| | | doc_1 | | doc_2 | decision | id |
|---|---|---|---|---|---|---|
| cases | authors | <ul><li>Xiaotong Li</li><li>Yixiao Ge</li><li>Kun Yi</li><li>Zixuan Hu</li><li>Ying Shan</li><li>Ling-Yu Duan</li></ul> | authors | <ul><li>Duan, Ling-Yu</li><li>Ge, Yixiao</li><li>Hu, Zixuan</li><li>Li, Xiaotong</li><li>Shan, Ying</li><li>Yi, Kun</li></ul> | DUPLICATES | 130 |
| | title | mc-BEiT: Multi-choice Discretization for Image BERT Pre-training | title | mc-BEiT: Multi-choice Discretization for Image BERT Pre-training | | |
| | publication_date | 2022-07-28 00:00:00 | publication_date | 2022-04-09 00:00:00 | | |
| | source | SupportedSources.INTERNET_ARCHIVE | source | SupportedSources.CORE | | |
| | journal | | journal | | | |
| | volume | | volume | | | |
| | doi | | doi | None | | |
| | urls | <ul><li>https://web.archive.org/web/20220729034617/https://arxiv.org/pdf/2203.15371v4.pdf</li></ul> | urls | <ul><li>http://arxiv.org/abs/2203.15371</li></ul> | | |
| | id | id9120167953335425519 | id | id-1434096295655901912 | | |
| | abstract | Image BERT pre-training with masked image modeling (MIM) becomes a popular practice to cope with self-supervised representation learning. A seminal work, BEiT, casts MIM as a classification task with a visual vocabulary, tokenizing the continuous visual signals into discrete vision tokens using a pre-learned dVAE. Despite a feasible solution, the improper discretization hinders further improvements of image pre-training. Since image discretization has no ground-truth answers, we believe that the masked patch should not be assigned with a unique token id even if a better tokenizer can be obtained. In this work, we introduce an improved BERT-style image pre-training method, namely mc-BEiT, which performs MIM proxy tasks towards eased and refined multi-choice training objectives. Specifically, the multi-choice supervision for the masked image patches is formed by the soft probability vectors of the discrete token ids, which are predicted by the off-the-shelf image tokenizer and further refined by high-level inter-patch perceptions resorting to the observation that similar patches should share their choices. Extensive experiments on classification, segmentation, and detection tasks demonstrate the superiority of our method, e.g., the pre-trained ViT-B achieves 84.1% top-1 fine-tuning accuracy on ImageNet-1K classification, 49.2% AP^b and 44.0% AP^m of object detection and instance segmentation on COCO, 50.8% mIOU on ADE20K semantic segmentation, outperforming the competitive counterparts. The code will be available at https://github.com/lixiaotong97/mc-BEiT. | abstract | Image BERT pre-training with masked image modeling (MIM) becomes a popular practice to cope with self-supervised representation learning. A seminal work, BEiT, casts MIM as a classification task with a visual vocabulary, tokenizing the continuous visual signals into discrete vision tokens using a pre-learned dVAE. Despite a feasible solution, the improper discretization hinders further improvements of image pre-training. Since image discretization has no ground-truth answers, we believe that the masked patch should not be assigned with a unique token id even if a better tokenizer can be obtained. In this work, we introduce an improved BERT-style image pre-training method, namely mc-BEiT, which performs MIM proxy tasks towards eased and refined multi-choice training objectives. Specifically, the multi-choice supervision for the masked image patches is formed by the soft probability vectors of the discrete token ids, which are predicted by the off-the-shelf image tokenizer and further refined by high-level inter-patch perceptions resorting to the observation that similar patches should share their choices. Extensive experiments on classification, segmentation, and detection tasks demonstrate the superiority of our method, e.g., the pre-trained ViT-B achieves 84.1% top-1 fine-tuning accuracy on ImageNet-1K classification, 50.8% mIOU on ADE20K semantic segmentation, 51.2% AP^b and 44.3% AP^m of object detection and instance segmentation on COCO, outperforming the competitive counterparts | | |
| | versions | | versions | | | |