

To make a prediction, a bagged ensemble must accumulate votes from all of its members. We refer to this process as **inference** in this context. So far, our description of bagging and dropout has not required that the model be explicitly probabilistic. Now, we assume that the model's role is to output a probability distribution. In the case of bagging, each model  $i$  produces a probability distribution  $p^{(i)}(y \mid \mathbf{x})$ . The prediction of the ensemble is given by the arithmetic mean of all of these distributions,

$$\frac{1}{k} \sum_{i=1}^k p^{(i)}(y \mid \mathbf{x}). \quad (7.52)$$

In the case of dropout, each sub-model defined by mask vector  $\boldsymbol{\mu}$  defines a probability distribution  $p(y \mid \mathbf{x}, \boldsymbol{\mu})$ . The arithmetic mean over all masks is given by

$$\sum_{\boldsymbol{\mu}} p(\boldsymbol{\mu}) p(y \mid \mathbf{x}, \boldsymbol{\mu}) \quad (7.53)$$

where  $p(\boldsymbol{\mu})$  is the probability distribution that was used to sample  $\boldsymbol{\mu}$  at training time.

Because this sum includes an exponential number of terms, it is intractable to evaluate except in cases where the structure of the model permits some form of simplification. So far, deep neural nets are not known to permit any tractable simplification. Instead, we can approximate the inference with sampling, by averaging together the output from many masks. Even 10-20 masks are often sufficient to obtain good performance.

However, there is an even better approach, that allows us to obtain a good approximation to the predictions of the entire ensemble, at the cost of only one forward propagation. To do so, we change to using the geometric mean rather than the arithmetic mean of the ensemble members' predicted distributions. [Warde-Farley \*et al.\* \(2014\)](#) present arguments and empirical evidence that the geometric mean performs comparably to the arithmetic mean in this context.

The geometric mean of multiple probability distributions is not guaranteed to be a probability distribution. To guarantee that the result is a probability distribution, we impose the requirement that none of the sub-models assigns probability 0 to any event, and we renormalize the resulting distribution. The unnormalized probability distribution defined directly by the geometric mean is given by

$$\tilde{p}_{\text{ensemble}}(y \mid \mathbf{x}) = \sqrt[d]{\prod_{\boldsymbol{\mu}} p(y \mid \mathbf{x}, \boldsymbol{\mu})} \quad (7.54)$$

where  $d$  is the number of units that may be dropped. Here we use a uniform distribution over  $\boldsymbol{\mu}$  to simplify the presentation, but non-uniform distributions are