

Supplemental Material: HashNet: Deep Learning to Hash by Continuation

Zhangjie Cao[†], Mingsheng Long[†], Jianmin Wang[†], and Philip S. Yu[‡]

[†]KLiss, MOE; NEL-BDS; TNList; School of Software, Tsinghua University, China

[‡]University of Illinois at Chicago, IL, USA

caozhangjie14@gmail.com {mingsheng, jimwang}@tsinghua.edu.cn psyu@uic.edu

1. Convergence Analysis

We briefly analyze that the continuation optimization in Algorithm 1 will decrease the loss of HashNet (4) in each stage and in each iteration until converging to HashNet with sign activation function that generates *exactly* binary codes.

Let $L_{ij} = w_{ij} (\log(1 + \exp(\alpha \langle \mathbf{h}_i, \mathbf{h}_j \rangle)) - \alpha s_{ij} \langle \mathbf{h}_i, \mathbf{h}_j \rangle)$ and $L = \sum_{s_{ij} \in \mathcal{S}} L_{ij}$, where $\mathbf{h}_i \in \{-1, +1\}^K$ are *binary* hash codes. Note that when optimizing HashNet by continuation in Algorithm 1, network activation in each stage t is $g = \tanh(\beta_t z)$, which is *continuous* in nature and will only become *binary* when convergence $\beta_t \rightarrow \infty$. Denote by $J_{ij} = w_{ij} (\log(1 + \exp(\alpha \langle \mathbf{g}_i, \mathbf{g}_j \rangle)) - \alpha s_{ij} \langle \mathbf{g}_i, \mathbf{g}_j \rangle)$ and $J = \sum_{s_{ij} \in \mathcal{S}} J_{ij}$ the true loss we optimize in Algorithm 1, where $\mathbf{g}_i \in \mathbb{R}^K$ and note that $\mathbf{h}_i = \text{sgn}(\mathbf{g}_i)$. We will show that HashNet loss $L(\mathbf{h})$ descends when minimizing $J(\mathbf{g})$.

Theorem 1. *The HashNet loss L will not change across stages t and $t+1$ with bandwidths switched from β_t to β_{t+1} .*

Proof. When the algorithm switches from stages t to $t+1$ with bandwidths changed from β_t to β_{t+1} , only the network activation is changed from $\tanh(\beta_t z)$ to $\tanh(\beta_{t+1} z)$ but its sign $h = \text{sgn}(\tanh(\beta_t z)) = \text{sgn}(\tanh(\beta_{t+1} z))$, i.e. the hash code, remains the same. Thus L is unchanged. \square

For each pair of binary codes $\mathbf{h}_i, \mathbf{h}_j$ and their continuous counterparts $\mathbf{g}_i, \mathbf{g}_j$, the derivative of J w.r.t. each bit k is

$$\frac{\partial J}{\partial g_{ik}} = w_{ij} \alpha \left(\frac{1}{1 + \exp(-\alpha \langle \mathbf{g}_i, \mathbf{g}_j \rangle)} - s_{ij} \right) g_{jk}, \quad (1)$$

where $k = 1, \dots, K$. The derivative of J w.r.t. \mathbf{g}_j can be defined similarly. Updating \mathbf{g}_i by SGD, the updated \mathbf{g}'_i is

$$\begin{aligned} g'_{ik} &= g_{ik} - \eta \frac{\partial J}{\partial g_{ik}} \\ &= g_{ik} - \eta w_{ij} \alpha \left(\frac{1}{1 + \exp(-\alpha \langle \mathbf{g}_i, \mathbf{g}_j \rangle)} - s_{ij} \right) g_{jk}, \end{aligned} \quad (2)$$

where η is the learning rate and \mathbf{g}'_j is computed similarly.

Lemma 1. *Denote by $\mathbf{h}_i = \text{sgn}(\mathbf{g}_i)$, $\mathbf{h}'_i = \text{sgn}(\mathbf{g}'_i)$, then*

$$\begin{cases} \langle \mathbf{h}'_i, \mathbf{h}'_j \rangle \geq \langle \mathbf{h}_i, \mathbf{h}_j \rangle, & s_{ij} = 1, \\ \langle \mathbf{h}'_i, \mathbf{h}'_j \rangle \leq \langle \mathbf{h}_i, \mathbf{h}_j \rangle, & s_{ij} = 0. \end{cases} \quad (3)$$

Proof. Since $\langle \mathbf{h}_i, \mathbf{h}_j \rangle = \sum_{k=1}^K h_{ik} h_{jk}$, Lemma 1 can be proved by verifying that $h'_{ik} h'_{jk} \geq h_{ik} h_{jk}$ if $s_{ij} = 1$ and $h'_{ik} h'_{jk} \leq h_{ik} h_{jk}$ if $s_{ij} = 0, \forall k = 1, 2, \dots, K$.

Case 1. $s_{ij} = 0$.

(1) If $g_{ik} < 0, g_{jk} > 0$, then $\frac{\partial J}{\partial g_{ik}} > 0, \frac{\partial J}{\partial g_{jk}} < 0$. Thus, $h'_{ik} \leq h_{ik} = -1, h'_{jk} \geq h_{jk} = 1$. And we have $h'_{ik} h'_{jk} = -1 = h_{ik} h_{jk}$.

(2) If $g_{ik} > 0, g_{jk} < 0$, then $\frac{\partial J}{\partial g_{ik}} < 0, \frac{\partial J}{\partial g_{jk}} > 0$. Thus, $h'_{ik} \geq h_{ik} = 1, h'_{jk} \leq h_{jk} = -1$. And we have $h'_{ik} h'_{jk} = -1 = h_{ik} h_{jk}$.

(3) If $g_{ik} < 0, g_{jk} < 0$, then $\frac{\partial J}{\partial g_{ik}} < 0, \frac{\partial J}{\partial g_{jk}} < 0$. Thus $h'_{ik} \geq h_{ik} = -1, h'_{jk} \geq h_{jk} = -1$. So h'_{ik} and h'_{jk} may be either +1 or -1 and we have $h'_{ik} h'_{jk} \leq 1 = h_{ik} h_{jk}$.

(4) If $g_{ik} > 0, g_{jk} > 0$, then $\frac{\partial J}{\partial g_{ik}} > 0, \frac{\partial J}{\partial g_{jk}} > 0$. Thus $h'_{ik} \leq h_{ik} = 1, h'_{jk} \leq h_{jk} = 1$. So h'_{ik} and h'_{jk} may be either +1 or -1 and we have $h'_{ik} h'_{jk} \leq 1 = h_{ik} h_{jk}$.

Case 2. $s_{ij} = 1$. It can be proved similarly as Case 1. \square

Theorem 2. *Loss L decreases when optimizing loss $J(\mathbf{g})$ by the stochastic gradient descent (SGD) within each stage.*

Proof. The gradient of loss L w.r.t. hash codes $\langle \mathbf{h}_i, \mathbf{h}_j \rangle$ is

$$\frac{\partial L}{\partial \langle \mathbf{h}_i, \mathbf{h}_j \rangle} = w_{ij} \alpha \left(\frac{1}{1 + \exp(-\alpha \langle \mathbf{h}_i, \mathbf{h}_j \rangle)} - s_{ij} \right). \quad (4)$$

We observe that

$$\begin{cases} \frac{\partial L}{\partial \langle \mathbf{h}_i, \mathbf{h}_j \rangle} \leq 0, & s_{ij} = 1, \\ \frac{\partial L}{\partial \langle \mathbf{h}_i, \mathbf{h}_j \rangle} \geq 0, & s_{ij} = 0. \end{cases} \quad (5)$$

By substituting Lemma 1: if $s_{ij} = 1$, then $\langle \mathbf{h}'_i, \mathbf{h}'_j \rangle \geq \langle \mathbf{h}_i, \mathbf{h}_j \rangle$, and thus $L(\mathbf{h}'_i, \mathbf{h}'_j) \leq L(\mathbf{h}_i, \mathbf{h}_j)$; if $s_{ij} = 0$, then $\langle \mathbf{h}'_i, \mathbf{h}'_j \rangle \leq \langle \mathbf{h}_i, \mathbf{h}_j \rangle$, and thus $L(\mathbf{h}'_i, \mathbf{h}'_j) \leq L(\mathbf{h}_i, \mathbf{h}_j)$. \square