

Mixture of Experts

Seungjin Choi

Department of Computer Science and Engineering
Pohang University of Science and Technology
77 Cheongam-ro, Nam-gu, Pohang 37673, Korea
seungjin@postech.ac.kr

Supervised Learning and Function Approximation

- ▶ The supervised learning involves estimating a function which relates inputs $\{\mathbf{x}_I\}$ to their associated targets $\{\mathbf{t}_I\}$, given a data set $\mathcal{D} = \{\mathbf{x}_I, \mathbf{t}_I\}$.
- ▶ We wish to learn either $p(\mathbf{t}|\mathbf{x})$ or $p(\mathbf{t}, \mathbf{x}) = p(\mathbf{t}|\mathbf{x})p(\mathbf{x})$ if we care to model the input distribution too.
- ▶ Once we have these probabilistic models, we can always choose a deterministic function through

$$f(\mathbf{x}) = \langle \mathbf{t} | \mathbf{x} \rangle,$$

or if the density function is nicely unimodal, the average could be a point of low probability and in such a case we might choose the maximum

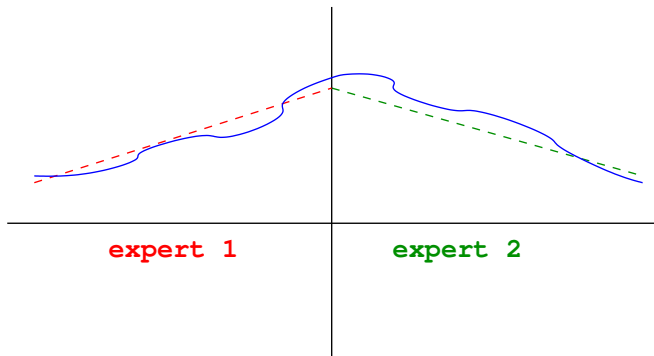
$$f(\mathbf{x}) = \arg \max_{\mathbf{t}} p(\mathbf{t} | \mathbf{x}).$$

Modeling Nonlinear Functions

- ▶ A powerful principle to model nonlinear functions is by approximating it through many **local linear functions**.
- ▶ The same strategy can be adopted for the density estimation problem.
- ▶ We employ many **simple local models** and combine them to approximate the density, very much in the same spirit as MoG.
- ▶ This leads to **mixture of experts**.

Mixture of Experts: A Model

- ▶ Model $p(\mathbf{t}|\mathbf{x})$ as $p(\mathbf{t}|\mathbf{x}) = \sum_{i=1}^M p(\mathbf{t}, i|\mathbf{x}) = \sum_{i=1}^M p(\mathbf{t}|i, \mathbf{x})p(i|\mathbf{x})$.
- ▶ The basic idea is to make local models to specialize on a certain task.
 - ▶ **Experts:** $p(\mathbf{t}|i, \mathbf{x})$ (local models).
 - ▶ **Gatekeepers:** $p(i|\mathbf{x})$ (what region experts should model, is regulated by the gatekeepers)



Models for Experts and Gates

- ▶ For the gates, we use the **softmax** functions,

$$p(i|\mathbf{x}) = g_i = \frac{e^{\xi_i}}{\sum_{j=1}^M e^{\xi_j}},$$

where $\xi_i = \mathbf{v}_i^\top \mathbf{x} + \mathbf{v}_i^0 = \tilde{\mathbf{v}}_i^\top \tilde{\mathbf{x}}$.

- ▶ Modelling experts is as follows. First we linearly combine the inputs

$$\boldsymbol{\mu}_i(\mathbf{x}) = \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^0 = \tilde{\mathbf{W}}_i \tilde{\mathbf{x}},$$

which will be used as mean values for a Gaussian distribution modelling the target data, i.e.,

$$p(\mathbf{t}|\mathbf{x}, i) = \mathcal{N}(\mathbf{t} | [\boldsymbol{\mu}_i(\mathbf{x}), \sigma^2 \mathbf{I}]),$$

where we treat the noise as a fixed parameter, but we could estimate it from the data.

Learn MoE

- ▶ How to learn MoE? \Rightarrow EM!
- ▶ In E-step, we compute responsibilities and calculate the expected complete-data log-likelihood.
- ▶ In M-step, we calculate the partial derivatives of the expected complete-data log-likelihood with respect to associated parameters, in order to derive updating rules for those parameters.

EM-MoE: E-Step

We compute responsibilities:

$$\begin{aligned} R_{in} &= p(i|\mathbf{x}_n, \mathbf{t}_n) \\ &= \frac{p(\mathbf{t}_n|\mathbf{x}_n, i) p(i|\mathbf{x}_n)}{\sum_{j=1}^M p(\mathbf{t}_n|\mathbf{x}_n, j) p(j|\mathbf{x}_n)}. \end{aligned}$$

The complete-data log-likelihood is

$$\begin{aligned} \mathcal{L}_c &= \sum_{n=1}^N \sum_{i=1}^M \log p(\mathbf{x}_n, \mathbf{t}_n, i) \\ &= \sum_{n=1}^N \sum_{i=1}^M \log [p(\mathbf{t}_n|\mathbf{x}_n, i) p(i|\mathbf{x}_n)]. \end{aligned}$$

Hence, the expected complete-data log-likelihood is

$$\langle \mathcal{L}_c \rangle = \sum_{n=1}^N \sum_{i=1}^M R_{in} [\log p(\mathbf{t}_n|\mathbf{x}_n, i) + \log p(i|\mathbf{x}_n)].$$

EM-MoE: M-Step

- ▶ Solving $\frac{\partial \langle \mathcal{L}_c \rangle}{\partial \tilde{\mathbf{W}}_i} = 0$ leads to $\tilde{\mathbf{W}}_i = \left(\sum_{n=1}^N R_{in} \mathbf{t}_n \tilde{\mathbf{x}}_n^\top \right) \left(\sum_{n=1}^N R_{in} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^\top \right)^{-1}$.
- ▶ We have $\frac{\partial \langle \mathcal{L}_c \rangle}{\partial \tilde{\mathbf{v}}_i} = \sum_{n=1}^N (R_{in} - g_{in}) \tilde{\mathbf{x}}_n$, which is hard to optimize since the softmax function depends on the parameters $\{\tilde{\mathbf{v}}_i\}$.
- ▶ We see that g_{in} is optimized to be as close as possible to the R_{in} , leading to

$$\log R_{in} \approx \tilde{\mathbf{v}}_i^\top \tilde{\mathbf{x}}_n - \underbrace{\log \left[\sum_j \exp \{ \tilde{\mathbf{v}}_j^\top \tilde{\mathbf{x}}_n \} \right]}_{\text{normalization constant}}.$$

- ▶ This is reduced to a simple **regression problem**. Minimizing $\sum_n \left(\log R_{in} - \tilde{\mathbf{v}}_i^\top \tilde{\mathbf{x}}_n \right)^2$, leads to

$$\tilde{\mathbf{v}}_i = \left(\sum_{n=1}^N \log [R_{in}] \tilde{\mathbf{x}}_n^\top \right) \left(\sum_{n=1}^N \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^\top \right)^{-1}.$$

Algorithm Outline: EM-MoE

- ▶ **E-Step:** Compute responsibilities

$$\begin{aligned} R_{in} &= p(i|\mathbf{x}_n, \mathbf{t}_n) \\ &= \frac{p(\mathbf{t}_n|\mathbf{x}_n, i) p(i|\mathbf{x}_n)}{\sum_{j=1}^M p(\mathbf{t}_n|\mathbf{x}_n, j) p(j|\mathbf{x}_n)}, \end{aligned}$$

where $p(\mathbf{t}_n|\mathbf{x}_n, i) = \mathcal{G}_{\mathbf{t}_n}[\boldsymbol{\mu}_i(\mathbf{x}_n), \sigma^2 \mathbf{I}]$ with $\boldsymbol{\mu}_i(\mathbf{x}_n) = \widetilde{\mathbf{W}}_i \widetilde{\mathbf{x}}_n$, and $p(i|\mathbf{x}_n) = g_{in} = \frac{e^{\xi_{in}}}{\sum_j e^{\xi_{jn}}}$ with $\xi_{in} = \widetilde{\mathbf{v}}_i^\top \widetilde{\mathbf{x}}_n$.

- ▶ **M-Step:** Update parameters

$$\begin{aligned} \widetilde{\mathbf{W}}_i &\leftarrow \left(\sum_{n=1}^N R_{in} \mathbf{t}_n \widetilde{\mathbf{x}}_n^\top \right) \left(\sum_{n=1}^N R_{in} \widetilde{\mathbf{x}}_n \widetilde{\mathbf{x}}_n^\top \right)^{-1}, \\ \widetilde{\mathbf{v}}_i &\leftarrow \left(\sum_{n=1}^N \log[R_{in}] \widetilde{\mathbf{x}}_n^\top \right) \left(\sum_{n=1}^N \widetilde{\mathbf{x}}_n \widetilde{\mathbf{x}}_n^\top \right)^{-1}. \end{aligned}$$