‌هایی است که برای بررسی یکسان بودن توزیع آماری دو مجموعه به کار می‌رود. در این بررسی فرض

صفر مبتنی بر اینکه توزیع داده‌ها یکنواخت است، در سطح خطای ۰/۰۵ که سطح معناداری (significance level) نام دارد، آزمایش می‌شود. اگر آماره آزمون یا همان مقدار P- (value) بزرگتر از یا مساوی با سطح خطای ۰/۰۵ باشد، در این صورت دلیلی برای رد فرض صفر وجود ندارد و توزیع داده‌ها یکنواخت است (فرجی، ۱۳۸۵). مقدار P در واقع اختلاف بین تابع توزیع تجربی هر یک از ویژگی‌های دو مجموعه است و برای اینکه توزیع آماری دو مجموعه داده یکسان باشد، باید مقدار این آماره برای همه ویژگی‌ها بزرگتر از سطح معناداری باشد. هر نمونه داده تولید شده در بخش ۲-۱ با ۳۵ ویژگی (۲۶ ویژگی مرتبط با پاسخ‌های مدل‌سازی پیشرو شامل ۱۳ مقدار مقاومت ویژه ظاهری و ۱۳ مقدار فاز به‌اضافه ۹ ویژگی مرتبط با پارامترهای فیزیکی و هندسی مدل زمین‌گرمایی انتخاب شده) معرفی می‌شود. جدول ۲ تعداد داده‌ها در سه مجموعه داده آموزش، اعتبارسنجی و آزمایش در نسبت‌های مختلف جداسازی تصادفی ۹۸/۱/۱، ۹۰/۵/۵ و ۷۰/۱۵/۱۵ و شکل ۳ نتیجه اعمال آزمون KS برای ویژگی‌های مختلف و نسبت‌های متفاوت جداسازی تصادفی داده‌ها را نشان می‌دهد.



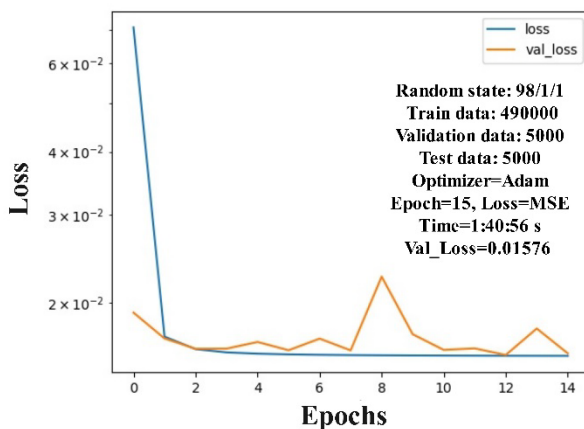
شکل ۳. نتیجه آزمون KS برای ارزیابی یکسان بودن توزیع آماری سه مجموعه داده آموزش، اعتبارسنجی و آزمایش با نسبت‌های مختلف جداسازی تصادفی ۹۸/۱/۱، ۹۰/۵/۵ و ۷۰/۱۵/۱۵ برای مجموعه داده ۵۰۰۰۰۰ تایی.

افزایش تعداد داده‌های اعتبارسنجی، اساساً توزیع آماری داده‌ها یکسان می‌شود. البته شایان ذکر است جداسازی داده‌ها به صورت تصادفی در برنامه نویسی، ارتباط مستقیم با عدد حالت تصادفی (Random State) دارد و با تغییر عدد حالت تصادفی ممکن است در نسبت‌های ذکر شده، در توزیع آماری داده‌ها در سه مجموعه داده تغییر ایجاد

همان‌طور که در شکل ۳ ملاحظه می‌شود در نسبت‌های ۹۸/۱/۱ و ۹۰/۵/۵ آزمون KS پایین‌تر از سطح معنی‌دار ۰/۰۵ قرار دارد و بیانگر این است که توزیع آماری بین داده‌های آموزش و آزمایش در نسبت ۹۸/۱/۱ و بین داده‌های آموزش و اعتبارسنجی در نسبت ۹۰/۵/۵، یکسان نیست. با کم شدن نسبت جداسازی داده‌ها (۷۰/۱۵/۱۵) و

نبودن توزیع آماری دو مجموعه داده آموزش و اعتبارسنجی، روند منحنی خطای آموزش، نزولی، اما روند منحنی خطای اعتبارسنجی، نامنظم و دارای نوسان است و مدل تعمیم‌پذیری مناسبی را ارائه نمی‌دهد.

شکل (۴) منحنی‌های خطای آموزش و اعتبارسنجی با معماری U-Net در جداسازی تصادفی با نسبت ۹۸/۱/۱ را نشان می‌دهد. در معماری U-Net از الگوریتم Adam و تنظیمات فرآیندهای گفته‌شده در بخش ۲-۳ استفاده شده‌است. همانگونه که مشاهده می‌شود، به دلیل یکسان



شکل ۴. رفتار نامنظم در منحنی خطای اعتبارسنجی به دلیل یکسان نبودن توزیع آماری داده‌های اعتبارسنجی و آموزش در نسبت تصادفی ۹۸/۱/۱.

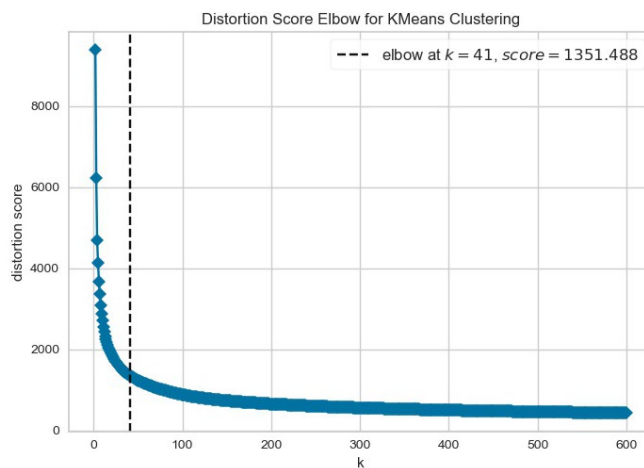
متفاوت نگه داشته می‌شوند. نقاط داده طوری به یک خوشه اختصاص پیدا می‌کنند که مجموع مجذور فاصله بین نقاط داده و مرکزوار خوشه (میانگین حسابی تمام نقاط داده متعلق به آن خوشه)، به‌عنوان معیار تشابه کمینه باشد. نحوه پیاده‌سازی الگوریتم به این ترتیب است: تعداد خوشه‌ها K مشخص می‌شود؛ با انتخاب تصادفی K نقطه داده برای مرکزوارها بدون جایگزینی، مرکزوارها در حالت اولیه قرار می‌گیرند؛ تا زمانی که بعد از آن هیچ تغییری در مرکزوارها ایجاد نشود، یعنی تخصیص نقاط داده به خوشه‌ها تغییر نکند، تکرارها ادامه پیدا می‌کنند؛ مجموع مجذور فاصله بین نقاط داده و همه مرکزوارها محاسبه می‌شوند؛ هر نقطه داده به نزدیک‌ترین خوشه یا مرکزوار اختصاص داده می‌شود؛ با محاسبه میانگین تمام نقاط داده‌ای که به هر خوشه تعلق دارند، مرکزوارها برای خوشه‌ها محاسبه می‌شوند (مک‌کوین، ۱۹۶۷). منحنی زانو یک روش گرافیکی برای یافتن مقدار بهینه K است. در این نمودار

رویکرد این مطالعه، اعمال روش‌های خوشه‌بندی بر داده‌ها و اختصاص درصد مشخصی از هر خوشه به سه مجموعه داده است. با این رویکرد، تعداد داده‌های مورد نیاز و زمان آموزش، کاهش و دقت پیش‌بینی افزایش می‌یابد. در ابتدا ۵۰۰۰۰۰ نمونه تولید شد (بخش ۲-۱) و هر نمونه داده با ۳۵ ویژگی، ۲۶ ویژگی مرتبط با پاسخ مدل‌سازی پیشرو یا ورودی و ۹ ویژگی مرتبط با پارامترهای مدل یا خروجی، معرفی می‌شود. الگوریتم K-means به‌عنوان یکی از پرکاربردترین الگوریتم‌های خوشه‌بندی بر داده‌های MT تولید شده اعمال و از نتایج خوشه‌بندی برای ایجاد و انتخاب مجموعه داده نمونه آموزشی استفاده شد. K-means یک الگوریتم مبتنی بر تکرار است که مجموعه داده را به زیرگروه‌ها یا خوشه‌های از قبل تعریف‌شده متمایز بدون همپوشانی تقسیم می‌کند که در آن هر نقطه داده فقط به یک گروه تعلق دارد. نقاط داده درون خوشه‌ای تا حد ممکن شبیه هم هستند و در عین حال خوشه‌ها تا حد ممکن

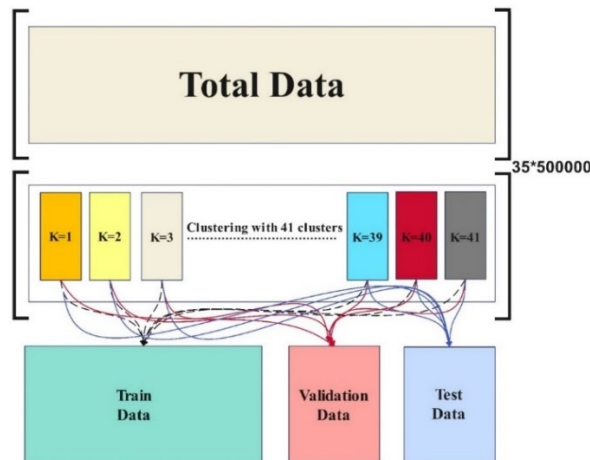
خوشه‌بندی می‌شوند و از هر خوشه با درصد‌های متفاوتی داده‌ها به سه مجموعه آموزش، اعتبارسنجی و آزمایش اختصاص داده می‌شوند (شکل ۶). هدف، دستیابی به بهترین نتیجه با کمترین تعداد داده است. جدول ۳ تعداد داده‌ها را در سه مجموعه آموزش، اعتبارسنجی و آموزش با رویکرد خوشه‌بندی نشان می‌دهد.

جدول ۳. تعداد داده‌ها در سه مجموعه آموزش، اعتبارسنجی و آزمایش با رویکرد خوشه‌بندی.

درصد‌های متفاوت از نسبت جداسازی ۹۸/۱/۱	تعداد داده آموزش	تعداد داده اعتبارسنجی	تعداد داده آزمایش
۹۸/۱/۱	۴۹۸۹۸۲	۴۹۸۰	۴۹۸۰
۱۰٪(۹۸/۱/۱)	۴۸۹۷۹	۴۸۲	۴۸۲
۲۰٪(۹۸/۱/۱)	۹۷۹۷۸	۹۸۱	۹۸۱
۳۰٪(۹۸/۱/۱)	۱۴۶۹۷۹	۱۴۷۸	۱۴۷۸
۴۰٪(۹۸/۱/۱)	۱۹۵۹۸۲	۱۹۸۱	۱۹۸۱
۵۰٪(۹۸/۱/۱)	۲۴۴۹۸۱	۲۴۷۹	۲۴۷۹



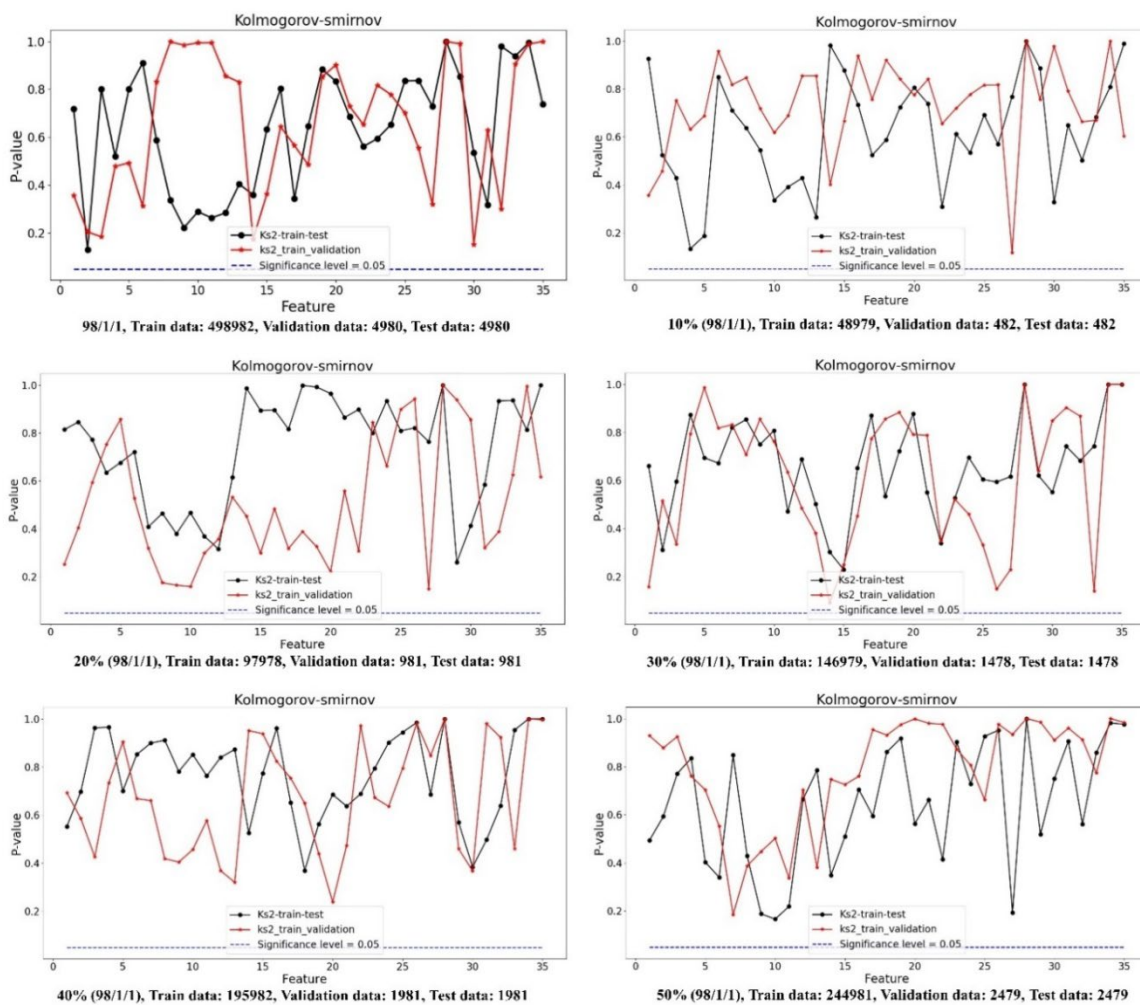
شکل ۵. منحنی مجموع مربعات فواصل درون خوشه‌ای بر حسب مقادیر مختلف K.



شکل ۶. پیاده‌سازی الگوریتم Kmeans برای خوشه‌بندی داده‌های آموزشی MT.

در همه حالت‌ها آماره آزمون KS بزرگتر از سطح معنادار ۰/۰۵ است و معرف یکسان بودن توزیع داده‌ها با نسبت‌های داده‌ای متفاوت است.

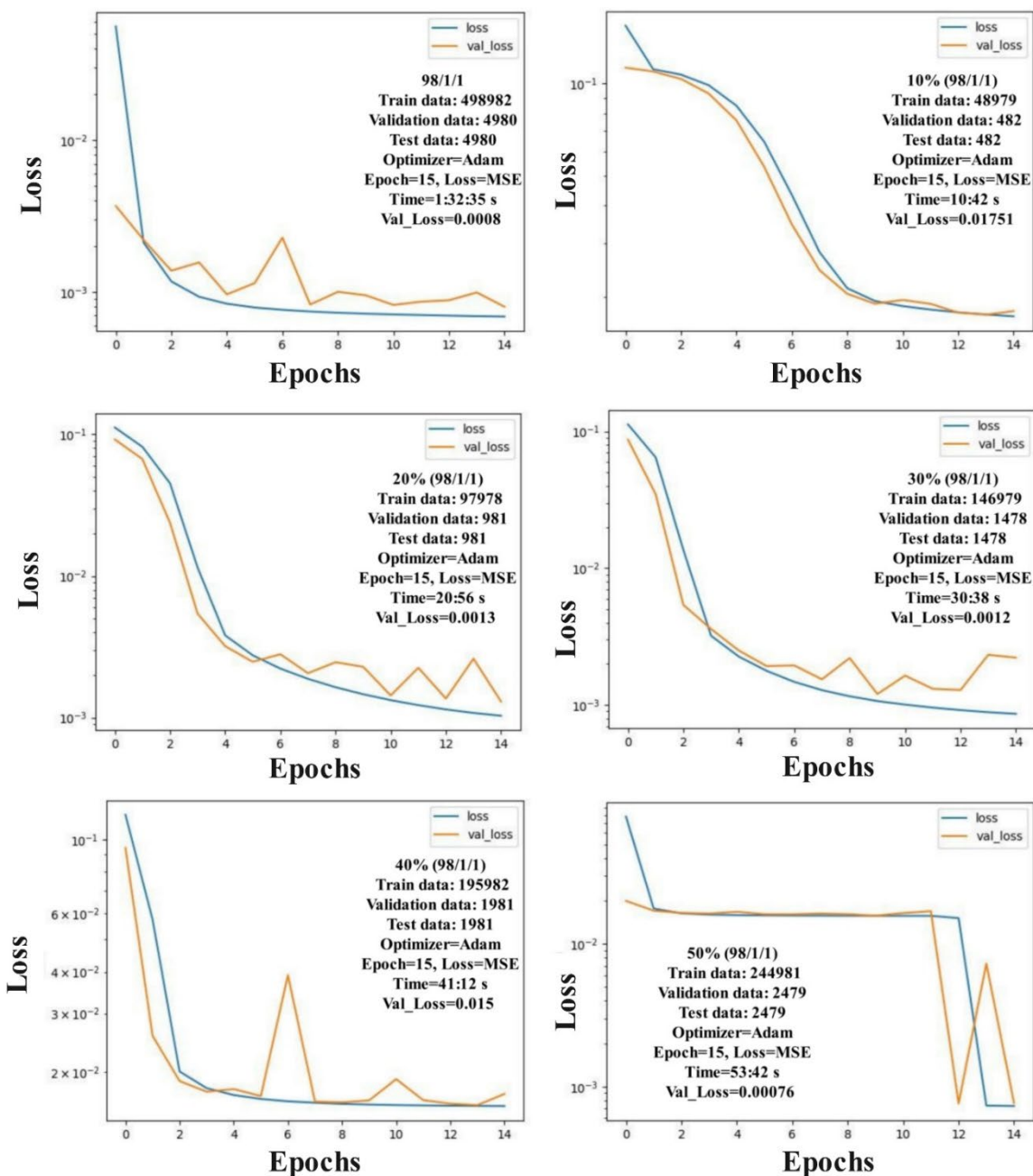
قبل از آموزش شبکه با داده‌های خوشه‌بندی شده، توزیع آماری سه مجموعه داده در حالت‌های متفاوت مقایسه می‌شوند. به این منظور از آزمون KS بهره گرفته می‌شود. نتایج در شکل ۷ نشان داده شده‌اند. چنانکه مشاهده می‌شود



شکل ۷. نتایج آزمون KS برحسب ویژگی برای درصد‌های متفاوت از کل داده‌ها در نسبت ۹۸/۱/۱.

می‌شوند. منحنی‌های خطای آموزش و اعتبارسنجی که در یادگیری عمیق پایش می‌شوند، در شکل ۸ نمایش داده شده‌اند. چنانکه مشاهده می‌شود در حالتی که ۵۰ درصد از داده‌های کل به کار گرفته شده‌اند، کمترین خطای یادگیری محقق می‌شود.

در مرحله بعد، داده‌های حاصل از خوشه‌بندی با درصد‌های متفاوت از مقدار کل و نسبت جداسازی ۹۸/۱/۱ با معماری U-Net و الگوریتم Adam و تنظیمات گفته شده در بخش ۲-۳ آموزش داده می‌شوند و نتایج با حالتی که کل داده‌ها با نسبت جداسازی ۹۸/۱/۱ وارد شده‌اند، مقایسه



شکل ۸. منحنی‌های خطای آموزش و اعتبارسنجی برای درصد‌های متفاوت از کل داده‌ها.

اعتبارسنجی، در تکراری که خطای اعتبارسنجی کمترین مقدار را دارد ماتریس وزن نوروها به عنوان بهترین مدل ذخیره می‌شود. نتایج وارون‌سازی داده‌های حاصل از خوشه‌بندی توسط معماری U-Net در جدول ۴ نشان می‌دهد که کمترین خطای آموزش در نسبت‌های ۹۸/۱/۱،

همانطور که در شکل ۸ ملاحظه می‌شود، در نسبت‌های ۹۸/۱/۱، ۱۰٪ (۹۸/۱/۱)، ۲۰٪ (۹۸/۱/۱)، ۵۰٪ (۹۸/۱/۱) و ۳۰٪ (۹۸/۱/۱) و ۴۰٪ (۹۸/۱/۱) برآزش بیش از حد اتفاق افتاده است که با فراخوانی تابع ذخیره بهترین وزن‌ها در پایش منحنی خطای

بوده است. در نسبت های ۹۸/۱/۱ و ۵۰٪(۹۸/۱/۱) کمینه خطای اعتبارسنجی نزدیک به هم هستند؛ بنابراین در آموزش شبکه در نسبت داده ای ۵۰٪(۹۸/۱/۱) با زمان آموزش کمتر و با تعداد داده کمتر، نتایج مشابه با نسبت داده ای ۹۸/۱/۱ حاصل خواهد شد.

۵۰٪(۹۸/۱/۱) و ۳۰٪(۹۸/۱/۱) اتفاق افتاده است. عملکرد یادگیری عمیق در نسبت های ۹۸/۱/۱ و ۵۰٪(۹۸/۱/۱) به صورت برازش خوب و در نسبت ۳۰٪(۹۸/۱/۱) به صورت بیش برازش است. در نسبت ۳۰٪(۹۸/۱/۱) ذخیره بهترین وزن ها با کمینه خطای اعتبارسنجی ۰/۰۰۱۲ در تکرار ۱۳

جدول ۴. نتایج وارون سازی داده های حاصل از خوشه بندی توسط معماری U-Net.

مدت زمان آموزش	شماره تکرار در ذخیره بهترین مدل	کمینه خطای اعتبارسنجی	کمینه خطای آموزش	عملکرد یادگیری	درصد های متفاوت از نسبت جداسازی ۹۸/۱/۱
۱ : ۳۲ : ۳۵ S	۱۵	۰/۰۰۰۸	۰/۰۰۰۷	برازش خوب	۹۸/۱/۱
۱۰ : ۴۲ S	۱۵	۰/۰۱۷۵۱	۰/۰۱۵	برازش خوب	۱۰٪(۹۸/۱/۱)
۲۰ : ۵۶ S	۱۵	۰/۰۰۱۳	۰/۰۰۱	برازش خوب	۲۰٪(۹۸/۱/۱)
۳۰ : ۳۸ S	۱۳	۰/۰۰۱۲	۰/۰۰۰۸	بیش برازش	۳۰٪(۹۸/۱/۱)
۴۱ : ۱۲ S	۱۴	۰/۰۱۵	۰/۰۱۴	بیش برازش	۴۰٪(۹۸/۱/۱)
۵۳ : ۴۲ S	۱۵	۰/۰۰۰۷۶	۰/۰۰۰۷	برازش خوب	۵۰٪(۹۸/۱/۱)

در شکل ۹ ترسیم شده است. میانه و دامنه میان چارکی نمودار جعبه ای با نسبت داده ای ۵۰٪(۹۸/۱/۱)، بیشترین نزدیکی را با نسبت داده ای ۹۸/۱/۱ دارد؛ ضمناً بیشترین تقارن هم در همین حالت مشاهده می شود، هرچند که کمترین تعداد داده پرت مربوط به نسبت ۲۰٪(۹۸/۱/۱) است.

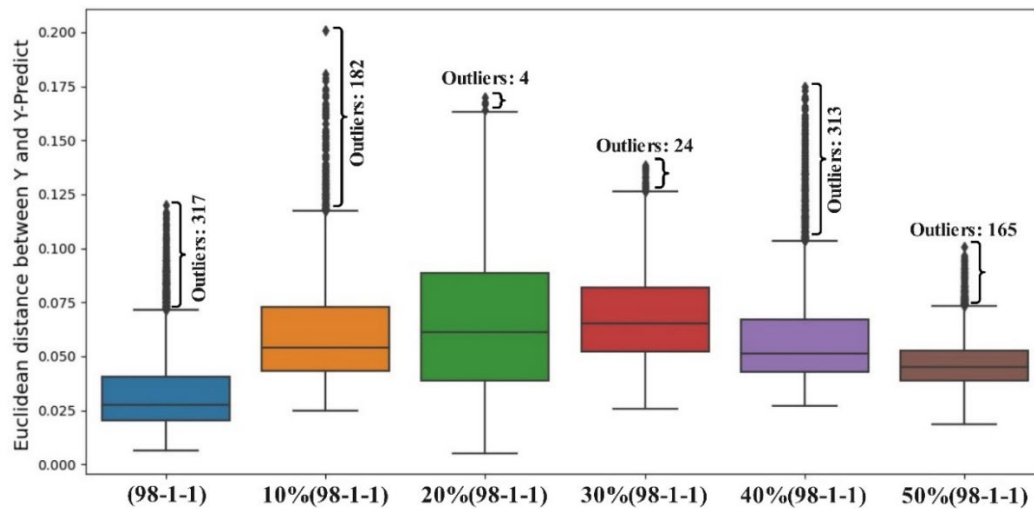
۴-۲ بررسی کمی نتایج وارون سازی

ضریب بهره وری ناش ساتکلیف (Nash-Sutcliffe efficiency: NSE) در سال ۱۹۷۰ به منظور ارزیابی کمی درجه انطباق سری های زمانی مشاهده و پیش بینی شده معرفی شد (ناش و ساتکلیف، ۱۹۷۰). از این معیار در اینجا برای ارزیابی میزان برازش بین منحنی های مقاومت ویژه ظاهری و فاز پیش بینی شده و واقعی استفاده شده است. ضریب NSE نشان می دهد که خط رگرسیون بین داده های مشاهده و پیش بینی شده تا چه حد به خط رگرسیون با شیب

۴ بحث و ارزیابی نتایج

۴-۱ بررسی کیفی نتایج وارون سازی

برای بررسی کیفی نتایج وارون سازی، از یک مجموعه داده آزمایشی یکسان (با ۴۹۸۰ نمونه داده) برای پیش بینی نتایج استفاده شد. از آنجایی که پیش بینی نتایج، نوعی رگرسیون دو متغیره است (مقاومت ویژه هر لایه همراه با ضخامت آن) به دلیل ماهیت متفاوت دو کمیت مقاومت ویژه و ضخامت، ابتدا مقادیر داده های پیش بینی و واقعی نرمال سازی شده و سپس از فاصله اقلیدوسی برای بررسی اختلاف مقادیر پیش بینی شده و واقعی استفاده می شود. نمودار جعبه ای، یک روش استاندارد برای نمایش توزیع یک مجموعه داده بر اساس پنج ویژگی عددی کمینه، چارک اول، میانه، چارک سوم و بیشینه آن نقاط داده است. بنابراین ابزار مفیدی برای نشان دادن مقادیر میانگین، میزان تقارن، مقادیر پرت و پراکندگی مجموعه داده ها محسوب می شود. فاصله اقلیدوسی در قالب نمودار جعبه ای برای شش نسبت مختلف



شکل ۹. نمودار جعبه‌ای فواصل اقلیدسی بین مقادیر واقعی و پیش‌بینی شده توسط شبکه عمیق.

مقادیر کوچک یا منفی معرف عملکرد ضعیف‌تر مدل هستند (واندرکلن و همکاران، ۲۰۱۸).

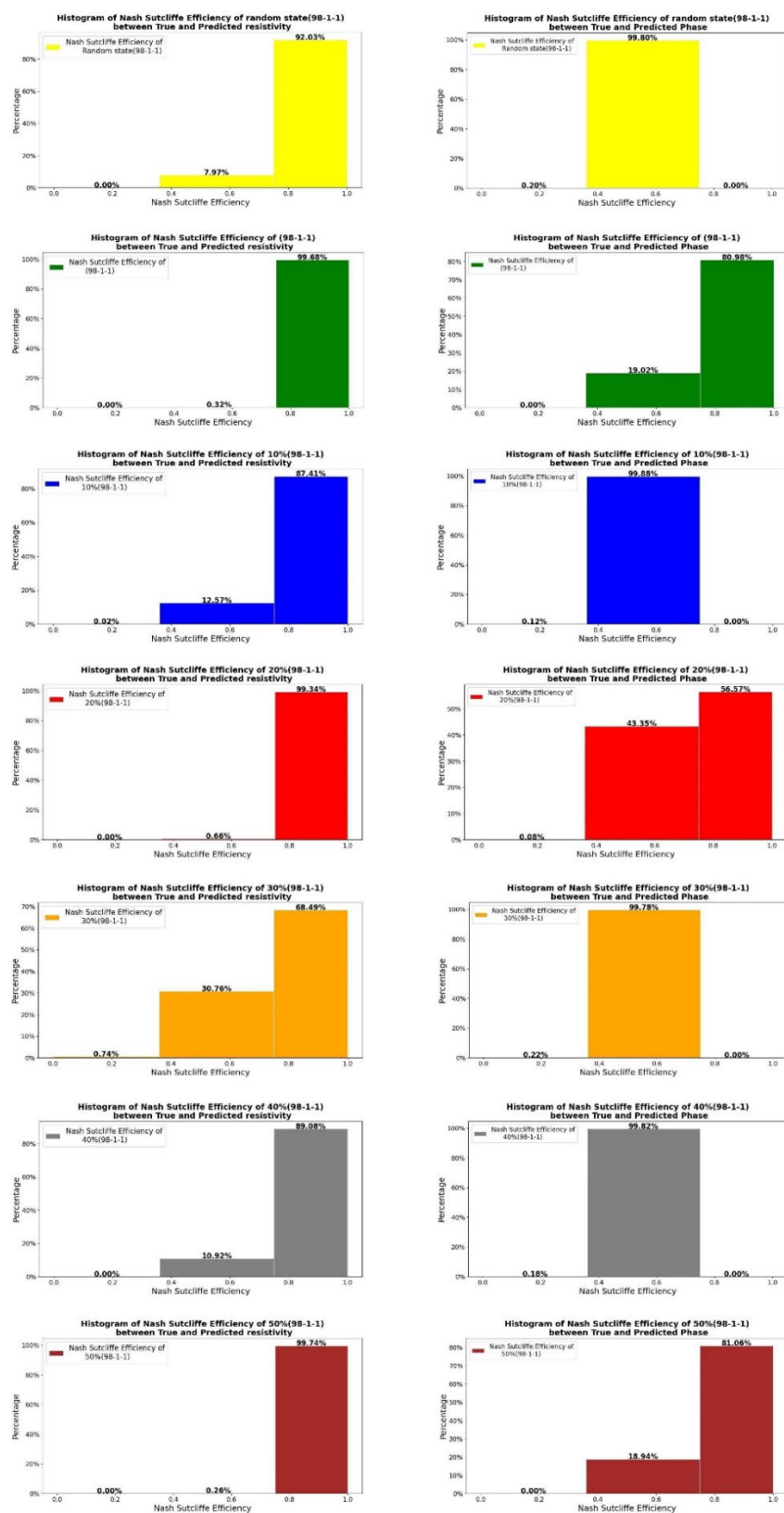
شکل (۱۰) هیستوگرام درصد فراوانی ضریب NSE برای منحنی‌های مقاومت ویژه ظاهری و فاز در جداسازی تصادفی و نسبت‌های مختلف داده‌ای در حالت خوشه‌بندی را نشان می‌دهد. در هیستوگرام منحنی‌های مقاومت ویژه ظاهری در جداسازی تصادفی (۹۸/۱/۱)، ۹۳/۰۳٪ نمونه‌ها دارای ضرائب بزرگتر از ۰/۷۵ و ۷/۹۳٪ نمونه‌ها دارای ضرائب بین ۰/۳۶ و ۰/۷۵ هستند؛ در هیستوگرام منحنی‌های فاز در همین حالت، ۹۹/۸٪ دارای ضرائب بین ۰/۳۶ و ۰/۷۵ و سایر نمونه‌ها دارای مقادیر کوچکتر از ۰/۷۵ هستند. در حالی که در هیستوگرام منحنی‌های مقاومت ویژه ظاهری در حالت خوشه‌بندی و در نسبت داده‌ای ۹۸/۱/۱، ۹۹/۶۸٪ نمونه‌ها ضرائب بزرگتر از ۰/۷۵ و در هیستوگرام منحنی‌های فاز در همین حالت، ۱۹/۰۲٪ نمونه‌ها ضرائب بین ۰/۳۶ و ۰/۷۵ و ۸۰/۹۹٪ نمونه‌ها ضرائب بزرگتر از ۰/۷۵ دارند. بنابراین نتایج وارون‌سازی با استفاده از داده‌های خوشه‌بندی شده، بسیار بهتر از نتایج وارون‌سازی با استفاده از داده‌هایی است که به شکل تصادفی جداسازی شده‌اند. با توجه به این شکل، از بین دیگر نسبت‌ها، تنها نسبتی که

یک نزدیک است. این ضریب بر اساس رابطه ۵ محاسبه می‌شود:

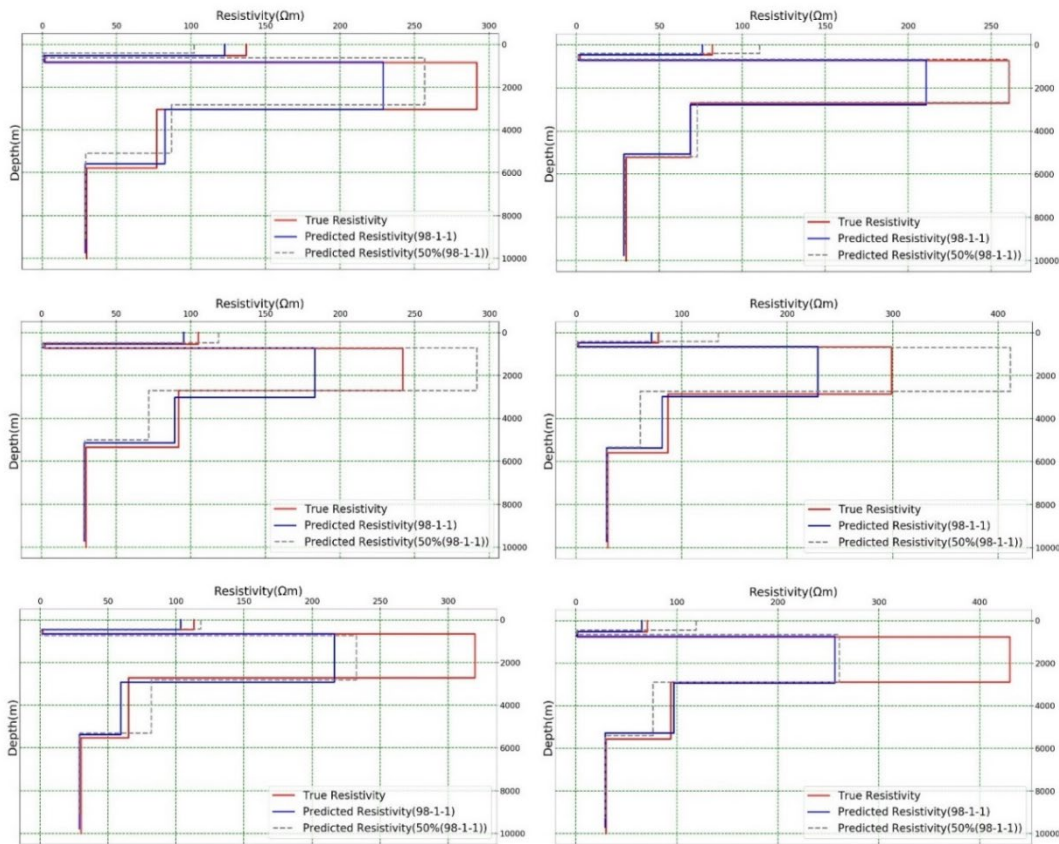
$$NSE_{\rho_a} = 1 - \frac{\left[\sum_i (\rho_{a_i}^{obs} - \rho_{a_i}^{pre})^2 \right]}{\left[\sum_i (\rho_{a_i}^{obs} - \rho^{mean})^2 \right]} ; NSE_{ph} \quad (5)$$

$$= 1 - \frac{\left[\sum_i (ph_i^{obs} - ph_i^{pre})^2 \right]}{\left[\sum_i (ph_i^{obs} - ph^{mean})^2 \right]}$$

که در آن $\rho_{a_i}^{obs}$ و ph_i^{obs} مقاومت ویژه ظاهری و فاز مشاهده شده، $\rho_{a_i}^{pre}$ و ph_i^{pre} مقاومت ویژه ظاهری و فاز پیش‌بینی شده و ρ^{mean} و ph^{mean} میانگین مقاومت ویژه ظاهری و فاز مشاهده شده هستند. در دسته بندی ضریب ناش ساتکلیف، مقادیر کمتر از ۰/۳۶، ضعیف، مقادیر بین ۰/۳۶ تا ۰/۷۵ به عنوان خوب و مقادیر بیشتر از ۰/۷۵ عالی در نظر گرفته می‌شوند. برای ارزیابی کمی نتایج وارون‌سازی، ضریب NSE منحنی‌های مقاومت ویژه ظاهری و فازی که با استفاده از مدل‌سازی پیشرو از روی نتایج پیش‌بینی شبکه به دست آمده‌اند و داده‌های واقعی برای نسبت‌های جداسازی مختلف محاسبه و درصد فراوانی آنها رسم می‌شود. مقادیر NSE از بی‌نهایت منفی تا یک متغیر هستند؛ مقادیر بزرگ‌تر نشان‌دهنده عملکرد بهتر و



شکل ۱۰. هیستوگرام درصد فراوانی ضرایب NSE برای منحنی‌های مقاومت ویژه ظاهری و فاز در حالت جداسازی تصادفی و نسبت داده‌ای با درصد‌های متفاوت در حالت خوشه‌بندی.



شکل ۱۱. نتایج وارون‌سازی شش نمونه اول از مجموعه داده آزمایش در دو نسبت داده‌ای (۹۸/۱/۱) و (۵۰/۹۸/۱/۱).

است. ضخامت‌های پیش‌بینی شده برای لایه‌ها دارای انطباق خوبی با مدل واقعی هستند. بیشترین اختلاف بین مقادیر پیش‌بینی شده و مدل واقعی، در پیش‌بینی مقاومت ویژه لایه سوم یا مخزن رخ داده‌است. این لایه بین دو لایه با مقاومت ویژه پایین‌تر قرار گرفته‌است و هرچه به بیشینه مقدار مقاومت ویژه خود نزدیک می‌شود، اختلاف مشاهده‌شده بزرگ‌تر می‌شود. با احتساب محدودیت ذاتی در توان تفکیک سونداژهای MT، می‌توان قابلیت مدل‌های یادگیری عمیق در غلبه بر بخشی از توان تفکیک محدود پاسخ‌های MT را بررسی کرد که در مطالعه انجام شده توسط رحمانی جوینانی و همکاران (۲۰۲۴) تا حدودی به آن پرداخته شده‌است. ضمناً می‌توان حدس زد که اگر تعداد فرکانس‌های به کاررفته در تولید داده بیشتر از مقدار

نتایج مشابه با نسبت داده‌ای (۹۸/۱/۱) دارد نسبت (۵۰/۹۸/۱/۱) است. در این نسبت داده‌ای، در هیستوگرام منحنی‌های مقاومت ویژه، ۹۹/۷۴٪ نمونه‌ها مقادیر بیشتر از ۰/۷۵ و در هیستوگرام منحنی‌های فاز ۱۸/۹۴٪ نمونه‌ها مقادیر بین ۰/۳۶ تا ۰/۷۵ و ۸۱/۰۶٪ نمونه‌ها مقادیر بیشتر از ۰/۷۵ دارند. بنابراین می‌توان نتیجه گرفت که با نیمی از داده‌ها می‌توان به نتایج حاصل از کل داده‌ها رسید و این امر باعث کاهش نیاز محاسباتی و زمان آموزش می‌شود.

شکل ۱۱ نتایج وارون‌سازی شبکه برای شش نمونه اول از مجموعه داده آزمایش را برای دو حجم ۱۰۰ و ۵۰ درصد از داده‌های جداسازی شده با نسبت ۹۸/۱/۱ نشان می‌دهد. نتایج در این دو حالت بسیار نزدیک به هم و همسو با مدل واقعی یعنی مدل زمین‌گرایی به کار رفته برای تولید داده

باعث کاهش نیاز محاسباتی و زمان آموزش می‌شود. این بررسی برای مدل ساده چندلایه و داده‌های تمیز اعمال شده است؛ با این وجود با اعمال روش بر داده‌های آلوده به نوفه و حاصل از مدل‌های پیچیده‌تر، مسلماً تعداد داده‌های آموزشی مورد نیاز و زمان آموزش افزایش می‌یابد و قطعاً جداسازی تصادفی داده‌ها رهیافت مناسبی نخواهد بود. بنابراین پیشنهاد می‌شود صرف نظر از ساده یا پیچیده بودن مدل، بعد از جداسازی داده‌ها آزمون KS روی سه مجموعه داده اعمال و از یکسان بودن توزیع آماری آنها اطمینان حاصل شود تا تعمیم پذیری مدل با مشکل مواجه نشود. در شرایطی که توزیع‌های آماری یکسانی رخ ندهند، خوشه‌بندی می‌تواند راهکار مناسبی برای یکسان‌سازی توزیع آماری سه مجموعه و کاهش تعداد داده‌های مورد نیاز برای آموزش باشد.

منابع

- خسروی، ن. ۱۳۸۵، آمار توصیفی و استنباطی. چاپ اول، انتشارات پوران پژوهش.
- Alali, A., Morgan, F.D., Coles, D., 2020. Novel approach for 1D resistivity inversion using the systematically determined optimum number of layers. *J Geol Geophys*, 9(6), 481.
- Caldwell, T.G., Bibby, H.M., Brown, C., 2004. The magnetotelluric phase tensor. *Geophys. J. Int.* 158, 457-469.
- Chen, J., Hoversten, G.M., Key, K., Nordquist, G., Cumming W., 2012. Stochastic inversion of magnetotelluric data using a sharp boundary parameterization and application to a geothermal site. *Geophysics*, 77(4), E265-E279.
- Comeau, M.J., Becken, M., Grayver, A.V., Käüfl, J.S., Kuvshinov, A.V., 2022. The geophysical signature of a continental intraplate volcanic system: from surface to mantle source. *Earth and Planetary Science Letters* 578, 117307.
- Constable, S.C., Parker, R.L., Constable, C.G.,

فعلی (۱۳) فرکانس در بازه ۱۰۰-۰/۰۱ هرتز) باشد، یادگیری شبکه و بنابراین انطباق بین مقادیر مقاومت ویژه پیش‌بینی شده و واقعی بهبود می‌یابد.

۵ نتیجه‌گیری

در این پژوهش برای غلبه بر دو چالش نیاز به تعداد داده‌های زیاد در آموزش شبکه‌های عصبی عمیق و عدم تعمیم‌پذیری شبکه آموزش دیده به سبب یکسان نبودن توزیع آماری سه مجموعه داده آموزش، اعتبارسنجی و آزمایش، رویکرد کاهش مبتنی بر خوشه‌بندی داده‌ها ارائه شده است. عموماً جداسازی تصادفی داده‌ها برای آموزش شبکه‌های عصبی عمیق توزیع آماری یکسانی را برای سه مجموعه ایجاد نمی‌کند و رویکرد مبتنی بر خوشه‌بندی جایگزین مناسبی برای جداسازی تصادفی داده‌ها است. برای حصول اطمینان از یکسان بودن توزیع آماری سه مجموعه می‌توان از آزمون KS استفاده کرد. در خوشه‌بندی به روش K-means پس از انتخاب تعداد بهینه خوشه‌ها، از هر خوشه سهمی در سه مجموعه داده وارد می‌شود. با این تدبیر چالش یکسان نبودن توزیع آماری سه مجموعه داده برطرف می‌گردد. برای کاهش تعداد داده‌ها و بار محاسباتی ناشی از آن، از ۴۱ خوشه بهینه با درصد‌های مختلف (۰٪، ۱۰٪، ۲۰٪، ۳۰٪، ۴۰٪ و ۵۰٪) از نسبت داده‌ای ۹۸/۱/۱ در قالب سه مجموعه داده آموزش، اعتبارسنجی و آزمایش، داده انتخاب می‌شود. از این داده‌ها برای آموزش یک شبکه یادگیری عمیق با معماری U-Net به منظور وارون‌سازی یک‌بعدی داده‌های مگنتوتلوریک و تعیین مشخصات یک مدل زمین‌گرمایی استفاده شده است. عملکرد شبکه آموزش دیده با معیارهای متنوعی سنجیده می‌شود. با بررسی نتایج وارون‌سازی به صورت کمی با کمک نمودار جعبه‌ای و کیفی با بهره‌گیری از معیار ارزیابی NSE می‌توان نتیجه گرفت که با نیمی از داده‌ها (۹۸/۱/۱) (۵۰٪) می‌توان به نتایج حاصل از اعمال کل داده‌ها (۹۸/۱/۱) رسید و این امر

1987. Occam's inversion: a practical algorithm for generating smooth models from electromagnetic data. *Geophysics*, 52, 289-300.
- Fischer, G., Schnegg, P.A., Peguiron, M., LeQuang, B.V., 1981. An analytic one-dimensional magnetotelluric inversion scheme. *Geophysical Journal of the Royal Astronomical Society*, 67, 257-278.
- Goodfellow. I., Bengio. Y., Courville. A., 2016. *Deep Learning*. London. The MIT Press.
- Junge, A., 2011. A concept for 1D inversion of MT data using phase tensor invariants. 24. Schmucker-Weidelt-Kolloquium Neustadt an der Weinstrasse, 19-23 September.
- Kim, Y., Nakata, N., 2018. Geophysical inversion versus machine learning inversion in inverse problems. *Leading Edge*, 894-901.
- Kingma, D.P., Ba, J.L., 2014. Adam: a method for stochastic optimization. In *International Conference on Learning Representations*.
- Kodinariya, T.M., Makwana, P.R., 2013. Review on determining number of Cluster in K-Means Clustering. *International Journal of Advance Research in Computer Science and Management Studies*, 1 (6), 90-95.
- Liao, X., Zhang, Z., Yan, Q., Shi, Z., Xu, K., Jia, D., 2022. Inversion of 1-D magnetotelluric data using CNN-LSTM hybrid network. *Arabian Journal of Geosciences* 15, 1430.
- Liao, X., Shi, Z., Zhang, Z., Yan, Q., Liu, P., 2022. 2D inversion of magnetotelluric data using deep learning technology. *Acta Geophysica* 70, 1047-1060.
- Liu, Z., Chen, H., Ren, Z., Tang, J., Xu, Z., Chen, Y., Liu, X., 2021. Deep learning audio magnetotellurics inversion using residual-based deep convolutional neural network. *Journal of Applied Geophysics*, 188, 104309.
- Liu, W., Wang, H., Xi, Z., Zhang, R., Huang, X., 2022. Physics-driven deep learning inversion with application to magnetotelluric. *Remote Sens* 14, 3218.
- MacQueen, J. B., 1967. Some Methods for classification and Analysis of Multivariate Observations. *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability*. Vol. 1. University of California Press. pp. 281-297.
- Miri, H., Habibian Dehkordi, B., Payrovian, G., 2021. Oil field imaging on the Sarab Anticline, southwest of Iran, using magnetotelluric data. *Journal of Petroleum Science and Engineering*, 202, 108497.
- Nair, V., and E. G. Hinton, 2010, Rectified linear units improve restricted Boltzmann machines: *Proceedings of the 27th International Conference on International Conference on Machine Learning*, 807-814.
- Nash J, Sutcliffe J (1970) River flow forecasting through conceptual models part I: A discussion of principles. *Journal of Hydrology*, 10, 282-290.
- Oh, S., Noh, K., Seol, S.J., Byun, J., 2020. Cooperative deep learning inversion of controlled-source electromagnetic data for salt delineation. *Geophysics* 85(4), E121-E137.
- Olaniyi Muraina., I., 2022. Ideal dataset splitting ratios in machine learning algorithms: general concerns for data scientists and data analysts. 7th International mardin artuklu scientific researches conference, 496-504, www.artuklukongresi.org
- Parker, R.L., Booker, J.R., 1996. Optional one-dimensional inversion and bounding of magnetotelluric apparent resistivity and phase measurements. *Physics of the Earth and Planetary Interiors* 98, 269-282.
- Puzyrev, V., 2019. Deep learning electromagnetic inversion with convolutional neural networks. *Geophys. J. Int* 218, 817-832.
- Rahmani Jevinani, M., Habibian Dehkordi, B., Ferguson, I.J., Rohban, M.H., 2024. Deep learning-based 1-D magnetotelluric inversion: performance comparison of architectures. *Earth Science Informatics* 17, 1663-1677.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional networks for

- biomedical image segmentation. In: International conference on medical image computing and computer-assisted intervention. Springer, 234-241.
- Segovia, M.J., Diaz, D., Selzak, K., Zuiga, F., 2021. Magnetotelluric study in the Los Lagos region (Chile) to investigate volcano-tectonic processes in the Southern Andes. *Earth, Planets and Space* 73(5).
- Shahriari, M., Pardo, D., Picon, A., Galdran, A., Del Ser, J., Torres-Verdin, C., 2020. A deep learning approach to the inversion of borehole resistivity measurements. *Computational Geosciences* 24, 971-994.
- Smith, J.T, Booker, J.R., 1988. Magnetotelluric inversion for minimum structure. *Geophysics* 53, 1565-1576.
- Vanderkelen, I.; Van Lipzig, N.P.M.; Thiery, W. Modelling the Water Balance of Lake Victoria (East Africa)-Part 2: Future Projections. *Hydrol. Earth Syst. Sci.* 2018, 22, 5527–5549.

Reduction of the data required for training deep learning models based on clustering of the data and its application in one-dimensional magnetotelluric inversion

Mehdi Rahmani Jevinani¹ and Banafsheh Habibian Dehkordi^{2*}

¹ Ph.D. Student, Institute of Geophysics, University of Tehran, Tehran, Iran

² Assistant Professor, Institute of Geophysics, University of Tehran, Tehran, Iran

(Received: 27 March 2024, Accepted: 25 September 2024)

Summary

Data-driven deep learning approaches have to deal with the challenge of generating large amounts of high-quality data, as well as the heavy computational cost and long training time imposed by it. Due to their ability to approximate complex nonlinear mapping functions, deep networks can be used effectively in geophysical inverse problems and better generalization can be achieved through deeper networks in many applications. In this research, an approach based on primary clustering of training data and assigning a certain percentage of each cluster to training, validation and test data has been used for data splitting. Kolmogorov Smirnov (KS) test has been applied to compare the distribution of three sets that are divided in this manner, and indicates that the training, validation and test data have the same distribution. A deep learning model based on modified U-Net architecture has been trained for one-dimensional inversion of magnetotelluric (MT) data, which is a highly non-linear regression problem. Supervised learning and back propagation error are used, and therefore, the inputs along with the corresponding outputs are given to the network in the form of training samples. For this purpose, a five-layer geoelectric model has been considered to simulate the conditions of a geothermal field. Using magnetotelluric forward modeling algorithm, the responses of this one-dimensional geoelectric model are analytically calculated in the frequency range of 0.01-100 Hz and in 13 frequencies that are uniformly distributed on a logarithmic scale, and total of 500000 sample data were generated. The thickness of the layers is variable and considered as part of the output. Pre-processing is done to scale the input and output variables before training and the network outputs are post-processed to be returned to the original interval. The mean square error (MSE) loss function and the Adam optimizer were used to train the network. Training is accomplished with a different amount of data separated by the mentioned method, and network performance is evaluated with some quantitative and qualitative criteria, including boxplots of Euclidean distance between true and predicted outputs and Nash Sutcliffe Efficiency coefficients. The trained network predicts the electrical resistivity and thickness of the layers from the new set of phase and apparent resistivity values. The results show that data splitting in this manner reduces the number of training data required to train the deep learning model by at least 50% without reducing the accuracy of the trained network. For noisy data and in more real scenarios, random separation is definitely not a suitable approach to form training, validation and test sets. In these conditions, the use of clustering is a suitable solution for equalizing the statistical distribution of the three sets and reducing the number of required data.

Keywords: Clustering, deep learning, inversion, magnetotelluric