

Глава 1

Полносвязные Сети

- \mathcal{N} — размер батча
- \mathcal{D} — размерность вектора признаков

1.1 Batch Normalization

Пусть \mathbf{x} — это вектор на входе слоя ($x \in \mathbb{R}^{\mathcal{N}}$). Тогда вектор \mathbf{y} на выходе слоя есть

$$\mathbf{y} = \gamma \cdot \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \varepsilon}} + \beta,$$

где

$$\mu = \frac{1}{\mathcal{N}} \sum_{i=1}^{\mathcal{N}} x_i, \quad \sigma^2 = \frac{1}{\mathcal{N}} \sum_{i=1}^{\mathcal{N}} (x_i - \mu)^2 = \frac{1}{\mathcal{N}} x_i^2 - \left(\frac{1}{\mathcal{N}} x_i \right)^2.$$

Пусть дана производная функции потерь по \mathbf{y} , т.е. $\partial \mathcal{L} / \partial \mathbf{y}$. Найдём производные функции потерь по \mathbf{x} , β , γ . Сначала найдём $\partial \mathcal{L} / \partial \mathbf{x}$:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \underbrace{\frac{\partial \mathbf{y}}{\partial \mathbf{x}}}_{\mathbb{R}^{\mathcal{N} \times \mathcal{N}}} \cdot \underbrace{\frac{\mathcal{L}}{\partial \mathbf{y}}}_{\mathbb{R}^{\mathcal{N}}}, \quad (1.1)$$

где

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_{\mathcal{N}}}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_{\mathcal{N}}}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_{\mathcal{N}}} & \frac{\partial y_2}{\partial x_{\mathcal{N}}} & \cdots & \frac{\partial y_{\mathcal{N}}}{\partial x_{\mathcal{N}}} \end{pmatrix} \quad (1.2)$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \gamma \frac{\partial \left(\frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \varepsilon}} \right)}{\partial \mathbf{x}}.$$

Рассмотрим производную

$$\begin{aligned} \frac{\partial y_i}{\partial x_j} &= \gamma \frac{\partial \left(\frac{x_i - \mu}{\sqrt{\sigma^2 + \varepsilon}} \right)}{\partial x_j} = \gamma \frac{\partial (x_i - \mu)}{\partial x_j} \frac{1}{\sqrt{\sigma^2 + \varepsilon}} - \frac{\gamma}{2} \frac{x_i - \mu}{(\sigma^2 + \varepsilon)^{3/2}} \frac{\partial (\sigma^2 + \varepsilon)}{\partial x_j} = \\ &= \gamma \left(\delta_{ij} - \frac{1}{\mathcal{N}} \right) \frac{1}{\sqrt{\sigma^2 + \varepsilon}} - \frac{\gamma}{2} \frac{x_i - \mu}{(\sigma^2 + \varepsilon)^{3/2}} \cdot 2 \left(\frac{x_j}{\mathcal{N}} - \frac{\mu}{\mathcal{N}} \right) = \gamma \left(\delta_{ij} - \frac{1}{\mathcal{N}} \right) \frac{1}{\sqrt{\sigma^2 + \varepsilon}} - \gamma \frac{(x_i - \mu)(x_j - \mu)}{\mathcal{N}(\sigma^2 + \varepsilon)^{3/2}}. \end{aligned}$$

Таким образом

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\gamma}{\sqrt{\sigma^2 + \varepsilon}} \left(\left(\mathbf{I} - \frac{1}{\mathcal{N}} \mathbf{E} \right) - \frac{\mathbf{C}}{\sigma^2 + \varepsilon} \right),$$

где $\mathbf{C} = \frac{1}{N}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T$, \mathbf{I} — единичная матрица размера $N \times N$, \mathbf{E} — матрица, полностью состоящая из единиц, размера $N \times N$. Тогда

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\gamma}{\sqrt{\sigma^2 + \varepsilon}} \left(\left(\mathbf{I} - \frac{1}{N} \mathbf{E} \right) - \frac{\mathbf{C}}{\sigma^2 + \varepsilon} \right) \frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \frac{\gamma}{\sqrt{\sigma^2 + \varepsilon}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{y}} - \mathbf{e} \cdot \frac{\overline{\partial \mathcal{L}}}{\partial \mathbf{y}} - \frac{\mathbf{x} - \boldsymbol{\mu}}{N(\sigma^2 + \varepsilon)} \cdot \left\langle \mathbf{x} - \boldsymbol{\mu}, \frac{\mathcal{L}}{\mathbf{y}} \right\rangle \right),$$

где \mathbf{e} — столбец из единиц размерности N , $\frac{\overline{\partial \mathcal{L}}}{\partial \mathbf{y}}$ — среднее значение элементов вектора $\frac{\partial \mathcal{L}}{\partial \mathbf{y}}$. Обозначив через \mathbf{z} “стандартизованный” вектор \mathbf{x} ($\mathbf{z} = (\mathbf{x} - \boldsymbol{\mu})/\sigma$), получим

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\gamma}{\sqrt{\sigma^2 + \varepsilon}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{y}} - \mathbf{e} \cdot \frac{\overline{\partial \mathcal{L}}}{\partial \mathbf{y}} - \frac{\mathbf{x} - \boldsymbol{\mu}}{N(\sigma^2 + \varepsilon)} \cdot \left\langle \mathbf{x} - \boldsymbol{\mu}, \frac{\mathcal{L}}{\mathbf{y}} \right\rangle \right) = \frac{\gamma}{\sqrt{\sigma^2 + \varepsilon}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{y}} - \mathbf{e} \cdot \frac{\overline{\partial \mathcal{L}}}{\partial \mathbf{y}} - \frac{\mathbf{z}}{N} \cdot \left\langle \mathbf{z}, \frac{\mathcal{L}}{\mathbf{y}} \right\rangle \right).$$

Теперь найдем производные $\partial \mathcal{L} / \partial \beta$ и $\partial \mathcal{L} / \partial \gamma$:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta} &= \frac{\partial \mathbf{y}}{\partial \beta} \frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \mathbf{I} \frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}}, \\ \frac{\partial \mathcal{L}}{\partial \gamma} &= \left(\frac{\partial \mathbf{y}}{\partial \gamma} \right)^T \frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \frac{\mathbf{x} - \boldsymbol{\mu}}{\sqrt{\sigma^2 + \varepsilon}} \end{aligned}$$

1.1.1 Реализация в Python

Представленная ниже реализация предполагает, что

- $\mathbf{X} \in \mathbb{R}^{N \times D}$ — матрица объектов-признаков
- $\mathbf{Z} \in \mathbb{R}^{N \times D}$ — “стандартизованная” вдоль оси 0 матрица \mathbf{X} , т.е. выход слоя `BatchNormalization`
- $\boldsymbol{\mu} \in \mathbb{R}^N$ — средние значения для признаков
- $\boldsymbol{\sigma} \in \mathbb{R}^N$ — стандартные отклонения для признаков

```
def batchnorm_backward_alt(output_grad, cache):
    X, Z, mu, sigma, gamma, beta, eps = cache
    N, D = X.shape
    X_grad = gamma / np.sqrt(sample_var + eps)[None, :] * \
        ((output_grad - np.mean(output_grad, axis=0)[None, :]) - \
         Z * np.mean(np.multiply(Z, output_grad), axis=0)[None, :])
    beta_grad = np.sum(output_grad, axis=0)
    gamma_grad = np.sum(np.multiply(X_norm, output_grad), axis=0)
    return X_grad, gamma_grad, beta_grad
```
