

# Эмпирические байесовские нейронные сети

Басов Дмитрий Константинович

## Аннотация

Данная статья посвящена применению техники эмпирического Байеса к байесовским нейронным сетям. Концептуально идея следующая:

1. Мы используем диагональное нормальное распределение для аппроксимации апостериорного распределения весов модели —  $q(W)$ .
2. Априорное распределение весов модели так же задаётся диагональным распределением с нулевым матожиданием —  $p(W)$ .
3. Используя вариационный вывод, мы приходим к ситуации, когда ELBO зависит от  $KL(q(W)||p(W))$ . Так как оба распределения являются нормальными, то KL дивергенция считается аналитически.
4. Мотивация следующего этапа была взята из RVM — взять дисперсию априорного распределения весов модели  $p(W)$  из данных. Там несложно берётся производная и всё получается красиво, кроме возможного деления 0/0. Но сделав замену переменных, от этой беды можно уйти.

Пункты 1–3 в принципе были описаны в статье [Weight Uncertainty in Neural Networks](#). А вот четвёртый пункт я ни в книгах, ни в статьях не находил.

## 1 Обозначения и сокращения

$N(\mu, \sigma^2)$  — нормальное распределение

$\mathbf{x} \odot \mathbf{y}$  — поэлементное произведение (произведение Адамара) векторов

$\mathcal{L}$  — Evidence Lower Bound (ELBO)

$KL(q||p) = \int q(\mathbf{Z}) \cdot \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z})} d\mathbf{Z}$  — дивергенция Кульбака–Лейблера

$\mathbf{x}$  — вектор признаков

$\mathbf{y}$  — вектор целевой переменной

$D$  — датасет — пары значений  $\{\mathbf{x}_i, \mathbf{y}_i\}$ , где  $i = 1, \dots, L$

$\mathbf{W}$  — веса модели — случайная величина размерности  $M$

$p(D|\mathbf{W}) = \prod_{i=1}^L p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{W})$  — правдоподобие (likelihood)

$p(\mathbf{W})$  — априорное распределение весов модели (prior)

$p(\mathbf{W}|D)$  — апостериорное распределение весов модели (posterior)

$p(D)$  — маргинальная вероятность датасета (evidence)

$p(\mathbf{W}, D) = p(D|\mathbf{W}) \cdot p(\mathbf{W}) = p(\mathbf{W}|D) \cdot p(D)$  — совместная вероятность весов модели и данных

$q(\mathbf{W}|\theta)$  — аппроксимация апостериорного распределения весов модели

$\theta$  — обучаемые параметры байесовской модели

## 2 Введение

В классическом машинном обучении делается следующее предположение: веса модели  $\mathbf{W}$  являются пусть и неизвестной, но фиксированной величиной. В этом случае можно получить точечную оценку весов модели согласно гипотезе максимального правдоподобия.

$$\mathbf{W}_{\text{ML}} = \underset{\mathbf{W}}{\operatorname{argmax}} p(D|\mathbf{W})$$

Тогда распределение  $p(\mathbf{y}|\mathbf{x}, D)$  аппроксимируется следующим образом:

$$p(\mathbf{y}|\mathbf{x}, D) \approx p(\mathbf{y}|\mathbf{W}_{\text{ML}}, \mathbf{x})$$

Однако это справедливо при условии, что количество объектов в датасете  $D$  сильно больше количества весов модели ( $L \gg M$ ). Если это не так, веса модели  $\mathbf{W}$  могут слишком сильно подстроиться под обучающую выборку  $D$ , что чревато переобучением.

Для борьбы с переобучением используется ряд приёмов (штрафы на норму весов, early stopping, dropout), однако для их настройки требуются вычислительные ресурсы и отложенные (не участвующие в обучении) выборки данных.

Альтернативным подходом к машинному обучению является нахождение апостериорного распределения весов модели  $p(\mathbf{W}|D)$  по теореме Байеса.

$$p(\mathbf{W}|D) = \frac{p(D|\mathbf{W}) \cdot p(\mathbf{W})}{\int p(D|\mathbf{W}) \cdot p(\mathbf{W}) d\mathbf{W}}$$

Байесовские модели машинного обучения устойчивы к переобучению, так как о размере обучающей выборки не делается никаких предположений. Однако использование байесовского машинного обучения сопряжено с двумя большими проблемами: выбором подходящего априорного распределения  $p(\mathbf{W})$  и вычислением апостериорного распределения  $p(\mathbf{W}|D)$ .

Неудачный выбор  $p(\mathbf{W})$  может сильно ухудшить качество модели, а расчёт апостериорного распределения  $p(\mathbf{W}|D)$  требует вычисления интеграла по всему пространству весов модели, что для нейронных сетей практически невозможно.

В данной статье предлагается следующий подход к решению этих проблем.

1. Для аппроксимации распределения  $p(\mathbf{W}|D)$  применяется техника вариационного вывода. В этом случае вводится распределение  $q(\mathbf{W}|\boldsymbol{\theta})$ , и задача сводится к максимизации нижней вариационной границы (ELBO)  $\mathcal{L}(q(\mathbf{W}|\boldsymbol{\theta}))$ .
2. Распределение  $q(\mathbf{W}|\boldsymbol{\theta})$  задаётся в виде нормального распределения с диагональной матрицей ковариации. То есть каждый вес модели определяется двумя числами, которые определяют его математическое ожидание и дисперсию. Таким образом, количество обучаемых параметров относительно классической нейронной сети возрастает в 2 раза.
3. Так как распределение  $q(\mathbf{W}|\boldsymbol{\theta})$  является нормальным, применяя трюк с репараметризацией, становится возможным использовать градиентные методы для максимизации  $\mathcal{L}(q(\mathbf{W}|\boldsymbol{\theta}))$ .
4. Априорное распределение весов модели  $p(\mathbf{W})$  задаётся в виде нормального распределения с нулевым матожиданием (из соображений симметрии) и диагональной матрицей ковариации. То есть мы вносим следующее априорное знание: веса модели находятся около нуля.
5. Для определения дисперсии априорного распределения весов модели  $p(\mathbf{W})$  применяется техника эмпирического Байеса. То есть значение дисперсии  $p(\mathbf{W})$  берётся из датасета  $D$ . Так как распределения  $q(\mathbf{W}|\boldsymbol{\theta})$  и  $p(\mathbf{W})$  являются нормальными, оптимальное значение дисперсии распределения  $p(\mathbf{W})$  определяется аналитически. Такой подход обладает большой универсальностью, однако из-за этого теряется теоретическая устойчивость к переобучению.
6. Вводятся новые параметры  $\boldsymbol{\gamma}$  и  $\boldsymbol{\rho}$ , через которые выражаются матожидание и дисперсия распределения  $q(\mathbf{W}|\boldsymbol{\theta})$ . Это нужно, чтобы избежать неопределённости деления  $\frac{0}{0}$ , которая может возникнуть из-за определения дисперсии распределения  $p(\mathbf{W})$  из данных.

Эксперименты показали, что предлагаемый подход сохраняет гибкость классических нейронных сетей и делает их устойчивыми к переобучению.

### 3 Постановка задачи

Задача машинного обучения с учителем в вероятностной постановке формулируется следующим образом: получить распределение вероятностей  $p(\mathbf{y}|\mathbf{x}, D)$  целевой переменной  $\mathbf{y}$  для неразмеченных  $\mathbf{x}$ , используя информацию из датасета  $D$ . В случае параметрических моделей, которыми являются нейронные сети, информация из датасета  $D$  кодируется посредством весов модели  $\mathbf{W}$ . Сделаем следующие преобразования:

$$p(\mathbf{y}|\mathbf{x}, D) = \int p(\mathbf{y}, \mathbf{W}|\mathbf{x}, D) d\mathbf{W} = \int p(\mathbf{y}|\mathbf{W}, \mathbf{x}, D) \cdot p(\mathbf{W}|\mathbf{x}, D) d\mathbf{W} = \int p(\mathbf{y}|\mathbf{W}, \mathbf{x}) \cdot p(\mathbf{W}|D) d\mathbf{W}$$

Пояснения:

- $p(\mathbf{y}|\mathbf{x}, D) = \int p(\mathbf{y}, \mathbf{W}|\mathbf{x}, D) d\mathbf{W}$ , так как для любых случайных величин  $\mathbf{a}$  и  $\mathbf{b}$  справедливо  $p(\mathbf{a}) = \int p(\mathbf{a}, \mathbf{b}) d\mathbf{b}$
- $p(\mathbf{y}, \mathbf{W}|\mathbf{x}, D) = p(\mathbf{y}|\mathbf{W}, \mathbf{x}, D) \cdot p(\mathbf{W}|\mathbf{x}, D)$ , так как для любых случайных величин  $\mathbf{a}$  и  $\mathbf{b}$  справедливо  $p(\mathbf{a}, \mathbf{b}) = p(\mathbf{a}|\mathbf{b}) \cdot p(\mathbf{b})$
- $p(\mathbf{y}|\mathbf{W}, \mathbf{x}, D) = p(\mathbf{y}|\mathbf{W}, \mathbf{x})$ , так как вся информация из датасета  $D$  отражена в весах  $\mathbf{W}$
- $p(\mathbf{W}|\mathbf{x}, D) = p(\mathbf{W}|D)$ , так как веса модели  $\mathbf{W}$  не зависят от неразмеченных  $\mathbf{x}$ , которых не было в датасете  $D$ .

Получим выражение для  $p(\mathbf{W}|D)$ , используя формулу Байеса:

$$p(\mathbf{W}|D) = \frac{p(\mathbf{W}, D)}{p(D)} = \frac{p(\mathbf{W}, D)}{\int p(\mathbf{W}, D) d\mathbf{W}} = \frac{p(D|\mathbf{W}) \cdot p(\mathbf{W})}{\int p(D|\mathbf{W}) \cdot p(\mathbf{W}) d\mathbf{W}}$$

Для аппроксимации распределения ответов модели можно воспользоваться методом Монте-Карло. Идея следующая: сэмплируем конечное количество весов  $\hat{\mathbf{W}}_1, \dots, \hat{\mathbf{W}}_T$  из распределения  $p(\mathbf{W}|D)$  и аппроксимируем распределение  $p(\mathbf{y}|\mathbf{x}, D)$  следующим образом:

$$p(\mathbf{y}|\mathbf{x}, D) \approx \frac{1}{T} \sum_{t=1}^T p(\mathbf{y}|\hat{\mathbf{W}}_t, \mathbf{x}), \text{ где } \hat{\mathbf{W}}_t \text{ — сэмпл весов модели из } p(\mathbf{W}|D)$$

Получить аналитическое решение интеграла  $\int p(D|\mathbf{W}) \cdot p(\mathbf{W}) d\mathbf{W}$  можно только в очень ограниченном числе случаев. Существует возможность сэмплировать из  $p(\mathbf{W}|D)$ , используя методы Монте-Карло для марковских цепей (MCMC). Однако для больших датасетов и большого числа весов это практически невозможно. Альтернативным подходом к решению такой задачи является вариационный вывод — аппроксимация распределения  $p(\mathbf{W}|D)$  распределением  $q(\mathbf{W}|\boldsymbol{\theta})$ , из которого сэмплировать намного проще.

### 4 Вариационный вывод

Идея вариационного вывода — сведение задачи байесовского вывода к задаче максимизации нижней вариационной границы (ELBO)  $\mathcal{L}$ , которая для распределения  $q(\mathbf{W}|\boldsymbol{\theta})$  записывается следующим образом:

$$\mathcal{L}(q(\mathbf{W}|\boldsymbol{\theta})) = \int q(\mathbf{W}|\boldsymbol{\theta}) \cdot \ln \frac{p(\mathbf{D}, \mathbf{W})}{q(\mathbf{W}|\boldsymbol{\theta})} d\mathbf{W}$$

Покажем мотивацию максимизации ELBO. Запишем выражение для  $KL(q(\mathbf{W}|\boldsymbol{\theta})||p(\mathbf{W}|D))$  и преобразуем его, используя тождество  $p(\mathbf{W}, D) = p(\mathbf{W}|D) \cdot p(D)$ :

$$\begin{aligned} KL(q(\mathbf{W}|\boldsymbol{\theta})||p(\mathbf{W}|D)) &= \int q(\mathbf{W}|\boldsymbol{\theta}) \cdot \ln \frac{q(\mathbf{W}|\boldsymbol{\theta})}{p(\mathbf{W}|D)} d\mathbf{W} = \int q(\mathbf{W}|\boldsymbol{\theta}) \cdot \ln \frac{p(D) \cdot q(\mathbf{W}|\boldsymbol{\theta})}{p(D, \mathbf{W})} d\mathbf{W} = \\ &= \ln p(D) \cdot \int q(\mathbf{W}|\boldsymbol{\theta}) d\mathbf{W} - \int q(\mathbf{W}|\boldsymbol{\theta}) \cdot \ln \frac{p(\mathbf{D}, \mathbf{W})}{q(\mathbf{W}|\boldsymbol{\theta})} d\mathbf{W} = \ln p(D) - \mathcal{L}(q(\mathbf{W}|\boldsymbol{\theta})) \end{aligned}$$

$\ln p(D)$  не зависит от  $\theta$ . Следовательно, максимизируя  $\mathcal{L}(q(\mathbf{W}|\theta))$ , мы минимизируем  $KL(q(\mathbf{W}|\theta)||p(\mathbf{W}|D))$ . Тем самым, при максимизации  $\mathcal{L}(q(\mathbf{W}|\theta))$  распределение весов модели  $q(\mathbf{W}|\theta)$  будет приближаться к апостериорному распределению весов модели  $p(\mathbf{W}|D)$ .

Преобразуем выражение для  $\mathcal{L}(q(\mathbf{W}|\theta))$ , используя тождество  $p(\mathbf{W}, D) = p(D|\mathbf{W}) \cdot p(\mathbf{W})$ :

$$\begin{aligned}\mathcal{L}(q(\mathbf{W}|\theta)) &= \int q(\mathbf{W}|\theta) \cdot \ln \frac{p(D, \mathbf{W})}{q(\mathbf{W}|\theta)} d\mathbf{W} = \int q(\mathbf{W}|\theta) \cdot \ln \frac{p(D|\mathbf{W}) \cdot p(\mathbf{W})}{q(\mathbf{W}|\theta)} d\mathbf{W} = \\ &= \int q(\mathbf{W}|\theta) \cdot \ln p(D|\mathbf{W}) d\mathbf{W} - \int q(\mathbf{W}|\theta) \cdot \ln \frac{q(\mathbf{W}|\theta)}{p(\mathbf{W})} d\mathbf{W} = \\ &= \int q(\mathbf{W}|\theta) \cdot \ln p(D|\mathbf{W}) d\mathbf{W} - KL(q(\mathbf{W}|\theta)||p(\mathbf{W}))\end{aligned}$$

Таким образом, задача байесовского вывода свелась к задаче максимизации  $\mathcal{L}(q(\mathbf{W}|\theta))$  по параметрам  $\theta$ .

При аппроксимации апостериорного распределения параметров модели  $p(\mathbf{W}|D)$  распределением  $q(\mathbf{W}|\theta)$  аппроксимация предсказательного распределения  $p(\mathbf{y}|\mathbf{x}, D)$  будет выглядеть следующим образом:

$$p(\mathbf{y}|\mathbf{x}, D) \approx \frac{1}{T} \sum_{t=1}^T p(\mathbf{y}|\hat{\mathbf{W}}_t, \mathbf{x}), \text{ где } \hat{\mathbf{W}}_t \text{ — сэмпл весов модели из } q(\mathbf{W}|\theta)$$

## 5 Задание функциональных форм распределений

Для дальнейшего вывода положим, что распределения  $p(\mathbf{W})$  и  $q(\mathbf{W}|\theta)$  являются нормальными с диагональными матрицами ковариации:

$$p(\mathbf{W}) = N(\mathbf{W}|\mathbf{0}, \text{diag}(\sigma_{\mathbf{p}(\mathbf{W})})^2), \text{ где } \sigma_{\mathbf{p}(\mathbf{W})} \text{ — вектор длины } M$$

$q(\mathbf{W}|\theta) = N(\mathbf{W}|\boldsymbol{\mu}, \text{diag}(\sigma_{\mathbf{q}(\mathbf{W})})^2)$ , где  $\boldsymbol{\mu}$  и  $\sigma_{\mathbf{q}(\mathbf{W})}$  — вектора длины  $M$ , которые вместе образуют вектор обучаемых параметров  $\theta$ .

Априорное распределение весов модели  $p(\mathbf{W})$  имеет нулевое математическое ожидание (из соображений симметрии), и среднеквадратическое отклонение  $\sigma_{\mathbf{p}(\mathbf{W})}$ . В классическом байесовском выводе параметр  $\sigma_{\mathbf{p}(\mathbf{W})}$  должен задаваться до начала обучения, то есть являться гиперпараметром. Однако мы можем воспользоваться техникой эмпирического Байеса, то есть определить параметр априорного распределения  $\sigma_{\mathbf{p}(\mathbf{W})}$  из данных.

Пусть  $\boldsymbol{\alpha} = \text{diag}(\sigma_{\mathbf{p}(\mathbf{W})})^{-2}$ . Тогда  $p(\mathbf{W}) = N(\mathbf{W}|\mathbf{0}, \boldsymbol{\alpha}^{-1})$ .

Так как распределения  $p(\mathbf{W})$  и  $q(\mathbf{W}|\theta)$  являются нормальными, то  $KL(q(\mathbf{W}|\theta)||p(\mathbf{W}))$  можно посчитать аналитически:

$$\begin{aligned}KL(q(\mathbf{W}|\theta)||p(\mathbf{W})) &= \frac{1}{2} \sum_{k=1}^M \left( \frac{\sigma_{q(W)_k}^2}{\sigma_{p(W)_k}^2} + \frac{\mu_k^2}{\sigma_{p(W)_k}^2} - \ln \frac{\sigma_{q(W)_k}^2}{\sigma_{p(W)_k}^2} - 1 \right) = \\ &= \frac{1}{2} \sum_{k=1}^M (\alpha_k (\sigma_{q(W)_k}^2 + \mu_k^2) - \ln (\alpha_k \cdot \sigma_{q(W)_k}^2) - 1)\end{aligned}$$

Так как в выражении  $\mathcal{L}(q(\mathbf{W}|\theta))$  интеграл  $\int q(\mathbf{W}|\theta) \cdot \ln p(D|\mathbf{W}) d\mathbf{W}$  не зависит от параметров распределения  $p(\mathbf{W})$ , то:

$$\frac{\partial \mathcal{L}(q(\mathbf{W}|\theta))}{\partial \alpha_k} = - \frac{\partial (KL(q(\mathbf{W}|\theta)||p(\mathbf{W})))}{\partial \alpha_k} = - \frac{1}{2} (\sigma_{q(W)_k}^2 + \mu_k^2 - \frac{1}{\alpha_k}) = - \frac{1}{2} (\sigma_{q(W)_k}^2 + \mu_k^2 - \sigma_{p(W)_k}^2)$$

Приравняв производную к нулю, получим:

$$\begin{aligned}- \frac{1}{2} (\sigma_{q(W)_k}^2 + \mu_k^2 - \sigma_{p(W)_k}^2) &= 0 \\ \sigma_{p(W)_k}^2 &= \sigma_{q(W)_k}^2 + \mu_k^2\end{aligned}$$

Подставив полученное выражение в  $KL(q(\mathbf{W}|\boldsymbol{\theta})||p(\mathbf{W}))$ , получим:

$$KL(q(\mathbf{W}|\boldsymbol{\theta})||p(\mathbf{W})) = \frac{1}{2} \sum_{k=1}^M \ln(1 + \frac{\mu_k^2}{\sigma_{q(W)_k}^2})$$

$$\mathcal{L}(q(\mathbf{W}|\boldsymbol{\mu}, \boldsymbol{\sigma}_q(\mathbf{W}))) = \int q(\mathbf{W}|\boldsymbol{\mu}, \boldsymbol{\sigma}_q(\mathbf{W})) \cdot \ln p(D|\mathbf{W}) d\mathbf{W} - \frac{1}{2} \sum_{k=1}^M \ln(1 + \frac{\mu_k^2}{\sigma_{q(W)_k}^2})$$

## 6 Замена переменных

При решении задачи оптимизации при попадании в такие области, где для какого-либо веса  $\sigma_{q(W)_k} = 0$  и  $\mu_k = 0$ , возникает неопределенность деления  $\frac{0}{0}$ . Чтобы избежать этой неопределенности, и чтобы  $\boldsymbol{\sigma}_q(\mathbf{W})$  была всегда положительна, сделаем следующую замену переменных:

$$\sigma_{q(W)_k} = \ln(1 + e^{\rho_k}) = \text{Softplus}(\rho_k)$$

$$\mu_k = \gamma_k \cdot \text{Softplus}(\rho_k)$$

Тогда:

$$KL(q(\mathbf{W}|\boldsymbol{\theta})||p(\mathbf{W})) = \frac{1}{2} \sum_{k=1}^M \ln(1 + \frac{\mu_k^2}{\sigma_{q(W)_k}^2}) = \frac{1}{2} \sum_{k=1}^M \ln(1 + \gamma_k^2)$$

Таким образом, задача сводится к минимизации следующей функции потерь:

$$\text{Loss}(\boldsymbol{\rho}, \boldsymbol{\gamma}) = -\frac{\mathcal{L}(q(\mathbf{W}|\boldsymbol{\theta}))}{L} = \int N(\mathbf{W}|\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}_q(\mathbf{W}))^2) \cdot NLL \, d\mathbf{W} + \frac{KL}{L}$$

где:

$$\boldsymbol{\sigma}_q(\mathbf{W}) = \text{Softplus}(\boldsymbol{\rho})$$

$$\boldsymbol{\mu} = \boldsymbol{\gamma} \cdot \text{Softplus}(\boldsymbol{\rho})$$

$$NLL = -\frac{1}{L} \sum_{i=1}^L \ln p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{W})$$

$$KL = \frac{1}{2} \sum_{k=1}^M \ln(1 + \gamma_k^2)$$

## 7 Алгоритм обучения

Задаем шаг градиентного спуска  $\alpha$  и инициализируем параметры распределения  $q(\mathbf{W}|\boldsymbol{\theta})$  —  $\boldsymbol{\rho}$  и  $\boldsymbol{\gamma}$ . Затем повторяем, пока не достигнем критерия остановки:

1.  $\boldsymbol{\sigma} \leftarrow \text{Softplus}(\boldsymbol{\rho})$  — расчёт среднеквадратических отклонений весов
2.  $\boldsymbol{\mu} \leftarrow \boldsymbol{\gamma} \odot \boldsymbol{\sigma}$  — расчёт математических ожиданий весов
3.  $\hat{\mathbf{W}} \leftarrow N(0, 1)$  — сэмплирование случайных весов
4.  $\hat{\mathbf{W}} \leftarrow \hat{\mathbf{W}} \odot \boldsymbol{\sigma} + \boldsymbol{\mu}$  — репараметризация
5.  $nll \leftarrow -\frac{1}{L} \sum_{i=1}^L \ln p(\mathbf{y}_i|\mathbf{x}_i, \hat{\mathbf{W}})$  — расчёт среднего отрицательного логарифма правдоподобия (возможна аппроксимация по батчам)
6.  $kl \leftarrow \frac{1}{2} \sum_{k=1}^M \ln(1 + \gamma_k^2)$  — расчёт KL-дивергенции
7.  $l \leftarrow nll + \frac{kl}{L}$  — расчёт функции потерь

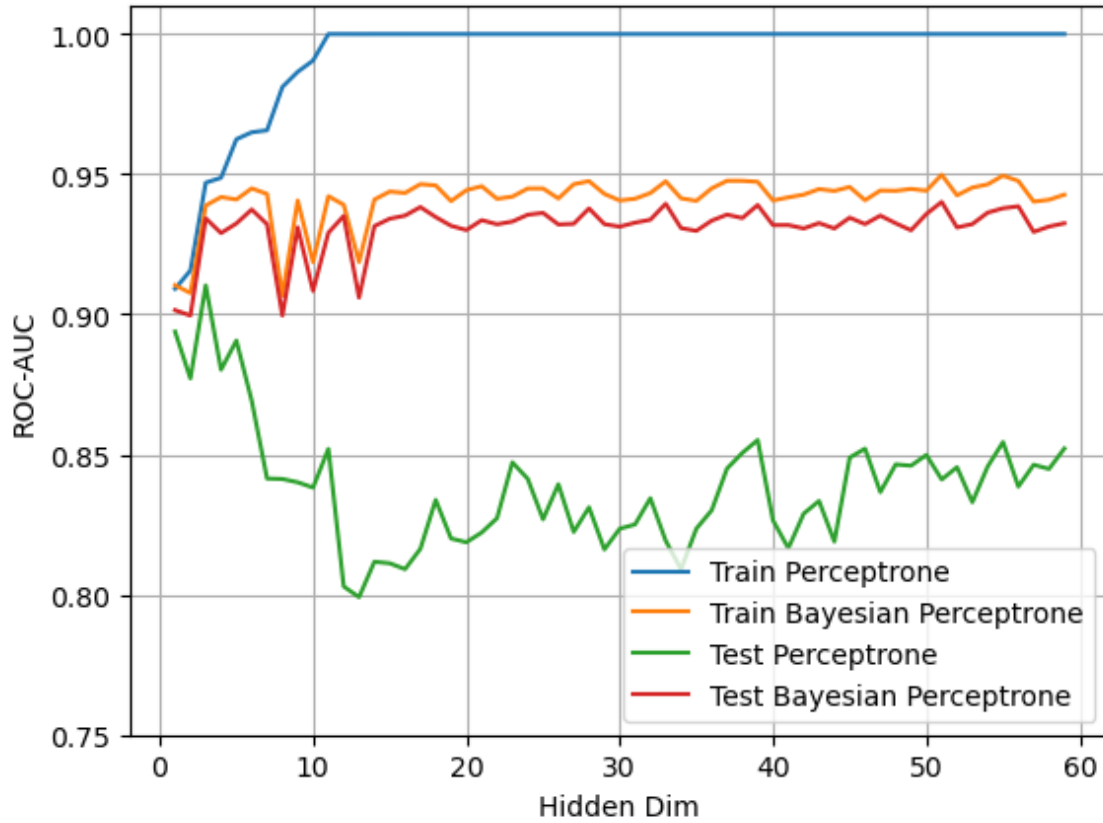


Рис. 1: Зависимость ROC-AUC от размерности скрытого состояния на тренировочных и тестовых данных

$$8. \rho \leftarrow \rho - \alpha \frac{\partial l}{\partial \rho} \text{ — обновление } \rho$$

$$9. \gamma \leftarrow \gamma - \alpha \frac{\partial l}{\partial \gamma} \text{ — обновление } \gamma$$

## 8 Эксперименты

Для проверки своей гипотезы я выбрал [Alzheimer's Disease Dataset](#). Данные были разбиты на тренировочную и тестовую часть в пропорции 80 на 20. В качестве архитектуры была выбрана полносвязная нейронная сеть с одним скрытым слоем и функцией активации ReLU. То есть:

$$z = \text{ReLU}(\text{matmul}(x, W_1))$$

$$y = \text{Sigmoid}(\text{matmul}(z, W_2))$$

Размерность скрытого состояния  $z$  варьировалась от 1 до 60. Для каждой размерности обучались 2 модели - классическая (без регуляризации) и байесовская. Для каждой модели производилась оценка ROC-AUC на тренировочной и тестовой выборках. На рисунке 1 представлены результаты экспериментов

## 9 Выводы

По результатам работы можно сделать следующие выводы:

- с ростом сложности модели байесовская нейронная сеть не переобучилась;
- значение ROC-AUC на тестовой выборке имеет очень высокую корреляцию со значением ROC-AUC на тренировочной выборке (0.97 по Пирсону). Следовательно, для подбора

гиперпараметров можно ориентироваться на метрики, полученные по тренировочной выборке. Это даёт нам возможность отказаться от деления на тренировочную и валидационную выборки для подбора гиперпараметров.

Так же стоит отметить, что данный подход переносится на другие архитектуры нейронных сетей (рекуррентные, свёрточные, трансформеры).

Имплементация данного подхода была выполнена с использованием PyTorch. Весь исходный код для проведения экспериментов размещён по адресу <https://github.com/dimabasow/bayesian-neural-networks>.