

CS 6.5: Theories of Deep Learning

Problem Sheet 4

Prof. Jared Tanner

November 16, 2020

Adversarial attacks for neural networks

Adversarial examples are intentionally designed optical illusions, where such inputs to learned models cause the model to make a mistake. Mathematically, given a point $\mathbf{x} \in \Omega$ drawn from class y , a scalar $\epsilon > 0$, and a metric d , we say that \mathbf{x} admits an adversarial example in the metric d if there exists a point $\mathbf{x}^* \in \Omega$ with $Class(\mathbf{x}^*) \neq y$, and $d(\mathbf{x}, \mathbf{x}^*) \leq \epsilon$. In practice d is chosen as ℓ^p -norms with ℓ^∞ being the most popular choice, which limits the absolute change that can be made to any one dimension of \mathbf{x} .

1. Task1: Write a short report summarizing the fast gradient sign method (FGSM) for adversarial attacks¹. Your report should be written in the format and style of a NIPS Proceedings, abridged to not exceed 2 pages. Latex style files and an exemplar template are provided on the course page, and are similar to last exercise.
2. Task2: One Layer Net: Consider the neural net defined as $\hat{y} = SM(\mathbf{W}\mathbf{x})$ trained with the cross-entropy loss $L(\mathbf{x}, y)$, where SM denotes softmax activation. Let \mathbf{x}^* be the adversarial image of \mathbf{x} resulting from FGSM attack with constant ϵ . Prove that $\forall \epsilon > 0$ we have $L(\mathbf{x}^*, y) \geq L(\mathbf{x}, y)$
3. Task3: Two Layer Net: Consider the neural net defined as $\hat{y} = SM(\mathbf{V}\sigma\mathbf{W}\mathbf{x})$ trained with the cross-entropy loss $L(\mathbf{x}, y)$, where \mathbf{V}, \mathbf{W} are weights, SM denotes softmax activation and σ is ReLU activation. Suppose every element of $\mathbf{W}\mathbf{x}$ is non-zero, if $\epsilon < \frac{|\mathbf{W}\mathbf{x}|_{min}}{\|\mathbf{W}\|_\infty}$, then prove that $L(\mathbf{x}^*, y) \geq L(\mathbf{x}, y)$, given the fact that for $j = 1, 2, \dots; sign(\mathbf{W}\mathbf{x})_j = sign(\mathbf{W}\mathbf{x}^*)_j$

¹<https://arxiv.org/pdf/1412.6572.pdf>

Solution:

1. The loss for a example \mathbf{x} and true class s can be expressed as:

$$\begin{aligned} L(\mathbf{x}, y) &= \text{crossentropy}(\text{softmax}(\mathbf{W}\mathbf{x}), y) \\ &= -\ln(\text{softmax}(\mathbf{W}\mathbf{x})_s) \\ &= -\ln\left[\frac{\exp(\mathbf{W}\mathbf{x})_s}{\exp(\mathbf{W}\mathbf{x})_1 + \exp(\mathbf{W}\mathbf{x})_2 + \dots + \exp(\mathbf{W}\mathbf{x})_k}\right] \end{aligned} \quad (1)$$

now each element of vector \mathbf{x}^* is expressed as:

$$\begin{aligned} x_i^* &= x_i + \epsilon \text{sign}\left(\frac{\partial L(\mathbf{x}, y)}{\partial x_i}\right) \\ &= x_i + \epsilon a_i \\ &= x_i + \epsilon \text{sign}\left(\sum_j^k \exp(\mathbf{W}\mathbf{x})_j w_{ji} - \left(\sum_j^k \exp(\mathbf{W}\mathbf{x})_j\right) w_{si}\right) \end{aligned} \quad (2)$$

Assuming the hypothesis is true we have to prove:

$$\begin{aligned} \frac{\exp(\mathbf{W}\mathbf{x})_s}{\sum_j^k \exp(\mathbf{W}\mathbf{x})_j} &\geq \frac{\exp(\mathbf{W}\mathbf{x}^*)_s}{\sum_j^k \exp(\mathbf{W}\mathbf{x}^*)_j} \\ \implies \frac{\exp(\mathbf{W}\mathbf{x}^*)_s}{\exp(\mathbf{W}\mathbf{x})_s} &\leq \sum_j^k \text{softmax}(\mathbf{W}\mathbf{x})_j \frac{\exp(\mathbf{W}\mathbf{x}^*)_j}{\exp(\mathbf{W}\mathbf{x})_j} \\ \implies \exp(\epsilon \mathbf{W}\mathbf{a})_s &\leq \sum_j^k \text{softmax}(\mathbf{W}\mathbf{x})_j \exp(\epsilon \mathbf{W}\mathbf{a})_j \end{aligned} \quad (3)$$

where $\mathbf{a} = [a_1 a_2 \dots]^T$. By property of softmax and Jensen's inequality the RHS can be lower bounded by:

$$RHS \geq \exp\left(\sum_j^k \epsilon \text{softmax}(\mathbf{W}\mathbf{x})_j (\mathbf{W}\mathbf{a})_j\right) \quad (4)$$

and hence we just need to prove

$$\begin{aligned} \sum_j^k \text{softmax}(\mathbf{W}\mathbf{x})_j (\mathbf{W}\mathbf{a})_j &\geq (\mathbf{W}\mathbf{a})_s \\ \implies \sum_j^k \exp(\mathbf{W}\mathbf{x})_j (\mathbf{W}\mathbf{a})_j - \left(\exp(\mathbf{W}\mathbf{x})_1 + \exp(\mathbf{W}\mathbf{x})_2 + \dots\right) (\mathbf{W}\mathbf{a})_s \\ &\geq 0 \end{aligned} \quad (5)$$

where the result follows from (2) and fact that $\mathbf{x} \text{sign}(\mathbf{x}) > 0$

2. Let $\mathbf{T} = \mathbf{V}\sigma\mathbf{W}$ i.e., $y = \mathbf{T}\mathbf{x}$ and define the following index set (and using property given in the problem):

$$A = \{i : \mathbf{W}\mathbf{x}_i > 0\} = \{i : \mathbf{W}\mathbf{x}_i^* > 0\}. \quad (6)$$

Here $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{V} \in \mathbb{R}^{k \times l}$ and $\mathbf{W} \in \mathbb{R}^{l \times n}$. Then we can express the operator \mathbf{T} as a linear operator:

$$\begin{aligned} (\mathbf{T}\mathbf{x})_j &= \sum_t^l v_{jt} \sigma(w_{t1}x_1 + w_{t2}x_2 + \dots + w_{tn}x_n) \\ &= \sum_{t \in A} v_{jt} (w_{t1}x_1 + w_{t2}x_2 + \dots + w_{tn}x_n) \end{aligned} \quad (7)$$

The loss for a example \mathbf{x} and true class s can be expressed as:

$$\begin{aligned} L(\mathbf{x}, y) &= \text{crossentropy}(\text{softmax}(\mathbf{T}\mathbf{x}), y) \\ &= -\ln(\text{softmax}(\mathbf{T}\mathbf{x})_s) \\ &= -\ln \left[\frac{\exp(\mathbf{T}\mathbf{x})_s}{\exp(\mathbf{T}\mathbf{x})_1 + \exp(\mathbf{T}\mathbf{x})_2 + \dots + \exp(\mathbf{T}\mathbf{x})_k} \right], \end{aligned} \quad (8)$$

which reduces to problem 1 with one linear layer.