

dropout as bagging an ensemble of models formed by including or excluding units. However, there is no need for this model averaging strategy to be based on inclusion and exclusion. In principle, any kind of random modification is admissible. In practice, we must choose modification families that neural networks are able to learn to resist. Ideally, we should also use model families that allow a fast approximate inference rule. We can think of any form of modification parametrized by a vector  $\boldsymbol{\mu}$  as training an ensemble consisting of  $p(y \mid \boldsymbol{x}, \boldsymbol{\mu})$  for all possible values of  $\boldsymbol{\mu}$ . There is no requirement that  $\boldsymbol{\mu}$  have a finite number of values. For example,  $\boldsymbol{\mu}$  can be real-valued. [Srivastava et al. \(2014\)](#) showed that multiplying the weights by  $\boldsymbol{\mu} \sim \mathcal{N}(\mathbf{1}, I)$  can outperform dropout based on binary masks. Because  $\mathbb{E}[\boldsymbol{\mu}] = \mathbf{1}$  the standard network automatically implements approximate inference in the ensemble, without needing any weight scaling.

So far we have described dropout purely as a means of performing efficient, approximate bagging. However, there is another view of dropout that goes further than this. Dropout trains not just a bagged ensemble of models, but an ensemble of models that share hidden units. This means each hidden unit must be able to perform well regardless of which other hidden units are in the model. Hidden units must be prepared to be swapped and interchanged between models. [Hinton et al. \(2012c\)](#) were inspired by an idea from biology: sexual reproduction, which involves swapping genes between two different organisms, creates evolutionary pressure for genes to become not just good, but to become readily swapped between different organisms. Such genes and such features are very robust to changes in their environment because they are not able to incorrectly adapt to unusual features of any one organism or model. Dropout thus regularizes each hidden unit to be not merely a good feature but a feature that is good in many contexts. [Warde-Farley et al. \(2014\)](#) compared dropout training to training of large ensembles and concluded that dropout offers additional improvements to generalization error beyond those obtained by ensembles of independent models.

It is important to understand that a large portion of the power of dropout arises from the fact that the masking noise is applied to the hidden units. This can be seen as a form of highly intelligent, adaptive destruction of the information content of the input rather than destruction of the raw values of the input. For example, if the model learns a hidden unit  $h_i$  that detects a face by finding the nose, then dropping  $h_i$  corresponds to erasing the information that there is a nose in the image. The model must learn another  $h_i$ , either that redundantly encodes the presence of a nose, or that detects the face by another feature, such as the mouth. Traditional noise injection techniques that add unstructured noise at the input are not able to randomly erase the information about a nose from an image of a face unless the magnitude of the noise is so great that nearly all of the information in