

Regularization method in Neural Network

Author : oilneck の人

1 tangent propagation

1.1 Tikhonov Regularization

入力ベクトル \mathbf{x}_n の変換 $\mathbf{s}(\mathbf{x}_n, \xi)$ を考える．ただし， $\mathbf{s}(\mathbf{x}_n, 0) = \mathbf{x}_n$ であり，例えば乱数ノイズを付加する入力の変換に対しては $\mathbf{s}(\mathbf{x}, \xi) = \mathbf{x} + \xi$ と変換される．正則化項係数 Ω を用いて誤差関数が $\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda\Omega$ (ただし λ は正則化の調節パラメータ) と表されるとき，次の正則化係数を考えたい．

$$\Omega = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left(\left. \frac{\partial y_{nk}}{\partial \xi} \right|_{\xi=0} \right)^2 \quad (1.1)$$

$y_{nk} = y_k |_{\mathbf{x}_n}$ であり， n 番目の訓練データに対してニューラルネットワークを順伝播した際， k 番目の出力ユニット値を意味する． $\Omega = \sum_n \Omega_n$ と記述するとき

$$\Omega_n = \frac{1}{2} \sum_{k=1}^K \left(\sum_{i=1}^D \underbrace{\frac{\partial y_{nk}}{\partial x_{ni}} \frac{\partial x_{ni}}{\partial \xi}}_{\tau_{ni}} \bigg|_{\xi=0} \right)^2 = \frac{1}{2} \sum_k \left(\sum_i J_{nki} \tau_{ni} \right)^2 \quad (1.2)$$

のように正則化項 (Tikhonov 正則化項) は表現される． τ は入力の接線情報を意味する．また， \sum_i は入力ユニットに関しての和の操作である．式 (1.2) 中に出てくる

$$J_{nki} = \left. \frac{\partial y_k}{\partial x_i} \right|_{\mathbf{x}_n} \quad (1.3)$$

はサンプル n の Jacobi 行列の (k, i) 成分を意味する．

1.2 Back propagation

以下に誤差逆伝播法のアルゴリズムをまとめる． $w^{(\ell)}$ は ℓ 番目の層における重み行列の成分を表す．

Error back-propagation

誤差関数が $E(\mathbf{w}) = \sum_n E_n(\mathbf{w})$ で与えられるとき，入力層 (第 0 層) と L 層ある出力ユニットから成るニューラルネットワークにおける誤差逆伝播法のアルゴリズムは以下で与えられる．

1. 活性化関数を $h(\cdot)$ と表すとき，次式を用いて順伝播を行う．

$$a_j^{(\ell)} = \sum_i w_{ji}^{(\ell)} z_i^{(\ell-1)} \quad z_j^{(\ell)} \equiv h(a_j^{(\ell)}) = h\left(\sum_i w_{ji}^{(\ell)} z_i^{(\ell-1)}\right) \quad (1.4)$$

2. 出力ユニットの誤差 $\delta^{(L)}$ 及びネットワーク内全ての隠れユニットの誤差 δ_j を計算する．

$$\begin{cases} \delta_j^{(\ell)} = h'(a_j^{(\ell)}) \sum_k w_{kj}^{(\ell+1)} \delta_k^{(\ell+1)} & [1 \leq \ell \leq L-1] \\ \delta_j^{(L)} = y_j - t_j \end{cases} \quad (1.5)$$

3. ユニット誤差を利用して誤差関数の各重みパラメータに関する微分を評価する．

$$\frac{\partial E_n}{\partial w_{ji}^{(\ell)}} = \delta_j^{(\ell)} z_i^{(\ell-1)} \quad (1.6)$$

式 (1.4)~(1.6) の公式は以下の導出によって得られる．まず，誤差関数が訓練集合の各データに対する誤差の和で表されると仮定する．

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \quad (1.7)$$

一般のフィードフォワードネットワークでは、それぞれのユニットで

$$a_j = \sum_i w_{ji} z_i, \quad z_j = h(a_j) \quad (1.8)$$

の形の出力 z_j が出現する．入力 z_i が出力ユニット a_j に影響を与え、その a_j が非線形活性化関数 $h(\cdot)$ を通して次の段の出力 z_j に信号を出す．ある特定のパターン E_n の重み w_{ji} に関する微分を考える． E_n は非線形活性化関数 $h(\cdot)$ の変数 a_j を介して w_{ji} に依存しているので

$$\frac{\partial E_n}{\partial w_{ji}} = \underbrace{\frac{\partial E_n}{\partial a_j}}_{\delta_j} \underbrace{\frac{\partial a_j}{\partial w_{ji}}}_{z_i} = \delta_j z_i \quad (1.9)$$

となる．ここで、式 (1.8) の線形和 a_j の結果より

$$\frac{\partial a_j}{\partial w_{ji}} = z_i$$

が成り立つことを用いた．また、誤差に関する便利な記法

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} \quad (1.10)$$

を導入した．次に式 (1.10) の誤差 δ_j に関して議論する．ロジスティックシグモイド活性化関数と交差エントロピー誤差関数の組み合わせや、多クラス交差エントロピー誤差関数とソフトマックス活性化関数の組み合わせのような、 $h(\cdot)$ に正準連結関数を用いたスキームでは

$$\delta_j = \frac{\partial E_n}{\partial a_j} = y_j - t_j \quad (1.11)$$

の性質が成り立つ．この性質は正準連結関数を $h(\cdot)$ に選んだ場合の一般的な結果である．上式は単純に誤差関数に 2 乗和を選んだ際、出力ユニット値が活性の線形出力に対して $y_j = a_j$ が成り立つことから、

$$\frac{\partial E_n}{\partial a_j} = \frac{\partial}{\partial a_j} \left[\frac{1}{2} \sum_{k=1}^K (y_k - t_k)^2 \right] = y_j - t_j$$

と導けばよい．次に隠れ層における δ_j の評価について考える．ユニット j につながっているユニット k を通して E_n 中にある a_j への影響があることを考えると、微分の連鎖則より

$$\frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} \quad (1.12)$$

と変形できる．ここで、式 (1.8) より a_j と a_k は z を経由して関係し

$$a_k = \sum_i w_{ki} h(a_i) \implies \frac{\partial a_k}{\partial a_j} = w_{kj} h'(a_j)$$

であることと、式 (1.10) の表記と合わせて、式 (1.12) は次の逆伝播の公式に変形される．

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k \quad (1.13)$$

これまでに得られた結果である式 (1.4)~(1.5) はベクトル形式で表すとより簡潔になる．

1. Forward propagation

$$\mathbf{a}^{(\ell)} = \hat{\mathbf{W}}^{(\ell)} \mathbf{z}^{(\ell-1)} \quad \mathbf{z}^{(\ell)} \equiv h\left(\mathbf{a}^{(\ell)}\right) = h\left(\hat{\mathbf{W}}^{(\ell)} \mathbf{z}^{(\ell-1)}\right) \quad (1.14)$$

2. Error Back propagation

$$\begin{cases} \boldsymbol{\delta}^{(\ell)} = h'\left(\mathbf{a}^{(\ell)}\right) \circ \left[\hat{\mathbf{W}}^{(\ell+1)}\right]^T \boldsymbol{\delta}^{(\ell+1)} & [1 \leq \ell \leq L-1] \\ \boldsymbol{\delta}^{(L)} = \mathbf{y} - \mathbf{t} \end{cases} \quad (1.15)$$

3. Evaluation of error gradient

$$\frac{\partial E_n}{\partial \hat{\mathbf{W}}^{(\ell)}} = \boldsymbol{\delta}^{(\ell)} \otimes \mathbf{z}^{(\ell-1)} \quad (1.16)$$

ただし、以上の表記で \circ はベクトルの対応する成分同士の積を意味するアダマール積であるとした。また、 \otimes は直積の記号を意味し、 $\mathbf{a} \otimes \mathbf{b} = \mathbf{ab}^T$ の略記である。直積の結果は行列となることに注意されたい。これは、物理でいうところの内積 $\langle a|b \rangle$ がスカラー出力を与えるのに対し、 $|a\rangle \langle b|$ は演算子として作用するのに似ている。その意味で直積の演算は外積 (内積の対義語) とも呼ばれる。python 言語を用いてプログラミングする際はこの直積を使った表式のほうが実装に便利である。演算は例えば numpy パッケージの `outer` 関数を用いればよい。式 (1.16) のようにひとたび誤差関数の重み勾配が求まると、後は以下に示す確率勾配降下法によって重みの更新を行い、この操作を適宜繰り返すことになる。

$$\hat{\mathbf{W}}_{(\text{new})}^{(\ell)} = \hat{\mathbf{W}}_{(\text{old})}^{(\ell)} - \eta \left. \frac{\partial E_n}{\partial \hat{\mathbf{W}}^{(\ell)}} \right|_{\hat{\mathbf{W}}_{(\text{old})}}$$

2 Regularization of NN with tangent propagation

ここでは、逆伝播アルゴリズムと正則化について考える。

$$\hat{\mathcal{G}} \equiv \sum_{m=1}^D \tau_{nm} \frac{\partial}{\partial x_{nm}} \quad (2.1)$$

というオペレータを導入すると、正則化項の式 (1.1) は

$$\Omega_n = \frac{1}{2} \sum_{k=1}^K \left(\hat{\mathcal{G}} y_{nk} \right)^2 \quad (2.2)$$

と簡潔に表記される。 x_{ni} はサンプル n における入力ユニット値 $x_i = z_i^{(0)}$ である。次に順伝播の更新式を導出する。以下のように $\hat{\mathcal{G}}$ オペレータを順伝播公式 (1.4) に作用させる。

$$\alpha_j^{(\ell)} \equiv \hat{\mathcal{G}} z_j^{(\ell)} \quad \beta_j^{(\ell)} \equiv \hat{\mathcal{G}} a_j^{(\ell)} \quad (2.3)$$

上の定義から、順伝播公式は

$$\alpha_j^{(\ell)} = h'\left(a_j^{(\ell)}\right) \hat{\mathcal{G}} a_j^{(\ell)} = h'\left(a_j^{(\ell)}\right) \beta_j^{(\ell)} \quad (2.4)$$

$$\beta_j^{(\ell)} = \hat{\mathcal{G}} \left(\sum_i w_{ji}^{(\ell)} z_i^{(\ell-1)} \right) = \sum_i w_{ji}^{(\ell)} \underbrace{\hat{\mathcal{G}} z_i^{(\ell-1)}}_{\alpha_i^{(\ell-1)}} \quad (2.5)$$

と変形され、特に入力層で τ の伝播は

$$\alpha_j^{(0)} = \hat{\mathcal{G}} x_{nj} = \sum_m \tau_{nm} \frac{\partial x_{nj}}{\partial x_{nm}} = \tau_{nj}$$

が得られる。以上をまとめると正則化順伝播の公式は次の表式となる。

$$\alpha_j^{(\ell)} = h'\left(a_j^{(\ell)}\right) \sum_i w_{ji}^{(\ell)} \alpha_i^{(\ell-1)} \equiv h'\left(a_j^{(\ell)}\right) \beta_j^{(\ell)} \quad \left[\alpha_i^{(0)} = \tau_{ni} \right] \quad (2.6)$$

2.1 Error gradient

式 (2.3) を用いると $y_k = z_k^{(L)}$ であることから

$$\Omega_n = \frac{1}{2} \sum_{k=1}^K [\hat{\mathcal{G}} y_{nk}]^2 = \frac{1}{2} \sum_{k=1}^K [\alpha_{nk}^{(L)}]^2 \quad (2.7)$$

と正則化項が別表式で表現できる．このとき，正則化項の重み勾配を調べる．

$$\begin{aligned} \frac{\partial \Omega_n}{\partial w_{rs}^{(\ell)}} &= \sum_k \left(\hat{\mathcal{G}} y_{nk} \right) \frac{\partial}{\partial w_{rs}^{(\ell)}} \left(\hat{\mathcal{G}} y_{nk} \right) \\ &= \sum_k \alpha_{nk}^{(L)} \hat{\mathcal{G}} \frac{\partial y_{nk}}{\partial w_{rs}^{(\ell)}} \\ &= \sum_k \alpha_{nk}^{(L)} \hat{\mathcal{G}} \left(\delta_{nkr}^{(\ell)} z_{ns}^{(\ell-1)} \right) \end{aligned} \quad (2.8)$$

上の式の 2 番目の等号では式 (1.9) と同様のロジックで

$$\frac{\partial y_{nk}}{\partial w_{rs}^{(\ell)}} = \frac{\partial y_{nk}}{\partial a_{nr}^{(\ell)}} \frac{\partial a_{nr}^{(\ell)}}{\partial w_{rs}^{(\ell)}} \equiv \delta_{nkr}^{(\ell)} z_{ns}^{(\ell-1)}$$

が成り立つことを用いた．誤差 δ_{nkr} の定義については後に述べる．式 (2.8) で $\hat{\mathcal{G}}$ は積の微分則によって δ と z のそれぞれに作用することと，式 (2.3) 及び，便利な表記

$$\hat{\mathcal{G}} \delta_{nkr}^{(\ell)} = \phi_{nkr}^{(\ell)} \quad (2.9)$$

を用いると，次の正則化項の重み勾配が与えられる．

$$\frac{\partial \Omega_n}{\partial w_{rs}^{(\ell)}} = \sum_k \alpha_{nk}^{(L)} \left\{ \phi_{nkr}^{(\ell)} z_{ns}^{(\ell-1)} + \delta_{nkr}^{(\ell)} \alpha_{ns}^{(\ell-1)} \right\} \quad (2.10)$$

2.2 Back propagation

ここで，誤差 δ の更新式についても考慮する．前述したように δ は出力ユニットの活性微分で定義される．式 (1.12) を真似ると

$$\begin{aligned} \delta_{nkr}^{(\ell)} &\equiv \frac{\partial y_{nk}}{\partial a_{nr}^{(\ell)}} = \sum_u \frac{\partial y_{nk}}{\partial a_{nu}^{(\ell+1)}} \frac{\partial a_{nu}^{(\ell+1)}}{\partial a_{nr}^{(\ell)}} \\ &= \sum_u \delta_{nku}^{(\ell+1)} \frac{\partial}{\partial a_{nr}^{(\ell)}} \left\{ \sum_s w_{us}^{(\ell+1)} h \left(a_{ns}^{(\ell)} \right) \right\} \end{aligned}$$

となるので，誤差 δ の更新式は以下ようになる．

$$\delta_{nkr}^{(\ell)} = h' \left(a_{nr}^{(\ell)} \right) \sum_u w_{ur}^{(\ell+1)} \delta_{nku}^{(\ell+1)} \quad (2.11)$$

また，式 (2.9) から ϕ に関する更新式は

$$\phi_{nkr}^{(\ell)} \equiv \hat{\mathcal{G}} \delta_{nkr}^{(\ell)} = \hat{\mathcal{G}} \left[h' \left(a_{nr}^{(\ell)} \right) \sum_u w_{ur}^{(\ell+1)} \delta_{nku}^{(\ell+1)} \right]$$

となるが，式 (2.3) から $\hat{\mathcal{G}} a_{nr}^{(\ell)} = \beta_{nr}^{(\ell)}$ が成り立つことと，式 (2.9) を用いて

$$\phi_{nkr}^{(\ell)} = h'' \left(a_{nr}^{(\ell)} \right) \beta_{nr}^{(\ell)} \sum_u w_{ur}^{(\ell+1)} \delta_{nku}^{(\ell+1)} + h' \left(a_{nr}^{(\ell)} \right) \sum_u w_{ur}^{(\ell+1)} \phi_{nku}^{(\ell+1)} \quad (2.12)$$

が与えられる．実際の正則化項の更新は次のステップを実行すればよい．

1. 式 (2.6) を用いて接線情報 τ の順伝播
2. 式 (2.11), (2.12) により δ, ϕ の更新
3. 式 (2.10) から誤差勾配を計算

3 Conjugate gradient method

3.1 formulation for conjugate gradient method

$E(\mathbf{w})$ を最小化する \mathbf{w} を考える．簡単化のため $E(\mathbf{w}^{(k)}) = E^{(k)}$ と表記する．現在の近似解が $\mathbf{w}^{(k)}$ であるとき，誤差関数 $E(\mathbf{w})$ の2次近似は次のように書ける．

$$E(\mathbf{w}) \simeq E^{(k)} + (\mathbf{w} - \mathbf{w}^{(k)})^T \nabla E^{(k)} + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(k)})^T \hat{\mathbf{H}}^{(k)} (\mathbf{w} - \mathbf{w}^{(k)}) \quad (3.1)$$

ここで，ヘッセ行列 $\hat{\mathbf{H}}^{(k)}$ の定義は $\nabla \nabla E^{(k)}$ とする．式 (3.1) の勾配は2次形式の微分を用いて次式のよう定まる．

$$\nabla E(\mathbf{w}) = \nabla E^{(k)} + \underbrace{\hat{\mathbf{H}}^{(k)} (\mathbf{w} - \mathbf{w}^{(k)})}_{\mathbf{p}^{(k)}} = 0 \quad (3.2)$$

ここで定義された $\mathbf{p}^{(k)} = \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}$ を共役勾配と呼ぶ．共役勾配は現在の地点から停留点の方向を向くベクトルであり，以下の性質を満たすことになる．

$$\nabla E^{(k)} \propto \hat{\mathbf{H}}^{(k)} \mathbf{p}^{(k)} \quad (3.3)$$

以降，表記の簡素化のため内積の定義を $\mathbf{a}^T \mathbf{b} \equiv \langle \mathbf{a}, \mathbf{b} \rangle$ と書くことにする．点 $\mathbf{w}^{(k)}$ (今は座標のように考えてよい) での，等高線の接線ベクトルを $\mathbf{t}^{(k)}$ とすると，接ベクトル $\mathbf{t}^{(k)}$ は誤差勾配 $\nabla E^{(k)}$ と直交するので，式 (3.3) を用いて

$$\langle \mathbf{t}^{(k)}, \hat{\mathbf{H}}^{(k)} \mathbf{p}^{(k)} \rangle = 0 \quad (3.4)$$

の関係式を満たすはずである．Fig.1 より共役勾配は勾配 $\nabla E^{(k)}$ と接線 $\mathbf{t}^{(k)}$ との線形結合で書けるので

$$\mathbf{p}^{(k)} = -\nabla E^{(k)} + \beta^{(k)} \mathbf{t}^{(k)} \quad (3.5)$$

となる．第一項にある勾配の前のマイナス符号は確率勾配降下法と対応させている．上の式は共役勾配を定める更新式として後で用いられる．式 (3.5) を式 (3.4) に代入して

$$-\langle \mathbf{t}^{(k)}, \hat{\mathbf{H}}^{(k)} \nabla E^{(k)} \rangle + \beta^{(k)} \langle \mathbf{t}^{(k)}, \hat{\mathbf{H}}^{(k)} \mathbf{t}^{(k)} \rangle = 0$$

となることから，調整パラメータ β の更新式を得る．

$$\beta^{(k)} = \frac{\langle \mathbf{t}^{(k)}, \hat{\mathbf{H}}^{(k)} \nabla E^{(k)} \rangle}{\langle \mathbf{t}^{(k)}, \hat{\mathbf{H}}^{(k)} \mathbf{t}^{(k)} \rangle} \quad (3.6)$$

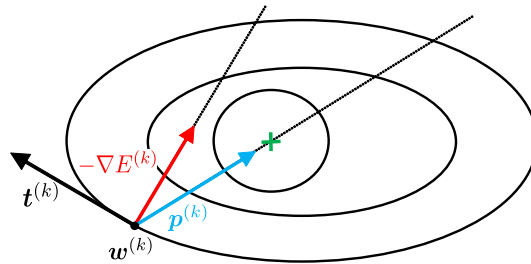


Fig 1: conjugate gradient direction

ここで，次の探索直線に関する定理を導入する．

Theorem 1

関数 $E(\mathbf{w})$ に対する勾配法の直線探索で定まる点では，その点を通る $E(\mathbf{w})$ の等値面 (等高線) が探索直線に接する．すなわち，次の直線探索の方向は直前の探索方向と直交する．

この定理は $\mathbf{p}^{(k)} = \mathbf{t}^{(k+1)}$ ，つまり $\mathbf{t}^{(k)} = \mathbf{p}^{(k-1)}$ を意味する．この関係を式 (3.5) と (3.6) に代入して次の

アルゴリズムを得る．下の3つの式をループ処理することにより \mathbf{w} の鞍点解を得ることができる．

$$\mathbf{p}^{(k)} = -\nabla E(\mathbf{w}^{(k)}) + \beta^{(k)} \mathbf{p}^{(k-1)} \quad [\beta^{(0)} = 0] \quad (3.7)$$

$$\beta^{(k)} = \frac{\langle \mathbf{p}^{(k-1)}, \hat{\mathbf{H}}^{(k)} \nabla E^{(k)} \rangle}{\langle \mathbf{p}^{(k-1)}, \hat{\mathbf{H}}^{(k)} \mathbf{p}^{(k-1)} \rangle} \quad (3.8)$$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \alpha^{(k)} \mathbf{p}^{(k)} \quad (3.9)$$

$\alpha^{(k)}$ の決め方については次のセクションで詳細を後述する．

3.2 Avoid Hessian calculation

式 (3.8) ではヘッセ行列を厳密に計算する必要があり，計算量の観点からあまり現実的でない．ヘッセ行列の計算を回避し β を更新する手法として Polak-Ribere の近似式がある．

$$\beta^{(k)} = \frac{\langle \nabla E^{(k)}, \nabla E^{(k)} - \nabla E^{(k-1)} \rangle}{\|\nabla E^{(k-1)}\|^2} \quad (3.10)$$

PR 式の導出を次に示す． $\mathbf{p}^{(k-1)} = \mathbf{w}^{(k)} - \mathbf{w}^{(k-1)} = \delta \mathbf{w}$ と置き，式 (3.4) から

$$\langle \mathbf{p}^{(k)} \hat{\mathbf{H}}^{(k)}, \mathbf{p}^{(k-1)} \rangle = \mathbf{p}^{(k)} \hat{\mathbf{H}}^{(k)} \delta \mathbf{w} = 0$$

が任意の $\delta \mathbf{w}$ について成り立つことから，

$$\mathbf{p}^{(k)} \hat{\mathbf{H}}^{(k)} = 0$$

となる．これに，式 (3.5) 及び， $\hat{\mathbf{H}}^{(k)} = \nabla \nabla E^{(k)} \simeq \delta \nabla E^{(k)}$ を代入して

$$\left\{ -\nabla E^{(k)} + \beta^{(k)} \mathbf{t}^{(k)} \right\} \left\{ \nabla E^{(k)} - \nabla E^{(k-1)} \right\} = 0$$

を得る．ここで， $\mathbf{t}^{(k)} = \mathbf{p}^{(k-1)} \simeq -\nabla E^{(k-1)}$ と近似し， $\nabla E^{(k-1)} \nabla E^{(k)} \simeq 0$ と落とすことにより

$$\beta^{(k)} = \frac{\langle \nabla E^{(k)}, \nabla E^{(k)} - \nabla E^{(k-1)} \rangle}{\|\nabla E^{(k-1)}\|^2}$$

が与えられる． $k \rightarrow \tau$ とラベルを張り替えると CG 法は次のようにまとめられる．

Conjugate Gradient method

$E^{(\tau)} \equiv E(\mathbf{w}^{(\tau)})$ と表記し，初期の重みベクトル $\mathbf{w}^{(0)}$ は適当に選ぶとする．

1. 調整パラメータ $\beta^{(\tau)}$ を以下の式に従って更新する．

$$\beta^{(\tau)} = \frac{\langle \nabla E^{(\tau)}, \nabla E^{(\tau)} - \nabla E^{(\tau-1)} \rangle}{\|\nabla E^{(\tau-1)}\|^2} \quad [\beta^{(0)} = 0] \quad (3.11)$$

2. 最急方向 $-\nabla E(\mathbf{w})$ を求め，共役方向 $\mathbf{p}^{(\tau)}$ を算出する．

$$\mathbf{p}^{(\tau)} = -\nabla E(\mathbf{w}^{(\tau)}) + \beta^{(\tau)} \mathbf{p}^{(\tau-1)} \quad (3.12)$$

3. 共役方向ベクトルを用いて重みの更新を行う．

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \alpha^{(\tau)} \mathbf{p}^{(\tau)} \quad (3.13)$$

- $\alpha^{(\tau)}$ の決め方

通常, ステップサイズとして知られる $\alpha^{(\tau)}$ のパラメータは次のように定義される.

$$\alpha^{(\tau)} = \operatorname{argmin}_{\alpha} E\left(\mathbf{w}^{(\tau)} + \alpha \mathbf{p}^{(\tau)}\right) \quad (3.14)$$

特に誤差関数 E が 2 次形式の時は式 (3.2) に $\mathbf{w} = \mathbf{w}^{(k)} + \alpha^{(k)} \mathbf{p}^{(k)}$ を代入して

$$\nabla E^{(k)} + \alpha^{(k)} \hat{\mathbf{H}}^{(k)} \mathbf{p}^{(k)} = 0$$

を満たす. 両辺に対して $\mathbf{p}^{(k)}$ との内積を取り, $k \rightarrow \tau$ とすると

$$\alpha^{(\tau)} = -\frac{\langle \mathbf{p}^{(\tau)}, \nabla E^{(\tau)} \rangle}{\langle \mathbf{p}^{(\tau)}, \hat{\mathbf{H}}^{(\tau)} \mathbf{p}^{(\tau)} \rangle} \quad (3.15)$$

が与えられる.

3.3 CG-algorithm

スケール化共役勾配法 (SCG 法) との比較のために CG 法のアルゴリズムをステップ化する.

- I. 初期ベクトル $\mathbf{w}^{(0)}$ を適当に選び, 初期勾配を計算する.

$$\mathbf{p}^{(0)} = \mathbf{r}^{(0)} = -\nabla E\left(\mathbf{w}^{(0)}\right) \quad : \tau = 0$$

- II. 次式を用いて 2 次情報の計算を行う.

$$\mathbf{s}^{(\tau)} = \nabla \nabla E\left(\mathbf{w}^{(\tau)}\right) \mathbf{p}^{(\tau)}$$

$$\delta^{(\tau)} = \langle \mathbf{p}^{(\tau)}, \mathbf{s}^{(\tau)} \rangle$$

- III. 共役勾配をもとにステップサイズ α を算出する.

$$\mu^{(\tau)} = \langle \mathbf{p}^{(\tau)}, \mathbf{r}^{(\tau)} \rangle$$

$$\alpha^{(\tau)} = \mu^{(\tau)} / \delta^{(\tau)}$$

- IV. 重みベクトルを更新し, 最急勾配を計算する.

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \alpha^{(\tau)} \mathbf{p}^{(\tau)}$$

$$\mathbf{r}^{(\tau+1)} = -\nabla E\left(\mathbf{w}^{(\tau+1)}\right)$$

- V. N をブレイクサイズとしたとき, 次の条件文を実行し共役勾配を再設定する.

if ($\tau \bmod N = 0$) : # setting as follow and restart Algorithm

$$\mathbf{p}^{(\tau+1)} = \mathbf{r}^{(\tau+1)}$$

elseif ($\tau \bmod N \neq 0$) : # create new conjugate direction

$$\beta^{(\tau+1)} = \left\{ \|\mathbf{r}^{(\tau+1)}\|^2 - \langle \mathbf{r}^{(\tau+1)}, \mathbf{r}^{(\tau)} \rangle \right\} / \mu^{(\tau)}$$

$$\mathbf{p}^{(\tau+1)} = \mathbf{r}^{(\tau+1)} + \beta^{(\tau+1)} \mathbf{p}^{(\tau)}$$

- VI. ループ中の収束条件をチェックする.

if ($|\mathbf{r}^{(\tau)}| \neq 0$) : $\tau = \tau + 1 \implies$ go to step II

else : return $\mathbf{w}^{(\tau+1)}$ # terminate

3.4 scaled CG method

I. 初期ベクトル $\mathbf{w}^{(0)}$ を適当に選び、初期勾配を計算する． $\tilde{\lambda}^{(0)} = 0$ とする．

$$\mathbf{p}^{(0)} = \mathbf{r}^{(0)} = -\nabla E(\mathbf{w}^{(0)}) \quad : \text{success} = \text{True}$$

II. $\text{success} = \text{True}$ のとき、次式を用いて 2 次情報の計算を行う．

$$\mathbf{s}^{(\tau)} = \left\{ \nabla E(\mathbf{w}^{(\tau)} + \sigma^{(\tau)} \mathbf{p}^{(\tau)}) - \nabla E(\mathbf{w}^{(\tau)}) \right\} / \sigma^{(\tau)} \quad \sigma^{(\tau)} = \frac{\sigma}{|\mathbf{p}^{(\tau)}|}$$

$$\delta^{(\tau)} = \langle \mathbf{p}^{(\tau)}, \mathbf{s}^{(\tau)} \rangle$$

III. $\delta^{(\tau)} := \delta^{(\tau)} + (\lambda^{(\tau)} - \tilde{\lambda}^{(\tau)}) |\mathbf{p}^{(\tau)}|^2$ を計算し $\delta^{(\tau)} \leq 0$ のときヘッセ行列を正定値にする．

$$\tilde{\lambda}^{(\tau)} = 2 \left(\lambda^{(\tau)} - \frac{\delta^{(\tau)}}{|\mathbf{p}^{(\tau)}|^2} \right)$$

$$\delta^{(\tau)} := -\delta^{(\tau)} + \lambda^{(\tau)} |\mathbf{p}^{(\tau)}|^2$$

$$\lambda^{(\tau)} := \tilde{\lambda}^{(\tau)}$$

IV. 共役勾配をもとにステップサイズ α を算出する．

$$\mu^{(\tau)} = \langle \mathbf{p}^{(\tau)}, \mathbf{r}^{(\tau)} \rangle$$

$$\alpha^{(\tau)} = \mu^{(\tau)} / \delta^{(\tau)}$$

V. 比較パラメータを算出する．

$$\Delta^{(\tau)} = \frac{2\delta^{(\tau)}}{|\mu^{(\tau)}|^2} \left[E(\mathbf{w}^{(\tau)}) - E(\mathbf{w}^{(\tau)} + \alpha^{(\tau)} \mathbf{p}^{(\tau)}) \right]$$

VI. $\Delta^{(\tau)} \geq 0$ のとき、次のプログラムを実行する

```

 $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \alpha^{(\tau)} \mathbf{p}^{(\tau)}$ 
 $\mathbf{r}^{(\tau+1)} = -\nabla E(\mathbf{w}^{(\tau+1)})$ 
 $\tilde{\lambda}^{(\tau)} = 0, \text{ success} = \text{True}$ 
if ( $\tau \bmod N = 0$ ) then restart algorithm:
     $\mathbf{p}^{(\tau+1)} = \mathbf{r}^{(\tau+1)}$ 
else:
     $\beta^{(\tau+1)} = \left\{ \|\mathbf{r}^{(\tau+1)}\|^2 - \langle \mathbf{r}^{(\tau+1)}, \mathbf{r}^{(\tau)} \rangle \right\} / \mu^{(\tau)}$ 
     $\mathbf{p}^{(\tau+1)} = \mathbf{r}^{(\tau+1)} + \beta^{(\tau+1)} \mathbf{p}^{(\tau)}$ 
if ( $\Delta^{(\tau)} \geq 0.75$ ) then reduce the scale:
     $\lambda^{(\tau)} = 0.25 \lambda^{(\tau)}$ 
else:
     $\tilde{\lambda}^{(\tau)} = \lambda^{(\tau)}, \text{ success} = \text{False}$ 

```

VII. $\Delta^{(\tau)} < 0.25$ のとき、 $\lambda^{(\tau)} = \lambda^{(\tau)} + \delta^{(\tau)} (1 - \Delta^{(\tau)}) / |\mathbf{p}^{(\tau)}|^2$ としスケール縮小を行う．

VIII. ループ中の収束条件をチェックする．

if ($|\mathbf{r}^{(\tau)}| \neq 0$) : $\tau = \tau + 1 \implies$ go to step II

else : return $\mathbf{w}^{(\tau+1)}$ # terminate

References

- [1] Pattern Recognition and Machine Learning (C.M.Bishop)
- [2] これなら分かる最適化数学 (金谷 健一)
- [3] 数理計画法による最適化 (北村 充)
- [4] Moller, M (1993), "A Scaled Conjugate Gradient Algorithm for Fast Supervised Learning"