

Анализ выбора априорного распределения для смеси экспертов *

А. В. Грабовой¹, В. В. Стрижов²

Аннотация: Данная работа посвящен анализу свойств смеси экспертов в зависимости от выбора априорного распределения. Предлагается проанализировать свойства моделей в случае, когда выбрано информативное и не информативное априорное распределения на веса параметров каждого эксперта. В качестве экспертов рассматриваются линейные модели, а в качестве гипермодели рассматривается нейросеть с функцией softmax на последнем слое. В качестве базовой задачи рассматривается задача поиска окружностей на картинке. Предполагается, что каждой картинке соответствует свой эксперт. В качестве данных рассматриваются синтетически сгенерированные окружности с разным уровнем шума. Предлагается сравнить устойчивость к шуму мультимоделей с априорными знаниями и без них.

Ключевые слова: смесь экспертов; байесовский выбор модели; априорное распределение.

DOI: 00.00000/0000000000000000

1 Введение

В настоящее время в разных задач анализа данных присутствует необходимость обработки поступающих данных локально.

Мультимодели показывают отличные результаты во многих задачах. Классическими методами являются беггинг и градиентный бустинг [1], случайный лес [2]. Более современный подход к мультимоделированию [3] предполагает, что вклад каждой модели в ответ должен зависеть от конкретного объекта. Смесь экспертов базируется на понятии шлюзовой функции, которая определяет значимость предсказания каждого

*Работа выполнена при поддержке РФФИ и правительства РФ.

¹Московский физико-технический институт, grabovoy.av@phystech.edu

²Московский физико-технический институт, strijov@ccas.ru

эксперта — отдельной модели. В качестве шлюзовой функции выступает: softmax-регрессия, процесс Дирихле [5], нейронная сеть [4] с softmax на последнем слое.

Несмотря на значимые успехи мультимodelей, они имеют ряд недостатков. Данные недостатки связаны с тем, что сходимость локальных моделей сильно зависит от начальной инициализации параметров. Для улучшения сходимости предлагается использовать априорные знания о данных.

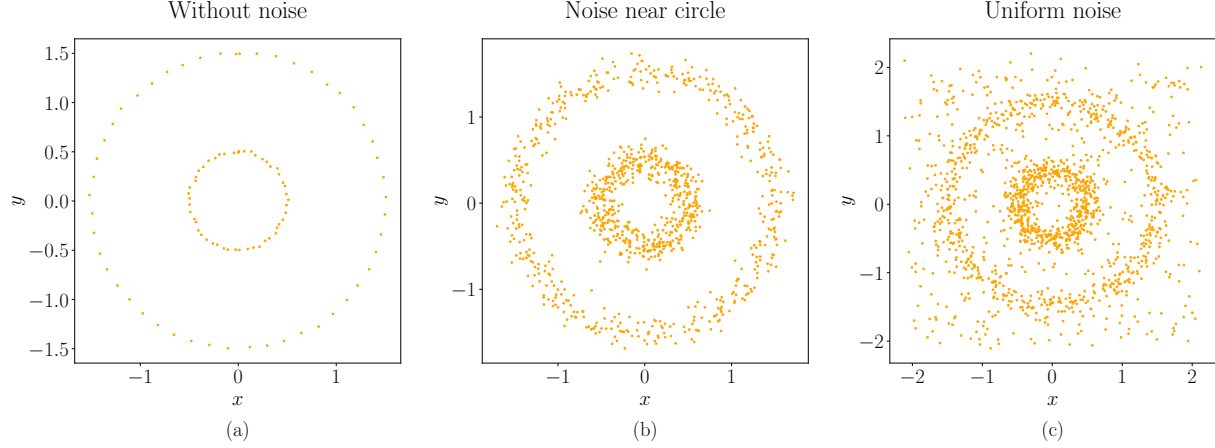


Рис. 1: Пример изображений с окружностями с разным уровнем шума: а) окружности без шума; б) окружности с зашумленным радиусом; в) окружности с зашумленным радиусом, а также с равномерным шумом по всему изображению

В данной работе предлагается исследовать, как влияют априорные знания на сходимость мультимodelей. Для данного исследования решается задача поиска окружностей на бинаризованной картинке — где каждый пиксель закрашен в один из двух цветов. Пример картинок показан на рис. 1. Предлагается рассмотреть как ведет себя модель с априорными знаниями и без них в случае картинок с разным уровнем шума. В данной работе в качестве отдельных экспертов рассматриваются линейные модели — каждая модель отвечает своей окружностью. В качестве шлюзовой функции рассматривается двухслойная нейронная сеть.

2 Постановка задачи нахождения параметров окружностей

Задано бинарное изображение:

$$\mathbf{M} \in \{0, 1\}^{m_1 \times m_2}, \quad (2.1)$$

где 0 отвечает белой точке — фону, 1 — черная точка изображения.

Используя изображения \mathbf{M} строится выборка \mathbf{C} , элементами которой являются координаты x_i, y_i черных точек на картинке:

$$\mathbf{C} \in \mathbb{R}^{N \times 2}, \quad (2.2)$$

где N — количество черных точек на изображении \mathbf{M} .

Пусть x_0, y_0 — центр окружности, которую нужно найти на бинарной картинке \mathbf{M} , а r ее радиус. Тогда элементы выборки \mathbf{C} должны удовлетворять уравнению окружности:

$$(x_i - x_0)^2 + (y_i - y_0)^2 = r^2. \quad (2.3)$$

Раскрыв скобки получим следующие уравнение:

$$(2x_0) \cdot x_i + (2y_0) \cdot y_i + (r^2 - x_0^2 - y_0^2) \cdot 1 = x_i^2 + y_i^2. \quad (2.4)$$

Получаем следующую задачу линейной регрессии для нахождения параметров окружности:

$$\mathbf{X}\mathbf{w} \approx \mathbf{Y}, \quad \mathbf{X} = \mathbf{C} \times \mathbf{1}, \quad \mathbf{Y} = \{x^2 + y^2 | \forall x, y \in \mathbf{C}\}, \quad (2.5)$$

где найденные параметры линейной регрессии $\mathbf{w} = \{w_1, w_2, w_3\}$ восстанавливают параметры окружности:

$$x_0 = \frac{w_1}{2}, \quad y_0 = \frac{w_2}{2}, \quad r = \sqrt{w_3 + x_0^2 + y_0^2}. \quad (2.6)$$

Данное решение позволяет искать параметры единственной окружности на рисунке. В случае, когда на картинке несколько окружностей предлагается использовать мультимодель, где в качестве каждой модели рассматривается единственная линейная модель, которая отвечает одной окружности на рисунке. В качестве мультимодели рассматривается смесь экспертов.

3 Постановка задачи смеси экспертов

Задана выборка:

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

где N — количество объектов в выборке, а n — размерность признакового пространства.

Определение 3.1. *Смесь экспертов — мультимодель, определяющая правдоподобие каждой π_k каждой модели \mathbf{f}_k на объекте \mathbf{x} на основе его признакового описания.*

$$\hat{\mathbf{f}} = \sum_{k=1}^K \pi_k \mathbf{f}_k, \quad \pi_k(\mathbf{x}, \mathbf{V}) : \mathbb{R}^{2 \times n} \rightarrow [0, 1], \quad \sum_{k=1}^K \pi_k(\mathbf{x}, \mathbf{V}) = 1 \quad (3.2)$$

где \mathbf{f} — мультимодель, а \mathbf{f}_k является некоторой моделью, π_k — параметрическая модель.

3.1 Общий случай

Правдоподобие модели:

$$p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) = \prod_{i=1}^N \prod_{k=1}^K [\pi_k(\mathbf{x}_i, \mathbf{V}) p_k(y_i|\mathbf{w}_k, \mathbf{x}_i)]^{z_{ik}} \prod_{k=1}^K p^k(\mathbf{w}_k|\mathbf{w}_k^0, \mathbf{A}_k) \quad (3.3)$$

Логарифм правдоподобия модели:

$$\begin{aligned} \log p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) &= \sum_{i=1}^N \sum_{k=1}^K z_{ik} [\log \pi_k(\mathbf{x}_i, \mathbf{V}) + \log p_k(y_i|\mathbf{w}_k, \mathbf{x}_i)] + \\ &+ \sum_{k=1}^K \log p^k(\mathbf{w}_k|\mathbf{w}_k^0, \mathbf{A}_k) \end{aligned} \quad (3.4)$$

E-step

Найдем $q(\mathbf{Z})$:

$$\begin{aligned} \log q(\mathbf{Z}) &= \mathbb{E}_{q/\mathbf{Z}} \log p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) \\ p(z_{ik} = 1) &= \frac{\exp(\log \pi_k(\mathbf{x}_i, \mathbf{V}) + \mathbb{E} \log p_k(y_i|\mathbf{w}_k, \mathbf{x}_i))}{\sum_{k'=1}^K \exp(\log \pi_{k'}(\mathbf{x}_i, \mathbf{V}) + \mathbb{E} \log p_{k'}(y_i|\mathbf{w}'_{k'}, \mathbf{A}'_k))} \end{aligned} \quad (3.5)$$

Найдем $q(\mathbf{W})$:

$$\begin{aligned} \log q(\mathbf{W}) &= \mathbb{E}_{q/\mathbf{W}} \log p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) \\ &= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E} z_{ik} [\log \pi_k(\mathbf{x}_i, \mathbf{V}) + \log p_k(y_i|\mathbf{w}_k, \mathbf{x}_i)] + \\ &+ \sum_{k=1}^K \log p^k(\mathbf{w}_k|\mathbf{w}_k^0, \mathbf{A}_k) \end{aligned} \quad (3.6)$$

M-step

$$\begin{aligned} \mathbb{E}_q p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) &= \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}) \\ \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}) &= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E} z_{ik} [\log \pi_k(\mathbf{x}_i, \mathbf{V}) + \mathbb{E} \log p_k(y_i|\mathbf{w}_k, \mathbf{x}_i)] \\ &+ \sum_{k=1}^K \mathbb{E} \log p_k(\mathbf{w}_k|\mathbf{w}_k^0, \mathbf{A}_k) \end{aligned} \quad (3.7)$$

Найдем \mathbf{A} из условия:

$$\frac{\partial \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A})}{\partial \mathbf{A}^{-1}} = 0 \quad (3.8)$$

Найдем \mathbf{V} из условия:

$$\frac{\partial \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A})}{\partial \mathbf{V}} = 0 \quad (3.9)$$

Найдем \mathbf{W}^0 из условия:

$$\frac{\partial \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A})}{\partial \mathbf{W}^0} = 0 \quad (3.10)$$

3.2 Случай линейной регрессионной модели

Рассмотрим следующий случай распределений:

1. $p_k(y_i | \mathbf{w}_k, \mathbf{x}_i) = \mathcal{N}(y_i | \mathbf{w}_k^\top \mathbf{x}_i, \beta^{-1})$
2. $p^k(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k) = \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k)$
3. $\pi(\mathbf{x}_i, \mathbf{V}) = \text{softmax}(\mathbf{F}(\mathbf{x}_i, \mathbf{V}))$, где $F : \mathbb{R}^n \times \mathbb{R}^V \rightarrow \mathbb{R}^K$ — нейросеть, V — число параметров нейросети.

Правдоподобие модели:

$$p(\mathbf{y}, \mathbf{W}, \mathbf{Z} | \mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) = \prod_{i=1}^N \prod_{k=1}^K [\pi_k(\mathbf{x}_i, \mathbf{V}) \mathcal{N}(y_i | \mathbf{w}_k, \mathbf{x}_i)]^{z_{ik}} \prod_{k=1}^K \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k) \quad (3.11)$$

Логарифм правдоподобия модели:

$$\begin{aligned} \log p(\mathbf{y}, \mathbf{W}, \mathbf{Z} | \mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) &= \sum_{i=1}^N \sum_{k=1}^K z_{ik} \left[\log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (y_i - \mathbf{w}_k^\top \mathbf{x}_i)^2 + \frac{1}{2} \log \frac{\beta}{2\pi} \right] + \\ &+ \sum_{k=1}^K \left[-\frac{1}{2} (\mathbf{w}_k - \mathbf{w}_k^0)^\top \mathbf{A}_k^{-1} (\mathbf{w}_k - \mathbf{w}_k^0) + \frac{1}{2} \log \det \mathbf{A}_k^{-1} - \frac{n}{2} \log 2\pi \right] \end{aligned} \quad (3.12)$$

E-step

Найдем $q(\mathbf{Z})$:

$$\begin{aligned} \log q(\mathbf{Z}) &= \mathbb{E}_{q(\mathbf{Z})} \log p(\mathbf{y}, \mathbf{W}, \mathbf{Z} | \mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) \propto \\ &\propto \sum_{i=1}^N \sum_{k=1}^K z_{ik} \left[\log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (y_i^2 - \mathbf{x}_i^\top \mathbb{E} \mathbf{w}_k + \mathbf{x}_i^\top \mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i) + \frac{1}{2} \log \frac{\beta}{2\pi} \right] \\ p(z_{ik} = 1) &= \frac{\exp(\log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (\mathbf{x}_i^\top \mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbb{T} \mathbf{w}_k))}{\sum_{k'=1}^K \exp(\log \pi_{k'}(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (\mathbf{x}_i^\top \mathbb{E} \mathbf{w}_{k'} \mathbf{w}_{k'}^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbb{T} \mathbf{w}_{k'}))} \end{aligned} \quad (3.13)$$

Найдем $q(\mathbf{W})$:

$$\begin{aligned}
\log q(\mathbf{W}) &= \mathbb{E}_{q/\mathbf{W}} \log p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) \\
&= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}_{z_{ik}} \left[\log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (y_i - \mathbf{w}_k^\top \mathbf{x}_i)^2 + \frac{1}{2} \log \frac{\beta}{2\pi} \right] + \\
&+ \sum_{k=1}^K \left[-\frac{1}{2} (\mathbf{w}_k - \mathbf{w}_k^0)^\top \mathbf{A}_k^{-1} (\mathbf{w}_k - \mathbf{w}_k^0) + \frac{1}{2} \log \det \mathbf{A}_k^{-1} - \frac{n}{2} \log 2\pi \right] \\
&= \sum_{k=1}^K \left[\mathbf{w}_k^\top \left(\mathbf{A}_k^{-1} \mathbf{w}_k^0 + \beta \sum_{i=1}^N \mathbf{x}_i y_i \mathbb{E}_{z_{ik}} \right) - \frac{1}{2} \mathbf{w}_k^\top \left(\mathbf{A}_k^{-1} + \beta \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{w}_k \right]
\end{aligned} \tag{3.14}$$

Получаем распределение параметров:

$$q(\mathbf{w}_k) = \mathcal{N}(\mathbf{w}_k | \mathbf{m}_k, \mathbf{B}_k), \tag{3.15}$$

где введены обозначения

$$\mathbf{m}_k = \mathbf{B}_k \left(\mathbf{A}_k^{-1} \mathbf{w}_k^0 + \beta \sum_{i=1}^N \mathbf{x}_i y_i \mathbb{E}_{z_{ik}} \right) \quad \mathbf{B}_k = \left(\mathbf{A}_k^{-1} + \beta \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \tag{3.16}$$

M-step

$$\begin{aligned}
\mathbb{E}_q \log p(\mathbf{y}, \mathbf{W}, \mathbf{Z}|\mathbf{X}, \mathbf{V}, \mathbf{W}^0, \mathbf{A}) &= \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}) \\
\mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}) &= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}_{z_{ik}} \left[\log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} \mathbb{E} (y_i - \mathbf{w}_k^\top \mathbf{x}_i)^2 + \frac{1}{2} \log \frac{\beta}{2\pi} \right] + \\
&+ \sum_{k=1}^K \left[-\frac{1}{2} \mathbb{E} (\mathbf{w}_k - \mathbf{w}_k^0)^\top \mathbf{A}_k^{-1} (\mathbf{w}_k - \mathbf{w}_k^0) + \frac{1}{2} \log \det \mathbf{A}_k^{-1} - \frac{n}{2} \log 2\pi \right]
\end{aligned} \tag{3.17}$$

Найдем \mathbf{A} из условия:

$$\frac{\partial \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A}_k)}{\partial \mathbf{A}^{-1}} = \frac{1}{2} \mathbf{A}_k - \frac{1}{2} \mathbb{E} (\mathbf{w}_k - \mathbf{w}_k^0) (\mathbf{w}_k - \mathbf{w}_k^0)^\top = 0 \tag{3.18}$$

Получаем \mathbf{A} :

$$\mathbf{A}_k = \mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top - 2 \mathbf{w}_k^0 \mathbb{E} \mathbf{w}_k^\top + \mathbf{w}_k^0 \mathbf{w}_k^{0\top} \tag{3.19}$$

Найдем \mathbf{V} :

Аналитически решение не ищется, поэтому воспользуемся градиентным спуском для максимизации правдоподобия модели:

$$\mathbf{V}^{j+1} = \mathbf{V}^j + \alpha \frac{\partial \mathcal{F}(\mathbf{W}, \mathbf{V}^j, \beta)}{\partial \mathbf{V}} \tag{3.20}$$

Найдем \mathbf{W}^0 из условия:

$$\frac{\partial \mathcal{F}(\mathbf{V}, \mathbf{W}^0, \mathbf{A})}{\partial \mathbf{w}_k^0} = -2\mathbf{A}_k^{-1} \mathbf{E} \mathbf{w}_k + 2\mathbf{A}_k^{-1} \mathbf{w}_k^0 = 0 \quad (3.21)$$

Получаем \mathbf{W}^0 :

$$\mathbf{w}_k^0 = \mathbf{E} \mathbf{w}_k \quad (3.22)$$

4 Вычислительный эксперимент

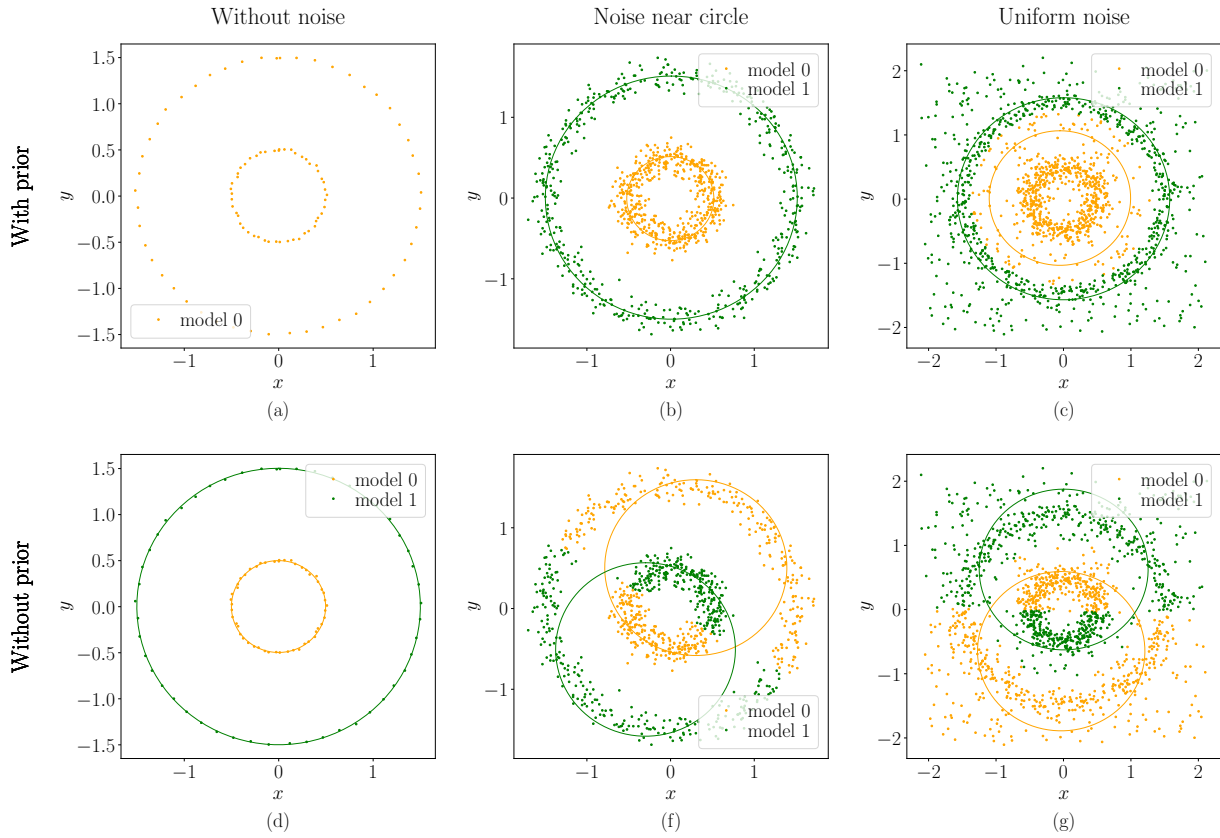


Рис. 2: Результат работы мультимодели в зависимости от априорных знаний и в зависимости от уровня шума: а) модель с априорными знаниями, окружности без шума; б) модель с априорными знаниями, окружности с зашумленным радиусом; в) модель с априорными знаниями, окружности с зашумленным радиусом, а также с равномерным шумом по всему изображению; д) модель без априорных знаний, окружности без шума; е) модель без априорных знаний, окружности с зашумленным радиусом; ф) модель без априорных знаний, окружности с зашумленным радиусом, а также с равномерным шумом по всему изображению

Вычислительный эксперимент проводится на синтетической выборке, которая получена при помощи генерации двух концентрических окружностей с разным уровнем шума.

Предлагается сравнить две разные постановки задачи смеси экспертов: в случае когда используются априорные знания об картинке и в случае когда априорные знания отсутствуют.

В качестве информативного априорного знания выступает, то, что предполагается что вектора параметров каждой модели имеют следующие распределения:

$$N(\mathbf{w}_1 | \mathbf{w}_1^0, \mathbf{I}), \quad N(\mathbf{w}_2 | \mathbf{w}_2^0, \mathbf{I}), \quad (4.1)$$

где $\mathbf{w}_1^0 = [0, 0, 0.1]$, $\mathbf{w}_2^0 = [0, 0, 5]$, что указывает, на то, что известно о концентричности окружностей, а также что у них радиусы различны.

Таблица 1: Результаты работы мультимodelей

Мультимodelь	Without Noise	Noise near circle	Uniform noise
With prior	99/100	97/100	95/100
Without prior	68/100	23/100	7/100

На рис. 2 показан случайный результаты работы мультимodelей с априорными знаниями и без них. На всех картинках обе модели работала 30 итераций. Так-как сходимость мультимodelей очень сильно зависит от начальной инициализации, также был проведен эксперимент с множественным запуском мультимodelей на одной и той же картинке. Обе модели на каждой картинке запускались по 100 раз. В таб. 1 показано сколько мультимodelей правильно отыскивали обе окружности на рисунке. Как видно мультимodelь с использованием априорных знаний является более стабильной, чем мультимodelь, которая не использует никаких априорных знаний.

5 Заключение

В данной работе проведено сравнение мультимodelей в случае, когда каждая модель имела априорные знания и в случае, когда априорных знаний не было. В качестве данных использовались изображения концентрических окружностей с разным уровнем шума. Для поиска окружностей использовались линейные модели. В качестве шлюзовой функции использовалась двухслойная нейросеть.

Как показано в эксперименте в случае, когда введены априорные знания на линейные модели, мультимodelь является более устойчивой к шуму. Также в случае задания априорных знаний моделей, мультимodelь менее зависит от начальной инициализации, что также позволяет сказать, что модель является более устойчивой к начальной инициализации.

В дальнейшем планируется улучшить мультимodelь при помощи задания априорных знаний на шлюзовую функцию.

Список литературы

- [1] *Chen Tianqi, Guestrin Carlos* XGBoost: A Scalable Tree Boosting System // KDD '16 Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2016.
- [2] *Chen Xi, Ishwaran Hemant* Random Forests for Genomic Data Analysis // Genomics. 2012. Issues. 99, No 6. pp. 323–329.
- [3] *Yuksel Seniha Esen, Wilson Joseph N., Gader Paul D* Twenty Years of Mixture of Experts // IEEE Transactions on Neural Networks and Learning Systems. 2012. Issues. 23, No 8. pp. 1177–1193.
- [4] *Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz* Outrageously large neural networks: the sparsely-gated mixture-of-experts layer // ICLR, 2017.
- [5] *Rasmussen Carl Edward, Ghahramani Zoubin* Infinite Mixtures of Gaussian Process Experts // Advances in Neural Information Processing Systems 14. 2002. pp. 881–888.