

Interpretabilidad y Métodos de Entrenamiento de Modelos de Deep Learning

Facundo Quiroga^{1,2}, Franco Ronchetti^{1,2}, Oscar Stanchi^{1,6}, Gastón Ríos^{1,3}, Pedro Dal Bianco^{1,3}, Santiago Ponte Ahon¹, Juan Seery¹, Tatiana Badaracco⁵, Federico Rabinovich⁵, Saif Khalid⁴, Hatem Rashwan⁴, Domenec Puig Valls⁴, Laura Lanzarini¹, Waldo Hasperué¹

¹ Instituto de Investigación en Informática LIDI, Facultad de Informática, Universidad Nacional de La Plata, La Plata, Argentina.*

² Comisión de Investigaciones Científicas de la Pcia. De Bs. As. (CICPBA)

³ Becario postgrado UNLP

⁴ Universitat Rovira i Virgili, Tarragona, España

⁵ Universidad de Buenos Aires, Buenos Aires, Argentina

⁶ Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET)

* Centro asociado de la Comisión de Investigaciones Científicas de la Pcia. De Bs. As. (CIC)

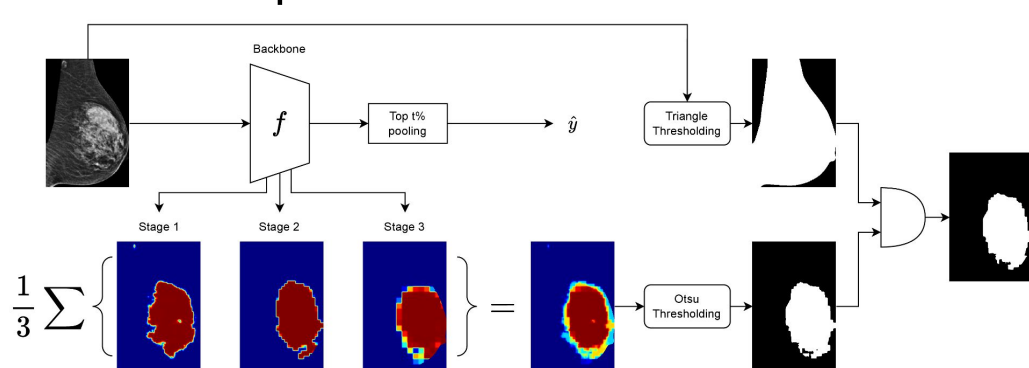
{fquiroga}@lidi.info.unlp.edu.ar

Contexto

Esta presentación corresponde a algunas de las tareas de investigación que se llevan a cabo en el III-LIDI en el marco del proyecto "Inteligencia de Datos. Técnicas y Modelos de Machine Learning" perteneciente al Programa de Incentivos (2023-2026).

Líneas de Investigación y Desarrollo

- Redes neuronales profundas, convolucionales y transformers.
- Herramientas de interpretabilidad para visión por computadora
- Interpretabilidad y entrenamiento no supervisado de Redes Convolucionales y ViT.
- Invarianza en redes neuronales, en particular de modelos ampliamente difundidos.



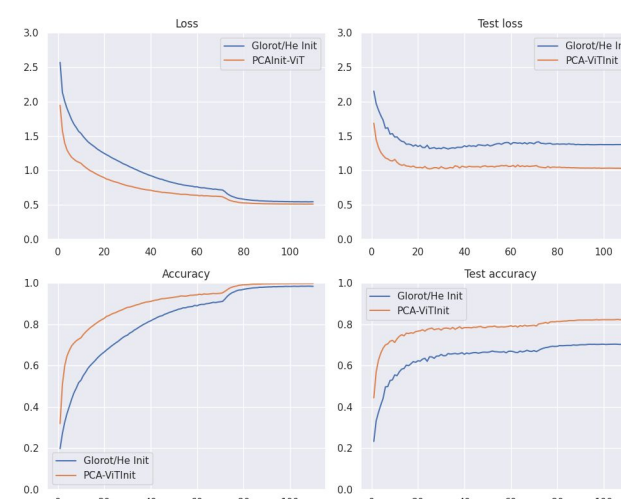
Método propuesto de generación de segmentaciones con etiquetas débiles. Se entrena un clasificador con etiquetas generales de clase (A vs D de la clasificación BI-RADS). Luego se generan máscaras mediante GradCAM, que se utilizan como etiquetas débiles

Formación de Recursos Humanos

- 3 Profesores y 1 JTP DE
- 2 Investigadores CIC-PBA
- 4 becarios de posgrado de la UNLP
- 3 Tesis de Doctorado y 2 de Maestría

Resultados Esperados y Obtenidos

- Interpretabilidad mediante invarianza de modelos EfficientNet y ViT
- Entrenamiento no supervisado de modelos EfficientNet y ViT para mejorar desempeño e interpretabilidad.
- Validación mediante métodos de interpretabilidad de un modelo de predicción de calidad de imágenes de retinopatía diabética basado en UNet.
- Validación de un método de entrenamiento semi supervisado para modelos de segmentación de mamografías basado en técnicas de interpretabilidad.



Los filtros de un modelo ViT inicializados mediante el método propuesto mejoran la velocidad de convergencia en algunos escenarios

Transformación Aplicada en el Entrenamiento

Brillo	0.35	0.55	0.36	0.63	0.71	0.75	0.89	0.86	0.82	0.62	0.65
Contraste	0.67	0.54	0.34	0.57	0.84	0.72	0.88	0.8	0.82	0.58	0.68
Escala grises	0.4	0.55	0.39	0.66	0.72	0.73	0.91	0.86	0.81	0.61	0.66
Inversión colores	0.43	0.58	0.37	0.64	0.76	0.74	0.86	0.85	0.8	0.61	0.66
Posterización	0.41	0.55	0.37	0.62	0.6	0.72	0.89	0.87	0.82	0.63	0.65
Solarización	0.44	0.56	0.36	0.63	0.81	0.62	0.82	0.78	0.79	0.6	0.64
Escala	0.43	0.56	0.37	0.64	0.75	0.74	0.88	0.86	0.8	0.61	0.67
Proyección	0.44	0.56	0.37	0.62	0.79	0.77	0.96	0.84	0.84	0.62	0.68
Rotación	0.41	0.57	0.36	0.62	0.77	0.75	0.97	0.9	0.81	0.64	0.68
Traslación	0.44	0.55	0.37	0.64	0.76	0.73	0.94	0.86	0.82	0.61	0.67
Identidad	0.43	0.57	0.36	0.63	0.81	0.76	0.95	0.89	0.79	0.61	0.68
Promedio											

Brillo
Contraste
Escala grises
Inversión colores
Posterización
Solarización
Escala
Proyección
Rotación
Traslación
Identidad
Promedio

Invarianza promedio de un modelo FasterViT variando la transformación aplicada en el entrenamiento (filas) y en la evaluación (columnas)