

[Deep Batch Active Learning by Diverse, Uncertain Gradient Lower Bounds]

Diversity ⊕ Uncertainty Approach

- Measure of Uncertainty: the gradient magnitude with respect to parameters in the final output layer, which is computed using the most likely label according to the model.
- Capturing Diversity: Collect a batch of examples where these gradients span a diverse set of directions.
 - ** K-means ⁺⁺ algorithm
 - ① captures magnitude of a candidate gradient
 - ② captures its distance from previously included points in the batch.

[Notation & Setting]

K : # of classes

$$[K] = \{1, 2, \dots, K\}$$

X : the instance space

y : the label space

D_x : the unlabeled data distribution

D_{yx} : the conditional distribution over labels given examples.

\mathbb{E}_D : Expectation under the data distribution D

h : a classifier that maps X to y (architectures are fixed in any given context)

$$l_{01}(h(x), y) = I(h(x) \neq y) \begin{cases} 1 & \text{if } h(x) \neq y \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{E}_D [l_{01}(h(x), y)] = \Pr_{(x,y) \sim D}(h(x) \neq y)$$

S : a set of labeled examples

\mathbb{E}_S : sample averages over S

$f(x; \theta) \in \mathbb{R}^K$: probability vector of scores assigned to candidate labels,
given examples X and parameters θ

$$h_\theta(x) = \operatorname{argmax}_{y \in [K]} f(x; \theta)_y$$

$$l_{CE}(p, y) = \sum_{i=1}^K I(y=i) \ln \frac{1}{p_i} = \ln \left(\frac{1}{p_y} \right)$$

$$\text{Loss} = \mathbb{E}_S [l_{CE}(f(x; \theta), y)]$$

[Algorithm]

Two main computations

- ① Gradient embedding computation
 - ② Sampling computation
- } \Rightarrow Computes
- $\hat{y}(x)$: label preferred by the current model
 - g_x : gradient of the loss on $(x, \hat{y}(x))$ with respect to the parameters of the last layer of the network.

<Gradient embedding computation>

- ✓ "Uncertain" if knowing the label induces a large gradient of the loss with respect to the model parameters and hence a large update to the model.
- ✓ Compute gradient as if the model's prediction on the example is true label.
- ✓ The gradient norm with respect to the last layer using the label

||

a lower bound on the gradient norm induced by any other label.

→ uncertainty

→ potential update direction upon receiving a label at an example.

<Sampling Step>

Want to secure both sample magnitude & batch diversity

- ✓ Sample using K-MEANS++ seeding algorithm

↳ Selects centroids by iteratively sampling points in proportion to their squared distances from the nearest centroid that has already been chosen.

→ Select diverse batch of high-magnitude samples

[Example: multi-class classification with softmax activations]

f : a neural network parameterized by θ w/ last non-linearity is softmax

$$\theta = (W, V)$$

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

W : $\theta_{\text{out}} = (W_1, \dots, W_k)^T \in \mathbb{R}^{K \times d}$ ← weights of the last layer

V : weights of all previous layers except the last layer.

$$f(x; \theta) = \sigma(W \cdot z(x; V))$$

z : non-linear function

$$p_i = f(x; \theta)_i$$

$$\begin{aligned} l_{CE}(f(x; \theta), y) &= l_{CE}(p_i, y) \\ &= \sum_{i=1}^K I(y=i) \log \frac{1}{p_i} \\ &= \sum_{i=1}^K I(y=i) \log \frac{1}{f(x; \theta)_i} \\ &= \sum_{i=1}^K I(y=i) \log \frac{1}{\sigma(W \cdot z(x; V))_i} \\ &= - \sum_{i=1}^K I(y=i) \log \sigma(W \cdot z(x; V))_i \\ &= - \sum_{i=1}^K I(y=i) \log \frac{e^{W \cdot z(x; V)}_i}{\sum_{j=1}^K e^{W \cdot z}_j} \\ &= - \sum_{i=1}^K I(y=i) \left\{ \log e^{W \cdot z(x; V)}_i - \log \sum_{j=1}^K e^{W \cdot z}_j \right\} \\ &= - \sum_{i=1}^K I(y=i) \left\{ W \cdot z(x; V)_i - \log \sum_{j=1}^K e^{W \cdot z}_j \right\} \\ &= - \sum_{i=1}^K I(y=i) \left\{ W \cdot z(x; V)_i - \frac{1}{K} \sum_{j=1}^K W \cdot z(x; V)_j \right\} \\ &= - W \cdot y \cdot z(x; V) + \log \sum_{j=1}^K e^{W \cdot y \cdot z(x; V)}_j \end{aligned}$$

$$g_x^y = \frac{\partial}{\partial W} l_{CE}(f(x; \theta), y) \text{ for a label } y$$

$$g_x = g_x^{\hat{y}}$$

$$\hat{y} = \operatorname{argmax}_{i \in [K]} p_i$$

gradient corresponding to label i :

$$(g_x)_i = \frac{\partial}{\partial W_i} l_{CE}(f(x; \theta), \hat{y}) = (p_i - I(\hat{y}=i)) z(x; v)$$

following observations:

- ✓ g_x is a scaling of $z(x; v)$
- ✓ g_x estimates the example's influence on the current model
- ✓ High confidence tends to have gradient embeddings of small magnitude.