# CSE4261: Neural Network and Deep Learning

Lecture: 24.06.2025

Sangeeta Biswas, Ph.D.
Associate Professor,
University of Rajshahi, Rajshahi-6205, Bangladesh

# Gradient Descent Algorithm for Parameter Updating

Gradient is calculated for cost function with respect to weights and biases for updating their values.

$$w^{[l]} = w^{[l]} - \alpha \frac{\partial C}{\partial w^{[l]}}$$

$$b^{[l]} = b^{[l]} - \alpha \frac{\partial C}{\partial b^{[l]}}$$

# Adversarial Attack by Fast Gradient Signed Method (FGSM)

- FGSM uses the gradients of the loss with respect to the input image to create a new image that maximises the loss.
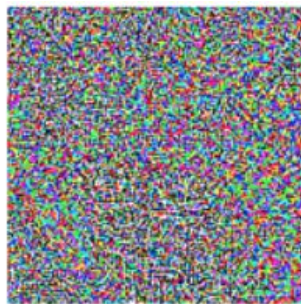- Adversarial image, $adv\_x = x + \epsilon * \text{sign}(\nabla_x J(\theta, x, y))$



$+.007 \times$      $=$

| $x$ | $\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$ | $\boldsymbol{x} +$ $\epsilon \text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$ |
|---|---|---|
| "panda" | "nematode" | "gibbon" |
| 57.7% confidence | 8.2% confidence | 99.3 % confidence |

# Grad-CAM (Gradient weighted Class Activation Mapping)



Input Image

Activation Maps

$A^k$

Any Network Layer Types Depending on Network Task

FC Layer

Softmax

C    Tiger Cat

$y^c$

(Prior to Softmax)

CNN

$\alpha_k^c A^k$

$\Sigma$

ReLU

Grad-CAM Heatmap

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$$

Global Average Pool Gradient Maps

Gradient Maps

Compute Gradients
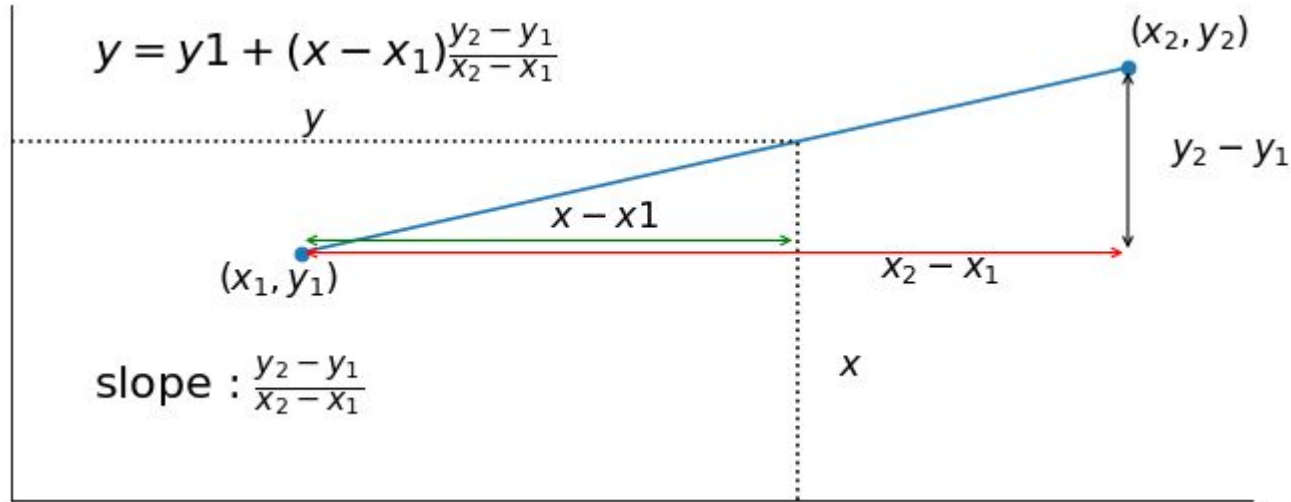
$\frac{\partial y^c}{\partial A^k}$

# Interpolation

- It is a method to estimate the value of a function at a point within the range of a known set of data points.
- Given a set of data points (x1, y1), ..., $(x_n, y_n)$, interpolation finds the value of y for a given x that lies within the range of $x_1$ and $x_n$
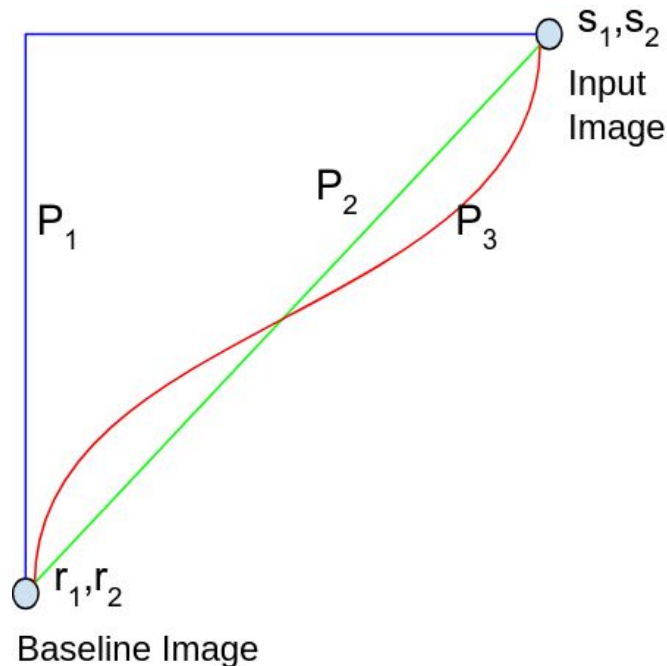
# Linear Interpolation

- It finds an intermediate value between two known values by assuming a linear relationship between them.

$$y = y1 + (x - x_1)\frac{y_2 - y_1}{x_2 - x_1}$$

$y$

$(x_2, y_2)$

$y_2 - y_1$

$x - x1$

$(x_1, y_1)$

$x_2 - x_1$

$x$

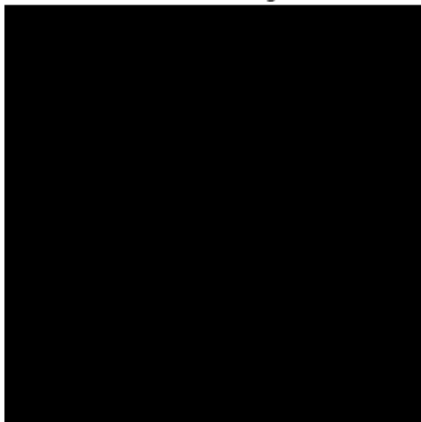slope : $\frac{y_2 - y_1}{x_2 - x_1}$

# Integrated Gradients (IG) Method as an XAI Technique

- In IG, a baseline image is used with the input image.
- Linear interpolation is used to approximate the integral of gradients along a path between a baseline and the input instance (i.e., path P2).
- The integral is approximated by summing the gradients at multiple interpolated points along this line.

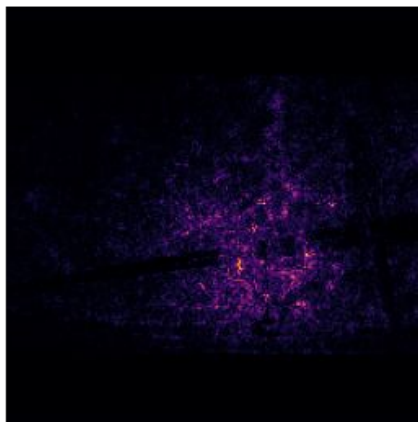# Integrated Gradients (IG) Method as an XAI Technique



Baseline Image | Original Image | IG Attribution Mask | Original + IG Attribution Mask Overlay
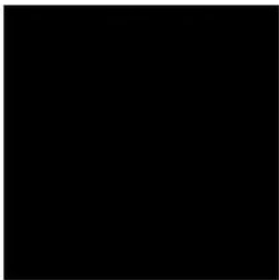
- Tutorial:
  https://www.tensorflow.org/tutorials/interpretability/integrated_gradients

# Example of Linear Interpolation

- Five-step interpolation between the baseline x' and the input image x

| alpha = 0.0 | alpha = 0.2 | alpha = 0.4 | alpha = 0.6 | alpha = 0.8 | alpha = 1.0 |



Baseline

Input image

# Formula of IG

$$IntegratedGrads_i^{approx}(x) ::= (x_i - x_i') \times \sum_{k=1}^{m} \frac{\partial F(x' + \frac{k}{m} \times (x - x'))}{\partial x_i} \times \frac{1}{m}$$

i = feature (individual pixel)

x = input (image tensor)

x' = baseline (image tensor)

k = scaled feature perturbation constant

m = number of steps in the Riemann sum approximation of the integral

$(x_i-x'_i)$= a term for the difference from the baseline

F = model prediction function

# Formula of IG

$$IntegratedGrads_i^{approx}(x) ::= (x_i - x_i') \times \sum_{k=1}^{m} \overbrace{\frac{\partial F(\text{interpolated images})}{\partial x_i}}^{\text{compute gradients}} \times \frac{1}{m}$$

$\frac{\partial F}{\partial x_i}$ = gradient

- The gradient tells us which pixels have the strongest effect on the model's predicted class probabilities.

# Gradient Computation for IG by Tensorflow

```
def compute_gradients(images, target_class_idx):

    with tf.GradientTape() as tape:

        tape.watch(images)

        logits = model(images)

        probs = tf.nn.softmax(logits, axis = -1)[:, target_class_idx]

    return tape.gradient(probs, images)
```

In Grad-CAM:
model.layers[-1].activation = None # Remove last layer's softmax

# Alternative Baseline



(a) Gaussian Baseline    (b) Blur Baseline    (c) Max Distance Baseline    (d) Uniform Baseline