

# (強化学習のための?) 深層学習メモ

菱沼 徹

2020 年 3 月 6 日

強化学習をやる際、深層学習についてだいたい次が分かればいいらしい [1] ので、それぞれメモしていく。

- architectures (MLP, vanilla RNN, LSTM, GRU, conv layers, resnets, attention mechanisms)
- common regularizers (weight decay, dropout)
- normalization (batch norm, layer norm, weight norm)
- optimizers (SGD, momentum SGD, Adam, others)
- reparameterization trick

表記の単純化のため、(できるだけ同じになるように努めるが) 記号の定義は各節のみにおいて有効とする。

## 1 入門: 多層パーセプトロン (MLP) を用いた回帰

### 1.1 2 層パーセプトロン

■2 層 NN 順伝播.  $n$  番目のデータに対して、入力  $\mathbf{x}_n \in \mathbb{R}^{H_0}$ , 出力  $\mathbf{y}_n \in \mathbb{R}^D$ , 重みパラメータ  $\mathbf{W}^{(1)} \in \mathbb{R}^{H_1 \times H_0}$ ,  $\mathbf{W}^{(2)} \in \mathbb{R}^{D \times H_1}$ , 活性化関数  $\phi(\cdot)$ , 正規ノイズ  $\epsilon_n \in \mathbb{R}^D$  として、次のモデルを考える。

$$y_{n,d} = \sum_{h_1=1}^{H_1} w_{d,h_1}^{(2)} \phi \left( \sum_{h_0=1}^{H_0} w_{h_1,h_0}^{(1)} x_{n,h_0} \right) + \epsilon_{n,d}.$$

あるいはベクトル表記で

$$\mathbf{y}_n = \mathbf{W}^{(2)} \phi(\mathbf{W} \mathbf{x}_n) + \epsilon_n$$

ここで、 $\phi(\cdot)$  は要素ごとに活性化関数を適用してベクトルにしたものを意味する。2 層だと上の書き方で十分だが、多層に拡張する際には次の表記の方が使いやすい。

$$\begin{aligned} \mathbf{z}_n^{(0)} &= \mathbf{x}_n \\ \mathbf{a}_n^{(1)} &= \mathbf{W}^{(1)} \mathbf{z}_n^{(0)} \\ \mathbf{z}_n^{(1)} &= \phi(\mathbf{a}_n^{(1)}) \\ \mathbf{a}_n^{(2)} &= \mathbf{W}^{(2)} \mathbf{z}_n^{(1)} \\ \mathbf{y}_n &= \mathbf{a}_n^{(2)} + \epsilon_n \end{aligned}$$

■2 層 NN 逆伝播. パラメータの集合を  $\mathbf{W}$ , 学習データ数を  $N$  とする。誤差関数を次で設計したとする。

$$\begin{aligned} E_n(\mathbf{W}) &= \frac{1}{2} \sum_{d=1}^D (y_{n,d} - a_{n,d}^{(2)})^2 \\ E(\mathbf{W}) &= \sum_{n=1}^N E_n(\mathbf{W}) \end{aligned}$$

この時、 $\phi'$  は活性化関数の微分として、次のように書ける。

$$\begin{aligned}
\frac{\partial E_n(\mathbf{W})}{\partial w_{d,h_1}^{(2)}} &= \frac{\partial E_n(\mathbf{W})}{\partial a_{n,d}^{(2)}} \frac{\partial a_{n,d}^{(2)}}{\partial w_{d,h_1}^{(2)}} = [y_{n,d} - a_{n,d}^{(2)}] [z_{n,h_1}^{(1)}] = [\delta_{n,d}^{(2)}] [z_{n,h_1}^{(1)}] = \delta_{n,d}^{(2)} z_{n,h_1}^{(1)} \\
\frac{\partial E_n(\mathbf{W})}{\partial w_{h_1,h_0}^{(1)}} &= \left\{ \sum_{d=1}^D \frac{\partial E_n(\mathbf{W})}{\partial a_{n,d}^{(2)}} \frac{\partial a_{n,d}^{(2)}}{\partial a_{n,h_1}^{(1)}} \right\} \frac{\partial a_{n,h_1}^{(1)}}{\partial w_{h_1,h_0}^{(1)}} = \left\{ \sum_{d=1}^D \delta_{n,d}^{(2)} \frac{\partial a_{n,d}^{(2)}}{\partial a_{n,h_1}^{(1)}} \right\} \frac{\partial a_{n,h_1}^{(1)}}{\partial w_{h_1,h_0}^{(1)}} \\
&= \left\{ \sum_{d=1}^D \delta_{n,d}^{(2)} \left[ \frac{\partial}{\partial a_{n,h_1}^{(1)}} \sum_{h'_1=1}^{H_1} w_{d,h'_1}^{(2)} z_{n,h'_1}^{(1)} \right] \right\} \left[ \frac{\partial}{\partial w_{h_1,h_0}^{(1)}} \sum_{h'_0=1}^{H_0} w_{h_1,h'_0}^{(1)} z_{n,h'_0}^{(0)} \right] \\
&= \left\{ \sum_{d=1}^D \delta_{n,d}^{(2)} \left[ w_{d,h_1}^{(2)} \frac{\partial z_{n,h_1}^{(1)}}{\partial a_{n,h_1}^{(1)}} \right] \right\} [z_{n,h_0}^{(0)}] = \left\{ \sum_{d=1}^D \delta_{n,d}^{(2)} [w_{d,h_1}^{(2)} \phi'(a_{n,h_1}^{(1)})] \right\} [z_{n,h_0}^{(0)}] \\
&= \left\{ \phi'(a_{n,h_1}^{(1)}) \sum_{d=1}^D \delta_{n,d}^{(2)} w_{d,h_1}^{(2)} \right\} [z_{n,h_0}^{(0)}] \\
&= [\delta_{n,h_1}^{(1)}] [z_{n,h_0}^{(0)}] = \delta_{n,h_1}^{(1)} z_{n,h_0}^{(0)}
\end{aligned}$$

■更新. 全てのパラメータに対する微分が上で計算できたので、勾配法で更新することができる。

$$w_{i,j}^{(\ell)} - \eta \frac{\partial E(\mathbf{W})}{\partial w_{i,j}^{(\ell)}} = w_{i,j}^{(\ell)} - \eta \sum_{n=1}^N \frac{\partial E_n(\mathbf{W})}{\partial w_{i,j}^{(\ell)}}$$

■まとめ.

- 順伝播で全ての  $\mathbf{a}$  と  $\mathbf{z}$  を計算.
- 逆伝播で全ての  $\delta$  を計算.
- 勾配法.

## 1.2 多層パーセプトロン

同様にして導出できる。

■MLP 順伝播.

$$\begin{aligned}
\mathbf{z}_n^{(\ell)} &= \begin{cases} \mathbf{x}_n & \text{if } \ell = 0 \\ \phi(\mathbf{a}_n^{(\ell)}) & \text{if } \ell \in \{1, \dots, L-1\} \end{cases} \\
\mathbf{a}_n^{(\ell)} &= \mathbf{W}^{(\ell)} \mathbf{z}_n^{(\ell-1)}, \quad \ell = 1, \dots, L \\
\mathbf{y}_n &= \mathbf{a}_n^{(L)} + \boldsymbol{\epsilon}_n
\end{aligned}$$

■MLP 逆伝播.

$$\begin{aligned}
\frac{\partial E_n(\mathbf{W})}{\partial w_{d,h_1}^{(L)}} &= \left[ \frac{\partial E_n(\mathbf{W})}{\partial a_{n,d}^{(L)}} \right] [z_{n,h_{L-1}}^{(L-1)}] = \delta_{n,d}^{(L)} z_{n,h_{L-1}}^{(L-1)} \\
\frac{\partial E_n(\mathbf{W})}{\partial w_{h_\ell,h_{\ell-1}}^{(\ell)}} &= \left[ \phi'(a_{n,h_\ell}^{(\ell)}) \sum_{h_{\ell+1}=1}^{H_{\ell+1}} \delta_{n,h_{\ell+1}}^{(\ell+1)} w_{h_{\ell+1},h_\ell}^{(\ell+1)} \right] [z_{n,h_{\ell-1}}^{(\ell-1)}] = \delta_{n,h_\ell}^{(\ell)} z_{n,h_{\ell-1}}^{(\ell-1)}, \quad \ell = 1, \dots, L-1
\end{aligned}$$

## 2 最適化 (optimization)

Adam [13] の導入を目標にして、書いていく。説明が載っている和書は、[24] が良いと思う。でも、英語で [19] を読むのが一番分かりやすいと思う。

### 2.1 復習：最急勾配法

変数を  $\mathbf{w}$ ，評価関数（滑らかとする）を  $E(\mathbf{w})$  とする。Taylor 展開により，

$$\begin{aligned} E(\mathbf{w} + \eta \mathbf{d}) &= E(\mathbf{w}) + \eta \mathbf{d}^T \left. \frac{\partial}{\partial \mathbf{w}} E(\mathbf{w}) \right|_{\mathbf{w}} + o(\eta) \\ &= E(\mathbf{w}) + \eta \mathbf{d}^T \nabla E(\mathbf{w}) + o(\eta) \end{aligned}$$

そのため、微小な係数  $\eta$  の高次項を無視して、 $\|\mathbf{d}\| = 1$  の下で  $E(\mathbf{w} + \eta \mathbf{d})$  を最小化する  $\mathbf{d}$  は、次を満たす。

$$\mathbf{d} = -\frac{1}{\|\nabla E(\mathbf{w})\|} \nabla E(\mathbf{w})$$

最急降下法の基本的アイデアは、上式に基づき次式で  $\mathbf{w}$  を更新していくことである。

$$\begin{aligned} \mathbf{d}^{(t)} &= -\frac{1}{\|\nabla E(\mathbf{w}^{(t)})\|} \nabla E(\mathbf{w}^{(t)}) = -\frac{1}{\left\| \left. \frac{\partial}{\partial \mathbf{w}} E(\mathbf{w}) \right|_{\mathbf{w}^{(t)}} \right\|} \left. \frac{\partial}{\partial \mathbf{w}} E(\mathbf{w}) \right|_{\mathbf{w}^{(t)}} \\ \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} + \eta^{(t)} \mathbf{d}^{(t)} \end{aligned}$$

最急降下法に限らず、 $\eta^{(t)}$  に相当する部分は学習係数やステップ幅と呼ばれ、様々な決め方が存在する。

### 2.2 確率的勾配降下 (SGD)

$n$  番目の入出力データ  $(\mathbf{x}_n, y)$ ，与えられたデータを、 $D = \{(\mathbf{x}_n, y)\}_{n=1}^N$  とする。回帰関数  $f(\mathbf{x}; \mathbf{w})$ ，パラメータ  $\mathbf{w}$  とする。fitting の誤差関数を次のように設計したとする。

$$\begin{aligned} E_n(\mathbf{w}) &= \frac{1}{2} (y_n - f(\mathbf{x}_n; \mathbf{w}))^2 \\ E(\mathbf{w}) &= \sum_{n=1}^N E_n(\mathbf{w}) \end{aligned}$$

やりたいことは、誤差関数を小さくするような  $\mathbf{w}$  を見つけることである。しかし、データの個数  $N$  が大きいほど、勾配  $\nabla E(\mathbf{w})$  を直接得る事そのものが難しくなる。

確率的勾配降下のアイデアは、「確率的に取り出した  $M < N$  となるデータの部分集合（ミニバッチと呼ばれる）に対して誤差関数を設定し、勾配を計算するための 1 回あたりの計算を減らす」というものである。具体的には、 $S$  をランダムに選択された  $M$  個のインデックスとして、次の誤差関数を扱う。

$$E_S(\mathbf{w}) = \frac{N}{M} \sum_{n \in S} E_n(\mathbf{w})$$

この期待値は  $E(\mathbf{w})$  と等価である（ように  $S$  がランダムに選ばれる）場合、元々やりたい  $E(\mathbf{w})$  の最適化をサンプル近似しながら解いているとラフに考えることができる。

学習係数のスケジューリングは、次の条件の下での  $\eta^{(t)}$  を用いれば、確率 1 で  $E(\mathbf{w})$  の停留点に収束する。

$$\sum_{t=1}^{\infty} \eta^{(t)} = \infty, \quad \sum_{t=1}^{\infty} [\eta^{(t)}]^2 < \infty$$

直感的には、 $\sum_{t=1}^{\infty} \eta^{(t)} = \infty$  で実行可能領域全体へ到達することを保証し、 $\sum_{t=1}^{\infty} [\eta^{(t)}]^2 < \infty$  で推定が収束することを保証している（和書なら [23] を参照）。これを満たすシンプルな方法は、正定数  $c$  に対して  $\eta^{(t)} = c/t$  である。

これ以降、バッチサイズが 1 である場合について書く。ステップ  $t$  におけるデータのインデックスを  $n^t$  とする。この時、誤差関数は  $E_{n^t}(\mathbf{w})$  である。

## 2.3 Momentum SGD

局所的に微分係数がゼロに近くなる場合（プラトー、平坦域とも）や、鞍点が存在する場合には、SGD の更新が遅くなってしまう。この問題の解決策の一つが慣性項の利用である。慣性係数  $\mu \in (0, 1]$  として、次のような更新を扱う。

$$\begin{aligned} m_i^{(t+1)} &= \mu m_i^{(t)} + \eta \left. \frac{\partial}{\partial w_i} E_{n^t}(\mathbf{w}) \right|_{\mathbf{w}^{(t)}} \\ w_i^{(t+1)} &= w_i^{(t)} - m_i^{(t+1)} \end{aligned}$$

この方法では、微分がゼロに近い場合でも過去の勾配の方向に基づいて更新が進むため、平坦域を抜けることができる。また、この手法は、ランダムサンプリングに起因するパラメータ勾配の振動を平滑化する特性を持っている。

## 2.4 学習係数の調整

AdaGrad [7] では、パラメータの要素ごとに異なる学習係数を設定する。

$$\begin{aligned} g_i &= \left. \frac{\partial}{\partial w_i} E(\mathbf{w}) \right|_{\mathbf{w}^{(t)}} \\ v_i^{(t+1)} &= v_i^{(t)} + g_i^2 \\ w_i^{(t+1)} &= w_i^{(t)} - \frac{\eta}{\sqrt{v_i^{(t+1)} + \epsilon_c}} g_i \end{aligned}$$

ここで、 $\epsilon_c$  はゼロ除算を防ぐための小さな定数である。

直感的な説明としては、大きな勾配を用いたパラメータ更新が続けば学習係数は早くなり、逆に小さな勾配を用いた更新が続けば現状の学習係数を維持してパラメータを更新する。（和書なら [24, 23] を参照）。

理論は（私には）難しいので省略（COLT で発表 → JMLR で完成版を発表という激強研究です）。

## 2.5 Adam

天下りの的に書くと、Adam アルゴリズムは以下である。

$$\begin{aligned}g_i &= \left. \frac{\partial}{\partial w_i} E(\mathbf{w}) \right|_{\mathbf{w}^{(t)}} \\m_i^{(t+1)} &= \beta_1 m_i^{(t)} + (1 - \beta_1) g_i \\v_i^{(t+1)} &= \beta_2 v_i^{(t)} + (1 - \beta_2) g_i^2 \\\hat{m}_i^{(t+1)} &= \frac{1}{1 - \beta_1^{m+1}} m_i^{(t+1)} \\\hat{v}_i^{(t+1)} &= \frac{1}{1 - \beta_2^{m+1}} v_i^{(t+1)} \\w_i^{(t+1)} &= w_i^{(t)} - \frac{\eta}{\sqrt{\hat{v}_i^{(t+1)} + \epsilon_c}} \hat{m}_i^{(t+1)}\end{aligned}$$

Adam は、いくつかの勾配法のアイデアを組み合わせたものになっている。

- 慣性項の利用： $\beta_1 m_i^{(t)} + (1 - \beta_1) g_i$  の部分。
- 学習係数のスケジューリング： $\frac{\eta}{\sqrt{\hat{v}_i^{(t+1)} + \epsilon_c}} \hat{m}_i^{(t+1)}$  の部分。
- 主に初期段階において大きくなるバイアスの修正： $\frac{1}{1 - \beta^{m+1}}$  の部分。

これ以上の説明は、[19] などを参照すること。

■ノート。 最近、オリジナル版の収束性に関する問題（そもそも非凸で確率的最適化やっているもので何であれ難しい）が提起され [18]，理論考察とそれに基づく変種が提案されたりしている（AdaBound [14] など）が，実用面・理論面の両方で adam をちゃんと超えた（と皆が納得できた）ものはまだ無い（という雰囲気がある）。

最近の強化学習実装（ICLR 2020 に通った研究）を眺めていると，Adam（あるいはそれより前の RMSProp など）を使っていれば今のところ文句は言われ無い（まあ，既存手法との比較をやる都合上，特に議論の対象にしたくない部分は既存手法の実験設定に合わせているだけなんだろうけど）。

### 3 再パラメータ化

変分ベイズをやる時に、解析的に得られない項がありサンプル近似をする必要がある。雑に言うと、再パラメータ化は、変数変換を用いてサンプル近似の性能を向上する方法である。

#### 3.1 復習：変分ベイズ

入出力データを  $(\mathbf{X}, \mathbf{Y})$ 、パラメータを  $\mathbf{w}$  とする。確率モデル  $p(\mathbf{Y}|\mathbf{w}, \mathbf{X})$  とする。事前分布  $p(\mathbf{w})$  として、Bayes 則に従って事後確率を次のように得る。

$$p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) = \frac{P(\mathbf{Y}|\mathbf{w}, \mathbf{X})P(\mathbf{w})}{P(\mathbf{Y}|\mathbf{X})} = \frac{P(\mathbf{Y}, \mathbf{w}|\mathbf{X})}{P(\mathbf{Y}|\mathbf{X})}$$

一般に、 $p(\mathbf{w}|\mathbf{Y}, \mathbf{X})$  は解析的な形にはならない。変分ベイズでは、 $q(\mathbf{w}; \boldsymbol{\xi})$  を用いて  $p(\mathbf{w}|\mathbf{Y}, \mathbf{X})$  を近似することを考える ( $\boldsymbol{\xi}$  は変分パラメータと呼ばれる)。近似の評価指標を KL ダイバージェンスとすると、

$$\begin{aligned} \text{KL}[q(\mathbf{w}; \boldsymbol{\xi})||p(\mathbf{w}|\mathbf{Y}, \mathbf{X})] &= \int q(\mathbf{w}; \boldsymbol{\xi}) \frac{\ln q(\mathbf{w}; \boldsymbol{\xi})}{\ln p(\mathbf{w}|\mathbf{Y}, \mathbf{X})} d\mathbf{w} = \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln q(\mathbf{w}; \boldsymbol{\xi}) - \ln p(\mathbf{w}|\mathbf{Y}, \mathbf{X})] d\mathbf{w} \\ &= \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln q(\mathbf{w}; \boldsymbol{\xi}) - \ln p(\mathbf{Y}, \mathbf{w}|\mathbf{X}) + \ln p(\mathbf{Y}|\mathbf{X})] d\mathbf{w} \\ &= \ln p(\mathbf{Y}|\mathbf{X}) - \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln p(\mathbf{Y}, \mathbf{w}|\mathbf{X}) - \ln q(\mathbf{w}; \boldsymbol{\xi})] d\mathbf{w} \\ &= \ln p(\mathbf{Y}|\mathbf{X}) - \mathcal{L}(\boldsymbol{\xi}) \end{aligned}$$

ここで、 $\mathcal{L}(\boldsymbol{\xi})$  は、変分下界 (ELBO) として知られる量である。 $\ln p(\mathbf{Y}|\mathbf{X})$  は  $\boldsymbol{\xi}$  に依存しないことに注意すれば、近似指標  $\text{KL}[q(\mathbf{w}; \boldsymbol{\xi})||p(\mathbf{w}|\mathbf{Y}, \mathbf{X})]$  を最小化する問題は、ELBO  $\mathcal{L}(\boldsymbol{\xi})$  を最大化する問題と等価である。ELBO は、次のように変形できる。

$$\begin{aligned} \mathcal{L}(\boldsymbol{\xi}) &= \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln p(\mathbf{Y}, \mathbf{w}|\mathbf{X}) - \ln q(\mathbf{w}; \boldsymbol{\xi})] d\mathbf{w} \\ &= \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X}) + \ln p(\mathbf{w}) - \ln q(\mathbf{w}; \boldsymbol{\xi})] d\mathbf{w} \\ &= \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X})] d\mathbf{w} - \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln q(\mathbf{w}; \boldsymbol{\xi}) - \ln p(\mathbf{w})] d\mathbf{w} \\ &= \int q(\mathbf{w}; \boldsymbol{\xi}) [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X})] d\mathbf{w} - \text{KL}[q(\mathbf{w}; \boldsymbol{\xi})||p(\mathbf{w})] \\ &= \mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \boldsymbol{\xi})} [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X})] - \text{KL}[q(\mathbf{w}; \boldsymbol{\xi})||p(\mathbf{w})] \end{aligned}$$

基本的には、ELBO の最大化問題を勾配法で解いていくことになるので、上式の各項の  $\boldsymbol{\xi}$  に関する微分が欲しいものである。ここで、 $\text{KL}[q(\mathbf{w}; \boldsymbol{\xi})||p(\mathbf{w})]$  の微分は、事前分布  $p(\mathbf{w})$  と変分分布  $q(\mathbf{w}; \boldsymbol{\xi})$  を  $\mathbf{w}$  の正規分布などで与えれば、解析的に計算することができる。一方で、 $\mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \boldsymbol{\xi})} [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X})]$  の微分は解析的には計算できない。

これを扱うシンプルな方法は、分布  $q(\mathbf{w}; \boldsymbol{\xi})$  から  $\mathbf{w}$  のサンプルを生成してサンプル近似することであるが、実用上は高い分散を生じてしまうことが知られている。

### 3.2 reparameterization trick [6]

与えられた  $q(\mathbf{w}; \boldsymbol{\xi})$  に対して,  $q(\mathbf{g}(\boldsymbol{\xi}; \epsilon); \boldsymbol{\xi})d\boldsymbol{\xi} = p(\epsilon)d\epsilon$  を満たす変数変換  $\mathbf{w} = \mathbf{g}(\boldsymbol{\xi}; \epsilon)$  と分布  $p(\epsilon)$  を定義できる場合には, 次が成り立つ (途中の  $(\cdot)_{i,j}$  は行列の成分表示,  $(\cdot)_j$  はベクトルの成分表示を指す).

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\xi}} \mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \boldsymbol{\xi})} [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X})] &= \frac{\partial}{\partial \boldsymbol{\xi}} \mathbb{E}_{\epsilon \sim p(\epsilon)} [\ln p(\mathbf{Y}|\mathbf{g}(\boldsymbol{\xi}; \epsilon), \mathbf{X})] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[ \frac{\partial}{\partial \boldsymbol{\xi}} \ln p(\mathbf{Y}|\mathbf{g}(\boldsymbol{\xi}; \epsilon), \mathbf{X}) \right] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[ \left( \frac{\partial \mathbf{g}(\boldsymbol{\xi}; \epsilon)}{\partial \boldsymbol{\xi}} \right) \left( \frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X}) \Big|_{\mathbf{w}=\mathbf{g}(\boldsymbol{\xi}; \epsilon)} \right) \right] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[ \left( \frac{\partial g_j(\boldsymbol{\xi}; \epsilon)}{\partial \xi_i} \right)_{i,j} \left( \frac{\partial}{\partial w_j} \ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X}) \Big|_{\mathbf{w}=\mathbf{g}(\boldsymbol{\xi}; \epsilon)} \right)_j \right] \end{aligned}$$

よって, 1 回のサンプル値  $\tilde{\epsilon} \sim p(\epsilon)$  によって近似すると,

$$\frac{\partial}{\partial \boldsymbol{\xi}} \mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \boldsymbol{\xi})} [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X})] \approx \frac{\partial}{\partial \boldsymbol{\xi}} \ln p(\mathbf{Y}|\mathbf{g}(\boldsymbol{\xi}; \tilde{\epsilon}), \mathbf{X}) = \left( \frac{\partial g_j(\boldsymbol{\xi}; \tilde{\epsilon})}{\partial \xi_i} \right)_{i,j} \left( \frac{\partial}{\partial w_j} \ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X}) \Big|_{\mathbf{w}=\mathbf{g}(\boldsymbol{\xi}; \tilde{\epsilon})} \right)_j$$

■具体例: 正規分布の場合.  $q(\mathbf{w}; \boldsymbol{\xi})$  として正規分布が与えられたとする (従って,  $\boldsymbol{\xi} = (\mu_\xi, \sigma_\xi^2)^T$  である).

$$q(w; \boldsymbol{\xi}) = \mathcal{N}(w|\mu_\xi, \sigma_\xi^2) = \frac{1}{\sqrt{2\pi\sigma_\xi^2}} \exp\left(-\frac{(w - \mu_\xi)^2}{2\sigma_\xi^2}\right)$$

ここで, 変数変換  $w = g(\boldsymbol{\xi}; \epsilon)$  と分布  $p(\epsilon)$  を次のように定義する.

$$\begin{aligned} g(\boldsymbol{\xi}; \epsilon) &= w = \mu_\xi + \sigma_\xi \epsilon \\ p(\epsilon) &= \mathcal{N}(0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\epsilon^2}{2}\right) \end{aligned}$$

$dw = \frac{dw}{d\epsilon} d\epsilon = \sigma_\xi d\epsilon$  であるため, 次が成り立つ.

$$\begin{aligned} q(w; \boldsymbol{\xi})dw &= \frac{1}{\sqrt{2\pi\sigma_\xi^2}} \exp\left(-\frac{((\mu_\xi + \sigma_\xi \epsilon) - \mu_\xi)^2}{2\sigma_\xi^2}\right) \sigma_\xi d\epsilon \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\epsilon^2}{2}\right) d\epsilon = p(\epsilon)d\epsilon \end{aligned}$$

よって, この変数変換と分布は,  $q(w; \boldsymbol{\xi})dw = p(\epsilon)d\epsilon$  を満たしている. そのため, 再パラメータ化が適用できて,

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\xi}} \mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \boldsymbol{\xi})} [\ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X})] &= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[ \left( \frac{\partial g(\boldsymbol{\xi}; \epsilon)}{\partial \boldsymbol{\xi}} \right) \left( \frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X}) \Big|_{\mathbf{w}=\mathbf{g}(\boldsymbol{\xi}; \epsilon)} \right) \right] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[ \left( \begin{array}{c} \frac{\partial g(\boldsymbol{\xi}; \epsilon)}{\partial \mu_\xi} \\ \frac{\partial g(\boldsymbol{\xi}; \epsilon)}{\partial \sigma_\xi} \end{array} \right) \left( \frac{\partial}{\partial w} \ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X}) \Big|_{\mathbf{w}=\mathbf{g}(\boldsymbol{\xi}; \epsilon)} \right) \right] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[ \left( \begin{array}{c} 1 \\ \epsilon \end{array} \right) \left( \frac{\partial}{\partial w} \ln p(\mathbf{Y}|\mathbf{w}, \mathbf{X}) \Big|_{\mathbf{w}=\mathbf{g}(\boldsymbol{\xi}; \epsilon)} \right) \right] \end{aligned}$$

■ノート. 再パラメータ化の拡張や, 効率化を与えるメカニズムの解析は, 現在でも研究トピックの一つである (例えば [22]).

## 4 正則化 (regularization)

### 4.1 weight decay

要するに L2 正則化であり, NN の分野ではこれを weight decay と呼んでいる. 過剰適合を防ぐ目的で,  $\mathbf{w}$  の大きさに対してペナルティを置いてパラメータを最適化する. L2 正則化では, 2 次のペナルティ項を用いる.

$$J(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

これは, 事前分布を  $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I})$  とする MAP 推定としても解釈できる (後述).

■MAP 推定. 次の回帰モデルを考える.

$$y = f(\mathbf{x}; \mathbf{w}) + \epsilon$$

ここで,  $\epsilon \sim \mathcal{N}(\epsilon|0, \sigma^2)$  とする.  $y$  の確率分布は, 次のように得られる.

$$\begin{aligned} p(y|\mathbf{x}, \mathbf{w}) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - f(\mathbf{x}; \mathbf{w}))^2}{2\sigma^2}\right) \\ \ln p(y|\mathbf{x}, \mathbf{w}) &= -\frac{(y - f(\mathbf{x}; \mathbf{w}))^2}{2\sigma^2} - \frac{1}{2} \ln(2\pi\sigma^2) \end{aligned}$$

したがって, 対数尤度関数は次のように得られる.

$$\ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - f(\mathbf{x}_n; \mathbf{w}))^2 - \frac{N}{2} \ln(2\pi\sigma^2)$$

また, 対数事前分布は次のように得られる.

$$\ln p(\mathbf{w}) = -\frac{\lambda}{2} \mathbf{w}^T \mathbf{w} - \frac{D}{2} \ln(2\pi\lambda^{-1})$$

ここで,  $D$  はパラメータ  $\mathbf{w}$  の次元である. 事後確率最大化は, 次のように書くことができる.

$$\begin{aligned} \mathbf{w}_{MAP} &= \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{y}, \mathbf{X}) \\ &= \arg \max_{\mathbf{w}} [\ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \ln p(\mathbf{w})] \\ &= \arg \max_{\mathbf{w}} \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - f(\mathbf{x}_n; \mathbf{w}))^2 - \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right] \\ &= \arg \min_{\mathbf{w}} \left[ \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - f(\mathbf{x}_n; \mathbf{w}))^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right] \end{aligned}$$

そのため, 誤差関数  $\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - f(\mathbf{x}_n; \mathbf{w}))^2$  と, ペナルティ  $\frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$  とした最適化問題と等価である.

### 4.2 dropout

#### 4.2.1 実装

■MLP 順伝播 (再掲).

$$\begin{aligned} \mathbf{z}_n^{(\ell)} &= \begin{cases} \mathbf{x}_n & \text{if } \ell = 0 \\ \phi(\mathbf{a}_n^{(\ell)}) & \text{if } \ell \in \{1, \dots, L-1\} \end{cases} \\ \mathbf{a}_n^{(\ell)} &= \mathbf{W}^{(\ell)} \mathbf{z}_n^{(\ell-1)}, \quad \ell = 1, \dots, L \\ \mathbf{y}_n &= \mathbf{a}_n^{(L)} + \epsilon_n \end{aligned}$$



dropout では、入力  $\mathbf{x}_n = \mathbf{z}_n^{(0)}$  と中間層  $\mathbf{z}_n^{(\ell)}$  に対してランダムにゼロを与えるような操作を行う。中間層の要素  $z_{n,i}^{(\ell)}$  にゼロを与えるかどうかを示すマスクを  $m_i^{(\ell)} \in \{0, 1\}$  とし、ベルヌーイ分布に従って得る。  $m_i^{(\ell)}$  のベクトルを  $\mathbf{m}^{(\ell)}$  とする。  $\mathbf{m}^{(\ell)}$  の成分が対角項に並んだ対角行列を、  $[\text{diag}\{\mathbf{m}^{(\ell)}\}]$  とする。この時、dropout を用いた NN 順伝播は次のように書ける。

■MLP 順伝播 with dropout.

$$\begin{aligned} \mathbf{z}_n^{(\ell)} &= \begin{cases} \mathbf{x}_n & \text{if } \ell = 0 \\ \phi(\mathbf{a}_n^{(\ell)}) & \text{if } \ell \in \{1, \dots, L-1\} \end{cases} \\ \tilde{\mathbf{z}}_n^{(\ell)} &= [\text{diag}\{\mathbf{m}^{(\ell)}\}] \mathbf{z}_n^{(\ell)}, \quad \ell = 1, \dots, L \\ \mathbf{a}_n^{(\ell)} &= \mathbf{W}^{(\ell)} \tilde{\mathbf{z}}_n^{(\ell-1)}, \quad \ell = 1, \dots, L \\ \mathbf{y}_n &= \mathbf{a}_n^{(L)} + \epsilon_n \end{aligned}$$

逆伝播も、  $\mathbf{z}_n^\ell$  を  $\tilde{\mathbf{z}}_n^\ell$  で置き換えて導出する。2層 NN の場合は、1.1 節の導出の途中から書いていくと、次のようになる。

$$\begin{aligned} \frac{\partial E_n(\mathbf{W})}{\partial w_{h_1, h_0}^{(1)}} &= \left\{ \sum_{d=1}^D \delta_{n,d}^{(2)} \left[ w_{d, h_1}^{(2)} \frac{\partial \tilde{z}_{n, h_1}^{(1)}}{\partial a_{n, h_1}^{(1)}} \right] \right\} [\tilde{z}_{n, h_0}^{(0)}] = \left\{ \sum_{d=1}^D \delta_{n,d}^{(2)} \left[ w_{d, h_1}^{(2)} m_{h_1}^{(1)} \frac{\partial z_{n, h_1}^{(1)}}{\partial a_{n, h_1}^{(1)}} \right] \right\} [m_{h_0}^{(0)} z_{n, h_0}^{(0)}] \\ &= \left\{ \sum_{d=1}^D \delta_{n,d}^{(2)} \left[ w_{d, h_1}^{(2)} m_{h_1}^{(1)} \phi'(a_{n, h_1}^{(1)}) \right] \right\} [m_{h_0}^{(0)} z_{n, h_0}^{(0)}] = \left\{ m_{h_1}^{(1)} \phi'(a_{n, h_1}^{(1)}) \sum_{d=1}^D \delta_{n,d}^{(2)} [w_{d, h_1}^{(2)}] \right\} [m_{h_0}^{(0)} z_{n, h_0}^{(0)}] \end{aligned}$$

同様に、逆伝播は次のように書ける。

■MLP 逆伝播（再掲）。

$$\begin{aligned} \frac{\partial E_n(\mathbf{W})}{\partial w_{d, h_1}^{(L)}} &= \left[ \frac{\partial E_n(\mathbf{W})}{\partial a_{n, d}^{(L)}} \right] [z_{n, h_{L-1}}^{(L-1)}] = \delta_{n, d}^{(L)} z_{n, h_{L-1}}^{(L-1)} \\ \frac{\partial E_n(\mathbf{W})}{\partial w_{h_\ell, h_{\ell-1}}^{(\ell)}} &= \left[ \phi'(a_{n, h_\ell}^{(\ell)}) \sum_{h_{\ell+1}=1}^{H_{\ell+1}} \delta_{n, h_{\ell+1}}^{(\ell+1)} w_{h_{\ell+1}, h_\ell}^{(\ell+1)} \right] [z_{n, h_{\ell-1}}^{(\ell-1)}] = \delta_{n, h_\ell}^{(\ell)} z_{n, h_{\ell-1}}^{(\ell-1)}, \quad \ell = 1, \dots, L-1 \end{aligned}$$

■MLP 逆伝播 with dropout.

$$\begin{aligned} \frac{\partial E_n(\mathbf{W})}{\partial w_{d, h_1}^{(L)}} &= \left[ \frac{\partial E_n(\mathbf{W})}{\partial a_{n, d}^{(L)}} \right] [m_{h_{L-1}}^{(L-1)} z_{n, h_{L-1}}^{(L-1)}] = m_{h_{L-1}}^{(L-1)} \times \delta_{n, d}^{(L)} z_{n, h_{L-1}}^{(L-1)} \\ \frac{\partial E_n(\mathbf{W})}{\partial w_{h_\ell, h_{\ell-1}}^{(\ell)}} &= \left[ m_{h_\ell}^{(\ell)} \phi'(a_{n, h_\ell}^{(\ell)}) \sum_{h_{\ell+1}=1}^{H_{\ell+1}} \delta_{n, h_{\ell+1}}^{(\ell+1)} w_{h_{\ell+1}, h_\ell}^{(\ell+1)} \right] [m_{h_{\ell-1}}^{(\ell-1)} z_{n, h_{\ell-1}}^{(\ell-1)}] = m_{h_\ell}^{(\ell)} m_{h_{\ell-1}}^{(\ell-1)} \times \delta_{n, h_\ell}^{(\ell)} z_{n, h_{\ell-1}}^{(\ell-1)}, \quad \ell = 1, \dots, L-1 \end{aligned}$$

■勾配 with dropout (and 正則化). dropout 付き順伝播における  $\mathbf{a}_n^{(\ell)}$  を求める式は、次のようにも書ける。

$$\mathbf{a}_n^{(\ell)} = \mathbf{W}^{(\ell)} [\text{diag}\{\mathbf{m}^{(\ell)}\}] \mathbf{z}_n^{(\ell-1)} = \tilde{\mathbf{W}}^{(\ell)} \mathbf{z}_n^{(\ell-1)}, \quad \ell = 1, \dots, L$$

$\tilde{\mathbf{W}}^{(\ell)}$  の集合を  $\tilde{\mathbf{W}}$  とする。  $\mathbf{W}$  と  $\mathbf{m}$  から  $\tilde{\mathbf{W}}$  を得る操作を、  $\tilde{\mathbf{W}} = \mathbf{g}(\mathbf{W}, \mathbf{m})$  と表す（具体的には、各層  $\ell$  で  $\tilde{\mathbf{W}}^{(\ell)} = \mathbf{W}^{(\ell)} [\text{diag}\{\mathbf{m}^{(\ell)}\}]$  をする）。  $\mathbf{a}_n^{(L)} = \mathbf{f}(\mathbf{x}_n; \tilde{\mathbf{W}}) = \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m}))$  と表記すると、回帰モデルは次のように書ける。

$$\mathbf{y}_n = \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m})) + \epsilon_n$$

ノイズ項の分布を正規分布  $\epsilon_n \sim \mathcal{N}(\epsilon_n | \mathbf{0}, \sigma^2 \mathbf{I})$  で仮定すると、対数尤度関数は

$$\ln p(\mathbf{y}_n | \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m}))) = -\frac{1}{2\sigma^2} (\mathbf{y}_n - \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m})))^T (\mathbf{y}_n - \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m}))) - \frac{1}{2} \ln(2\pi\sigma^2)$$

ミニバッチのデータ数を  $M$ 、インデックス集合を  $\mathcal{S}$  とする。評価関数（最小化したい）を、二乗誤差プラス正則化項とする。

$$\begin{aligned} J(\mathbf{W}) &= \frac{1}{2M} \sum_{n \in \mathcal{S}} (\mathbf{y}_n - \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m})))^T (\mathbf{y}_n - \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m}))) + \frac{1}{2} \sum_{\ell=1}^L \lambda_{\ell} \|\mathbf{W}^{(\ell)}\| \\ &= - \left\{ \sigma^2 \sum_{n \in \mathcal{S}} \ln p(\mathbf{y}_n | \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m}))) - \frac{1}{2} \sum_{\ell=1}^L \lambda_{\ell} \|\mathbf{W}^{(\ell)}\| \right\} + \text{const.} \end{aligned}$$

勾配を取ると、

$$\nabla_{\mathbf{W}} J(\mathbf{W}) = - \left\{ \sigma^2 \sum_{n \in \mathcal{S}} \nabla_{\mathbf{W}} \ln p(\mathbf{y}_n | \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\mathbf{W}, \mathbf{m}))) - \frac{1}{2} \sum_{\ell=1}^L \lambda_{\ell} \nabla_{\mathbf{W}} \|\mathbf{W}^{(\ell)}\| \right\}$$

なお、マスク  $\mathbf{m}$  は確率的に得られていることに注意する。

#### 4.2.2 変分ベイズとしての解釈 [8]

■勾配 with dropout の記号の書き換え。 前小節で導出した勾配の記号を、 $\mathbf{W} \rightarrow \boldsymbol{\xi}$  と  $\mathbf{m} \rightarrow \tilde{\epsilon}$  と書き替える。

$$\nabla_{\boldsymbol{\xi}} J(\boldsymbol{\xi}) = - \left\{ \sigma^2 \sum_{n \in \mathcal{S}} \nabla_{\boldsymbol{\xi}} \ln p(\mathbf{y}_n | \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\boldsymbol{\xi}, \tilde{\epsilon}))) - \frac{1}{2} \sum_{\ell=1}^L \lambda_{\ell} \nabla_{\boldsymbol{\xi}} \|\boldsymbol{\xi}^{(\ell)}\| \right\} \quad (1)$$

■変分ベイズ with 再パラメータ化勾配（3 節の導出をまとめると得られる）。

$$\begin{aligned} \nabla_{\boldsymbol{\xi}} \mathcal{L}(\boldsymbol{\xi}) &= \nabla_{\boldsymbol{\xi}} \mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \boldsymbol{\xi})} [\ln p(\mathbf{Y} | \mathbf{w}, \mathbf{X})] - \nabla_{\boldsymbol{\xi}} \text{KL}[q(\mathbf{w}; \boldsymbol{\xi}) || p(\mathbf{w})] \\ &\approx \nabla_{\boldsymbol{\xi}} \ln p(\mathbf{Y} | \mathbf{g}(\boldsymbol{\xi}, \tilde{\epsilon}), \mathbf{X}) - \nabla_{\boldsymbol{\xi}} \text{KL}[q(\mathbf{w}; \boldsymbol{\xi}) || p(\mathbf{w})] \end{aligned}$$

ここで、 $\mathbf{w} = \mathbf{g}(\boldsymbol{\xi}; \epsilon)$  と  $p(\epsilon)$  は再パラメータ化の条件を満たしていて、また  $\tilde{\epsilon}$  は  $p(\epsilon)$  から得られたサンプルである。尤度関数を  $p(\mathbf{Y} | \mathbf{w}, \mathbf{X}) = \prod_{n \in \mathcal{S}} p(\mathbf{y}_n | \mathbf{f}(\mathbf{x}_n; \mathbf{w}))$  と書くとき、

$$\nabla_{\boldsymbol{\xi}} \mathcal{L}(\boldsymbol{\xi}) \approx \sum_{n \in \mathcal{S}} \nabla_{\boldsymbol{\xi}} \ln p(\mathbf{y}_n | \mathbf{f}(\mathbf{x}_n; \mathbf{g}(\boldsymbol{\xi}, \tilde{\epsilon}))) - \nabla_{\boldsymbol{\xi}} \text{KL}[q(\mathbf{w}; \boldsymbol{\xi}) || p(\mathbf{w})] \quad (2)$$

■比較。 定数倍を無視して式 (1) と (2) を見比べると、第 1 項についてはそれぞれ同じ形で表現できる。そのため、式 (1) の第 2 項（正則化項の勾配←データに依存しない）と、式 (2) の第 2 項（事前分布・変分分布間の KL ダイバージェンス←データに依存しない）を同じ形で表現できるなら、dropout は変分ベイズの枠組みの枠組みで解釈できる。これは、特定の  $q(\mathbf{w}; \boldsymbol{\xi})$  と  $p(\mathbf{w})$  を選ぶことにより成立させることができる（条件は [8] とその関連研究を参照）。

これらをまとめると、

- dropout におけるパラメータ  $\mathbf{W}$  は、変分ベイズにおけるハイパラメータ  $\boldsymbol{\xi}$  に相当する。
- dropout におけるマスク  $\mathbf{m}$  のサンプリングは、再パラメータ化における  $\epsilon$  のサンプリングに相当する。
- dropout の背後には、暗に確率変数（式 (2) における  $\mathbf{w}$  相当）が存在しているが、周辺化されているので表には出てこない。

一般に、確率変数の周辺化は、確率変数の点推定と比較して overfitting に対してロバストである（その代償に計算量をより多く必要とする）と理解される（例：NN における重みパラメータの点推定よりも、Bayesian NN における重みパラメータの周辺化の方がロバストである）。上述の議論より、dropout はある種の周辺化を暗にやっていると解釈することができて、それが dropout において経験的に得られているロバスト性の理由だと考えられている。

## 5 正規化 (normalization)

### 5.1 batch normalization [12]

多層 NN のパラメータを更新する場合、前の層のパラメータの変化によって次の層に対する入力の変動が大きく変化してしまい、学習が不安定になりやすい。正規化のアイデアは、入力を標準正規（ガウス）分布を使って変換して用いることにより、学習を安定化させようというものである。

$n$  番目のデータの  $\ell$  層目の中間変数を、 $\mathbf{z}_n^{(\ell)} = (z_{n,1}^{(\ell)}, z_{n,2}^{(\ell)}, \dots)^T$  と表記する。ミニバッチに含まれるデータのインデックスの集合を  $\mathcal{S}$  とする。具体的には、次の変換を用いる。

$$\begin{aligned}\mu_d^{(\ell)} &= \frac{1}{m} \sum_{n \in \mathcal{S}} z_{n,d}^{(\ell)} \\ \sigma_d^{2(\ell)} &= \frac{1}{m} \sum_{n \in \mathcal{S}} (z_{n,d}^{(\ell)} - \mu_d^{(\ell)})^2 \\ \hat{z}_{n,d}^{(\ell)} &= \frac{z_{n,d}^{(\ell)} - \mu_d^{(\ell)}}{\sqrt{\sigma_d^{2(\ell)} + \epsilon_c}} \\ \tilde{z}_{n,d}^{(\ell)} &= \gamma_d^{(\ell)} \hat{z}_{n,d}^{(\ell)} + \beta_d^{(\ell)}\end{aligned}$$

ここで、 $\gamma_d$  と  $\beta_d$  は学習可能な係数である。

■ノート。 バッチ正規化が収束と汎化を改善する事、つまり学習安定化／効率化と overfitting 低減化をしている事は、多くの研究において経験的には得られている（だからみんな使っている）。また、理解の試みもなされている [15]。強化学習でも、汎化性能を改善するみたいな話（要検討）がある [5]。

### 5.2 layer normalization [2]

別の空間で正規化する。[3]（ICLR 2020 で興味深いと言われつつ検証不足で reject されたもの）を見てると、強化学習において効率的になるかは微妙な感じ？

### 5.3 weight normalization [20, 15]

別の空間で正規化する。いくつかの強化学習実装（例えば [4, 11]）でも見かける。

## 6 雑感と文献メモ

- 各項目について, [25] が和書でだいたい網羅してくれているので, まずはこれの該当箇所を読めばいいと思う.
- 活性化関数は, ここ最近の強化学習実装 (例えば SAC [9]) ではほぼ ReLU 一択なのが現状. ちょっと古い強化学習実装 (例えば TRPO [21] とか) だと, tanh とかも見かける (ReLU vs tanh の比較は [10] を見ると良い). 強化学習に限らず, swish [17] とか mish [16] とか提案されてはいるものの, 実用面は ReLU 一択で決着がついた印象.
- アーキテクチャに関するのは, (1) mujoco など状態が関節角なら MLP, (2) 画像を状態とするなら CNN, (3) 完全状態観測ではなく履歴が欲しいなら (あるいは動的な特徴量が欲しいなら) RNN や LSTM, を使っている印象. 最新の強化学習でもまだ最近の深層学習アーキテクチャを使いこなしていない (あるいは手が回っていない) 雰囲気なので, MLP さえ押さえておけば我々にとっては十分だと思う.

## 参考文献

- [1] Joshua Achiam. Spinning Up in Deep Reinforcement Learning. *GitHub repository*, 2018. <https://github.com/openai/spinningup>.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] Aditya Bhatt, Max Argus, Artemij Amiranashvili, and Thomas Brox. Crossnorm: Normalization for off-policy td reinforcement learning. *arXiv preprint arXiv:1902.05605*, 2019.
- [4] Ignasi Clavera, Jonas Rothfuss, John Schulman, Yasuhiro Fujita, Tamim Asfour, and Pieter Abbeel. Model-based reinforcement learning via meta-policy optimization. In *Conference on Robot Learning*, pages 617–629, 2018.
- [5] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In *International Conference on Machine Learning*, pages 1282–1289, 2019.
- [6] P Kingma Diederik and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- [7] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(Jul):2121–2159, 2011.
- [8] Yarın Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.
- [9] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870, 2018.
- [10] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *AAAI Conference on Artificial Intelligence*, 2018.
- [11] Zhewei Huang, Wen Heng, and Shuchang Zhou. Learning to paint with model-based deep reinforcement learning. In *The IEEE International Conference on Computer Vision*, 2019.
- [12] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by

- reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015.
- [13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [14] Liangchen Luo, Yuanhao Xiong, and Yan Liu. Adaptive gradient methods with dynamic bound of learning rate. In *International Conference on Learning Representations*, 2019.
- [15] Ping Luo, Xinjiang Wang, Wenqi Shao, and Zhanglin Peng. Towards understanding regularization in batch normalization. In *International Conference on Learning Representations*, 2019.
- [16] Diganta Misra. Mish: A self regularized non-monotonic neural activation function. *arXiv preprint arXiv:1908.08681*, 2019.
- [17] Quoc V. Le Prajit Ramachandran, Barret Zoph. Searching for activation functions, 2018.
- [18] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- [19] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [20] Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in neural information processing systems*, pages 901–909, 2016.
- [21] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897, 2015.
- [22] Ming Xu, Matias Quiroz, Robert Kohn, and Scott A Sisson. Variance reduction properties of the reparameterization trick. In *International Conference on Artificial Intelligence and Statistics*, pages 2711–2720, 2019.
- [23] 鈴木 大慈. 確率的最適化 = *Stochastic optimization*. MLP 機械学習プロフェッショナルシリーズ. 講談社, 2015.
- [24] 原田 達也. 画像認識 = *Image recognition*. MLP 機械学習プロフェッショナルシリーズ. 講談社, 2017.
- [25] 須山 敦志. ベイズ深層学習 = *Bayesian deep learning*. MLP 機械学習プロフェッショナルシリーズ. 講談社, 2019.