
Disentangled Latent Representations with CapsNet

Chuanquan (Charles) Shu
Massachusetts Institute of Technology
Cambridge, MA 02139
cggz21@mit.edu

Abstract

1 A distinguishing feature of Capsule Neural Network (CapsNet) is that the input
2 and output neurons of a Caps Layer are vectors. The orientation of a neuron
3 vector encodes features of a particular entity, but the semantic meanings of such
4 orientation are entangled and thus uninterpretable. Therefore, in this project, I
5 seek to use a training procedure that enables CapsNet to learn disentangled and
6 semantically interpretable latent representations. Demonstrating the feasibility
7 of this concept is the primary goal, hence MNIST dataset is used and only two
8 features (digit rotation and scale) are experimented.

9
10 In order to achieve this goal, a deliberate training procedure is implemented. The
11 first part creates supporting training batches by creating Rotation Batches (only
12 rotation transformation), Scale Batches (only scaling transformation) and Standard
13 Batches (no transformation). The second part ensures that only one group of
14 variable(s)/dimension(s) is allowed to learn at each step by imposing restrictions
15 on neuron activations and gradients flow during the forward and backward stages,
16 respectively.

17
18 Under the new method, almost all 10 digits unanimously choose the first and
19 second dimensions of their encodings to explain rotation and scaling, respectively.
20 Under the original method, the best explaining dimension for each feature is
21 rather arbitrary. Moreover, regressing the true rotation and scaling against the
22 best explaining dimension under the original methods has a mean R^2 of just
23 0.45, suggesting the lack of disentangled representations for rotation and scaling.
24 Under the new methods, the mean R^2 for rotation and scaling are 0.63 and
25 0.83, respectively. The significant improvements show that the first and second
26 dimensions have truly learned to explain the rotation and scaling of a given image.

27
28 It is worth mentioning that the 0.63 R^2 for rotation, although better than the old
29 method counterpart, is not impressive. This could be largely attributed to the
30 sub-optimal quality of the Standard Batches. The images in the Standard Batches
31 are supposed to have no rotation and scaling so irrelevant dimensions in the latent
32 representations will not be exposed to any variance in these two features during
33 training. Nevertheless, MNIST does have substantial innate variance in rotation,
34 leading to sub-optimal Standard Batches. This undermines the constraint on the first
35 dimension to explain rotation alone and hence resulting in a less powerful rotation
36 encoding. It also explains why the scale encoding achieves great performance
37 and the rotation encoding does not, suggesting that rotation encoding could be
38 as powerful if presented with better Standard Batches (ones with less variance in
39 rotation).