

深層学習 day1 レポート

導入

・確認テスト

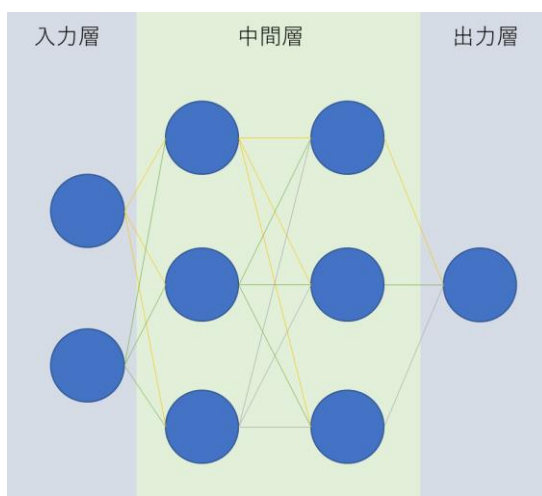
[ディープラーニングは、結局何をやろうとしているか 2行以内で述べよ。]

万能近似定理によりニューラルネットワークでの複雑な関数を表現し、入力データから目的関数となる出力を算出する。

[次の中のどの値の最適化が最終目的か]

③重み[W]、④バイアス[b]の最適化

[入力層：2 ノード 1 層、中間層：3 ノード 2 層、出力層：1 ノード 1 層]



Section1：入力層～中間層

入力層：ニューラルネットワークが最初に情報を受け取る層となる。

ニューラルネットワークに何かしらのインプットを行う層であり、各ノードで入力より数値を受け取る。(生き物では電気信号で情報を受け取る)

中間層：入力層から情報を受け継ぎ、さまざまな計算を行う層。

中間層が多いほど複雑な分析が可能（万能近似定理）

中間層は隠れ層と呼ばれる場合もある。

出力層：入力層、中間層で重みをかけ、活性化関数で処理された値が示される出力層。

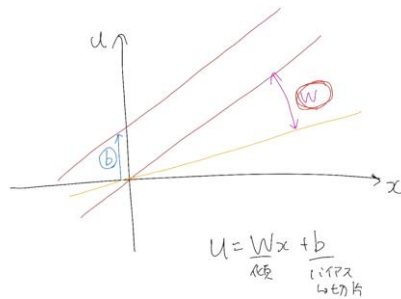
重み：入力層の重要度を表現するもの。

入力の値をいい感じにシャッフルし出力するため重みを調整する。

バイアス： w と x だけでは値をシフトすることができない。(傾きを表現している。)

バイアスを用いることでシフトすることができる。

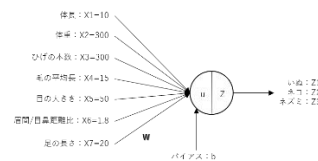
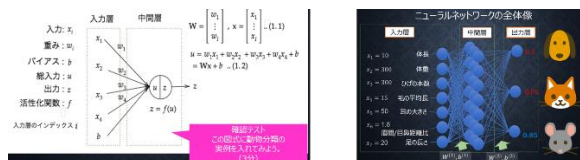
参考資料：[Neural Network Console](#) (SONY)



傾きと切片が決まれば、一次関数では関数を決めることができる。

ニューラルネットワークの学習ではこのwとbを学習することになる。

・確認テスト



$$u = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$

$$= Wx + b \quad \text{.. (1.2)}$$

➤ z

確認テスト
この数式をPythonで書け。(2分)

解答: `u = np.dot(W, x) + b`

ソースコード

- 1_1_forward_propagation.html
- 1_1_forward_propagation.ipynb
- 1_1_forward_propagation.py
- 1_1_forward_propagation_after.html
- 1_1_forward_propagation_after.ipynb
- 1_1_forward_propagation_after.py
- 1_2_back_propagation.html
- 1_2_back_propagation.ipynb
- 1_2_back_propagation.py

確認テスト
1-1のファイルから
中間層の出力を定義しているソースを抜き出せ。(2分)

解答:

```
# 2層の総入力
u2 = np.dot(z1, W2) + b2

# 2層の総出力
z2 = functions.relu(u2)
```

・実習演習

順伝播（単層・複数ユニット）

```
1 # 順伝播（単層・複数ユニット）
2
3 # 重み
4 W = np.array([
5     [0.1, 0.2, 0.3, 0],
6     [0.2, 0.3, 0.4, 0.5],
7     [0.3, 0.4, 0.5, 1],
8 ])
9
10 ## 試してみよう_配列の初期化
11 #W = np.zeros((4,3))
12 #W = np.ones((4,3))
13 W = np.random.rand(3,4)
14 #W = np.random.randint(5, size=(4,3))
15
16 print_vec("重み", W)
17
18 # バイアス
19 b = np.array([0.1, 0.2, 0.3])
20 print_vec("バイアス", b)
21
22 # 入力値
23 x = np.array([1.0, 5.0, 2.0, -1.0])
24 print_vec("入力", x)
25
26 # 総入力
27 u = np.dot(W, x) + b
28 print_vec("総入力", u)
29
30 # 中間層出力
31 z = functions.sigmoid(u)
32 print_vec("中間層出力", z)
33
```

支給ファイルの試してみようでは、W が 4x3 で記入されていたためエラーとなった。
W=np.random.rand(3,4)へ変更し、実行を行った。

```
*** 重み ***
[[0.233063  0.02587027 0.31105677 0.96633896]
 [0.95506073 0.34413852 0.85488792 0.64611252]
 [0.44716423 0.5759112  0.24649931 0.46086534]]
shape: (3, 4)

*** バイアス ***
[0.1 0.2 0.3]
shape: (3, )

*** 入力 ***
[ 1.  5.  2. -1.]
shape: (4, )

*** 総入力 ***
[0.11818892 3.93941664 3.65885349]
shape: (3, )

*** 中間層出力 ***
[0.52951288 0.98091188 0.97488498]
shape: (3, )
```

無事実行完了。

W と X で行列演算できる形へすることは重要と感じた。

Section2：活性化関数

活性化関数を使うことでニューラルネットワークの幅が広がる。

活性化関数は非線形の関数となっているため。

中間層用の活性化関数

- ReLU 関数
- シグモイド（ロジスティック）関数
- ステップ関数（ほぼ使われることはない。パーセプトロンで利用された関数）

出力層用の活性化関数

- ソフトマックス関数
- 恒等写像
- シグモイド（ロジスティック）関数

活性化関数に微分を適用できるようになり、

・確認テスト

[線形と非線形の違いを図にかいて簡易に説明せよ。]

線形関数

非線形関数



※ 非線形関数は加法性・斉次性を満たさない（講師）

[配布されたソースコードより該当する箇所を抜き出せ。]

$z_3 = f(u) \dots (1.5)$

配布されたソースコードより
該当する箇所を抜き出せ。（3分）

```
46 # 1層の総出力
47 z1 = functions.relu(u1)
48
49 # 2層の総入力
50 u2 = np.dot(z1, W2) + b2
51
52 # 2層の総出力
53 z2 = functions.relu(u2)
```

ここでは ReLU 関数を使用

※提供別ファイルで活性化関数は定義（common ファイル）

・実習演習

順伝播（3層・複数ユニット）

活性化関数を ReLU 関数→シグモイド関数へ変更して出力の変化を確認してみる。

ReLU 関数

```
*** 出力 ***
[85.46151642 42.8777085 35.42770326 63.55549018]
shape: (4,)
```

出力合計: 227.3224183500664

シグモイド関数

```
# 1層の総出力
# z1 = functions.relu(u1)
z1 = functions.sigmoid(u1)

# 2層の総入力
u2 = np.dot(z1, W2) + b2

# 2層の総出力
#z2 = functions.relu(u2)
z2 = functions.sigmoid(u2)

*** 出力 ***
[3.01742837 2.32818806 3.23545458 3.25864259]
shape: (4, )

出力合計: 11.839713600998046
```

ReLU では 227.32 だった結果が、シグモイド関数では 11.83 となる。
これは入力層・中間層の出力が活性化関数により制限がかかったためと考える。

・その他

活性化関数について調べたところ多くの関数が活用されているとのこと。

参考に下記に記載する。

恒等関数／線形関数

ステップ関数

シグモイド関数

tanh 関数

ReLU

ソフトプラス関数

Leaky ReLU

PReLU／Parametric ReLU

ELU

SELU

Swish 関数

Mish 関数

ソフトマックス関数

Section3：出力層

出力層：入力層、中間層で重みをかけ、活性化関数で処理された値が示される出力層。

私たち人間がほしい情報へ変換した出力を行ってくれる層となる。

動物分類のような問題では、イヌの確率やネコの確率などで出力する必要がある。

誤差関数：ニューラルネットワーク全体を学習させるために重要な関数となる。

訓練データとして入力データとそれに対する正解データをペアで用意し、
 訓練データの入力データよりニューラルネットワークで予測を行う。
 この結果を正解データと比べることで答え合わせを行う。
 この際に使用する関数が誤差関数となる。
 具体的には下記に示す二乗和誤差などがある。

$$E_n(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^J (y_j - d_j)^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{d}\|^2$$

出力層の中間層との違い

中間層 : しきい値の前後で信号の強弱を調整

出力層 : 信号の大きさ (比率) はそのままに変換

出力層と中間層で利用される活性化関数が異なる

全結合NN - 出力層の種類			
	回帰	二値分類	多クラス分類
活性化関数	恒等写像 $f(u) = u$	シグモイド関数 $f(u) = \frac{1}{1 + e^{-u}}$	ソフトマックス関数 $f(\mathbf{i}, \mathbf{u}) = \frac{e^{u_i}}{\sum_{k=1}^K e^{u_k}}$
誤差関数	二乗誤差	交差エントロピー	
【訓練データサンプルあたりの誤差】 $E_n(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^J (y_i - d_i)^2$.. 二乗誤差 $E_n(\mathbf{w}) = - \sum_{i=1}^J d_i \log y_i$.. 交差エントロピー		【学習サイクルあたりの誤差】 $E(\mathbf{w}) = \sum_{n=1}^N E_n$	

・確認テスト

[なぜ、引き算でなく二乗するか述べて]

誤差はプラス、マイナス両方発生するため、単純に引き算すると誤差を表現できないのでは。
 二乗することで正の数にそろえることができる。

[下式の 1/2 はどういう意味を持つか述べて]

ニューラルネットワークの学習の際にバックプロパゲーションを行うが、
 その時に微分が必要となり、誤差関数の二乗により 2 が発生するため 1/2 を用いている。

[①~③の数式に該当するソースコードを示し、一行ずつ処理の説明をせよ。]

```
def softmax(x):
    ① if x.ndim == 2:
        x = x.T
        x = x - np.max(x, axis=0)
        y = np.exp(x) / np.sum(np.exp(x), axis=0)
        return y.T ②
    x = x - np.max(x) # オーバーフロー対策
    return np.exp(x) / np.sum(np.exp(x)) ③
```

[①~②の数式に該当するソースコードを示し、一行ずつ処理の説明をせよ。]

```
①
def cross_entropy_error(d, y):
    if y.ndim == 1:
        d = d.reshape(1, d.size)
        y = y.reshape(1, y.size)
        # 教師データがone-hot-vectorの場合、正解ラベルのインデックスに変換
        if d.size == y.size:
            d = d.argmax(axis=1)
        batch_size = y.shape[0]
    return -np.sum(np.log(y[np.arange(batch_size), d] + 1e-7)) / batch_size
    ②
```

1e-7 は np.log(0)は計算できないため,0 とならないためにつけている。

・実習演習

回帰 (2-3-2 ネットワーク)

誤差関数を二乗誤差とクロスエントロピー誤差で比較してみる。

```
54 # 入力値
55 x = np.array([1., 2., 3.])
56 network = init_network()
57 y, z1 = forward(network, x)
58 # 目標出力
59 d = np.array([2., 4.])
60 # 誤差
61 loss = functions.mean_squared_error(d, y)
62 loss_2 = functions.cross_entropy_error(d, y)
63
64 ## 表示
65 print("\n##### 結果表示 #####")
66 print_vec("中間層出力", z1)
67 print_vec("出力", y)
68 print_vec("訓練データ", d)
69 print_vec("二乗誤差", loss)
70 print_vec("クロスエントロピー誤差", loss_2)
```

*** 二乗誤差 ***
3613.87373486424
shape: ()

*** クロスエントロピー誤差 ***
-4.469187970953645
shape: ()

やはり二乗誤差とクロスエントロピー誤差では値が全く異なる。

・その他 (参考: [AVILEN AI Trend](#) より)

線形回帰の損失関数: 平均二乗誤差(mean squared error)

$$MSE(y_i, \hat{y}_i) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

y_i は実値、 \hat{y}_i は予測値を指す。

回帰問題に用いる損失関数:

平均二乗誤差(mean squared error)、

$$MSE(y_i, \hat{y}_i) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

y_i は実値、 \hat{y}_i は予測値を指す。

平均二乗誤差の性質として外れ値に対して敏感である。

外れ値を含むデータに適用すると、予測結果が不安定となる。(外れ値に引っ張られる)

平均絶対誤差 / Mean Absolute Error

$$MAE(y_i, \hat{y}_i) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

y_i は実値、 \hat{y}_i は予測値を指す。

平均絶対誤差の性質として外れ値に強い

平均二乗対数誤差 / Mean Squared Logarithmic Error

$$MSLE(y_i, \hat{y}_i) = \frac{1}{n} \sum_{i=1}^n \{\log(1 + y_i) - \log(1 + \hat{y}_i)\}^2$$

y_i は実値、 \hat{y}_i は予測値を指す。

平均二乗対数誤差を用いたモデルは予測が実値を上回りやすくなるという傾向がある。

交差エントロピー誤差 / cross entropy error

$$E = - \sum_k t_k \log y_k$$

交差エントロピー誤差は実際のカテゴリに対する予測確率のみを評価する。

異なる予測に対しては 0 を掛ける。

Section4：勾配降下法

ニューラルネットワークを学習させるための手法となる。

学習を通じて誤差を最小にするネットワークを作成するが、そのパラメータ更新に利用。

・学習率

学習率が大きすぎた場合、最小値にいつまでもたどり着かず発散してしまう。

学習率が小さい場合発散することはないが、収束するまでに時間がかかってしまう。

また、極小解から抜け出せないことがある。

勾配降下法の学習率の決定、収束性向上のためのアルゴリズムについて複数の論文が公開され、よく利用されている。

- ・ Momentum
- ・ AdaGrad
- ・ Adadelta
- ・ Adam

※上記アルゴリズム詳細は Day2 で触れることになる。

誤差関数の値をより小さくする方向へ学習のサイクルをエポックと呼ぶ。

1 エポックでは 1 回の学習となる。

確率的勾配降下法（SGD）

勾配降下法では全サンプルの平均誤差で学習したが、SGD ではランダムに抽出したサンプルの誤差で学習する。下記メリットがある。

- ・ データが冗長な場合の計算コストの軽減
- ・ 望まない局所極小解に収束するリスクの軽減
- ・ オンライン学習ができる（バッチ学習の逆となる。）

ミニバッチ勾配降下法

SGD ではランダムに抽出したサンプルでの学習となるが、ミニバッチ勾配降下法ではランダムに分割したデータ群に属するサンプル平均誤差により学習する方法となる。

確率的勾配降下法のメリットを損なわず、計算機の計算資源を有効利用できる

→ CPU を利用したスレッド並列化や GPU を利用した SIMD 並列化のメリットがある。

誤差勾配の計算・・・微分

プログラムで微小な数値を生成し擬似的に微分を計算する一般的な手法 ← 計算量大

→ 誤差逆伝播法を利用する

・確認テスト

[該当するソースコードを探してみよう。]

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \varepsilon \nabla E$$

$$\nabla E = \frac{\partial E}{\partial \mathbf{w}} = \left[\frac{\partial E}{\partial w_1} \cdots \frac{\partial E}{\partial w_M} \right]$$

```
# パラメータに勾配適用
for key in ('W1', 'W2', 'b1', 'b2'):
    network[key] -= learning_rate * grad[key]
```

```
grad = backward(x, d, z1, y)
# パラメータに勾配適用
```

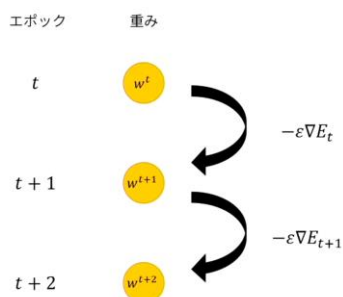
[オンライン学習とは何か]

学習毎にニューラルネットワークのパラメータの更新をする学習方法。

バッチ学習は学習データをまとめ平均誤差などを使ってパラメータの更新を行う方法。

[この数式の意味を図に書いて説明せよ。]

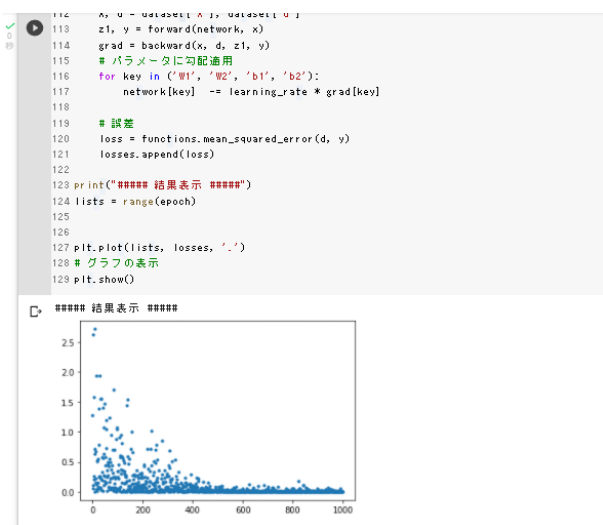
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \varepsilon \nabla E_t$$



各学習毎にネットワーク（重み）を更新している。

・実習演習

1_3_stochastic_gradient_descent.ipynb



ここでは確率的勾配効果法による学習を実施する。

当たり前ではあるが、学習を繰り返すことで誤差が小さくなっている。

Section5：誤差逆伝播法

Deep Learning の学習においては、誤差関数の値がゼロになる（近づく）ようにパラメータを決めたい。この際に利用するのが誤差逆伝播法となる。

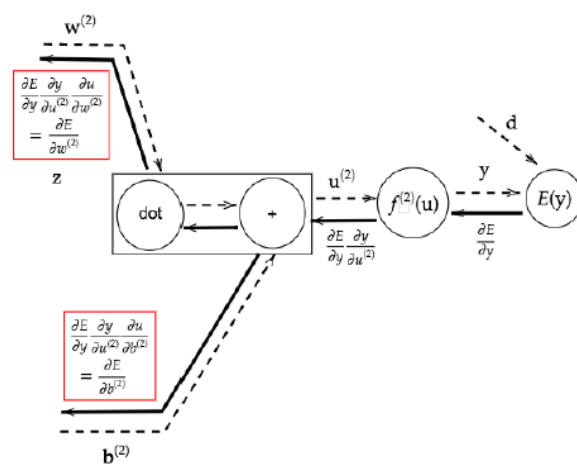
【誤差逆伝播法】

算出された誤差を、出力層側から順に微分し、前の層前の層へと伝播。

最小限の計算で各パラメータでの微分値を解析的に計算する手法

誤差勾配の計算では、誤差逆伝播法を用いる。

誤差から微分を逆算することで、不要な再帰的計算（計算量大）を避けて微分を算出できる



誤差逆伝播法の計算式

微分の連鎖律を用い計算量を減らしている。

$$E(y) = \frac{1}{2} \sum_{j=1}^I (y_j - d_j)^2 = \frac{1}{2} \|y - d\|^2 \quad : \text{誤差関数} = \text{二乗誤差関数}$$

$$y = u^{(L)} \quad : \text{出力層の活性化関数} = \text{恒等写像}$$

$$u^{(l)} = w^{(l)} z^{(l-1)} + b^{(l)} \quad : \text{総入力} の \text{計算}$$

$$\frac{\partial E}{\partial w_{ji}^{(2)}} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial u} \frac{\partial u}{\partial w_{ji}^{(2)}}$$

$$\frac{\partial E(y)}{\partial y} = \frac{\partial}{\partial y} \frac{1}{2} \|y - d\|^2 = y - d$$

$$\frac{\partial y(u)}{\partial u} = \frac{\partial u}{\partial u} = 1$$

・確認テスト

[既に行った計算結果を保持しているソースコードを抽出せよ。]

```
66     ## 試してみよう
67     delta1 = np.dot(delta2, W2.T) * functions.d_sigmoid(z1)
```

$\frac{\partial E}{\partial y}$
`delta2 = functions.d_mean_squared_error(d, y)`

$\frac{\partial E}{\partial y} \frac{\partial y}{\partial u}$
`delta1 = np.dot(delta2, W2.T) * functions.d_sigmoid(z1)`

$\frac{\partial E}{\partial y} \frac{\partial y}{\partial u} \frac{\partial u}{\partial w_{ji}^{(2)}}$
`grad['W1'] = np.dot(x.T, delta1)`

2つの空欄に該当するソースコードを探せ (3分)

※ここで用いられるz1は以下のコードで生成される

`z1, y = forward(network, x)`

・実習演習

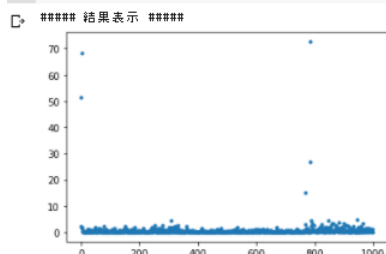
1_3_stochastic_gradient_descent.ipynb

学習率を 0.07 から 0.5 へ変更した場合を下記に示す。

```

98 losses = []
99 # 学習率
100 #learning_rate = 0.07
101 learning_rate = 0.5
102
103 # 抽出数
104 epoch = 1000
105
106 # パラメータの初期化
107 network = init_network()
108 # データのランダム抽出
109 random_datasets = np.random.choice(data_sets, epoch)
110
111 # 勾配降下の繰り返し
112 for dataset in random_datasets:
113     x, d = dataset['x'], dataset['d']
114     z1, y = forward(network, x)
115     grad = backward(x, d, z1, y)
116     # パラメータに勾配適用
117     for key in ('W1', 'W2', 'b1', 'b2'):
118         network[key] -= learning_rate * grad[key]
119
120 # 誤差
121 loss = functions.mean_squared_error(d, y)
122 losses.append(loss)
123
124 print("#### 結果表示 ####")
125 lists = range(epoch)
126
127
128 plt.plot(lists, losses, '.')
129 # グラフの表示
130 plt.show()

```



Section4：勾配降下法で実施した際に比べ誤差が収束せず、800 エポック前後で誤差増加が発生した。このことから誤差逆伝播法での学習時の学習率は重要であることがわかる。