

تکلیف دوم درس شبکه های عصبی

وحید ملکی، پوریا دادستان

۱۴۰۴ آذر

سوال ۱

توضیح دهید چرا شبکه های عصبی اغلب به عنوان "جعبه سیاه" شناخته می شوند و این مسئله چه چالش هایی برای قابلیت تفسیر و شفافیت مدل ایجاد می کند و راه های رفع این مشکل چیست؟

۱. چرا شبکه های عصبی "جعبه سیاه" (Black Box) هستند؟

اصطلاح "جعبه سیاه" در هوش مصنوعی به مدل های اشاره دارد که فرآیند تصمیم گیری درونی آنها برای انسان به سادگی قابل درک نیست. دلایل اصلی این ویژگی در شبکه های عصبی عبارتند از:

- پیچیدگی ساختاری و ابعاد بالا: شبکه های عصبی عمیق (Deep Neural Networks) اغلب شامل میلیون ها یا حتی میلیاردها پارامتر (وزن ها و بایاس ها) هستند. ردیابی اینکه چگونه یک ورودی خاص از طریق این اتصالات متعدد به یک خروجی منجر می شود، عملاً برای انسان غیرممکن است.

- غیرخطی بودن: استفاده از توابع فعال ساز غیرخطی (مانند Sigmoid ReLU یا Tanh) در لایه های پنهان باعث می شود که رابطه بین ورودی و خروجی یک نگاشت ریاضی بسیار پیچیده باشد که نمی توان آن را با یک رابطه خطی ساده توضیح داد.

- غایش ویژگی های انتزاعی: لایه های میانی شبکه، ویژگی هایی را استخراج می کنند که اغلب برای انسان معنای بصری یا منطقی مستقیمی ندارند (مثلاً ترکیبی از پیکسل ها که نه خط است و نه دایره، اما برای شبکه معنی دار است).

۲. چالش های ناشی از عدم تفسیر پذیری

عدم شفافیت در عملکرد شبکه های عصبی چالش های جدی ای را به وجود می آورد، از جمله:

- عدم اعتماد (Lack of Trust): در حوزه های حساس مانند پزشکی (تشخیص سرطان) یا مالی (تایید وام)، کاربران و متخصصان نیاز دارند بدانند چرا مدل یک تصمیم خاص را گرفته است. بدون توضیح منطقی، اعتماد به سیستم دشوار است.

- عیب یابی دشوار (Debugging Difficulty): وقتی یک شبکه عصبی اشتباه می کند، تشخیص اینکه کدام لایه یا کدام نورون باعث خطا شده است بسیار دشوار است، که فرآیند بهبود مدل را کند می کند.

• سوگیری و تبعیض (Bias and Fairness): مدل‌ها ممکن است سوگیری‌های موجود در داده‌های آموزشی را یاد بگیرند (مثلاً تبعیض نژادی یا جنسیتی). در یک جعبه سیاه، کشف این سوگیری‌ها قبل از وقوع فاجعه سخت است.

• مسائل قانونی و مقرراتی: قوانین مانند GDPR در اروپا بر "حق توضیح" (Right to Explanation) تاکید دارند. استفاده از مدل‌های جعبه سیاه در تصمیم‌گیری‌هایی که بر زندگی افراد تاثیر می‌گذارد، ممکن است با موانع قانونی رو به رو شود.

۳. راه‌های رفع مشکل: هوش مصنوعی توضیح‌پذیر (XAI)

برای مقابله با این چالش‌ها، حوزه "هوش مصنوعی توضیح‌پذیر" (Explainable AI) راهکارهایی را ارائه می‌دهد:

• روش‌های مستقل از مدل (Model-Agnostic Methods): تکنیک‌هایی مانند LIME و SHAP سعی می‌کنند رفتار مدل پیچیده را با تقریب زدن آن به یک مدل ساده‌تر (مثل رگرسیون خطی) در همسایگی یک نمونه خاص توضیح دهند. این روش‌ها نشان می‌دهند کدام ویژگی‌های ورودی بیشترین تاثیر را بر تصمیم نهایی داشته‌اند.

• بصری‌سازی (Visualization): در شبکه‌های کاتولوشنی، (CNN) می‌توان از روش‌هایی مانند Grad-CAM یا Maps استفاده کرد تا مشخص شود شبکه به کجا تصویر توجه کرده است.

• مکانیزم توجه (Attention Mechanisms): در مدل‌های پردازش زبان (مانند ترنسفورمرها)، مکانیزم توجه نشان می‌دهد که مدل هنگام تولید خروجی، روی کدام کلمات ورودی تمرکز کرده است که خود نوعی تفسیر‌پذیری ذاتی ایجاد می‌کند.

• استفاده از مدل‌های تفسیر‌پذیر: در مواردی که شفافیت اولویت بالاتری نسبت به دقت نهایی دارد، می‌توان به جای شبکه‌های عصبی عمیق از مدل‌های ساده‌تر درخت تصمیم (Decision Trees) یا رگرسیون لجستیک استفاده کرد.