

Нейронные сети. Краткий курс.

Лекция 5

Обучение сети RBF как плохо обусловленная задача.

Описанная в предыдущей лекции процедура обучения нейронной сети на основе радиальных базисных функций (далее RBF-сеть) не всегда работает. Множество обучающих точек может быть очень большим, матрица интерполяции Φ может оказаться плохо обусловленной, множество примеров может неудачно покрывать пространство и все это может привести к плохой обобщающей способности сети.

Здесь мы сталкиваемся с типом задач, которые называются *некорректными или плохо обусловленными*. Чтобы почувствовать, что значит плохо обусловленная задача, рассмотрим для примера простую систему линейных уравнений

$$\begin{aligned}x_1 + 10 x_2 &= 11 \\10 x_1 + 101 x_2 &= 111\end{aligned}$$

Или то же в матричной форме

$$\begin{pmatrix} 1 & 10 \\ 10 & 101 \end{pmatrix} \times \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 111 \end{pmatrix}$$

Решением этой системы является точка на плоскости с координатами $x_1=1, x_2=1$. Теперь представим, что в результате измерительной неточности при определении правой части обусловленной шумами система слегка видоизменилась:

$$\begin{pmatrix} 1 & 10 \\ 10 & 101 \end{pmatrix} \times \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 11.1 \\ 111 \end{pmatrix}.$$

И решение этой «испорченной» шумами системы уже $x_1=11.1, x_2=0$! Довольно далеко от решения «точной» системы. В матричной форме эта простая задачка представлена чтобы вызвать ассоциации с системой (4.12) из предыдущей лекции.

Итак рассмотрим неизвестное отображение $f: x \in X \rightarrow y \in Y, X \subset \mathbb{R}^{m_0}, Y \subset \mathbb{R}$, то есть отображающее область X в диапазон Y , что собственно и призвана делать рассматриваемая нами RBF-сеть. Задача реконструкции отображения f считается хорошо обусловленной (или корректной) если выполнены следующие три условия.

1. *Существование.* Для любого входного вектора $\mathbf{x} \in X$ существует выходное значение $y = f(\mathbf{x})$, где $y \in Y$.
2. *Однозначность.* Для любой пары векторов $\mathbf{x}, \mathbf{t} \in X$ равенство $f(\mathbf{x}) = f(\mathbf{t})$ выполняется тогда и только тогда, когда $\mathbf{x} = \mathbf{t}$.
3. *Непрерывность.* Отображение считается непрерывным, если для любого $\epsilon > 0$ существует $\delta(\epsilon) > 0$, такое, что из условия $\rho_x(\mathbf{x}, \mathbf{t}) < \epsilon$ вытекает условие $\rho_y(f(\mathbf{x}), f(\mathbf{t})) < \delta$, где ρ - расстояние в соответствующих пространствах. Свойство непрерывности еще называют *устойчивостью*.

Если какое-то из этих условий не выполнено, задача считается плохо обусловленной. В обучающем наборе данных информации может оказаться недостаточно для реконструкции отображения, шумы могут оказаться достаточно большими, чтобы вывести результат отображения из диапазона Y , что будет означать нарушение первого условия. Собственно сама задача может оказаться таковой, что является неустойчивой относительно возмущений в выходных данных. Рассмотрим далее как с помощью методов *регуляризации* плохо обусловленную задачу сделать хорошо обусловленной. Обычно теорию регуляризации связывают с именем Тихонова А.Н., который в 1963 г. предложил новый метод названный регуляризацией и предназначенный для решения плохо обусловленных задач. Главная идея регуляризации заключается в стабилизации решения с помощью некоторой вспомогательной неотрицательной функции, которая несет в себе априорную информацию о решении. Наиболее общей формой априорной информации является предположение о гладкости искомого отображения.

Пусть для создания RBF-сети (или реконструкции отображения) имеется множество пар данных (\mathbf{x}_i, d_i) доступных для аппроксимации, где $\mathbf{x}_i \in \mathbb{R}^{m_0}$, $i = 1, 2, \dots, N$ и $y_i \in \mathbb{R}$, $i = 1, 2, \dots, N$. Функцию аппроксимации обозначим как $F(\mathbf{x})$, где для простоты обозначений опущен вектор весов \mathbf{w} . В теории регуляризации Тихонова функция ошибок, которую нужно минимизировать для определения параметров искомого отображения модифицируется следующим образом:

$$E(F) = \frac{1}{2} \sum_{i=1}^N (d_i - y_i)^2 + \frac{1}{2} \lambda \|\mathbf{D}F\| = E_s(F) + \lambda E_c(F) . \quad (5.1)$$

Здесь первое слагаемое это *слагаемое стандартной ошибки* описывающее расстояние между желаемым откликом и фактическим выходным сигналом, а второе слагаемое это *слагаемое регуляризации* в котором \mathbf{D} - линейный дифференциальный оператор, называемый иногда *стабилизатором* решения, содержащий априорную информацию о задаче. Параметр λ - действительное положительное число, называемое *параметром регуляризации*. Обозначим аргминимум функционала Тихонова $E(F)$ как $F_\lambda(\mathbf{x})$. Параметр регуляризации λ можно рассматривать как индикатор достаточности имеющегося набора данных для определения решения $F_\lambda(\mathbf{x})$. Крайний случай $\lambda \rightarrow 0$ означает, что задача является хорошо обусловленной и имеет решение целиком зависящее от набора учебных данных.

Другой крайний случай большого параметра регуляризации означает, что самого априорного ограничения на гладкость определяемого оператором D достаточно для определения решения или что учебные данные содержат недостаточно информации для определения отображения. В практических приложениях параметр регуляризации подбирается между этими крайними случаями. Можно рассматривать функцию $E_c(F)$ как *функцию штрафа за сложность модели*, влияние которой определяется параметром регуляризации.

Итак, принцип регуляризации заключается в следующем. Необходимо найти функцию $F(x)$ минимизирующую функционал Тихонова (5.1). Чтобы минимизировать функционал, нужно задать правило оценки его дифференциала. Для этого используем *дифференциал Фреше*:

$$dE(F, h) = \frac{d}{d\beta} E(F(x) + \beta h(x))_{\beta=0}, \quad (5.2)$$

где $h(x)$ - фиксированная функция. Необходимым условием того, чтобы функция $F(x)$ была экстремумом функционала $E(F)$ является равенство

$$dE(F, h) = dE_s(F, h) + \lambda dE_c(F, h) = 0, \quad (5.3)$$

где dE_s и dE_c - соответствующие дифференциалы Фреше слагаемых стандартной ошибки и регуляризации. Дифференцирование в определении дифференциала Фреше осуществляется по обычным правилам. Рассмотрим первое слагаемое в (5.3):

$$\begin{aligned} dE_s(F, h) &= \frac{1}{2} \frac{d}{d\beta} \sum_{i=1}^N (d_i - F(x_i) - \beta h(x_i))^2_{\beta=0} = \\ &= - \sum_{i=1}^N (d_i - F(x_i) - \beta h(x_i)) h(x_i)_{\beta=0} = - \sum_{i=1}^N (d_i - F(x_i)) h(x_i) \end{aligned} \quad (5.4)$$

Нам понадобится следующая теорема Ритца о представлении из функционального анализа. Пусть f - линейный ограниченный функционал в Гильбертовом пространстве \mathbf{H} . Тогда существует единственная функция $h_0 \in \mathbf{H}$ такая что $f = (h, h_0)$ для всех $h \in \mathbf{H}$, при этом $\|f\|_{\mathbf{H}^*} = \|h_0\|_{\mathbf{H}}$, где \mathbf{H}^* - пространство сопряженное пространству \mathbf{H} . В связи с этим можно переписать (5.4) в виде скалярного произведения в функциональном пространстве:

$$dE_s(F, h) = - \left(h, \sum_{i=1}^N (d_i - F(x_i)) \delta(x - x_i) \right), \quad (5.5)$$

где $\delta(x - x_i)$ - дельта функция Дирака. Пользуясь подобным представлением, запишем выражение для дифференциала Фреше для dE_c :

$$\begin{aligned}
dE_c(F, h) &= \frac{d}{d\beta} E_c(F(\mathbf{x}) + \beta h(\mathbf{x}))_{\beta=0} = \frac{1}{2} \frac{d}{d\beta} \int_{\mathbb{R}^{m_0}} (\mathbf{D}[F(\mathbf{x}) + \beta h(\mathbf{x})])^2 d\mathbf{x}_{\beta=0} = \\
&= \int_{\mathbb{R}^{m_0}} \mathbf{D}[F(\mathbf{x}) + \beta h(\mathbf{x})] \mathbf{D}h(\mathbf{x}) d\mathbf{x}_{\beta=0} = \int_{\mathbb{R}^{m_0}} \mathbf{D}F(\mathbf{x}) \mathbf{D}h(\mathbf{x}) d\mathbf{x} = \\
&= (\mathbf{D}F(\mathbf{x}), \mathbf{D}h(\mathbf{x})).
\end{aligned} \tag{5.6}$$

Для данного оператора можно найти такой однозначно определенный оператор $\tilde{\mathbf{D}}$ называемый *сопряженным*, такой что для любых функций $u(\mathbf{x})$ и $v(\mathbf{x})$ удовлетворяющих соответствующим граничным условиям, можно записать

$$\int_{\mathbb{R}^{m_0}} u(\mathbf{x}) \mathbf{D}v(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^{m_0}} v(\mathbf{x}) \tilde{\mathbf{D}}u(\mathbf{x}) d\mathbf{x} . \tag{5.7}$$

Это тождество называют тождеством Грина. Если теперь переобозначить $u(\mathbf{x}) = \mathbf{D}F(\mathbf{x})$ и $\tilde{\mathbf{D}}v(\mathbf{x}) = \mathbf{D}h(\mathbf{x})$, то можно записать следующее выражение:

$$dE_c(F, h) = \int_{\mathbb{R}^{m_0}} \mathbf{D}F(\mathbf{x}) \mathbf{D}h(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^{m_0}} h(\mathbf{x}) \tilde{\mathbf{D}}\mathbf{D}F(\mathbf{x}) d\mathbf{x} = (h, \tilde{\mathbf{D}}\mathbf{D}F(\mathbf{x})) \tag{5.8}$$

Если теперь переписать условие экстремума функционала Тихонова в виде $\frac{1}{\lambda} dE_s + dE_c = 0$ и подставить туда выражения для дифференциалов Фреше, то получим

$$-\frac{1}{\lambda} \left(h(\mathbf{x}), \sum_{i=1}^N (d_i - F(\mathbf{x})) \delta(\mathbf{x} - \mathbf{x}_i) \right) + (h(\mathbf{x}), \tilde{\mathbf{D}}\mathbf{D}F(\mathbf{x})) = 0 . \tag{5.9}$$

Условие (5.9) можно переписать в виде одного скалярного произведения:

$$\left(h(\mathbf{x}), \tilde{\mathbf{D}}\mathbf{D}F(\mathbf{x}) - \frac{1}{\lambda} \sum_{i=1}^N (d_i - F(\mathbf{x})) \delta(\mathbf{x} - \mathbf{x}_i) \right) = 0 . \tag{5.10}$$

Условие (5.10) будет выполняться для произвольной функции $h(\mathbf{x})$ если

$$\tilde{\mathbf{D}}\mathbf{D}F(\mathbf{x}) = \frac{1}{\lambda} \sum_{i=1}^N (d_i - F(\mathbf{x})) \delta(\mathbf{x} - \mathbf{x}_i) - \tag{5.11}$$

уравнение *Эйлера-Лагранжа* для функционала Тихонова $E(F)$. Оно является необходимым условием существования экстремума для функции $F_\lambda(\mathbf{x})$. Уравнение (5.11) это дифференциальное уравнение в частных производных для функции аппроксимации F .

Для заданного линейного дифференциального оператора L , зададимся функцией $G(\mathbf{x}, \xi)$ обладающей следующими свойствами:

- 1). Для фиксированного значения ξ функция $G(x, \xi)$ удовлетворяет граничным условиям.
- 2). Во всех точках, кроме $x = \xi$, все производные функции $G(x, \xi)$ при фиксированном ξ непрерывны, а их количество определяется порядком L .
- 3). $G(x, \xi)$ как функция x удовлетворяет уравнению в частных производных

$$L G(x, \xi) = \delta(x - \xi) . \quad (5.12)$$

Такая функция называется *функцией Грина* для оператора дифференцирования L . Она играет ту же роль для L , что и обратная матрица в матричном исчислении.

Пусть $\phi(x)$ - непрерывная или кусочно-непрерывная функция. Тогда функция $F(x)$, определенная как

$$F(x) = \int_{\mathbb{R}^m} G(x, \xi) \phi(\xi) d\xi \quad (5.13)$$

является решением дифференциального уравнения

$$L F(x) = \phi(x) . \quad (5.14)$$

Проверка этого утверждения выполняется непосредственно

$$L F(x) = L \int_{\mathbb{R}^m} G(x, \xi) \phi(x) d\xi = \int_{\mathbb{R}^m} L G(x, \xi) \phi(x) d\xi = \int_{\mathbb{R}^m} \delta(x - \xi) \phi(x) d\xi = \phi(x) . \quad (5.15)$$

Перестановка интеграла и оператора дифференцирования возможно потому, что дифференцирование и интегрирование производятся по разным переменным, дифференцирование по x , а интегрирование по ξ .

Вернемся к уравнению Эйлера-Лагранжа. В нем в качестве оператора L выступает оператор $\tilde{D}D$, а в качестве функции ϕ функция

$$\phi(\xi) = \frac{1}{\lambda} \sum_{i=1}^N (d_i - F(x_i)) \delta(\xi - x_i) . \quad (5.16)$$

Тогда запись $F(x)$ через функцию Грина выглядит следующим образом:

$$\begin{aligned} F_\lambda(x) &= \int_{\mathbb{R}^m} G(x, \xi) \left\{ \frac{1}{\lambda} \sum_{i=1}^N (d_i - F(x_i)) \delta(\xi - x_i) \right\} = \\ &= \frac{1}{\lambda} \sum_{i=1}^N (d_i - F(x_i)) \int_{\mathbb{R}^m} G(x, \xi) \delta(\xi - x_i) d\xi = \\ &= \frac{1}{\lambda} \sum_{i=1}^N (d_i - F(x)) G(x, x_i) \end{aligned} \quad (5.17)$$

А это означает, что решение задачи регуляризации является линейной суперпозицией N функций Грина с весами $\frac{1}{\lambda}(d_i - F(\mathbf{x}_i))$. Обозначив $w_i = \frac{1}{\lambda}(d_i - F(\mathbf{x}_i))$ можно записать (5.17) как

$$F_\lambda(\mathbf{x}) = \sum_{i=1}^N w_i G(\mathbf{x}, \mathbf{x}_i) . \quad (5.18)$$

Вычисляя $F_\lambda(\mathbf{x}_j)$ в точках $\mathbf{x}_j, j=1,2,\dots,N$ можно получить систему уравнений

$$F_\lambda(\mathbf{x}_j) = \sum_{i=1}^N w_i G(\mathbf{x}_j, \mathbf{x}_i) \quad (5.19)$$

$$j=1,2,\dots,N$$

или то же самое в матричной форме

$$\mathbf{F}_\lambda = \mathbf{G} \mathbf{w} , \quad (5.20)$$

где $\mathbf{F}_\lambda = (F_\lambda(\mathbf{x}_1), F_\lambda(\mathbf{x}_2), \dots, F_\lambda(\mathbf{x}_N))^T$, $\mathbf{w} = \frac{1}{\lambda}(\mathbf{d} - \mathbf{F}_\lambda)$, $\mathbf{d} = (d_1, d_2, \dots, d_N)^T$,

$$\mathbf{G} = \begin{pmatrix} G(\mathbf{x}_1, \mathbf{x}_1) & \dots & G(\mathbf{x}_1, \mathbf{x}_N) \\ \dots & \dots & \dots \\ G(\mathbf{x}_N, \mathbf{x}_1) & \dots & G(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix} .$$

Используя выражение для вектора весов \mathbf{w} , переписав его как $\mathbf{F}_\lambda + \lambda \mathbf{w} = \mathbf{d}$, затем подставив вместо \mathbf{F}_λ правую часть (5.20), можно получить следующее уравнение для определения весов:

$$(\mathbf{G} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{d} , \quad (5.21)$$

где \mathbf{I} - единичная матрица. Матрица \mathbf{G} называется *матрицей Грина*, она аналогична матрице интерполяции Φ и так же является симметричной в силу того, что мы определили оператор $L = \tilde{D}D$ как самосопряженный, а это соответствует симметрии функции Грина, то есть $G(\mathbf{x}_j, \mathbf{x}_i) = G(\mathbf{x}_i, \mathbf{x}_j)$ для всех $i, j=1,2,\dots,N$. Симметрия \mathbf{G} еще означает, что $\mathbf{G}^T = \mathbf{G}$. На практике всегда можно подобрать λ , так чтобы матрица $(\mathbf{G} + \lambda \mathbf{I})$ была достаточно хорошей для ее обращения, следовательно система (5.21) имеет единственное решение

$$\mathbf{w} = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{d} . \quad (5.22)$$

Итак, выбрав оператор дифференцирования \mathbf{D} и имея для него набор функций Грина можно получить набор весов \mathbf{w} для заданного λ . Решение задачи регуляризации задается разложением (5.18) в котором функция Грина соответствует самосопряженному оператору $\tilde{\mathbf{D}}\mathbf{D}$. Характеристики функции Грина зависят от выбранной формы оператора \mathbf{D} и выбранных граничных условий. Количество функций Грина в таком решении равно числу примеров обучения. Чтобы функция Грина оказалась еще и радиальной функцией нужно, чтобы оператор \mathbf{D} был инвариантен к преобразованием связанным с переносом координат (тогда $G(\mathbf{x}, \mathbf{x}_i) = G(\mathbf{x} - \mathbf{x}_i)$) и инвариантен к поворотам. В этом случае $G(\mathbf{x}, \mathbf{x}_i) = G(\|\mathbf{x} - \mathbf{x}_i\|)$. В этом случае решение

$$F_\lambda(\mathbf{x}) = \sum_{i=1}^N w_i G(\|\mathbf{x} - \mathbf{x}_i\|) \quad (5.23)$$

очень похоже на решение задачи интерполяции, однако оно регуляризовано и будет совпадать с интерполяционным только в случае $\lambda = 0$.

Такая функция Грина, которая одновременно является радиальной представляет практический интерес для построения RBF-сетей. Примером такой функции является многомерная функция Гаусса:

$$G(\mathbf{x}, \mathbf{x}_i) = \exp\left(-\frac{1}{2\sigma_i^2} \|\mathbf{x} - \mathbf{x}_i\|^2\right) . \quad (5.24)$$

Такую функцию Грина определяет оператор

$$\mathbf{L} = \sum_{n=0}^{\infty} (-1)^n \alpha_n \nabla^{2n} , \quad (5.25)$$

где $\alpha_n = \frac{\sigma_i^{2n}}{n! 2^n}$, ∇^{2n} - n -кратное применение оператора Лапласа в \mathbb{R}^{m_0}

$$\nabla^2 = \frac{\partial^2}{\partial x_1^2} + \frac{\partial^2}{\partial x_2^2} + \dots + \frac{\partial^2}{\partial x_{m_0}^2} . \quad (5.26)$$