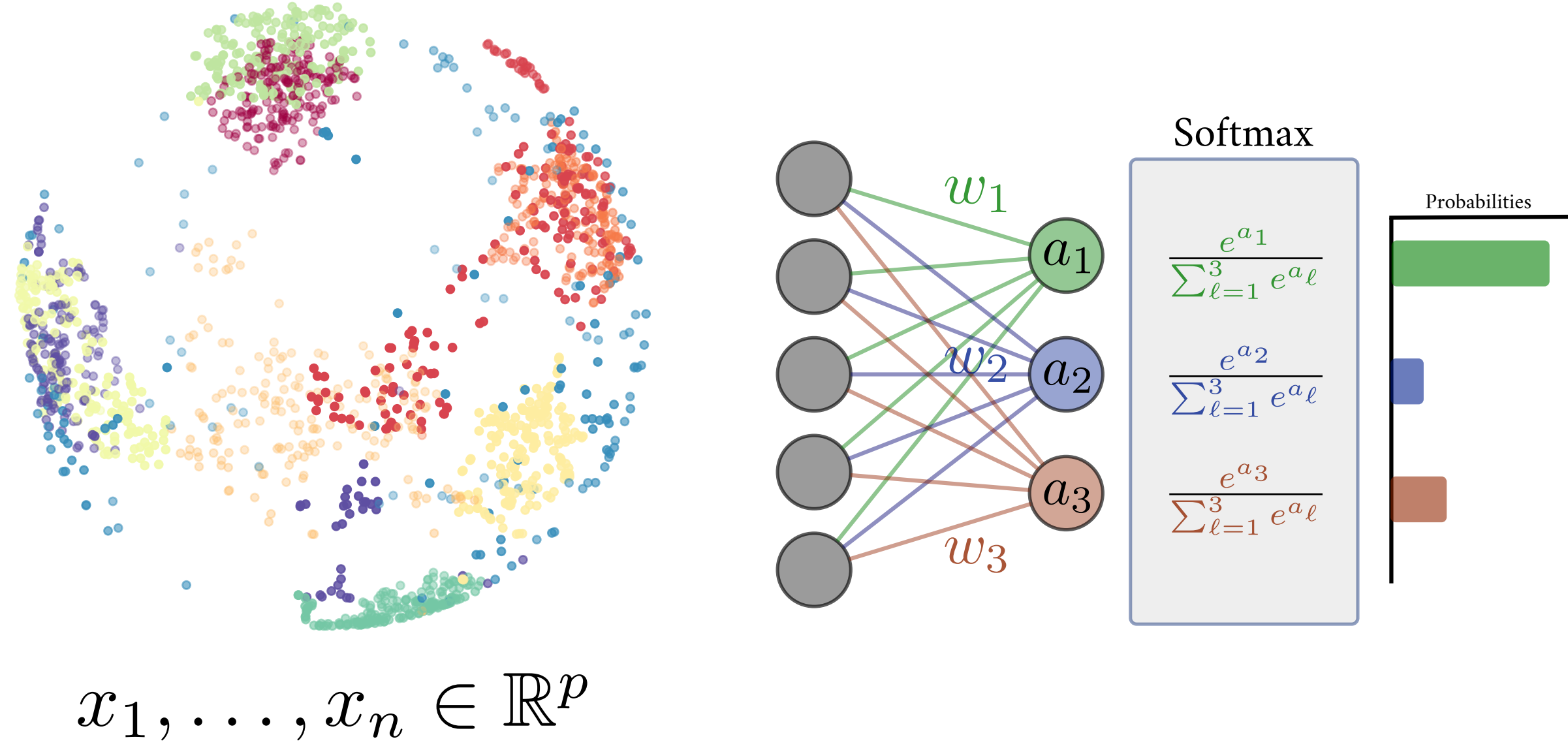


## Abstract



- **Theoretically:** Asymptotic analysis of the Softmax classifier suggests that weights correlate with the class-wise means of the input features.
- **Empirically:** The theoretical observation extends to the Softmax layer of feed-forward neural networks trained for classification tasks.
- **Consequences for transfer learning:** a simple initialization procedure of the Softmax weights is proposed based on the theoretical findings.

## Settings

We consider  $n$  data points  $\mathbf{d}_1, \dots, \mathbf{d}_m$  with their corresponding labels  $\ell_1, \dots, \ell_n$  distributed in  $k$  different classes  $\mathcal{C}_1, \dots, \mathcal{C}_k$ . Denote

$$\mathbf{x} = \varphi \circ \phi(\mathbf{d}; \Theta) \in \mathbb{R}^p$$

where

- $\phi$  is implemented by a deep CNN model parameterized by  $\Theta$ .
- $\varphi$  is the activation function at the representation layer.

The final class prediction is given by a classifier function  $\psi : \mathbb{R}^p \rightarrow \mathbb{R}^k$  as

$$\arg \max \psi(\mathbf{x}) \quad \text{with} \quad \psi(\mathbf{x}) = \text{softmax}(\mathbf{W}^\top \mathbf{x}).$$

Suppose that the following statistics exist and well defined

$$\mathbf{m}_\ell = \mathbb{E}_{\mathbf{d} \in \mathcal{C}_\ell}[\mathbf{x}] \quad \mathbf{C}_\ell = \mathbb{E}_{\mathbf{d} \in \mathcal{C}_\ell}[(\mathbf{x} - \mathbf{m}_\ell)(\mathbf{x} - \mathbf{m}_\ell)^\top]$$

## Softmax Classifier

Minimize:

$$\mathcal{L}(\mathbf{w}_1, \dots, \mathbf{w}_k) = -\frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^k y_{i\ell} \log p_{i\ell}$$

$$p_{i\ell} = \frac{\exp(\mathbf{w}_\ell^\top \mathbf{x}_i)}{\sum_{j=1}^k \exp(\mathbf{w}_j^\top \mathbf{x}_i)}, \quad y_\ell^{(i)} = \alpha_{c(i)} \frac{|\delta_{\ell, c(i)} - \varepsilon|}{1 + (k-2)\varepsilon}$$

where  $\alpha_{c(i)}$  and  $\varepsilon > 0$  are hyper-parameters,  $c(i)$  returns the class index of the  $i$ -th datum and  $\delta_{i,j}$  stands for the Kronecker delta.

The classical labels are recovered by setting  $\alpha_{c(i)} = 1$  and  $\varepsilon = 0$ .

## Assumption on statistical model

For  $\mathbf{x}_i \in \mathcal{C}_\ell$ , assume  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{m}_\ell, \mathbf{C}_\ell)$ .

Denote  $\pi_\ell = \lim_n \frac{|\mathcal{C}_\ell|}{n}$  the proportion of class  $\mathcal{C}_\ell$ .

## Expression of the gradients of the Softmax class-weight vectors

$$\nabla_{\mathbf{w}_\ell} \mathcal{L} = \frac{1}{n} \sum_{i=1}^n \alpha_{c(i)} \left( \frac{e^{\mathbf{w}_\ell^\top \mathbf{x}_i}}{\sum_{j=1}^k e^{\mathbf{w}_j^\top \mathbf{x}_i}} - \frac{|\delta_{\ell, c(i)} - \varepsilon|}{1 + (k-2)\varepsilon} \right) \mathbf{x}_i \equiv \frac{1}{n} \sum_{i=1}^n \alpha_{c(i)} f_{\ell, i}(\mathbf{w}_\ell^\top \mathbf{x}_i) \mathbf{x}_i.$$

## Asymptotic Softmax class-weight vectors

Let  $\bar{\mathbf{w}}_1, \dots, \bar{\mathbf{w}}_k$  be the deterministic vectors satisfying  $\mathbb{E}[\nabla_{\bar{\mathbf{w}}_\ell} \mathcal{L}] = 0$ . Thus, each  $\bar{\mathbf{w}}_\ell$  satisfies the implicit equation

$$\bar{\mathbf{w}}_\ell = \left( \sum_{j=1}^k \alpha_j \pi_j \mathbb{E}_j[f'_{\ell, i}(\bar{\mathbf{w}}_\ell^\top \mathbf{x}_i)] \mathbf{C}_j \right)^{-1} \left( \sum_{j=1}^k \alpha_j \pi_j \mathbb{E}_j[f_{\ell, i}(\bar{\mathbf{w}}_\ell^\top \mathbf{x}_i)] \mathbf{m}_j \right)$$

where the notation  $\mathbb{E}_j[g(\mathbf{x}_i)] \equiv \mathbb{E}[g(\mathbf{x}_i) | \mathbf{x}_i \in \mathcal{C}_j]$  for some  $g : \mathbb{R} \rightarrow \mathbb{R}$ .

Setting the parameters  $\alpha_j = (k\pi_j)^{-1}$  results in the class-weight vectors becoming independent of the proportion  $\pi_j$ :

$$\bar{\mathbf{w}}_\ell = \left( \sum_{j=1}^k \mathbb{E}_j[f'_{\ell, i}(\bar{\mathbf{w}}_\ell^\top \mathbf{x}_i)] \mathbf{C}_j \right)^{-1} \left( \sum_{j=1}^k \mathbb{E}_j[f_{\ell, i}(\bar{\mathbf{w}}_\ell^\top \mathbf{x}_i)] \mathbf{m}_j \right)$$

## Near optimal representations (NOR)

Let  $\varepsilon > 0$ , assume

- $\mathbf{m}_i^\top \mathbf{m}_j = \delta_{i,j} \mu_1 + (1 - \delta_{i,j}) \mu_2$  with  $\mu_1 = \mathcal{O}(1)$  and  $\mu_2 = \mathcal{O}(p^{-\varepsilon})$ .
- $\mathbf{C}_j = \sigma_{1,j}^2 \mathbf{I}_p + \sigma_{2,j}^2 (\mathbf{1}_p \mathbf{1}_p^\top - \mathbf{I}_p)$  with  $\sigma_{1,j}^2 = \mathcal{O}(p^{-\varepsilon})$  and  $\sigma_{2,j}^2 = \mathcal{O}(p^{-1-3\varepsilon})$ .

Notably, these conditions ensure (as  $p \rightarrow \infty$ ) that the between-class means are asymptotically orthogonal (maximize the between-class variance) and the within class covariances asymptotically isotropic (independent features).

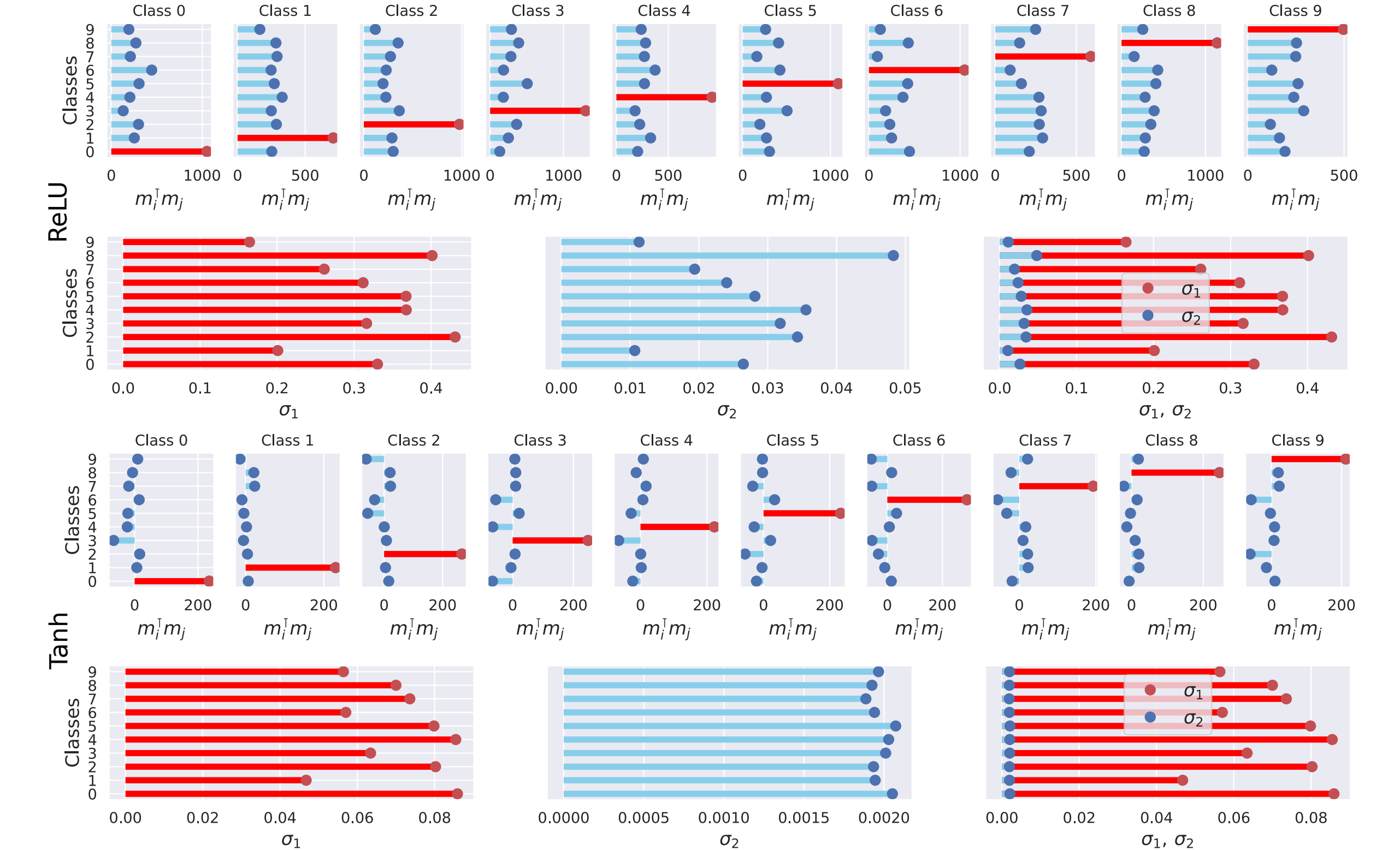
## Asymptotic Softmax weight vectors under NOR assumption

For near optimal representations, for sufficiently large  $p$  and letting  $\varepsilon \rightarrow 0$  in the expression of the generalized labels  $y_\ell^{(i)}$ , the class-weight vectors are asymptotically proportional to the centred class-wise means as

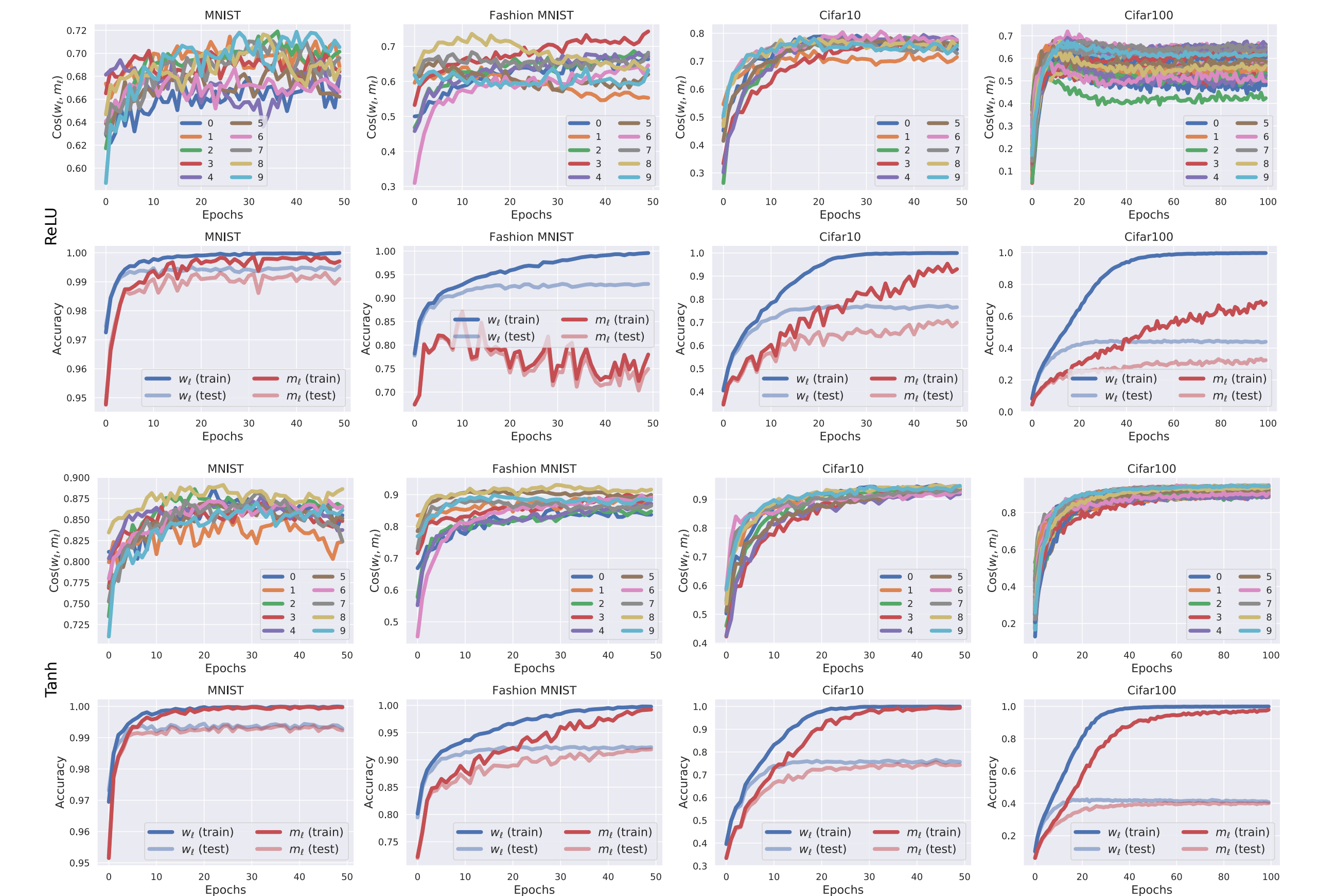
$$\bar{\mathbf{w}}_\ell = \frac{\gamma_\ell k e^{-\kappa \mu_1}}{1 + (k-1) e^{-\kappa \mu_1}} \left( \mathbf{m}_\ell - \frac{1}{k} \sum_{j=1}^k \mathbf{m}_j \right) + \mathcal{O}_{\|\cdot\|}(1) \quad (1)$$

for some constant  $\kappa > 0$  and  $\gamma_\ell = \left( \sum_{j=1}^k \sigma_{1,j}^2 \mathbb{E}_j[-f'_{\ell, i}(\bar{\mathbf{w}}_\ell^\top \mathbf{x}_i)] \right)^{-1} \geq \frac{4}{k \sigma_{1, \max}^2}$ .

## NOR in practice



## Correlation between $\bar{\mathbf{w}}_\ell$ and $\bar{\mathbf{m}}_\ell$



## Conclusion

Our findings suggest three main procedures for efficient transfer learning: (i) use of symmetric representation activations to ensure the near-optimal representations assumption; (ii) source model selection without training the Softmax layer; (iii) initialization procedure which accelerates the training of the Softmax layer as the target domain gets closer to the source domain.