

# Prior distribution choices for a mixture of experts \*

A. V. Grabovoy<sup>1</sup>, V. V. Strijov<sup>2</sup>

**Abstract:** The paper investigates a mixture of expert models. The mixture of experts is a combination of experts, local approximation model, and a gate function, which weighs these experts and forms their ensemble. In this work, each expert is a linear model. The gate function is a neural network with softmax on the last layer. The paper analyzes various prior distributions for each expert. The authors propose a method that takes into account the relationship between prior distributions of different experts. The EM algorithm optimises both parameters of the local models and parameters of the gate function. As an application problem, the paper solves a problem of shape recognition on images. Each expert fits one circle in an image and recovers its parameters: the coordinates of the center and the radius. The computational experiment uses synthetic and real data to test the proposed method. The real data is a human eye image from the iris detection problem.

**Keywords:** mixture of Experts; bayesian model selection; prior distribution.

## 1 Introduction

The paper studies the problem of mixture of experts construction. The mixture of experts is a multimodel. The mixture of experts uses a gate function to weight predictions of each expert. It is a combination of weighed local models to approximate a dataset. The weighting coefficients depend on objects in the dataset. Examples of multimodel are bagging, gradient boosting [1] and random forest [2]. The paper [3] suggests that the contribution of each expert to the answer depends on the object from the dataset.

The main problem of multimodel construction is high dependence of the resulted ensemble on the initial value of the parameters. The authors propose to use the probability

---

\*This research was supported by RFBR (project ???) and NTI (project ???).

<sup>1</sup>Moscow Institute of Physics and Technology, grabovoy.av@phystech.edu

<sup>2</sup>Moscow Institute of Physics and Technology, Dorodnicyn Computing Centre, Federal Research Center “Computer Science and Control” of the Russian Academy of Sciences, strijov@phystech.edu

approach to find optimal gate function parameters and local parameters, and proposes to use various prior distributions of the parameters to improve stability of the multimodel. The paper introduces a method, which use dependence between prior distributions to improve the multimodel quality.

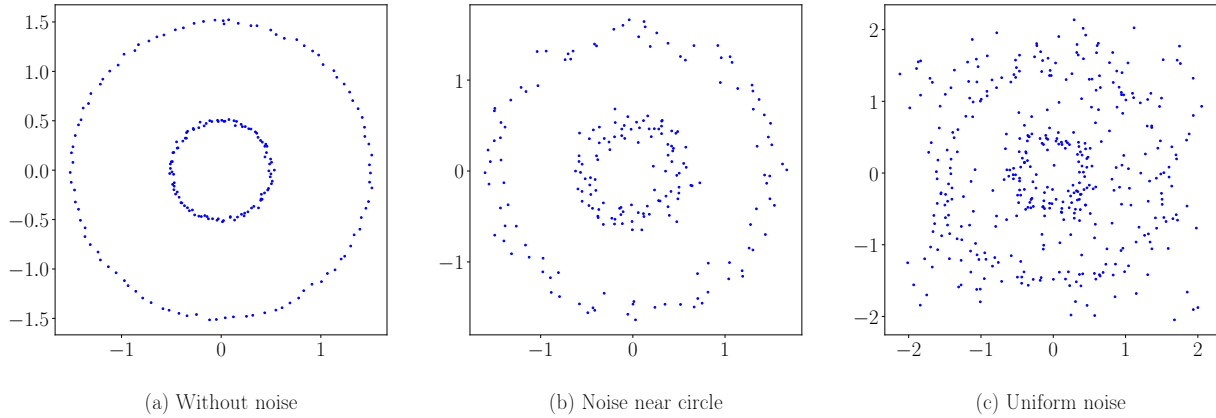


Figure 1: An example of circles with different noise levels: (a) circle without noise; (b) noisy radius circle; (c) noisy radius circle, and uniform noisy on image

As a practical example, the circles approximation on a binarized image are solved. Examples of images are shown in fig. 1. In this paper, each expert is a linear model. A gate function is a two-layer fully connected neural network.

**Related work.** Many papers on a mixture of experts are devoted to gateway function selection: softmax, the Dirichlet process [4], neural networks [5] with softmax function on the last layer. Some papers are devoted to the choice of the expert type. Papers [6, 7] analyze linear models as experts. Papers [8, 9] analyze SVMs as experts. The paper [3] contains an overview of various gate function and various expert types.

The method of mixture of experts has many applications. Papers [10, 11, 12] use the mixture of experts in time series forecasting. Paper [13] uses the mixture of experts in the task of recognizing handwritten numbers. Papers [14, 15, 16, 17] are devoted to methods of text and speech recognition by using the mixture of experts. Paper [18] analyzes the mixture of experts to recognize three-dimensional human movements. Paper [19] is devoted to a review of the study results on the iris detection in the image. The methods of highlighting the borders of the iris and pupil are described in papers [20, 21].

## 2 Problem statement of circle parameters estimation

There given binary image

$$\mathbf{M} \in \{0, 1\}^{m_1 \times m_2},$$

where 1 is a black pixel, an image foreground, and 0 is a white pixel, the image background. An image example is shown in fig. 1. The image  $\mathbf{M}$  is mapped to a set of coordinates  $\mathbf{C} = \{x_i, y_i\}_{i=1}^N$ . The pair of coordinates  $(x_i, y_i)$  is a black pixel in  $\mathbf{M}$ :

$$\mathbf{C} \in \mathbb{R}^{N \times 2},$$

where  $N$  is the number of black pixels.

Let  $(x_0, y_0)$  be the center of the circle, and  $r$  is radius of the circle. The coordinates  $(x_i, y_i) \in \mathbf{C}$  is a circle locus of points defined by

$$(x_i - x_0)^2 + (y_i - y_0)^2 = r^2.$$

Expand brackets:

$$(2x_0) \cdot x_i + (2y_0) \cdot y_i + (r^2 - x_0^2 - y_0^2) \cdot 1 = x_i^2 + y_i^2. \quad (2.1)$$

Rewrite equation (2.1) to set the linear regressions problem for all points in the dataset:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|, \quad \mathbf{X} = [\mathbf{C}, \mathbf{1}], \quad \mathbf{y} = [x_1^2 + y_1^2, x_2^2 + y_2^2, \dots, x_N^2 + y_N^2]^\top. \quad (2.2)$$

The parameters  $\hat{\mathbf{w}} = [w_1, w_2, w_3]^\top$  reconstruct the circle parameters  $x_0, y_0, r$ :

$$x_0 = \frac{w_1}{2}, \quad y_0 = \frac{w_2}{2}, \quad r = \sqrt{w_3 + x_0^2 + y_0^2}.$$

The solution of problem (2.2) reconstructs the circle parameters only if the number of circles in an image is equal to one. The authors propose to use the multimodel for the image, which consists of several circles. The multimodel is an ensemble of the linear models. Each linear model approximates only one circle in the image. In this paper, multimodel is a mixture of experts.

### 3 Problem statement of building a mixture of experts

Generalize one-circle approximation (2.2) problem to the case of several circles. Each circle is a local model. The data for this case is

$$\mathbf{X} \in \mathbb{R}^{N \times n}, \quad \mathbf{y} \in \mathbb{R}^N \quad (3.1)$$

where  $N$  is the sample size and  $n$  is the number of features. In this paper,  $n$  is equal to 3.

**Definition 3.1.** A model  $\mathbf{g}$  is a local model on dataset  $\mathbf{X}$  if  $\mathbf{g}$  approximates some no-empty subset  $\mathbf{X}' \subset \mathbf{X}$ .

**Definition 3.2.** Call the multimodel  $\mathbf{f}$  a mixture of experts

$$\mathbf{f} = \sum_{k=1}^K \pi_k \mathbf{g}_k(\mathbf{w}_k), \quad \pi_k(\mathbf{x}, \mathbf{V}) : \mathbb{R}^{n \times |\mathbf{V}|} \rightarrow [0, 1], \quad \sum_{k=1}^K \pi_k(\mathbf{x}, \mathbf{V}) = 1, \quad (3.2)$$

where  $\mathbf{g}_k$  is a local model,  $\pi_k$  is a gate function, vector  $\mathbf{w}_k$  is some parameters of the local model and  $\mathbf{V}$  is some parameters of the gate function.

This paper asserts the local model be linear model. The gate function is the two-layer fully connected neural network

$$\mathbf{g}_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x}, \quad \boldsymbol{\pi}(\mathbf{x}, \mathbf{V}) = \text{softmax}(\mathbf{V}_1^\top \boldsymbol{\sigma}(\mathbf{V}_2^\top \mathbf{x})), \quad (3.3)$$

where  $\mathbf{V} = \{\mathbf{V}_1, \mathbf{V}_2\}$  is a set of the gate function parameters.

The paper proposes to use a probabilistic approach to describe a mixture of experts. Let  $\mathbf{y}$  be a random variable with density function  $p(\mathbf{y}|\mathbf{X})$ . Let density  $p(\mathbf{y}|\mathbf{X}, \mathbf{f})$  approximate truth density  $p(\mathbf{y}|\mathbf{X})$ :

$$p(\mathbf{y}|\mathbf{X}, \mathbf{f}) = \prod_{i=1}^N \left( \sum_{k=1}^K \pi_k p_k(y_i | \mathbf{g}_k(\mathbf{x}_i)) \right), \quad (3.4)$$

where  $\mathbf{f}$  is the mixture of experts and  $\mathbf{g}_k, \boldsymbol{\pi}$  are defined by (??).

Suppose that  $\mathbf{w}_k$  is the random variable with density function  $p^k(\mathbf{w}_k)$ , and get joint probability distribution of the target and the parameters:

$$p(\mathbf{y}, \mathbf{W}|\mathbf{X}, \mathbf{V}) = \prod_{k=1}^K p^k(\mathbf{w}_k) \prod_{i=1}^N \left( \sum_{k=1}^K \pi_k p_k(y_i | \mathbf{w}_k, \mathbf{x}_i) \right), \quad (3.5)$$

where  $\mathbf{W} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K\}$ . The optimal parameters is delivered by the evidence maximisation:

$$\hat{\mathbf{V}} = \arg \max_{\mathbf{V}} p(\mathbf{y}|\mathbf{X}, \mathbf{V}).$$

## 4 Probabilistic statement of mixture of expert

To build the mixture of experts (??, ??), set the following probabilistic statement for the dataset (??):

- 1) the likelihood  $p_k(y_i | \mathbf{w}_k, \mathbf{x}_i) = \mathcal{N}(y_i | \mathbf{w}_k^\top \mathbf{x}_i, \beta^{-1})$ , where parameter  $\beta$  is the noise level,
- 2) the prior distribution for the parameters  $p^k(\mathbf{w}_k) = \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k)$ , where  $\mathbf{w}_k^0$  is a vector of size  $n \times 1$  and  $\mathbf{A}_k$  is a covariance matrix,
- 3) the prior regularisation  $p(\boldsymbol{\varepsilon}_{k,k'} | \boldsymbol{\Xi}) = \mathcal{N}(\boldsymbol{\varepsilon}_{k,k'} | \mathbf{0}, \boldsymbol{\Xi})$ , where  $\boldsymbol{\Xi}$  is a covariance matrix and  $\boldsymbol{\varepsilon}_{k,k'} = \mathbf{w}_k^0 - \mathbf{w}_{k'}^0$ .

The assumption 2) is the prior distribution for the parameters  $\mathbf{w}_k$ . It sets some restrictions to the model. For example, if  $\mathbf{w}_k^0$  is  $[0,0,1]$ , then the  $k$ -th local model fits with the most probability a circle with parameters  $x_0 = 0, y_0 = 0, r = 1$ .

The assumption 3) is a prior regularisation. It sets restrictions on the prior distribution parameters. For example, if  $\text{diag}(\boldsymbol{\Xi}) = [0.001, 0.001, 1]$ , then the centers of different circles, which corresponds to different local models are equal.

Combining assumptions 1), 2), 3) and equation (??), obtain the new likelihood:

$$p(\mathbf{y}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) = \prod_{i=1}^N \left( \sum_{k=1}^K \pi_k \mathcal{N}(y_i | \mathbf{w}_k^\top \mathbf{x}_i, \beta^{-1}) \right) \cdot \prod_{k=1}^K \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k) \cdot \prod_{k,k'=1}^K \mathcal{N}(\boldsymbol{\varepsilon}_{k,k'} | \mathbf{0}, \mathbf{\Xi}), \quad (4.1)$$

where  $\mathbf{A} = \{\mathbf{A}_1, \dots, \mathbf{A}_K\}$ .

Introduce a binary matrix  $\mathbf{Z}$ . Its element  $z_{ik}$  is equal to 1 if and only if the  $i$ -th object corresponds to the  $k$ -th local model. Using the binary matrix  $\mathbf{Z}$  in (??) take logarithm and obtain

$$\begin{aligned} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) = & \\ = \sum_{i=1}^N \sum_{k=1}^K z_{ik} & \left[ \log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (y_i - \mathbf{w}_k^\top \mathbf{x}_i)^2 + \frac{1}{2} \log \frac{\beta}{2\pi} \right] + \\ + \sum_{k=1}^K & \left[ -\frac{1}{2} (\mathbf{w}_k - \mathbf{w}_k^0)^\top \mathbf{A}_k^{-1} (\mathbf{w}_k - \mathbf{w}_k^0) + \frac{1}{2} \log \det \mathbf{A}_k^{-1} - \frac{n}{2} \log 2\pi \right] + \\ + \sum_{k=1}^K \sum_{k'=1}^K & \left[ -\frac{1}{2} (\mathbf{w}_k^0 - \mathbf{w}_{k'}^0)^\top \mathbf{\Xi}^{-1} (\mathbf{w}_k^0 - \mathbf{w}_{k'}^0) + \frac{1}{2} \log \det \mathbf{\Xi} - \frac{n}{2} \log 2\pi \right]. \end{aligned} \quad (4.2)$$

Set the new optimisation problem to optimise the evidence function. The function is obtained by integrating equation (??) over parameters  $\mathbf{W}, \mathbf{Z}$ :

$$\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta = \arg \max_{\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta} \int_{\mathbf{W}, \mathbf{Z}} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) d\mathbf{W} d\mathbf{Z}. \quad (4.3)$$

## 5 EM–algorithm as a solver of optimisation problem

Let  $q(\mathbf{W}, \mathbf{Z})$  be some density distribution of parameters  $\mathbf{W}, \mathbf{Z}$ . Rewrite the evidence:

$$\begin{aligned}
\log p(\mathbf{y}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) &= \int_{\mathbf{W}, \mathbf{Z}} q(\mathbf{W}, \mathbf{Z}) \log p(\mathbf{y}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) d\mathbf{W} d\mathbf{Z} = \\
&= \int_{\mathbf{W}, \mathbf{Z}} q(\mathbf{W}, \mathbf{Z}) \log \frac{p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta)}{p(\mathbf{W}, \mathbf{Z}|\mathbf{y}, \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta)} d\mathbf{W} d\mathbf{Z} = \\
&= \int_{\mathbf{W}, \mathbf{Z}} q(\mathbf{W}, \mathbf{Z}) \log \frac{p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) q(\mathbf{W}, \mathbf{Z})}{p(\mathbf{W}, \mathbf{Z}|\mathbf{y}, \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) q(\mathbf{W}, \mathbf{Z})} d\mathbf{W} d\mathbf{Z} = \\
&= \int_{\mathbf{W}, \mathbf{Z}} q(\mathbf{W}, \mathbf{Z}) \frac{p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta)}{q(\mathbf{W}, \mathbf{Z})} d\mathbf{W} d\mathbf{Z} + \\
&+ \int_{\mathbf{W}, \mathbf{Z}} q(\mathbf{W}, \mathbf{Z}) \frac{q(\mathbf{W}, \mathbf{Z})}{p(\mathbf{W}, \mathbf{Z}|\mathbf{y}, \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta)} d\mathbf{W} d\mathbf{Z} = \\
&= \mathcal{L}(q, \mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta) + D_{KL}(q(\mathbf{W}, \mathbf{Z}) || p(\mathbf{W}, \mathbf{Z}|\mathbf{y}, \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta))
\end{aligned} \tag{5.1}$$

Using (??), we get lower bound of evidence:

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) \geq \mathcal{L}(q, \mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta),$$

where  $\mathcal{L}(q, \mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta)$  is called lower bound of evidence.

Use the expectation-maximization [22, 23] algorithm to find the solution to optimisation problem (??). The EM-algorithm instead of  $\log p(\mathbf{y}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta)$  optimizes the lower bound  $\mathcal{L}(q, \mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta)$ .

**E-step.** E-step solves the optimisation problem:

$$\mathcal{L}(q, \mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta) \rightarrow \max_{q(\mathbf{W}, \mathbf{Z})},$$

where  $\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta$  are fixed.

Let the joint distribution  $q(\mathbf{Z}, \mathbf{W})$  satisfies the assumption of independence  $q(\mathbf{Z}, \mathbf{W}) = q(\mathbf{Z})q(\mathbf{W})$  [23]. The symbol  $\propto$  means that both sides are equal to up to an additive constant. First, find the distribution  $q(\mathbf{Z})$ :

$$\begin{aligned}
\log q(\mathbf{Z}) &= E_{q/\mathbf{Z}} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W}|\mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) \propto \\
&\propto \sum_{i=1}^N \sum_{k=1}^K z_{ik} \left[ \log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (y_i^2 - \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k + \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i) + \frac{1}{2} \log \frac{\beta}{2\pi} \right] \\
p(z_{ik} = 1) &= \frac{\exp(\log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (\mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k))}{\sum_{k'=1}^K \exp(\log \pi_{k'}(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (\mathbf{x}_i^\top \mathbf{E} \mathbf{w}_{k'} \mathbf{w}_{k'}^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_{k'}))}.
\end{aligned} \tag{5.2}$$

Using (??), we get that distribution  $q(z_{ik})$  is the Bernoulli distribution with probability  $z_{ik}$  from equation (??). Second, find the distribution  $q(\mathbf{W})$ :

$$\begin{aligned}
\log q(\mathbf{W}) &= \mathbb{E}_{q/\mathbf{W}} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) \propto \\
&\propto \sum_{i=1}^N \sum_{k=1}^K \mathbb{E} z_{ik} \left[ \log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (y_i - \mathbf{w}_k^\top \mathbf{x}_i)^2 + \frac{1}{2} \log \frac{\beta}{2\pi} \right] + \\
&+ \sum_{k=1}^K \left[ -\frac{1}{2} (\mathbf{w}_k - \mathbf{w}_k^0)^\top \mathbf{A}_k^{-1} (\mathbf{w}_k - \mathbf{w}_k^0) + \frac{1}{2} \log \det \mathbf{A}_k^{-1} - \frac{n}{2} \log 2\pi \right] \\
&\propto \sum_{k=1}^K \left[ \mathbf{w}_k^\top \left( \mathbf{A}_k^{-1} \mathbf{w}_k^0 + \beta \sum_{i=1}^N \mathbf{x}_i y_i \mathbb{E} z_{ik} \right) - \frac{1}{2} \mathbf{w}_k^\top \left( \mathbf{A}_k^{-1} + \beta \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{w}_k \right].
\end{aligned} \tag{5.3}$$

Using (??), we get that distribution  $q(\mathbf{w}_k)$  is the normal distribution with mean  $\mathbf{m}_k$  and covariance matrix  $\mathbf{B}_k$ :

$$\mathbf{m}_k = \mathbf{B}_k \left( \mathbf{A}_k^{-1} \mathbf{w}_k^0 + \beta \sum_{i=1}^N \mathbf{x}_i y_i \mathbb{E} z_{ik} \right), \quad \mathbf{B}_k = \left( \mathbf{A}_k^{-1} + \beta \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \mathbb{E} z_{ik} \right)^{-1}.$$

**M-step.** M-step solves the optimisation problem:

$$\mathcal{L}(q, \mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta) \rightarrow \max_{\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta},$$

where  $q(\mathbf{W}, \mathbf{Z})$  are fixed density function. The distribution  $q(\mathbf{Z}, \mathbf{W})$  is fixed, while the lower bound  $\mathcal{L}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta)$  is maximized with respect to the parameters  $\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta$ :

$$\begin{aligned}
\mathcal{L}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta) &= \mathbb{E}_q \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \mathbf{\Xi}, \beta) = \\
&= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E} z_{ik} \left[ \log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} \mathbb{E} (y_i - \mathbf{w}_k^\top \mathbf{x}_i)^2 + \frac{1}{2} \log \frac{\beta}{2\pi} \right] + \\
&+ \sum_{k=1}^K \left[ -\frac{1}{2} \mathbb{E} (\mathbf{w}_k - \mathbf{w}_k^0)^\top \mathbf{A}_k^{-1} (\mathbf{w}_k - \mathbf{w}_k^0) + \frac{1}{2} \log \det \mathbf{A}_k^{-1} - \frac{n}{2} \log 2\pi \right] + \\
&+ \sum_{k=1}^K \sum_{k'=1}^K \left[ -\frac{1}{2} (\mathbf{w}_k^0 - \mathbf{w}_{k'}^0)^\top \mathbf{\Xi}^{-1} (\mathbf{w}_k^0 - \mathbf{w}_{k'}^0) + \frac{1}{2} \log \det \mathbf{\Xi} - \frac{n}{2} \log 2\pi \right].
\end{aligned} \tag{5.4}$$

First, to find the optimal parameters  $\mathbf{V}$ , use the gradient optimisation algorithm. It converges to some local maximum. Second, using (??), we get optimal value for  $\mathbf{A}_k$

$$\begin{aligned}
\frac{\partial \mathcal{L}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta)}{\partial \mathbf{A}_k^{-1}} &= \frac{1}{2} \mathbf{A}_k - \frac{1}{2} \mathbb{E} (\mathbf{w}_k - \mathbf{w}_k^0) (\mathbf{w}_k - \mathbf{w}_k^0)^\top = 0, \\
\mathbf{A}_k &= \mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top - \mathbf{w}_k^0 \mathbb{E} \mathbf{w}_k^\top - \mathbb{E} \mathbf{w}_k \mathbf{w}_k^{0\top} + \mathbf{w}_k^0 \mathbf{w}_k^{0\top}.
\end{aligned}$$

Similarly, we get optimal value for  $\beta$  and for  $\mathbf{w}_k^0$

$$\begin{aligned}
\frac{\partial \mathcal{L}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta)}{\partial \beta} &= \sum_{k=1}^K \sum_{i=1}^N \left( \frac{1}{\beta} \mathbf{E} z_{ik} - \frac{1}{2} \mathbf{E} z_{ik} [y_i^2 - 2y_i \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k + \mathbf{x}_i^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i] \right) = 0, \\
\frac{1}{\beta} &= \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K [y_i^2 - 2y_i \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k + \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i] \mathbf{E} z_{ik}. \\
\frac{\partial \mathcal{L}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}, \beta)}{\partial \mathbf{w}_k^0} &= \mathbf{A}_k^{-1} (\mathbf{E} \mathbf{w}_k - \mathbf{w}_k^0) + \Xi \sum_{k'=1}^K [\mathbf{w}_{k'}^0 - \mathbf{w}_k^0] = 0, \\
\mathbf{w}_k^0 &= [\mathbf{A}_k^{-1} + (K-1) \Xi]^{-1} \left( \mathbf{A}_k^{-1} \mathbf{E} \mathbf{w}_k + \Xi \sum_{k'=1, k' \neq k}^K \mathbf{w}_{k'}^0 \right).
\end{aligned} \tag{5.5}$$

Formulas (5.5) define an iterative procedure, which convergence to some local maximum of the optimisation problem (5.5).

## 6 Computational experiment

The computational experiment analyzes quality of various multimodels for circle approximation. The experiment analyzes next multimodels: multimodel  $\mathbf{f}_1$  without prior distribution for the parameters, multimodel  $\mathbf{f}_2$  with prior distribution (6.1) for the parameters and multimodel  $\mathbf{f}_3$  with prior regularisation. The approximation quality of model  $\mathbf{f}_i$  is

$$\mathcal{S}_{\mathbf{f}_i} = \sum_{k=1}^K (x_0^k - x_{\text{pr}}^k)^2 + (y_0^k - y_{\text{pr}}^k)^2 + (r^k - r_{\text{pr}}^k)^2, \tag{6.1}$$

where  $x_0^k, y_0^k, r^k$  are the true center and radius for  $k$ -th circle,  $x_{\text{pr}}^k, y_{\text{pr}}^k, r_{\text{pr}}^k$  are the predicted center and radius for  $k$ -th circle.

To compare models with various prior distribution, use the log-likelihood without any priors (6.1). The prior distribution for the parameters in the experiment is

$$p^1(\mathbf{w}_1) \sim \mathcal{N}(\mathbf{w}_1^0, \mathbf{I}), \quad p^2(\mathbf{w}_2) \sim \mathcal{N}(\mathbf{w}_2^0, \mathbf{I}), \tag{6.2}$$

where  $\mathbf{w}_1^0 = [0, 0, 0.1]$ ,  $\mathbf{w}_2^0 = [0, 0, 2]$ .

**Synthetic data with various types of noise in the image.** The computational experiment compares quality of mixture of experts  $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$  on a synthetic dataset. The synthetic data were generated as two concentric circles with different noise levels. Synthetic 1 is the image without any noises, Synthetic 2 is the image with a noise radius and the Synthetic 3 is an image with uniform noise. The fig. 2 shows results of multimodels  $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$ . All multimodels run 50 iterations of the EM-algorithm. Multimodels  $\mathbf{f}_2, \mathbf{f}_3$  approximate circles better than multimodel  $\mathbf{f}_1$ . Tab. 1 shows the approximation quality (6.1) for all multimodels.



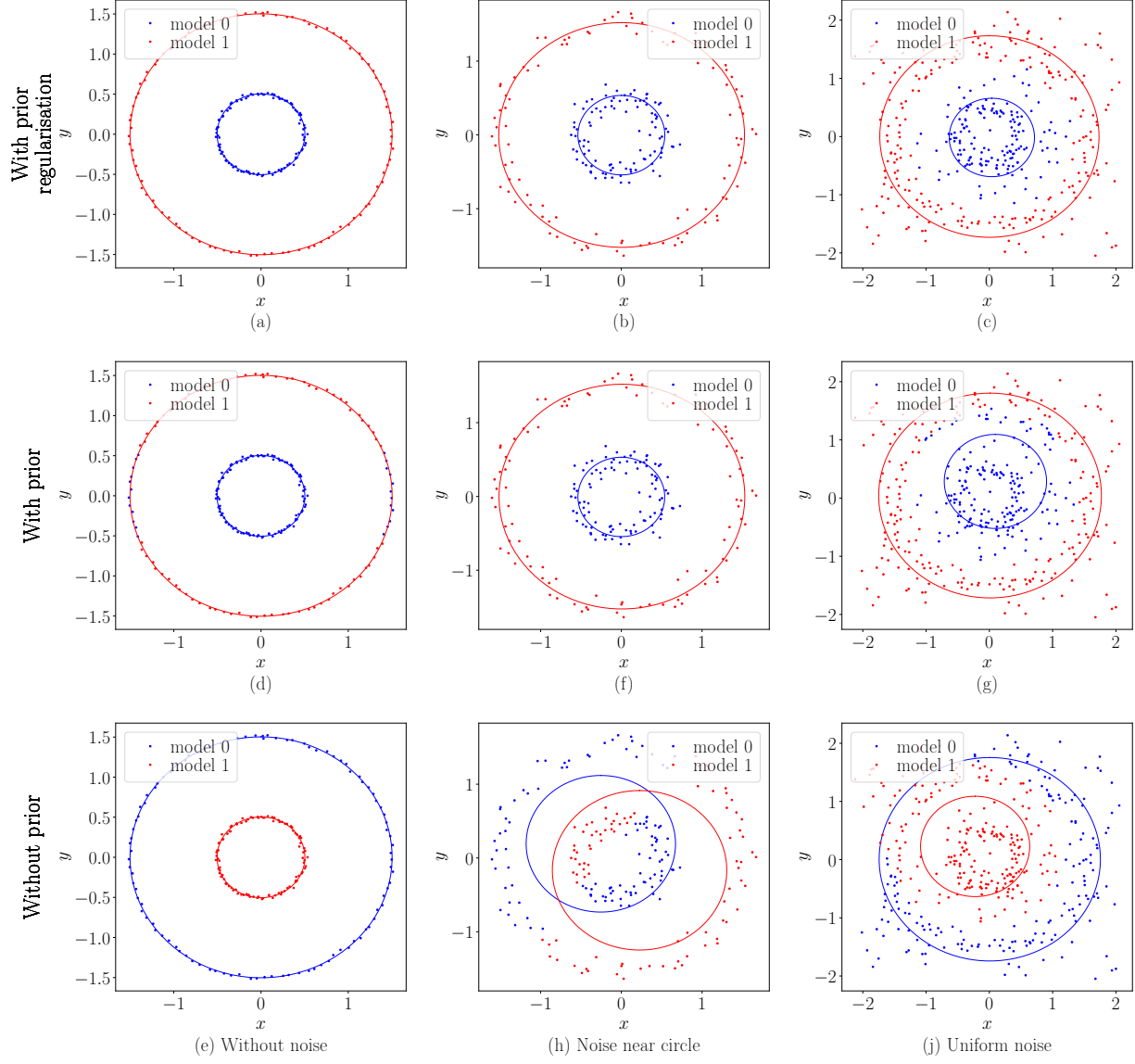


Figure 2: The multimodel depends on different prior distribution and depends on different noise levels: (a)–(c) multimodel with prior regularization; (d)–(g) multimodel with simple prior; (e)–(j) multimodel without any priors.

Table 1: The approximation quality (??) for all multimodels

Dataset	$\mathcal{S}_{f_1}$	$\mathcal{S}_{f_2}$	$\mathcal{S}_{f_3}$
Synthetic 1	$10^{-5}$	$10^{-5}$	$10^{-5}$
Synthetic 2	0.6	$10^{-3}$	$10^{-3}$
Synthetic 3	0.6	$10^{-3}$	$10^{-3}$

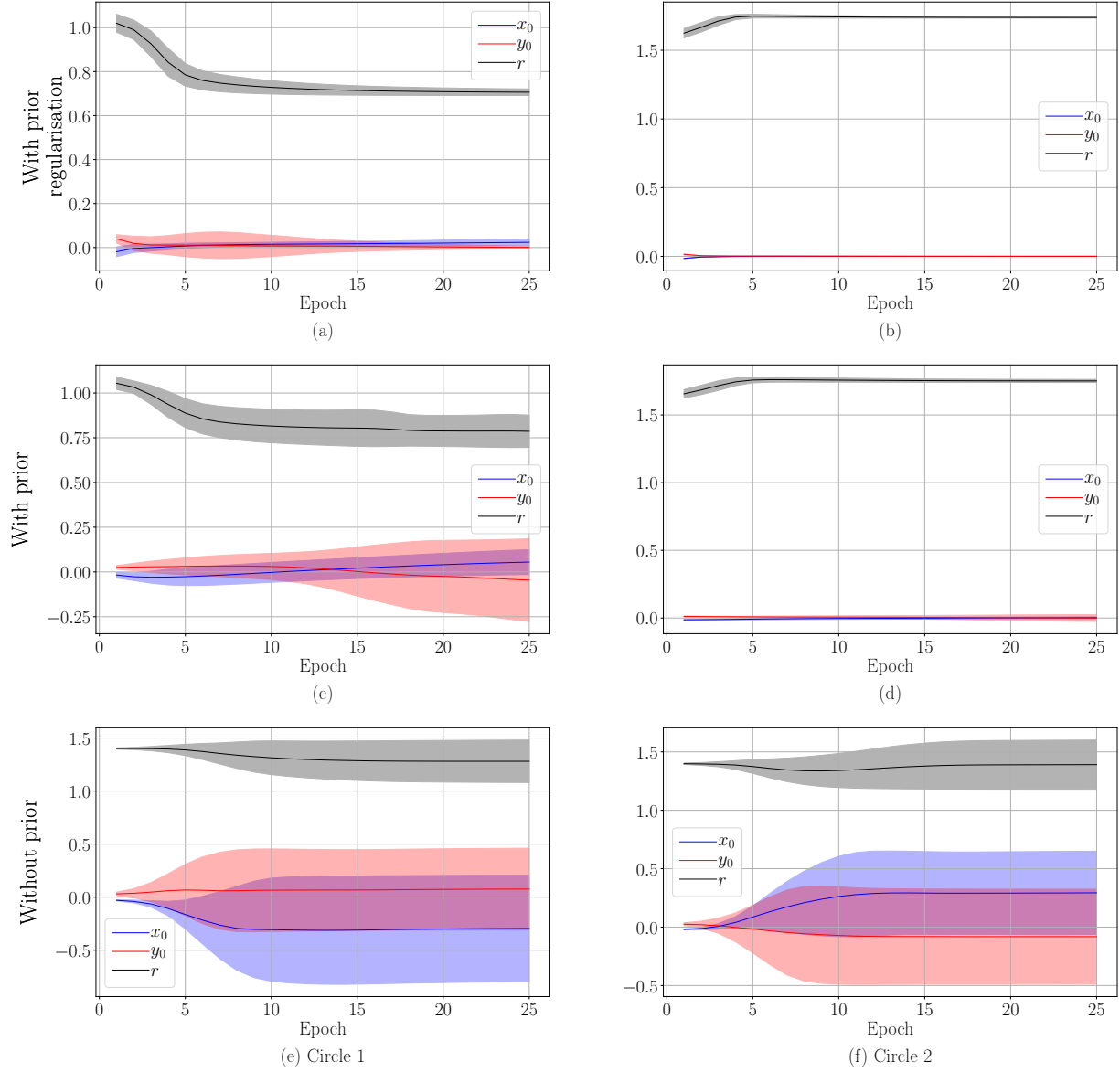


Figure 3: The dependence of center and radius on the iteration number: (a)–(b) multimodel with prior regularization; (c)–(d) multimodel with simple prior; (e)–(f) multimodel without any priors.

**Analysis of convergence on synthetic data.** The experiment analyzes multimodels  $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$  during the EM-algorithm convergence. The multimodels were analysed on the synthetic dataset Synthetic 3.

Fig. 3 shows the dependence of the radius and center on the EM-algorithm iteration number. The multimodel  $\mathbf{f}_2$  with prior distribution for the parameters approximates circles better than multimodel  $\mathbf{f}_1$  without any priors. The multimodel  $\mathbf{f}_3$  with prior regularisation approximates circles more stable than multimodel  $\mathbf{f}_2$ .

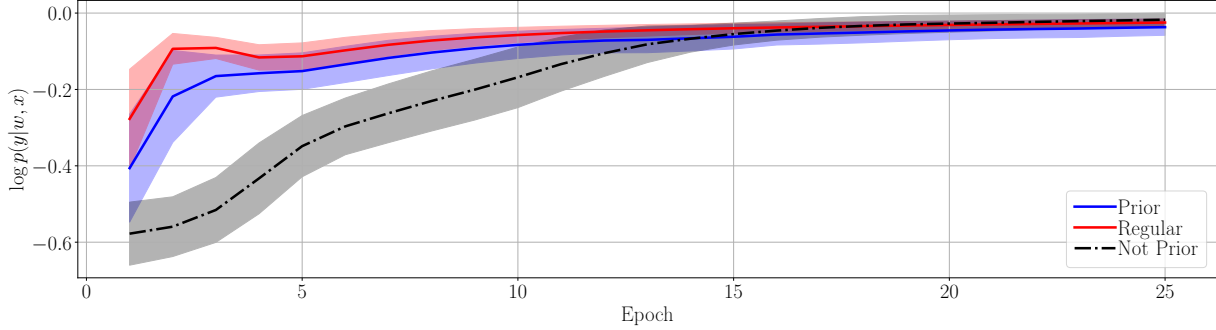


Figure 4: The dependence of log-likelihood (??) on the iteration number.

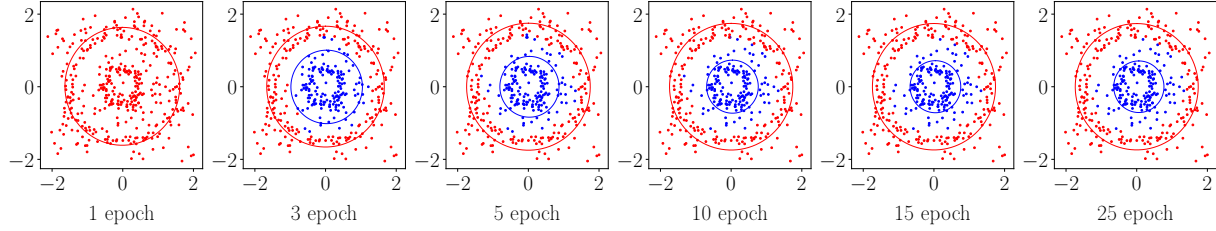


Figure 5: Visualization of convergence for the multimodel with prior regularisation.

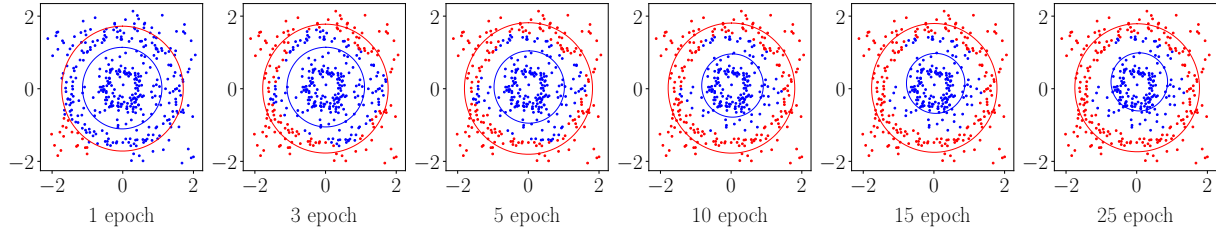


Figure 6: Visualization of convergence for the multimodel with simple prior.

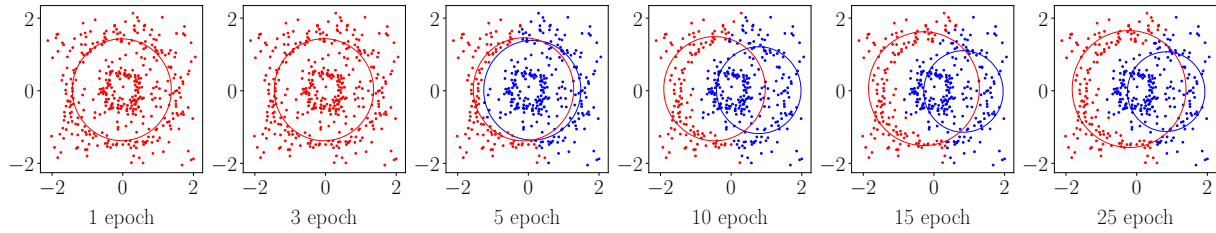


Figure 7: Visualization of convergence for the multimodel without any priors.

Fig. 4 shows the dependence of log-likelihood (??) on the EM-algorithm iteration number. The log-likelihood of multimodels  $\mathbf{f}_2, \mathbf{f}_3$  is growing faster than multimodel  $\mathbf{f}_1$ . After the 20-th iteration all three multimodels have the same log-likelihood.

Fig. 5-7 show learning process for multimodels  $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$ . Fig. 7 shows multimodel  $\mathbf{f}_1$ ,

which does not approximate circles correctly. Fig. 5-6 show multimodels  $\mathbf{f}_2, \mathbf{f}_3$ , which approximate circles correctly.

The experiment shows that multimodels  $\mathbf{f}_2, \mathbf{f}_3$  with prior distribution for the parameters approximate circles better than multimodel  $\mathbf{f}_1$  without any priors.

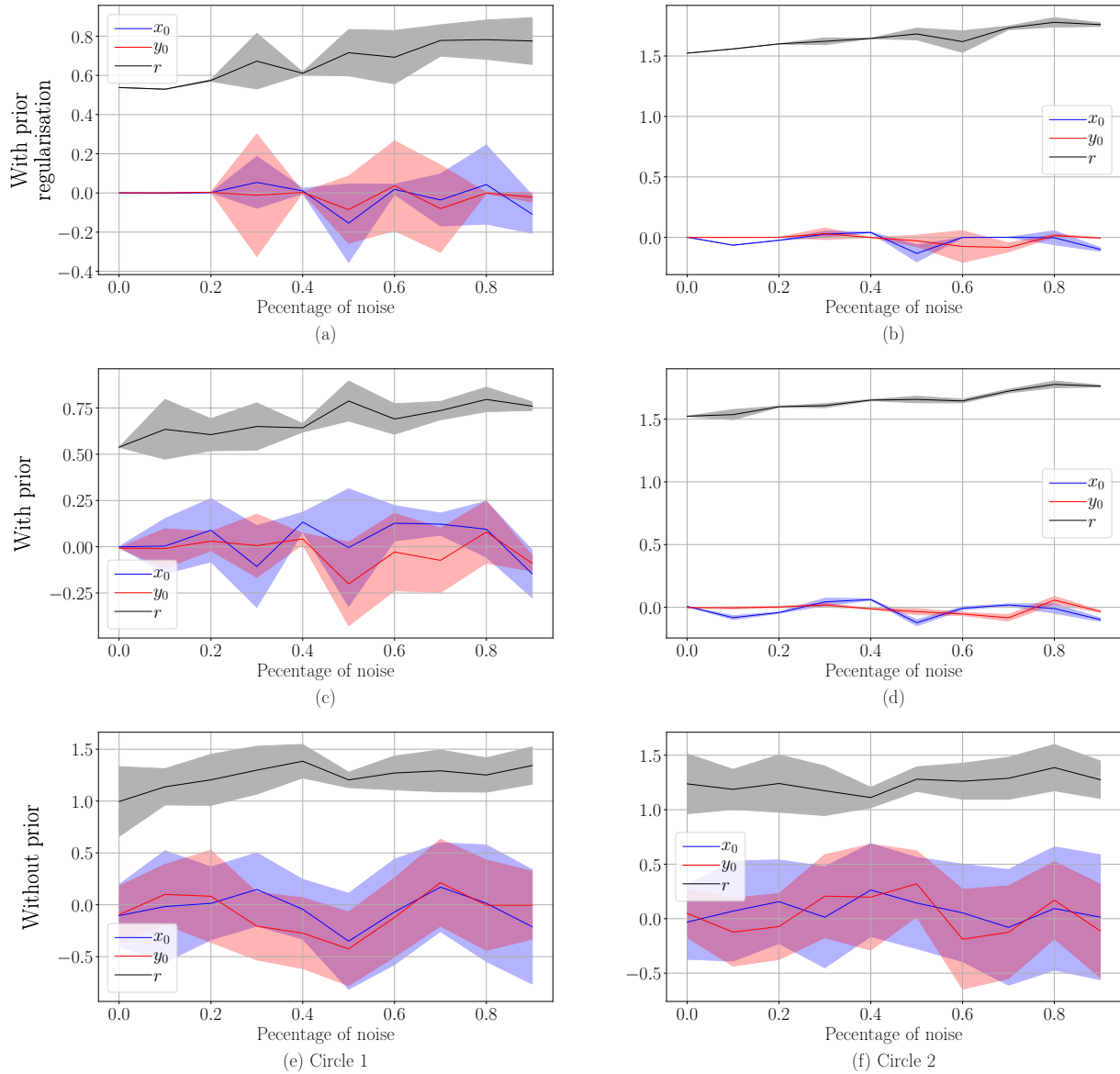


Figure 8: The dependence of center and radius on the noise level: (a)–(b) multimodel with prior regularization; (c)–(d) multimodel with simple prior; (e)–(f) multimodel without any priors.

**Multimodels analysis depending on the noise level.** The experiment analyzes dependence of multimodels  $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$  on the noise level. The multimodels were analysed on

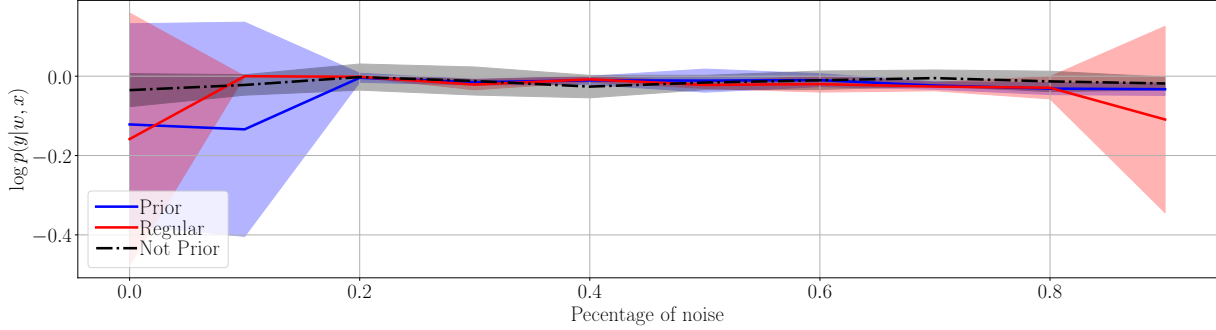


Figure 9: The dependence of log-likelihood (??) on the noise level.

the synthetic dataset Synthetic 1, with adding various noise level. The lowest noise level is equal to 0, when number of noise samples is equal to 0. The highest noise level is equal to 1, when number of noise samples is equal to number of point in circles. Fig. 8 shows the dependence of center and radius on the noise level. It shows that circle radius increases with increasing noise level. The multimodels  $\mathbf{f}_2, \mathbf{f}_3$  approximate circles center correctly, but the multimodel  $\mathbf{f}_3$  is more stable. Fig. 9 shows the dependence of log-likelihood (??) on the noise level. It shows that log-likelihood (??) is the same for all multimodels, but the fig. 8 shows that the approximation quality (??) depends on multimodels. This part of the experiment shows that multimodel  $\mathbf{f}_3$  with prior regularisation is the most stable model.

**Real data.** This part of the experiment analyzes multimodels  $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$  on the real data. Fig. 10 shows the result of different multimodels. The multimodel  $\mathbf{f}_1$  approximates incorrectly one of the circles. The multimodels  $\mathbf{f}_2, \mathbf{f}_3$  approximate both circles correctly.

Fig. 11-13 show learning process for multimodels  $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$ . Fig. 11 shows multimodel  $\mathbf{f}_1$ . Fig. 12 shows multimodel  $\mathbf{f}_2$ . Fig. 13 shows multimodel  $\mathbf{f}_3$ .

This part of experiment shows that multimodels  $\mathbf{f}_2, \mathbf{f}_3$  approximate circles better than multimodel  $\mathbf{f}_1$  even for the real images.

## 7 Conclusion

The paper compares multimodels with various prior distributions of model parameters. The computational experiment hold on the concentric circles with different noises. The linear models were used to approximate the circles in the image. The gate function is the two-layer fully connected neural network. The experiment compares the multimodel with the prior distribution and without it. The multimodel with prior distribution is more accurate in comparison to the multimodel without it. Another experiment compares different types of regularization. The experiment showed that multimodel with regularization is more stable. The experiment shows that all models in this article are sensitive to outlines. To solve this problem, the aauthors proposed to use a local model, which approximate noise.

The future work plannes to improve the multimodel by adding a prior distribution for

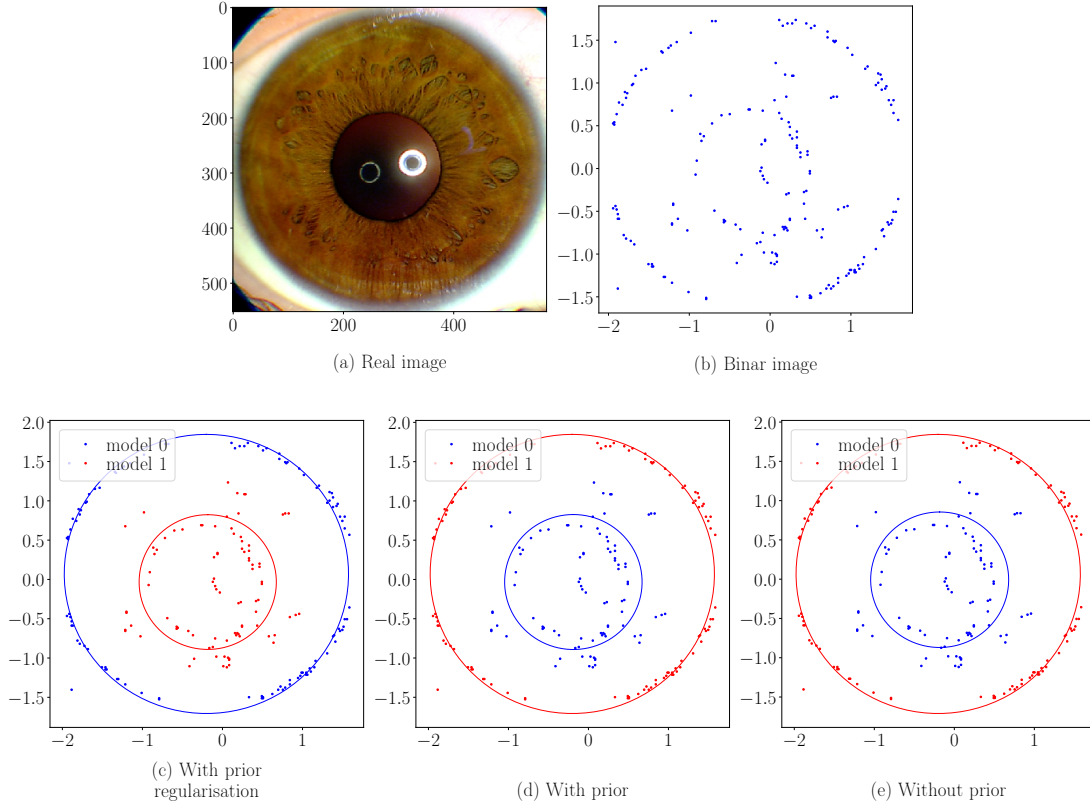


Figure 10: Different multimodels on the real image: (a) source image; (b) binarize image; (c) multimodel without any priors; (d) multimodel with simple prior; (e) multimodel with prior regularisation.

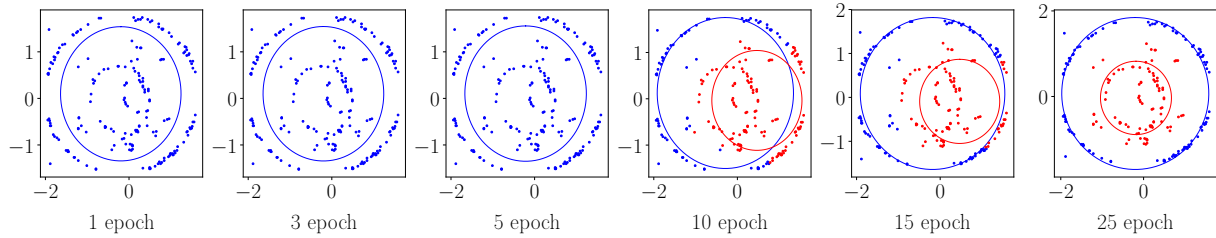


Figure 11: Visualization of convergence for the multimodel without any priors.

the gate function parameters. It plans to add a local model that approximates noise in the data. It assumes that the probability of noise is low.

## References

- [1] *Tianqi C., Carlos G.* XGBoost: A Scalable Tree Boosting System // Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data

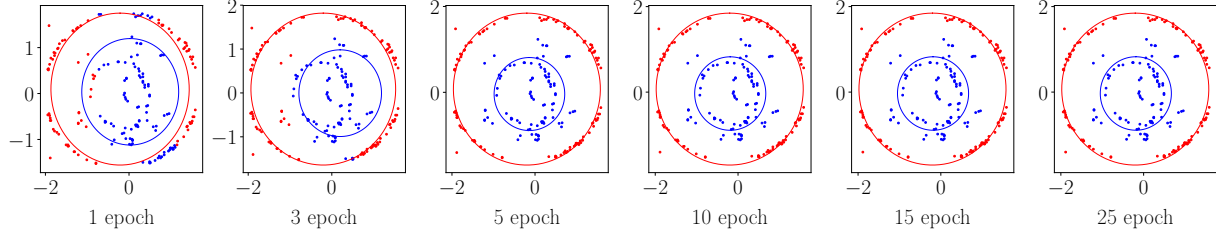


Figure 12: Visualization of convergence for the multimodel with simple prior.

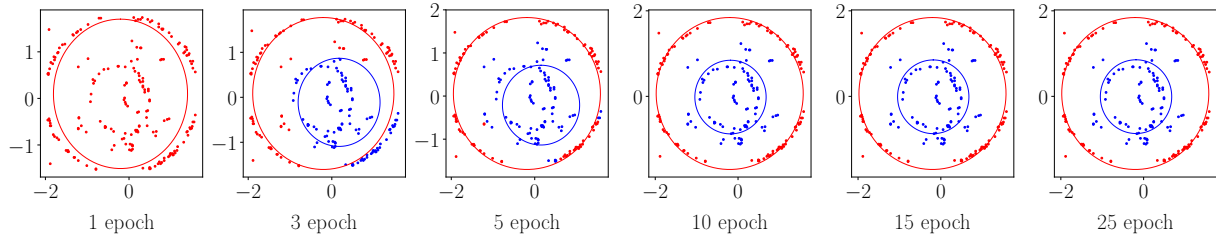


Figure 13: Visualization of convergence for the multimodel with prior regularisation.

Mining. 2016.

- [2] *Xi C., Hemant I.* Random Forests for Genomic Data Analysis // Genomics. 2012. Issues. 99. No. 6. P. 323–329.
- [3] *Esen Y.S., Wilson J., Gader P.D.* Twenty Years of Mixture of Experts // IEEE Transactions on Neural Networks and Learning Systems. 2012. Issues. 23. No. 8. P. 1177–1193.
- [4] *Rasmussen C.E., Ghahramani Z.* Infinite Mixtures of Gaussian Process Experts // Advances in Neural Information Processing Systems 14. 2002. P. 881–888.
- [5] *Shazeer N., Mirhoseini A., Maziarz K.* Outrageously large neural networks: the sparsely-gated mixture-of-experts layer // International Conference on Learning Representations. 2017.
- [6] *Jordan M.I.* Hierarchical mixtures of experts and the EM algorithm // Neural Comput. 1994. Vol. 6, No. 2. P. 181–214.
- [7] *Jordan M.I., Jacobs R.A.* Hierarchies of adaptive experts // In Advances in Neural Information Processing Systems. 1991. P. 985–992.
- [8] *Lima C., Coelho A., Zuben F.J.* Hybridizing mixtures of experts with support vector machines: Investigation into nonlinear dynamic systems identification // Inf. Sci. 2007. Vol. 177. No. 10. P. 2049–2074.
- [9] *Cao L.* Support vector machines experts for time series forecasting // Neurocomputing. 2003. Vol. 51. P. 321–339.

- [10] *Yumlu M. S., Gorgen F. S., Okay N.* Financial time series prediction using mixture of experts // In Proc. 18th Int. Symp. Comput. Inf. Sci. 2003. P. 553–560.
- [11] *Cheung Y. M., Leung W. M., Xu L.* Application of mixture of experts model to financial time series forecasting // On Proc. Int. Conf. Neural Netw. Signal Process. 1995. P. 1–4.
- [12] *Weigend A. S., Shi S.* Predicting daily probability distributions of S&P500 returns // J. Forecast. 2000. Vol. 19. No. 4. P. 375–392.
- [13] *Ebrahimipour R., Moradian M. R., Esmkhani A., Jafarlou F. M.* Recognition of Persian handwritten digits using characterization loci and mixture of experts // J. Digital Content Technol. Appl. 2009. Vol. 3. No. 3. P. 42–46.
- [14] *Estabrooks A., Japkowicz N.* A mixture-of-experts framework for text classification // In Proc. Workshop Comput. Natural Lang. Learn., Assoc. Comput. Linguist. 2001. P. 1–8.
- [15] *Mossavat S., Amft O., Petkov Vries B., Kleijn W.* A Bayesian hierarchical mixture of experts approach to estimate speech quality // In Proc. 2nd Int. Workshop Qual. Multimedia Exper. 2010. P. 200–205.
- [16] *Peng F., Jacobs R. A., Tanner M. A.* Bayesian inference in mixtures-of-experts and hierarchical mixtures-of-experts models with an application to speech recognition // J. Amer. Stat. Assoc. 1996. Vol. 91. No. 435. P. 953–960.
- [17] *Tuerk A.* The state based mixture of experts HMM with applications to the recognition of spontaneous speech. Ph.D. thesis. Cambridge: Univ. Cambridge, 2001.
- [18] *Sminchisescu C., Kanaujia A., Metaxas D.* Discriminative density propagation for visual tracking // IEEE Trans. Pattern Anal. Mach. Intell. 2007. Vol. 29. No. 11. P. 2030–2044.
- [19] *Bowyer K., Hollingsworth K., Flynn P.* A Survey of Iris Biometrics Research: 2008–2010.
- [20] *Matveev I.* Detection of iris in image by interrelated maxima of brightness gradient projections // Appl. Comput. Math. 2010. Vol. 9. No. 2. P. 252–257.
- [21] *Matveev I., Simonenko I.* Detecting precise iris boundaries by circular shortest path method // Pattern Recognition and Image Analysis. 2014. Vol. 24. P. 304–309.
- [22] *Dempster A. P., Laird N. M., Rubin D. B.* Maximum Likelihood from Incomplete Data via the EM Algorithm // Journal of the Royal Statistical Society. Series B (Methodological). 1977. Vol. 39. No. 1 P. 1–38.
- [23] *Bishop C.* Pattern Recognition and Machine Learning. Berlin: Springer, 2006. P. 758.