

## روش‌های بصری سازی Feature Map ها در CNN

(1) Activation Maximization: یکی از اولین تکنیک‌های است که برای درک آنچه توسط یک فیلتر خاص در CNN فعال می‌شود، استفاده می‌شود. این روش با تولید یک تصویر که فعالیت یک واحد خاص را به حداکثر رسانده شروع می‌شود. این روش اغلب به استفاده از روش‌های بهینه‌سازی انجام می‌شود که تصویر ورودی را تغییر می‌دهند تا پاسخ خاص فیلتر را بیشینه کنند.

(2) Feature Maps Visualization: این تکنیک در واقع نمایش مستقیم نقشه‌های ویژگی است که در آن؟ می‌توانیم ببینیم که هر فیلتر چگونه به بخش‌های مختلف تصویر واکنش نشان می‌دهد. نقشه‌های ویژگی به ما نشان می‌دهند که کدام ویژگی‌ها در تصویر به هر لایه شبکه فعال می‌شوند و به ما کمک می‌کنند که بفهمیم شبکه چگونه تصمیم‌گیری می‌کند.

## (3) Backpropagation & Deconvolution

این روش با بازسازی تصاویر از ویژگی‌های میانی شبکه، تصمیم‌گیری‌های آن را قابل فهم و تفسیر می‌کند. در این روش فعالسازی‌های حاصل از یک لایه مشخص در شبکه‌های CNN مورد استفاده قرار می‌گیرند و سپس فراکنش معکوس انجام می‌شود تا بتوان ورودی‌هایی که به این فعالسازها منجر شده‌اند را بازسازی کرد. با این روش می‌توان دید که کدام بخش‌های تصویر در تصمیم‌گیری‌های مدل تأثیرگذار بوده‌اند.

مراحل اجرای Deconvolution:

ابتدا یک لایه از مدل CNN انتخاب می‌شود که می‌خواهیم فعال سازهای آن را بررسی کنیم. (انتخاب لایه)

در مرحله بعدی یک تصویر ورودی از شبکه عبور داده می‌شود تا فعالسازی‌های لایه مورد نظر بدست آید. (Forward Pass)

در مرحله بعدی معکوس کردن فعالسازی‌هاست. با استفاده از روش‌های مثل *Transposed Conv.* فعالسازی به سمت ورودی شبکه معکوس می‌شوند، به این معنی که سعی می‌کنیم تصویری را بسازیم که به این فعالسازها منجر شده است. بازسازی می‌شود.

مرحله نهایی مرحله تحلیل و بررسی تصاویر بازسازی شده است تا متوجه شویم کدام ویژگی‌های تصویری در تصمیم‌گیری‌های شبکه تأثیر داشته‌اند.

این مدل یک روش اصلاح شده از شکل معکوس  $Back\ propagation$  است که تنها دانش های مثبت از فعال سازی را در نظر بگیرد. این روش برابر شناسایی دقیق تر قسمت های تصویر که بیشترین تأثیر را بر تصمیمات شبکه دارند، استفاده می شود.

## Grad-CAM و CAM (5)

$CAM$  (Class Activation Maps) به معنای نقشه های فعال سازی کلاس است که امکان مشاهده این را می دهد که کدام نواحی از تصویر به طور خاص در تصمیم گیری های یک شبکه عصبی (CNN) دخالت دارند. این تکنیک برابر مدل هایی که شامل لایه های استخراج و سپس لایه های گلوبال - اوریج پولینگ (GAP) می باشند، به خوبی کار می کند.  $CAM$  می تواند نقشه های حرارتی (Heat Map) تولید کند که نشان دهنده میزان تأثیر هر بخش از تصویر بر تصمیم گیری های کلاس خاص است.

$Grad-CAM$  یا  $Gradient-weighted - CAM$  یک نسخه پیشرفته تر و محسوس تر از  $CAM$  است که با استفاده از گرادیان برای محاسبه نسبت به نقشه های ویژگی هر لایه استخراج، امکان تولید نقشه های حرارتی را می دهد که نشان می دهد کدام قسمت های تصویر بیشترین تأثیر را برای تشخیص یک کلاس خاص دارند.

## Integrated Gradients (6)

این روش یکی از روش های محبوب برای تفسیر مدل های یادگیری عمیق است که بر پایه تشخیص اهمیت هر ویژگی ورودی (مثل پیکسل های تصویر) در تصمیم گیری نهایی مدل استفاده می شود. این تکنیک با محاسبه گرادیان خروجی نسبت به ورودی و ادغام این گرادیان بر روی یک مسیر از یک ورودی پایه (مثلاً تصویری که همه پیکسل های آن صفر هستند) به تصویر واقعی انجام می شود.

## نقشه های حساسیت (Saliency Maps) (7)

نقشه های حساسیت تکنیک هایی هستند که برای بررسی سازی و تحلیل شبکه های عصبی عمیق بکار می روند. این نقشه ها نشان می دهند که کدام نواحی از یک تصویر ورودی برای تصمیم گیری یک شبکه عصبی مهم تر هستند و عبارت دیگر یک شبکه هنگام طبقه بندی یک تصویر به چه قسمت هایی از تصویر توجه بیشتری دارد.

این روش یک تکنیک تفسیرپذیری مدل‌های یادگیری ماشین است که به طرز مستقل از نوع مدل، توضیحات محلی (Local) برای تصمیمات مدل فراهم می‌کند. این روش با ایجاد نمونه‌های محلی از داده‌ها و آموزش یک مدل ساده روی این نمونه‌ها، ویژگی‌هایی که بیشترین تأثیر را بر تصمیم مدل اصلی دارند، شناسایی می‌کند.