

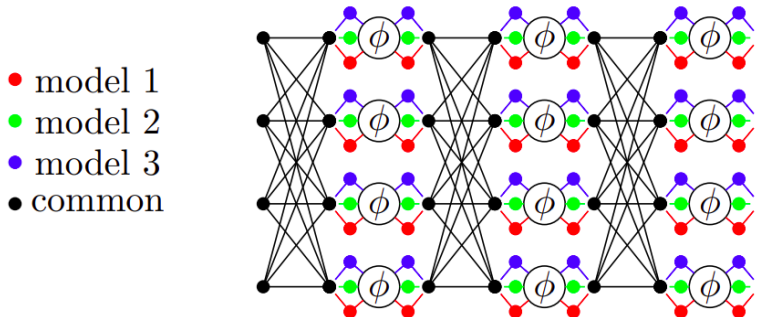
Embedded Ensembles: Infinite Width Limit and Operating Regimes

Шокоров Вячеслав Александрович

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем

Москва,
2022 г.

Embedded Ensembles scheme



All models in the BatchEnsemble have common fully-connected (or convolutional) weights (colored black), and a small number of pre- and post- activation modulations (colored red, green or blue) that differ for each model. For each model, only the respective family of modulations (i.e., red, green or blue) is active on a forward pass.

Standard fully-connected layer:

$$z_j^l = \frac{1}{\sqrt{N_{l-1}}} \sum_{i=1}^{N_{l-1}} x_i^{l-1} W_{ij}^l + b_j^l, \quad x_j^l = \phi(z_j^l)$$

BatchEnsemble:

$$\begin{cases} z_{\alpha j}^l = \frac{1}{\sqrt{N_{l-1}}} \sum_{i=1}^{N_{l-1}} x_{\alpha i}^{l-1} W_{ij}^l + b_j^l, \\ x_{\alpha j}^l = u_{\alpha j}^l \phi(v_{\alpha j}^l z_{\alpha j}^l), \end{cases}$$

where $u_{\alpha} = \{u_{\alpha j}^l, v_{\alpha j}^l\}$ - modulating weights.

$w = \{W_{ij}^l, b_j^l\}$ - shared weights.

Embedded Ensembles: MC dropout ensembles

BatchEnsemble:

$$\begin{cases} z_{\alpha j}^l = \frac{1}{\sqrt{N_{l-1}}} \sum_{i=1}^{N_{l-1}} x_{\alpha i}^{l-1} W_{ij}^l + b_j^l, \\ x_{\alpha j}^l = u_{\alpha j}^l \phi(v_{\alpha j}^l z_{\alpha j}^l), \end{cases}$$

where $u_{\alpha} = \{u_{\alpha j}^l, v_{\alpha j}^l\}$ - modulating weights.

$w = \{W_{ij}^l, b_j^l\}$ - shared weights.

MC dropout ensemble - modification of BatchEnsemble:

$$\begin{cases} z_{\alpha j}^l = \frac{1}{\sqrt{N_{l-1}}} \sum_{i=1}^{N_{l-1}} x_{\alpha i}^{l-1} W_{ij}^l + b_j^l, \\ x_{\alpha j}^l = u_{\alpha j}^l \phi(z_{\alpha j}^l), \end{cases}$$

where $u_{\alpha j}$ have a Bernoulli 0-1 distribution.

$w = \{W_{ij}^l, b_j^l\}$ - shared weights.

$$\Delta \mathbf{u}_\alpha = -\mu \frac{\partial L_\alpha(\mathbf{w}, \mathbf{u}_\alpha)}{\partial \mathbf{u}_\alpha}$$

$$\Delta \mathbf{w} = -\mu \mathbf{w} \frac{\gamma(M)}{M} \sum_{\alpha=1}^M \frac{\partial L_\alpha(\mathbf{w}, \mathbf{u}_\alpha)}{\partial \mathbf{w}}$$

μ_w, μ_u - different learning rates. Authors propose to control accumulation of gradients by means of scaling factor $\gamma(M)$. In the sequel we will argue that the natural choice for scaling factor is either $\gamma(M) = 1$ or $\gamma(M) = M$ depending on whether the ensemble is in independent or collective regime.

Independence of two different models can be achieved if:

- independent initialization
- have dynamic independence

Dynamic independence:

$$\Delta f_{\alpha}(x) = \sum_{\beta=1}^M \frac{\partial f_{\alpha}(x)}{\partial \mathbf{w}} \Delta_{\beta} \mathbf{w} \propto - \sum_{\beta=1}^M \frac{\partial f_{\alpha}(x)}{\partial \mathbf{w}} \frac{\partial \mathcal{L}_{\beta}}{\partial \mathbf{w}}$$

So, dynamic independence $\Leftrightarrow \frac{\partial f_{\alpha}(x)}{\partial \mathbf{w}} \frac{\partial \mathcal{L}_{\beta}}{\partial \mathbf{w}} = 0$

Neural Tangent Kernel (NTK) for single network

$$\frac{df(x)}{dt} = -\frac{1}{B} \sum_{b=1}^B \Theta(x, x_b, w) \frac{\partial L(f(x_b), y_b)}{\partial f(x_b)},$$
$$\Theta(x, x', w) = \frac{\partial f(w, x)}{\partial w} \frac{\partial f(w, x')}{\partial w}.$$

Here $\Theta(x, x', w)$ is the Neural Tangent Kernel of network $f(x, w)$.
If $N \rightarrow \infty$, the following three statements are hold:

- (A1) The network output $f(x)$ is a draw from a Gaussian Process (GP) at initialization.
- (A2) The NTK $\Theta(x, x')$ converges to a non-random deterministic value at initialization.
- (A3) The NTK $\Theta(x, x')$ stays constant during training.

[Jacot et al., 2018]

Neural Tangent Kernel (NTK) for Embedded Ensemble

$$\frac{df_{\alpha}(x)}{dt} = -\frac{1}{B} \sum_{b,\beta} \Theta_{\alpha\beta}(x, x_b) \frac{\partial L(f_{\beta}(x_b), y_b)}{\partial f_{\beta}(x_b)},$$

$$\Theta_{\alpha\beta}(x, x') = \frac{\gamma(M)}{M} \Theta_{\alpha\beta}^{\text{com}}(x, x') + \delta_{\alpha} \Theta_{\alpha}^{\text{ind}}(x, x')$$

$$\Theta_{\alpha\beta}^{\text{com}}(x, x') = \frac{\partial f(w, u_{\alpha}, x)}{\partial w} \frac{\partial f(w, w_{\beta} x')}{\partial w}.$$

$$\Theta_{\alpha}^{\text{ind}}(x, x') = \frac{\partial f(w, u_{\alpha}, x)}{\partial u_{\alpha}} \frac{\partial f(w, w_{\alpha} x')}{\partial u_{\alpha}}.$$

Theorem 1 (Outputs f_α).

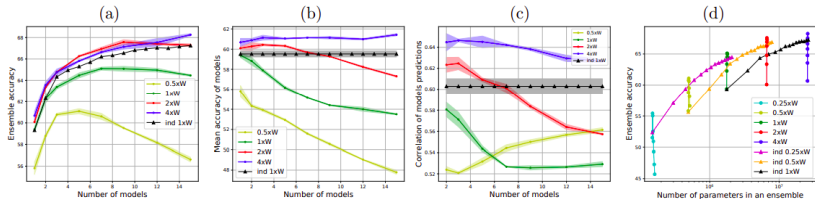
- 1 **(Gaussianity)** Consider a BatchEnsemble or MC dropout ensemble at initialization. Then in the sequential infinite-width limit $N_l \rightarrow \infty$ the collection of ensemble model outputs $f_\alpha(x)$ converge (in law) to a zero mean Gaussian Process (GP).
- 2 **(Independence)** If $U_1^L \equiv \mathbb{E}[u_{\alpha j}^L] = 0$, then the GP covariance $\mathbb{E}[f_\alpha(x)f_\beta(x')] = 0$ for all x, x' and different ensemble models $\alpha \neq \beta$.
- 3 **(Breakdown of independence)** Let the activation function $\phi \in \mathcal{S}$ and $U_1^L \equiv \mathbb{E}[u_{\alpha j}^L] \neq 0$. Then $\mathbb{E}[f_\alpha(x)f_\beta(x')] > 0$ for all $\alpha \neq \beta$ and all (resp., all linearly independent) pairs of non-zero inputs x, x' for network depths $L > 1$.

\mathcal{S} is family of non-negative non-decreasing locally Lipschitz functions.

Theorem 2 (NTK $\Theta_{\alpha\beta}$).

- ① **(Determinacy)** Consider a BatchEnsemble or MC dropout ensemble at initialization. Then in the sequential limit $N_I \rightarrow \infty$ the ensemble NTK $\Theta_{\alpha\beta}(\mathbf{x}, \mathbf{x}')$ converges to a deterministic value.
- ② **(Dynamic independence)** If $U_1^L = 0$, then the ensemble NTK $\Theta_{\alpha\beta}(\mathbf{x}, \mathbf{x}') = 0$ for all \mathbf{x}, \mathbf{x}' and different ensemble models $\alpha \neq \beta$.
- ③ **(Breakdown of dynamic independence)** If $U_1^L \neq 0$ and the activation function $\phi \in \mathcal{S}$, then $\Theta_{\alpha\beta}(\mathbf{x}, \mathbf{x}') > 0$ for all $\alpha \neq \beta$ and all pairs of inputs \mathbf{x}, \mathbf{x}' with scalar product $\mathbf{x}^T \mathbf{x}' > 0$.

Experiments



Performance of BatchEnsembles in independent regime and usual independent ensembles for different model widths and numbers of models. Each curve corresponds to a particular model width. The respective value (0.25x – 4x) is the number of neurons in each layer relative to the baseline network. (a), (b), (c): From left to right: Accuracy of ensemble predictions on test set, mean accuracy of individual EE models on test set, mean accuracy of individual EE models on train set, . (d): Scatter plot comparing various ensembles with respect to accuracy and the absolute number of parameters.