

the Monte Carlo approximation. It appears that the optimal choice of inference approximation is problem-dependent.

[Srivastava et al. \(2014\)](#) showed that dropout is more effective than other standard computationally inexpensive regularizers, such as weight decay, filter norm constraints and sparse activity regularization. Dropout may also be combined with other forms of regularization to yield a further improvement.

One advantage of dropout is that it is very computationally cheap. Using dropout during training requires only  $O(n)$  computation per example per update, to generate  $n$  random binary numbers and multiply them by the state. Depending on the implementation, it may also require  $O(n)$  memory to store these binary numbers until the back-propagation stage. Running inference in the trained model has the same cost per-example as if dropout were not used, though we must pay the cost of dividing the weights by 2 once before beginning to run inference on examples.

Another significant advantage of dropout is that it does not significantly limit the type of model or training procedure that can be used. It works well with nearly any model that uses a distributed representation and can be trained with stochastic gradient descent. This includes feedforward neural networks, probabilistic models such as restricted Boltzmann machines ([Srivastava et al., 2014](#)), and recurrent neural networks ([Bayer and Osendorfer, 2014](#); [Pascanu et al., 2014a](#)). Many other regularization strategies of comparable power impose more severe restrictions on the architecture of the model.

Though the cost per-step of applying dropout to a specific model is negligible, the cost of using dropout in a complete system can be significant. Because dropout is a regularization technique, it reduces the effective capacity of a model. To offset this effect, we must increase the size of the model. Typically the optimal validation set error is much lower when using dropout, but this comes at the cost of a much larger model and many more iterations of the training algorithm. For very large datasets, regularization confers little reduction in generalization error. In these cases, the computational cost of using dropout and larger models may outweigh the benefit of regularization.

When extremely few labeled training examples are available, dropout is less effective. Bayesian neural networks ([Neal, 1996](#)) outperform dropout on the Alternative Splicing Dataset ([Xiong et al., 2011](#)) where fewer than 5,000 examples are available ([Srivastava et al., 2014](#)). When additional unlabeled data is available, unsupervised feature learning can gain an advantage over dropout.

[Wager et al. \(2013\)](#) showed that, when applied to linear regression, dropout is equivalent to  $L^2$  weight decay, with a different weight decay coefficient for