

# 自然言語処理 —CNN—

<https://satoyoshiharu.github.io/nlp/>

# 100本ノック第9章とCNNの位置づけ

- 100本ノック課題集第9章は、RNN、CNN、Transformerを扱っている。
- CNNは、現在、画像処理系では基本となる技術である。100本ノックの課題では、自然言語処理へ応用する事例が載せられている。



# 自然言語処理 CNN

[解説動画](#)



# Convolution層 (CNN)

特徴を抽出する(2次元データの場合)

1	2	3	4	5
6	7	8	9	0
1	2	3	4	5
6	7	8	9	0
1	2	3	4	5

入力

\*

0	1	0
0	1	0
0	1	0

フィルター



	11	14	17	
	16	19	22	
	11	14	17	

出力



# Convolution層

1	0	2	1	3	0	4	5
6	0	7	1	8	0	9	0
1	0	2	1	3	0	4	5
6	7	8	9	0			
1	2	3	4	5			

入力の窓内の値  
と、フィルターの値  
の積和をとっていく



	11	14	17	
	16	19	22	
	11	14	17	



# Convolution層

1	20	31	40	5
6	70	81	90	0
1	20	31	40	5
6	7	8	9	0
1	2	3	4	5

入力の窓内の値  
と、フィルターの値  
の積和をとっていく



	11	14	17	
	16	19	22	
	11	14	17	



# Convolution層

1	2	3 <sup>0</sup>	4 <sup>1</sup>	5 <sup>0</sup>
6	7	8 <sup>0</sup>	9 <sup>1</sup>	0 <sup>0</sup>
1	2	3 <sup>0</sup>	4 <sup>1</sup>	5 <sup>0</sup>
6	7	8	9	0
1	2	3	4	5

入力の窓内の値  
と、フィルターの値  
の積和をとっていく



	11	14	17	
	16	19	22	
	11	14	17	



# Convolution層

1	2	3	4	5
60	71	80	9	0
10	21	30	4	5
60	71	80	9	0
1	2	3	4	5

入力の窓内の値  
と、フィルターの値  
の積和をとっていく

	11	14	17	
	16	19	22	
	11	14	17	

以下省略...





# フィルタの各要素が学習パラメータ

1	2	3	4	5
6	7	8	9	0
1	2	3	4	5
6	7	8	9	0
1	2	3	4	5

入力

\*

$k_{11}$	$k_{12}$	$k_{13}$
$k_{21}$	$k_{22}$	$k_{23}$
$k_{31}$	$k_{32}$	$k_{33}$

フィルター

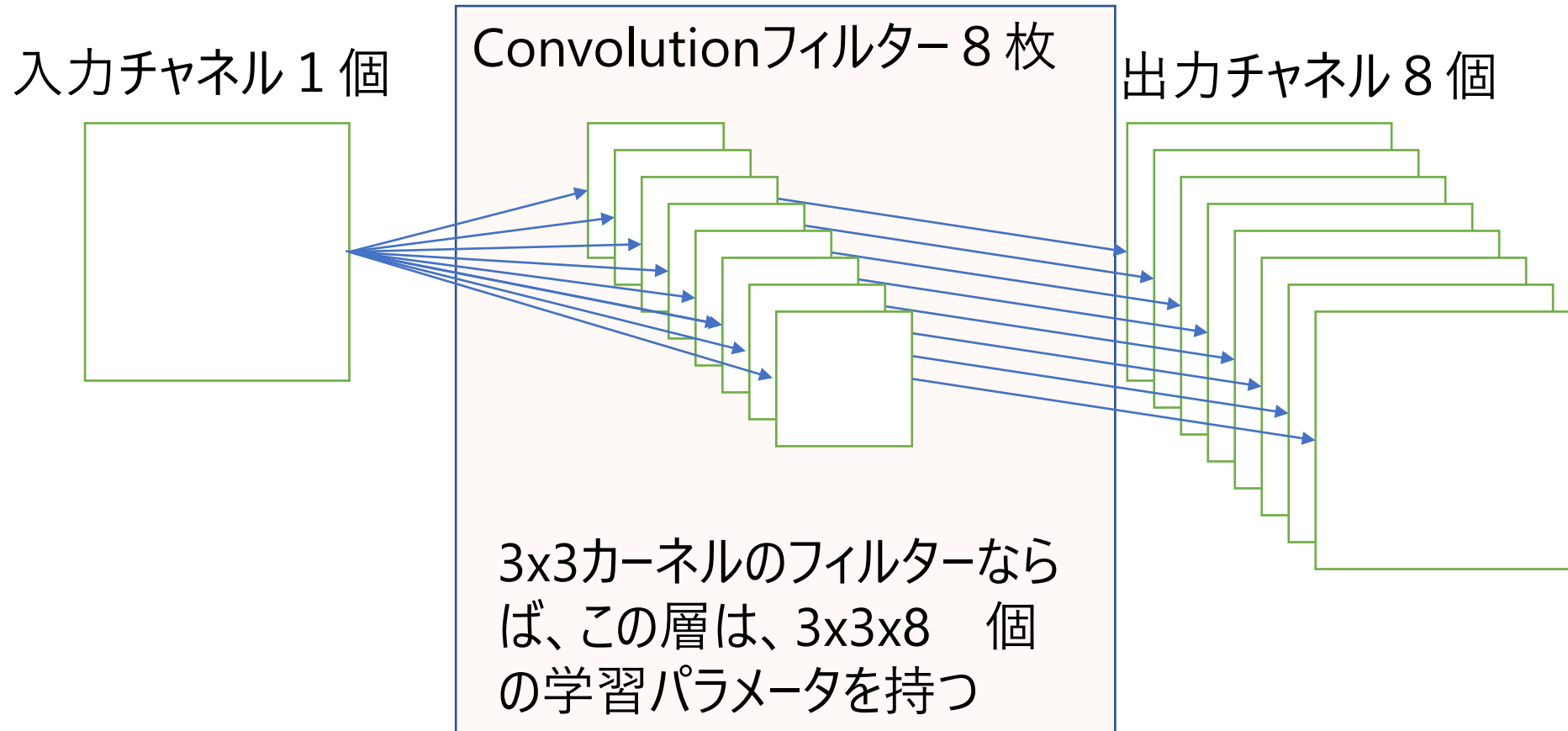


	11	14	17	
	16	19	22	
	11	14	17	

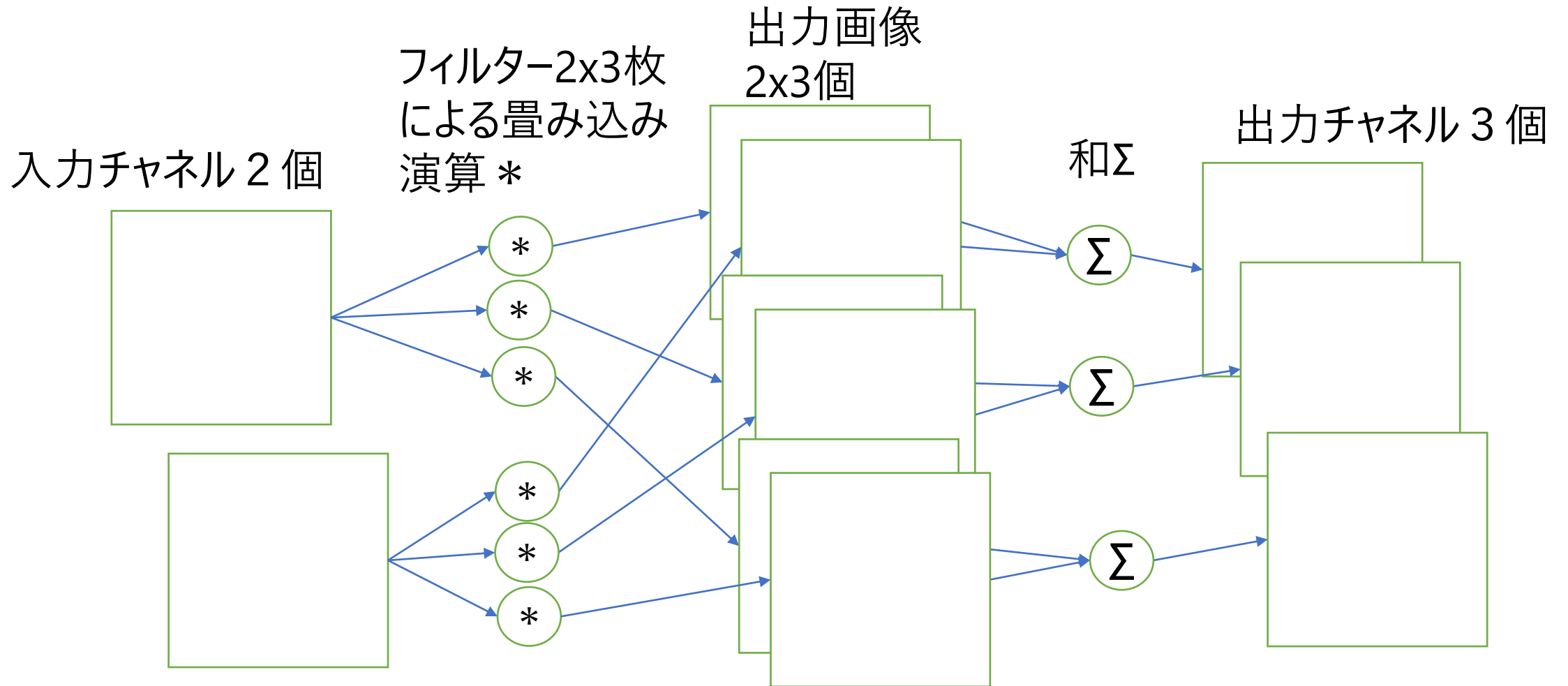
出力



# Convolution層の入出力チャンネルとパラメータ

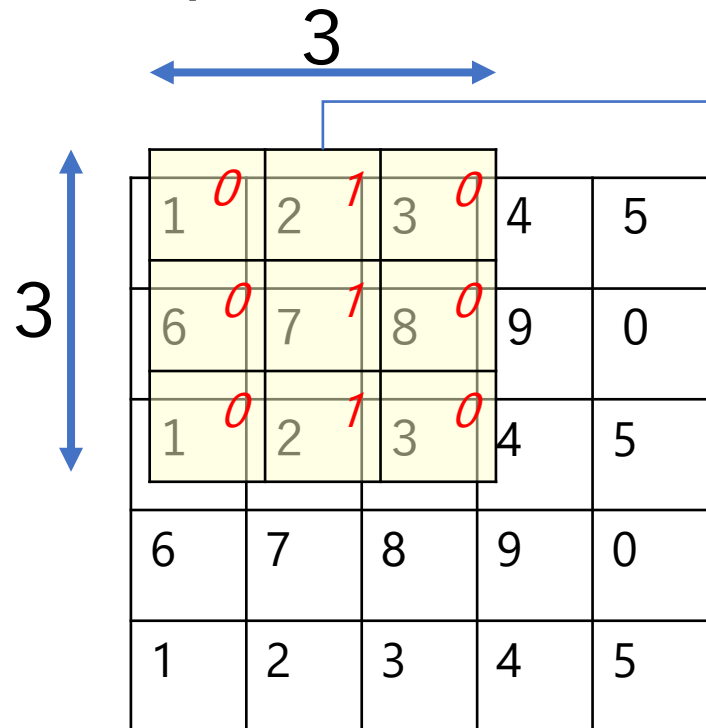


# 入力チャンネルが複数の場合



# カーネルサイズ

3x3のカーネル



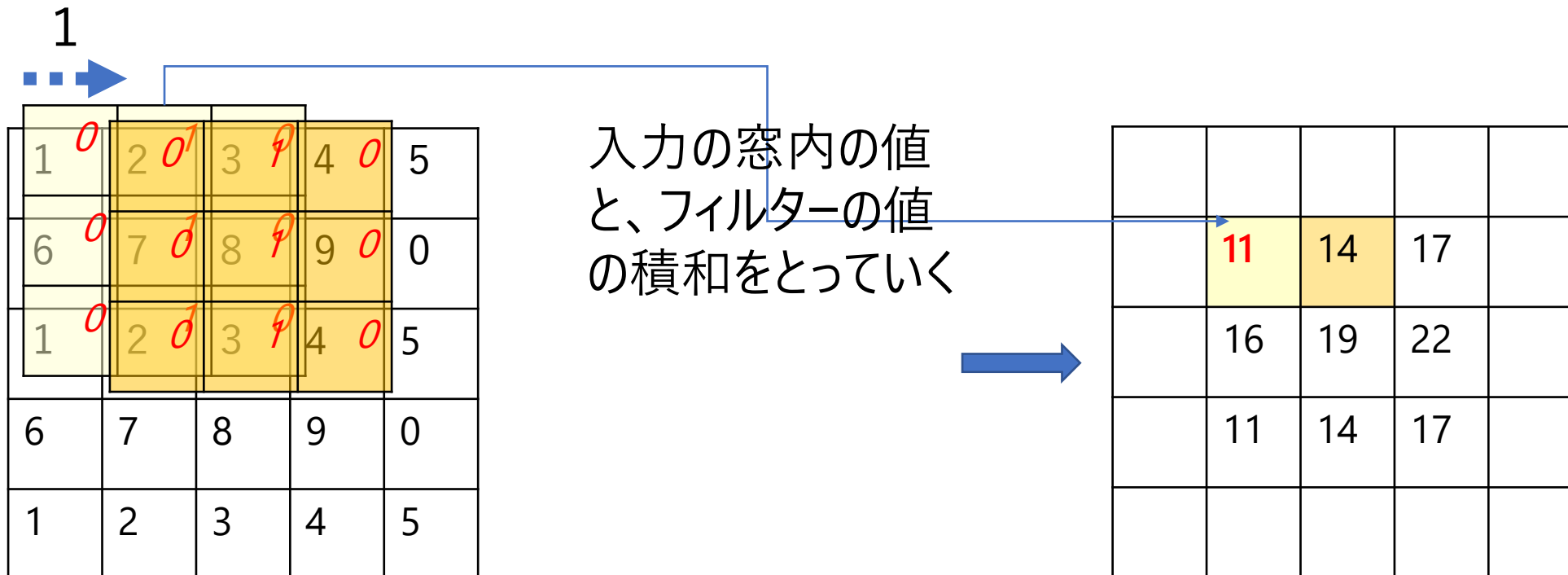
入力の窓内の値  
と、フィルターの値  
の積和をとっていく



カーネルサイズが大きいほど、グローバルに観察できるが、学習パラメータが増える



# ストライド（歩幅、移動幅）



ストライドが小さいほどきめ細かに観察するが、計算量が増える。



# パディング

1	0	2	1	3	0	4	5
6	0	7	1	8	0	9	0
1	0	2	1	3	0	4	5
6	7	8	9	0			
1	2	3	4	5			

入力の窓内の値  
と、フィルターの値  
の積和をとっていく



11	14	17
16	19	22
11	14	17

Convolutionはデフォルトで出力は小さくなる。



# パディング

0	1	0			
0	1	2	3	4	5
0	6	7	8	9	0
	1	2	3	4	5
	6	7	8	9	0
	1	2	3	4	5

入力の窓内の値  
と、フィルターの値  
の積和をとっていく



7	9	11	13	5
8	11	14	17	10
13	16	19	22	5
8	11	14	17	10
7	9	11	13	5

外境界上のConvolutionに関し、入力の外側に0とかを想定して計算し、出力を同じサイズにする。



# ディレイション(拡張、伸長)

ディレイション=1

1	0	2	1	3	0	4	5
6	0	7	1	8	0	9	0
1	0	2	1	3	0	4	5
6	7	8	9	0			
1	2	3	4	5			

入力の窓内の値  
と、フィルターの値  
の積和をとっていく

11	14	17
16	19	22
11	14	17

ディレイション=2

1	0	2	3	1	4	5	0	
6	7	8	9	0				
1	0	2	3	1	4	5	0	
6	7	8	9	0				
1	0	2	3	0	4	-	5	0

入力の窓内の値  
と、フィルターの値  
の積和をとっていく

6

よりグローバルに入力を観察して、高速に出力を得る



# CNNの貢献と今

- CNNはDeep Learningの発展を支えた。
- 今、Self-AttentionがCNNに影響を及ぼしている。



自然言語処理  
MaxPooling

[解説動画](#)



# MaxPooling

1	2	3	4
6	7	8	9
1	2	3	4
6	7	8	9

入力



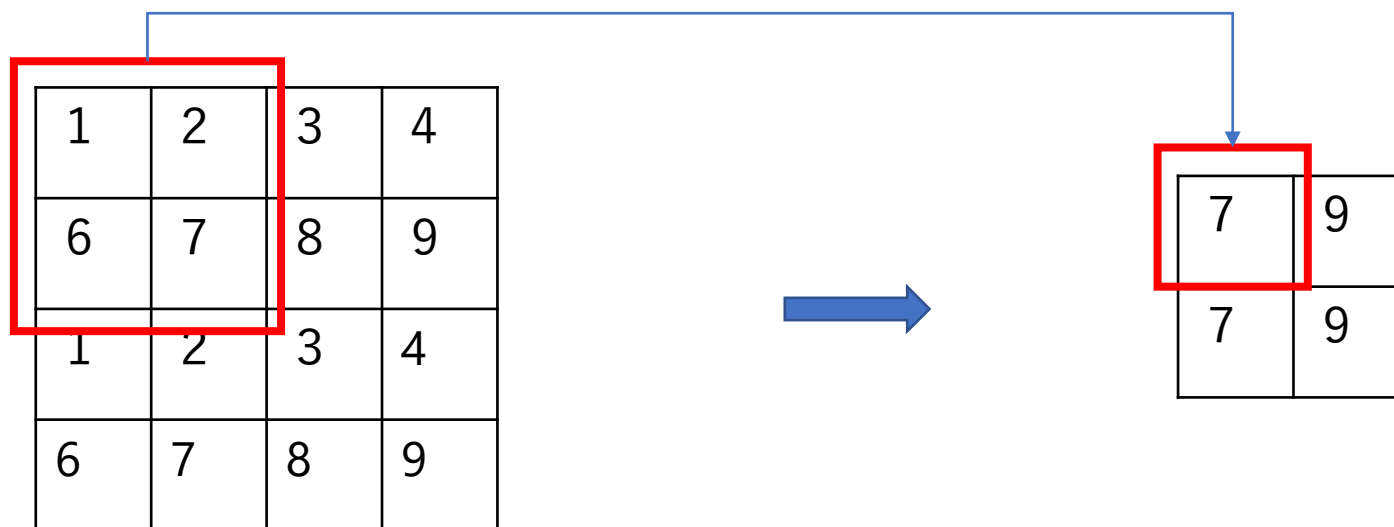
7	9
7	9

出力



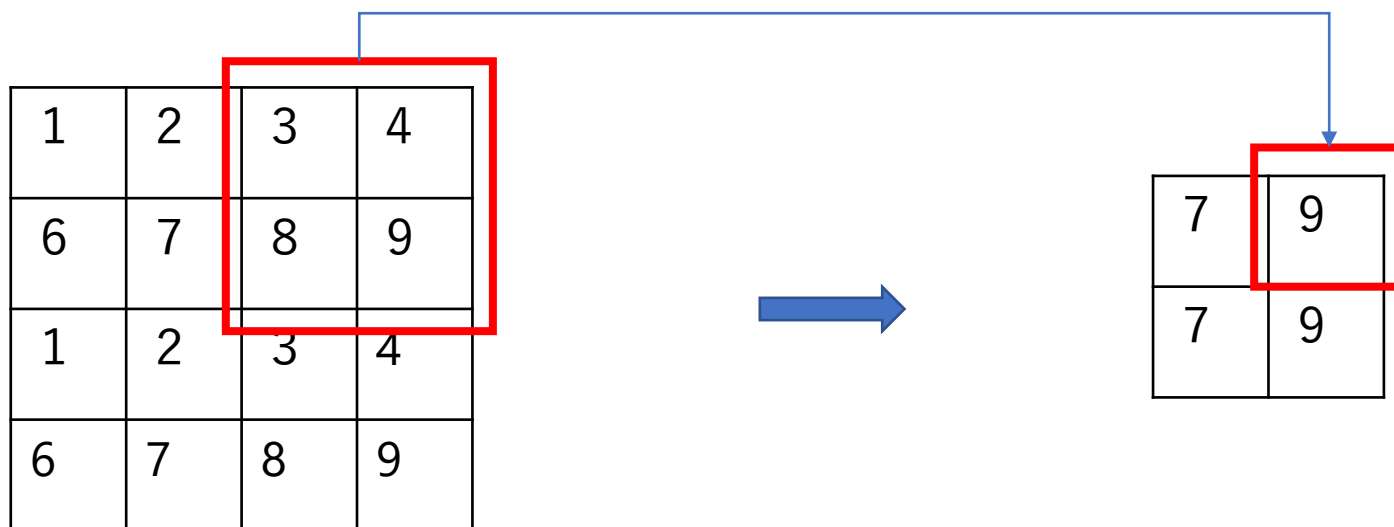
# MaxPooling

最も大きな値を拾っていく



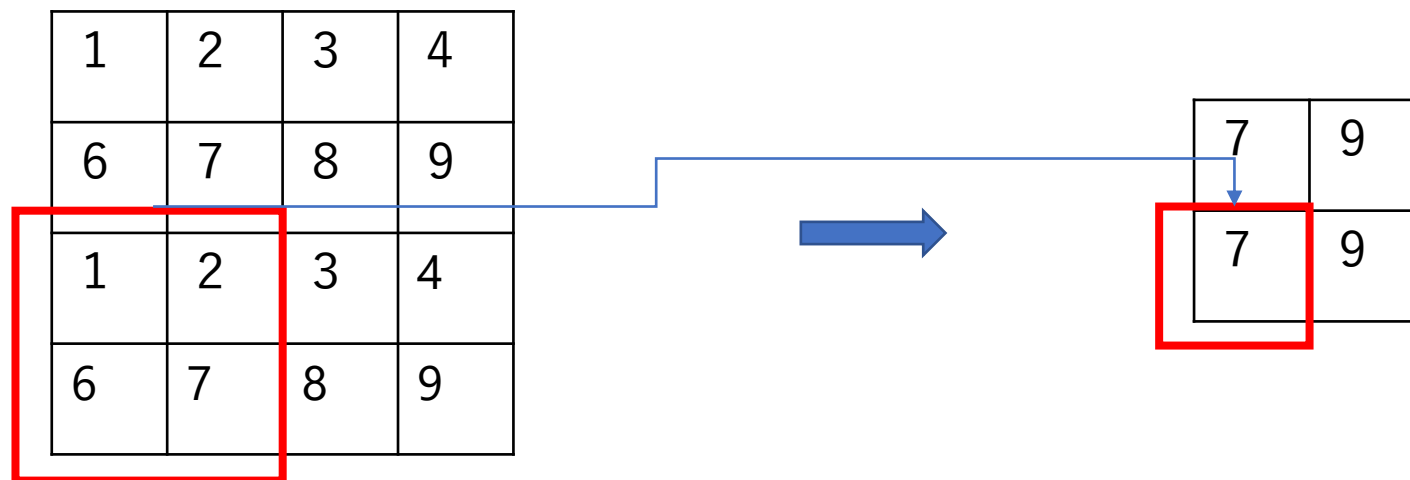
# MaxPooling

最も大きな値を拾っていく



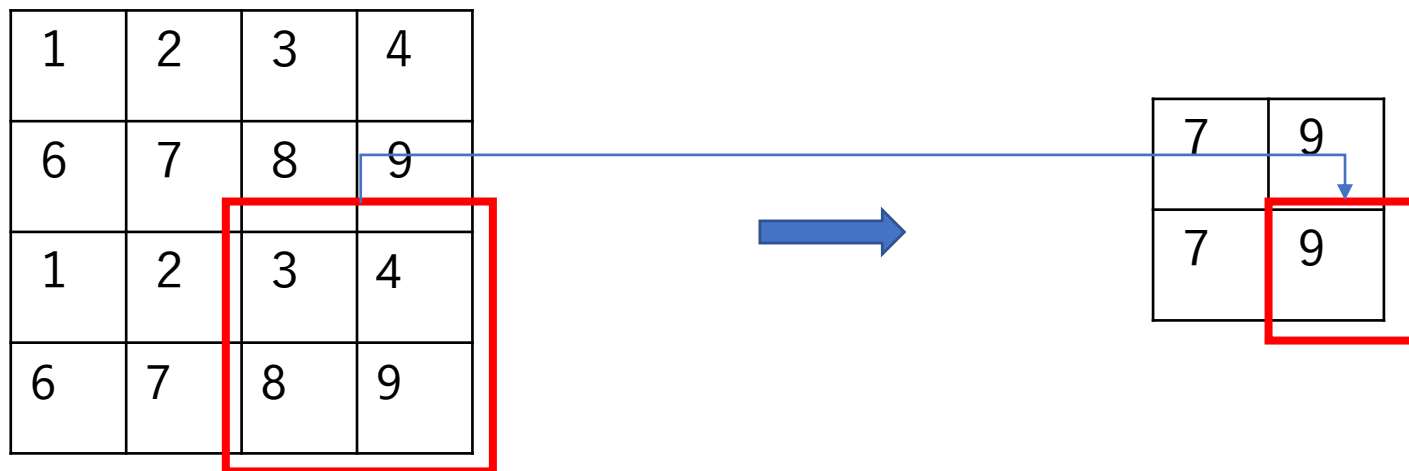
# MaxPooling

最も大きな値を拾っていく



# MaxPooling

最も大きな値を拾っていく



# MaxPooling

- 広い範囲の情報を集約する。
- データを小さくして後続の計算量を小さくする。



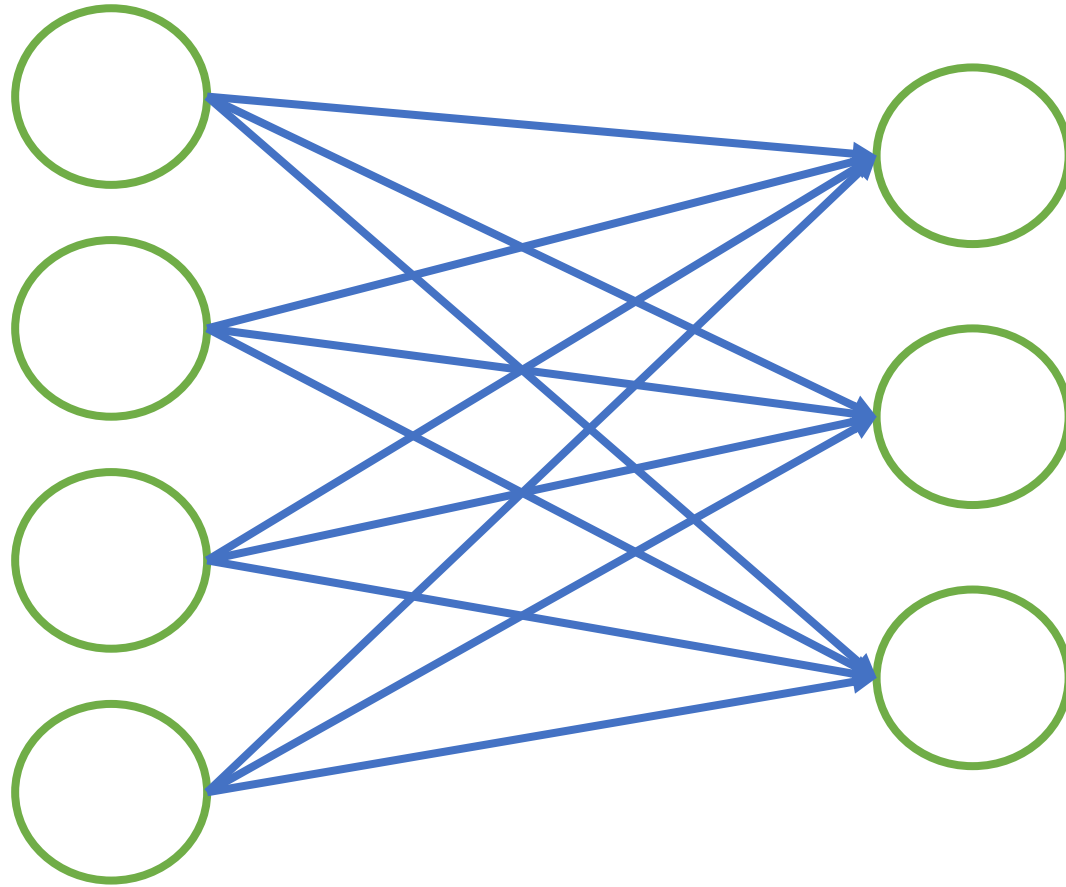


# 自然言語処理 Dropout

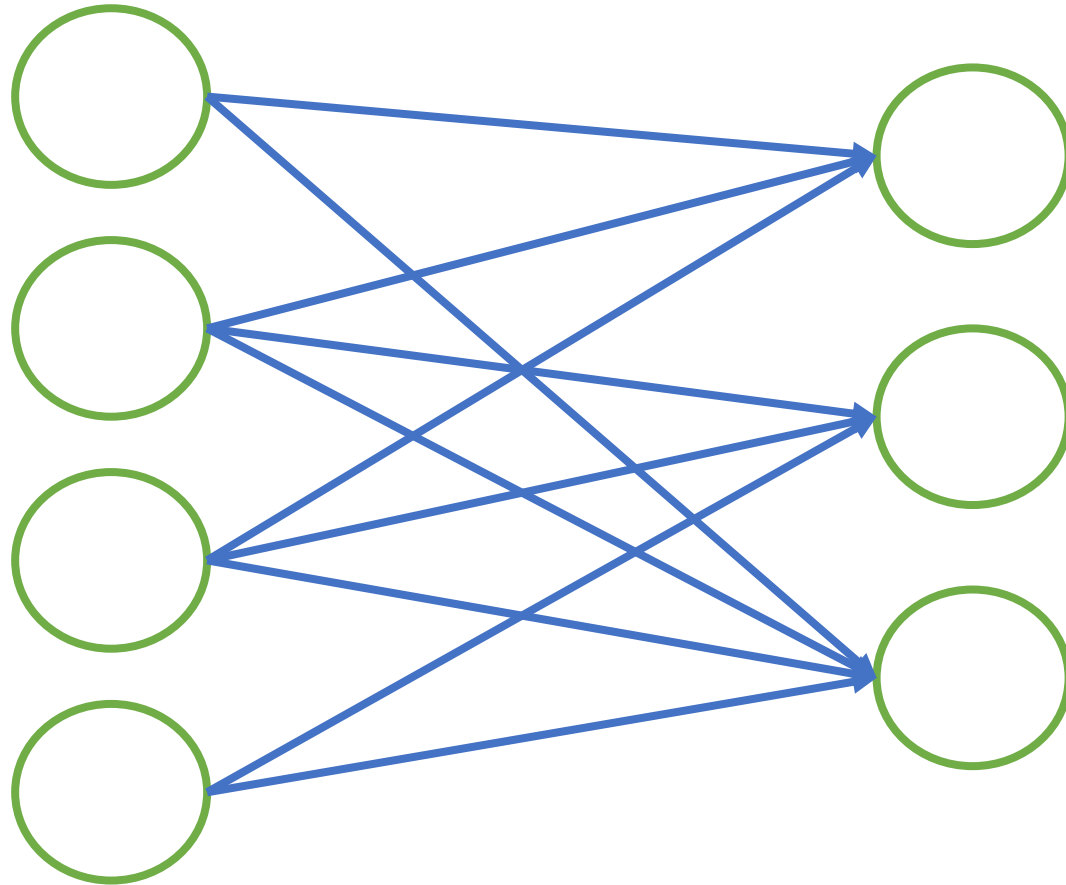
[解説動画](#)



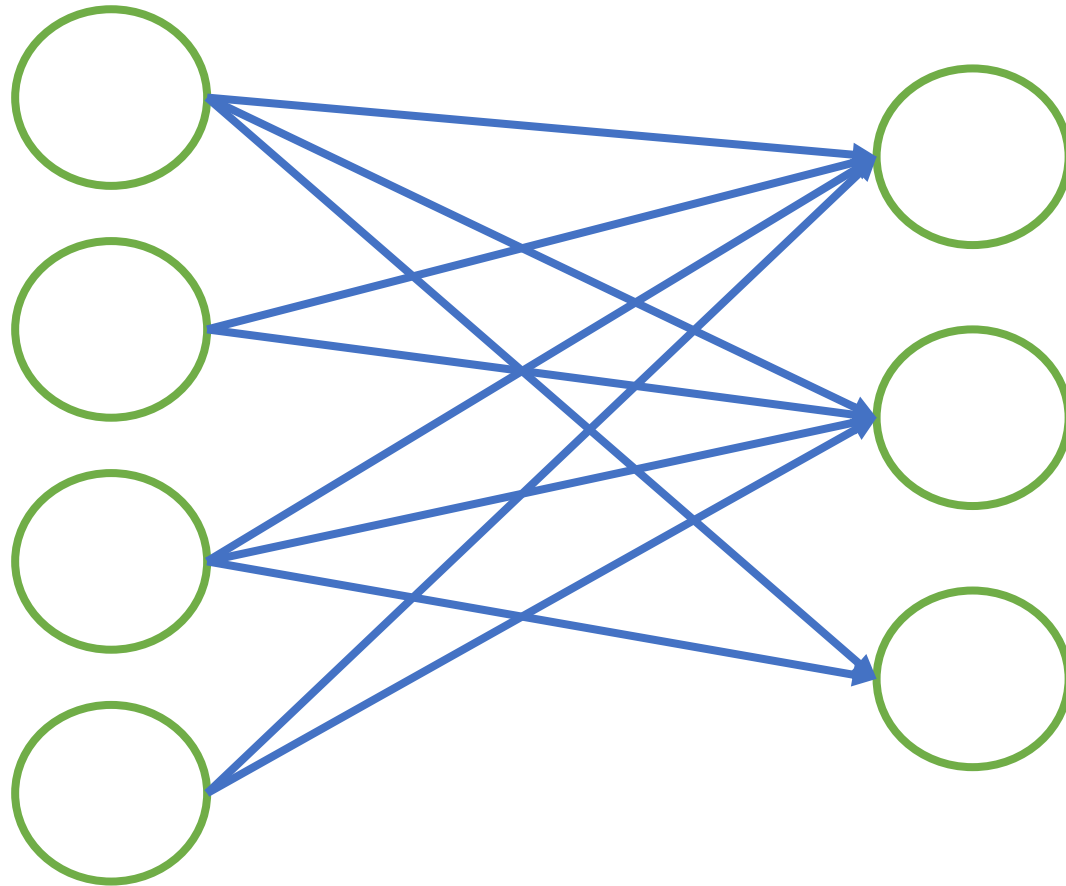
# Dropout



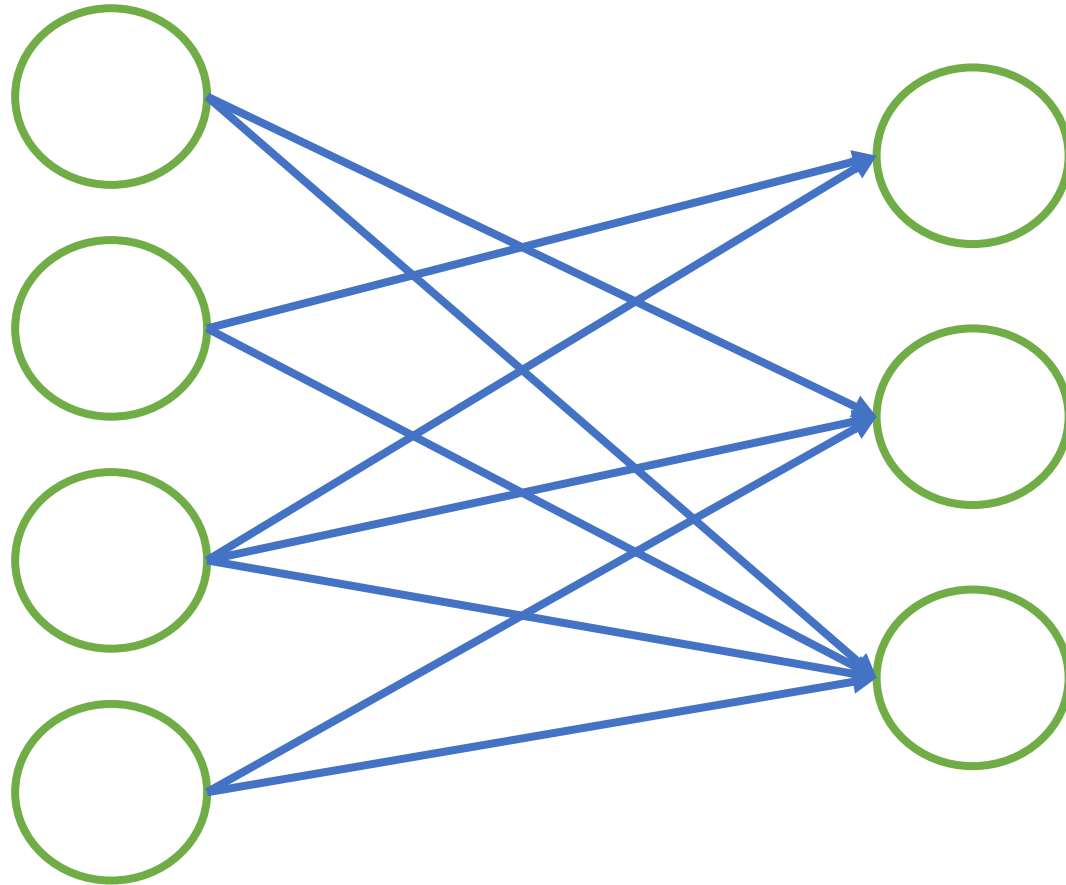
# Dropout



# Dropout



# Dropout



# Dropout

仮想的にサブネットワークを構成しそれらの多数決で出力を得る、感じにすることで、般化能力を高める、と言われる（経験的）。



# 自然言語処理 Batch Normalization

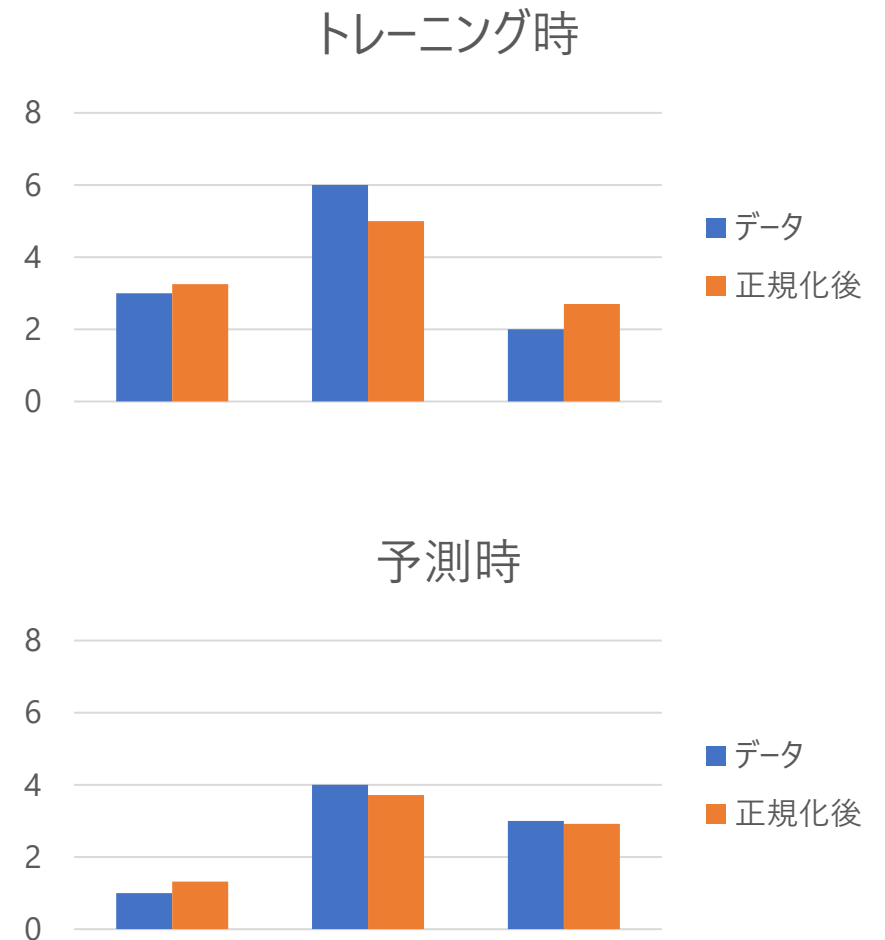
[解説動画](#)



# Normalizationとは？

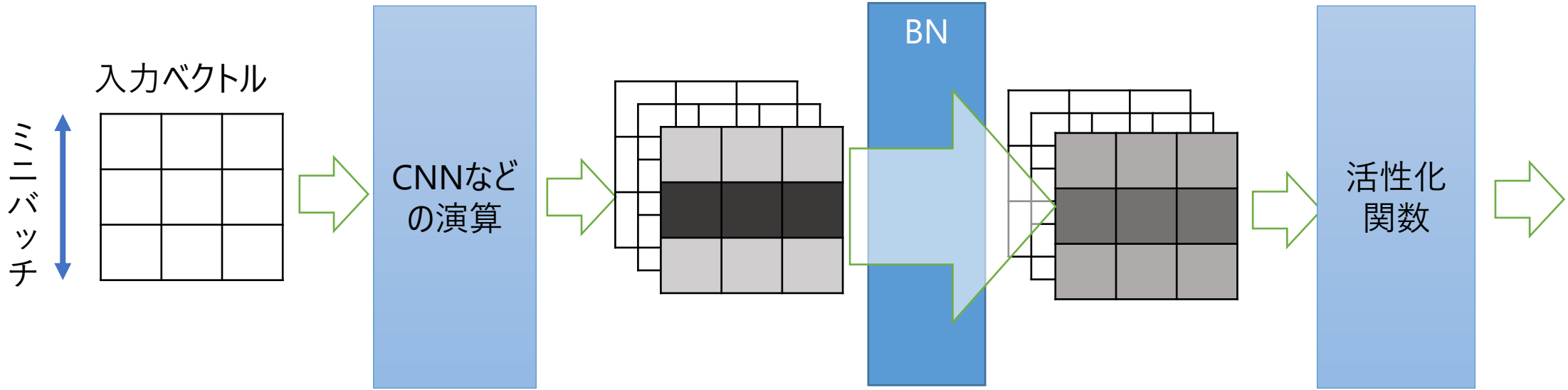
- 入力データや推論中間データが、トレーニング時と予測時とで、分布が大きく異なると困る。
- ある範囲のデータに関して、分布の偏り（凸凹）を小さくする。

$$\mu = \frac{1}{m} \sum_i H_{i,:}$$
$$\sigma = \sqrt{\delta + \frac{1}{m} \sum_i (H - \mu)_i^2}$$
$$H' = \frac{H - \mu}{\sigma}$$





# Batch Normalization



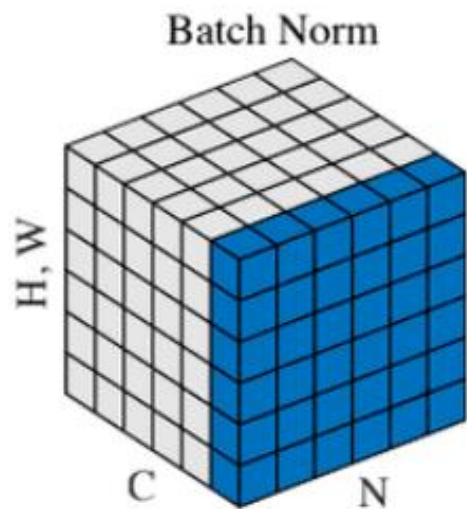
# Batch Normalization

- データのばらつきを抑えることで、学習が高速化し、過学習にも効果がある、と言われる。

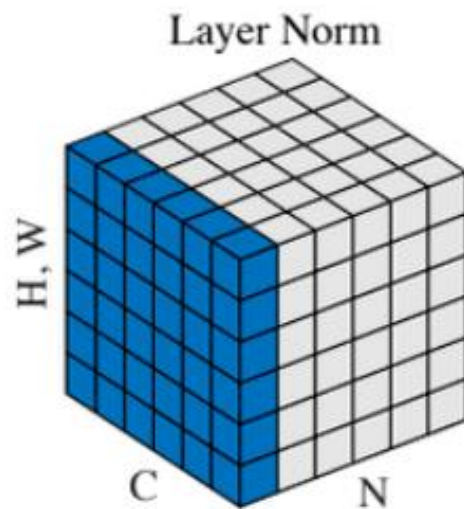


# データの平均と分散をとって補正する

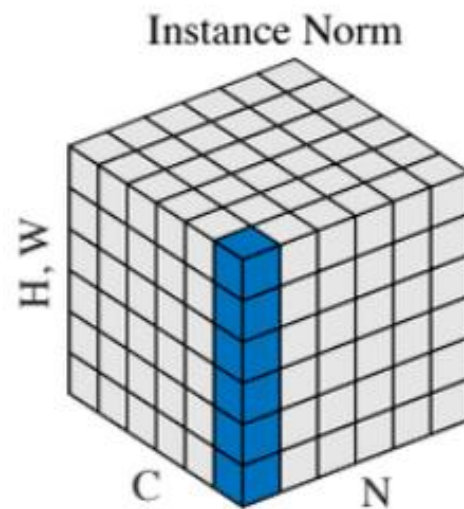
以下、入力データをシリアライズし (H,W)、チャンネル (C)、ミニバッチ (N) との3次元立体で示す



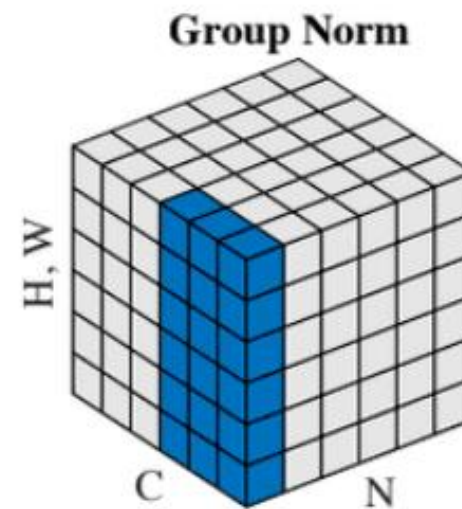
ミニバッチ内の正規化



チャンネル内の正規化



1 入力データ内の正規化



特定のチャンネルグループ内の正規化



MNISTサンプルコード

# MNISTサンプル

- 自然言語処理ではなく、画像処理になりますが、ディープラーニングの典型的な流れと構成要素の概略を把握するために、いいサンプルですので、見ておきます（ニューラルネットの入門書で最終章に用いられていたりする）。
- これは、CNNを利用しています。自然言語処理でもCNNを使うことがあります（例：100本ノック課題86）。

# 課題MNISTサンプルの理解

- Mnist\_sample.ipynbをダウンロードして、Google Colaboratoryで動かしてみてください。
- 最初のコードは、PyTorchのサンプル集の中の <https://github.com/pytorch/examples/tree/master/mnist>（もとのPyTorchのコードは、煩雑になりがちなデータ回りの処理もクラス化されているため、Mnistの主要な処理構造が見えやすいものとなっています。）を、教育用に単純化したものです。後続のコード解説、コードのコメント、PyTorchのマニュアルを参照しながら、このコードを咀嚼して下さい。
- そのあとには、参考のため、Convolution結果を表示したり、モデルを変えたコードを入れてあります。

# MNistサンプルコードを読むための 追加情報

# PyTorchマニュアルの読み方

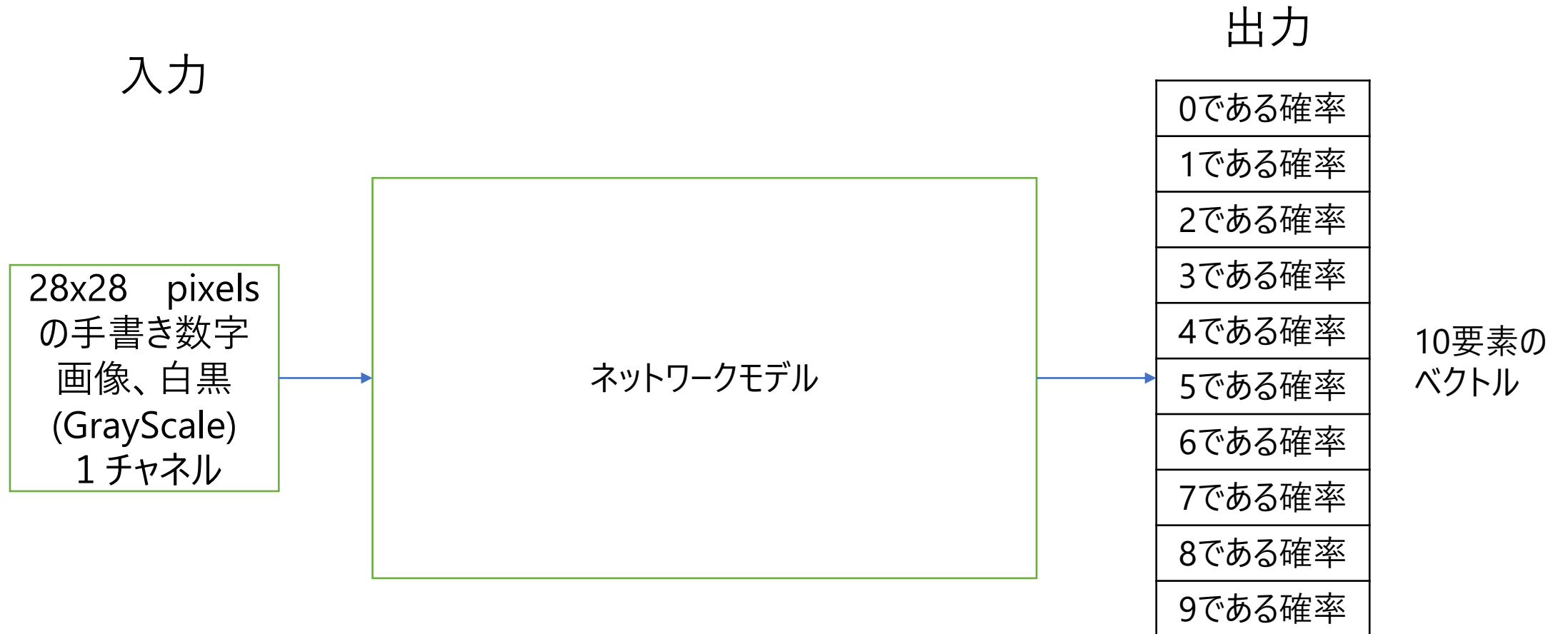
- <https://pytorch.org/docs/stable/index.html> へ行ってください。
- 左上に検索ボックスがあるので、そこにわからない書き方をタイプしてください。
- 各ページは英語です。Chromeの場合、該当ページ上で、右クリックで、「日本語に翻訳」を選んでください。



# torch.nnとtorch.nn.functionalの違い

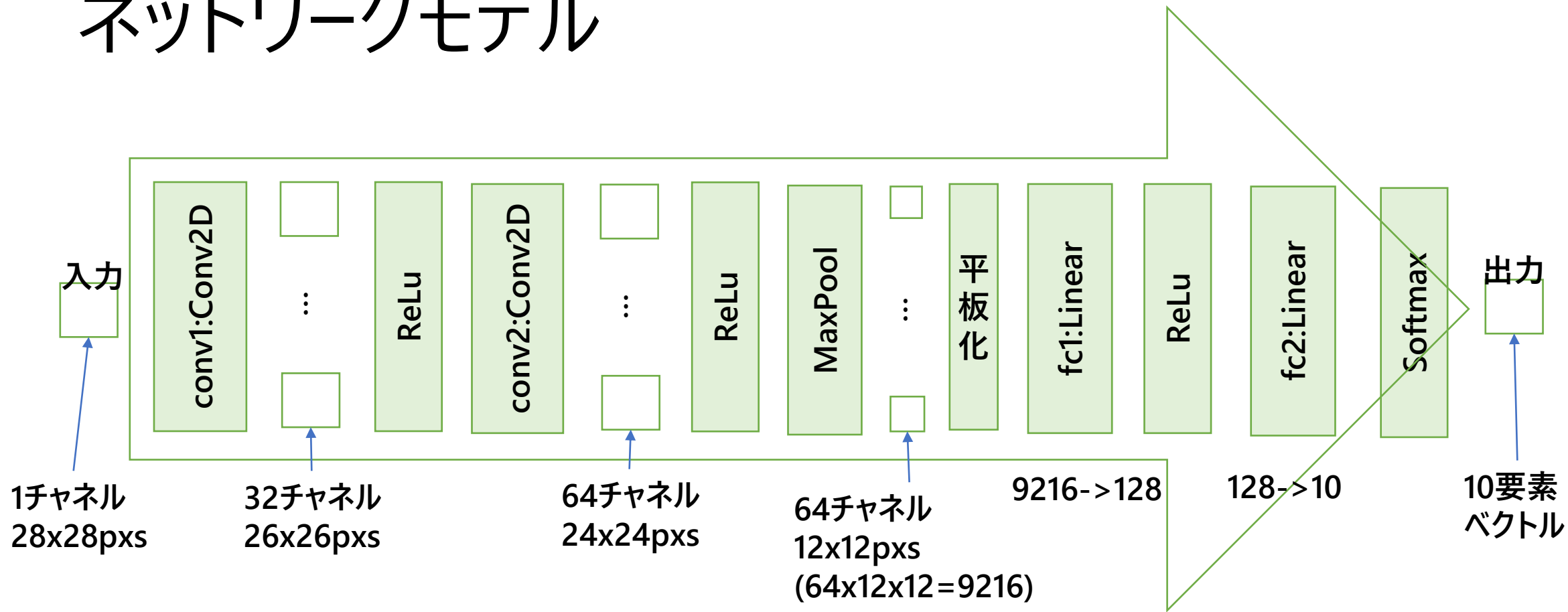
- nn モジュール
  - クラス定義
  - 例えば、ちょっと複雑な初期設定したい場合、オブジェクトを作っておく。
    - `self.conv1=nn.Conv2d(1,32,3,1)`
      - nn パッケージ内の Conv2dクラスのオブジェクトを作成し、初期設定する。
    - `x=self.conv1(x)`
      - conv1オブジェクトを入力データに対し作用させる。Pythonの呼び出し可能オブジェクトの機能を使い、`_call_`している。
- functional モジュール
  - nn クラスの機能を、関数として利用できるようにしたもの。
  - 例えば
    - `x=F.relu(x)`
    - `x=F.max_pool2d(x, 2)`
      - nn モジュールの MaxPool2dのオブジェクトを作成し、引数で初期化して、入力データに作用させる。

# Mnist サンプルの入出力

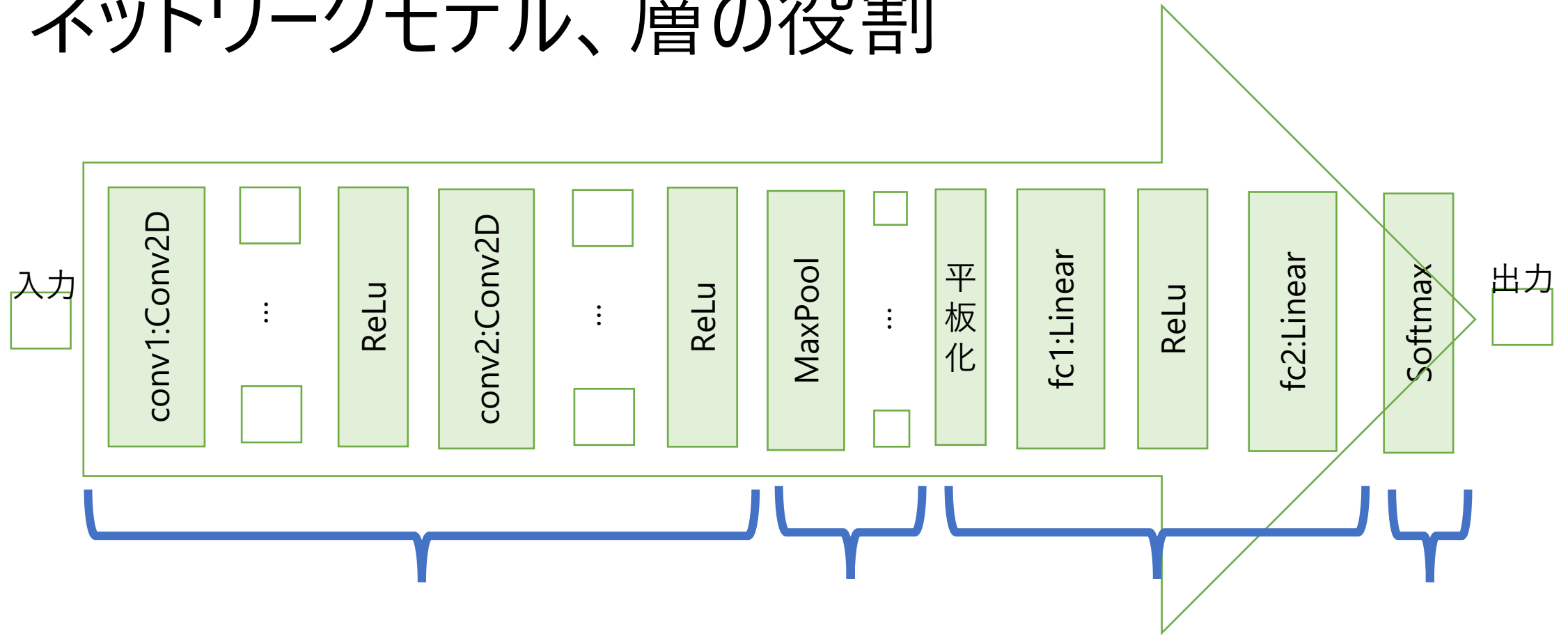


Train/Testは、データローダーを介して  
64個ずつのミニバッチ単位で計算

# ネットワークモデル



# ネットワークモデル、層の役割



CNN2層で画像の何らかの特徴を抽出

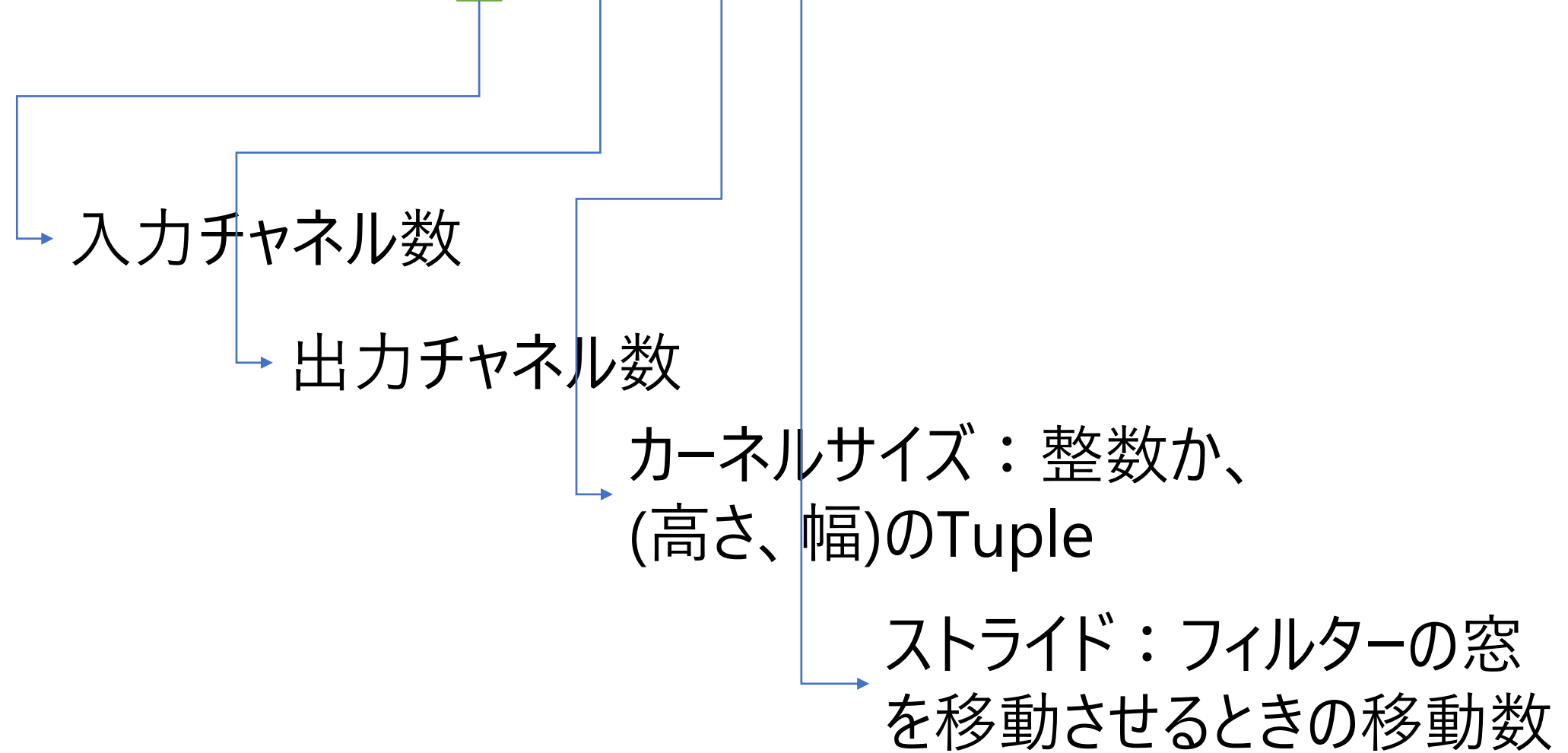
サイズを小さくし、グローバルに特徴をまとめる

全結合2層で情報の取捨選択・集約

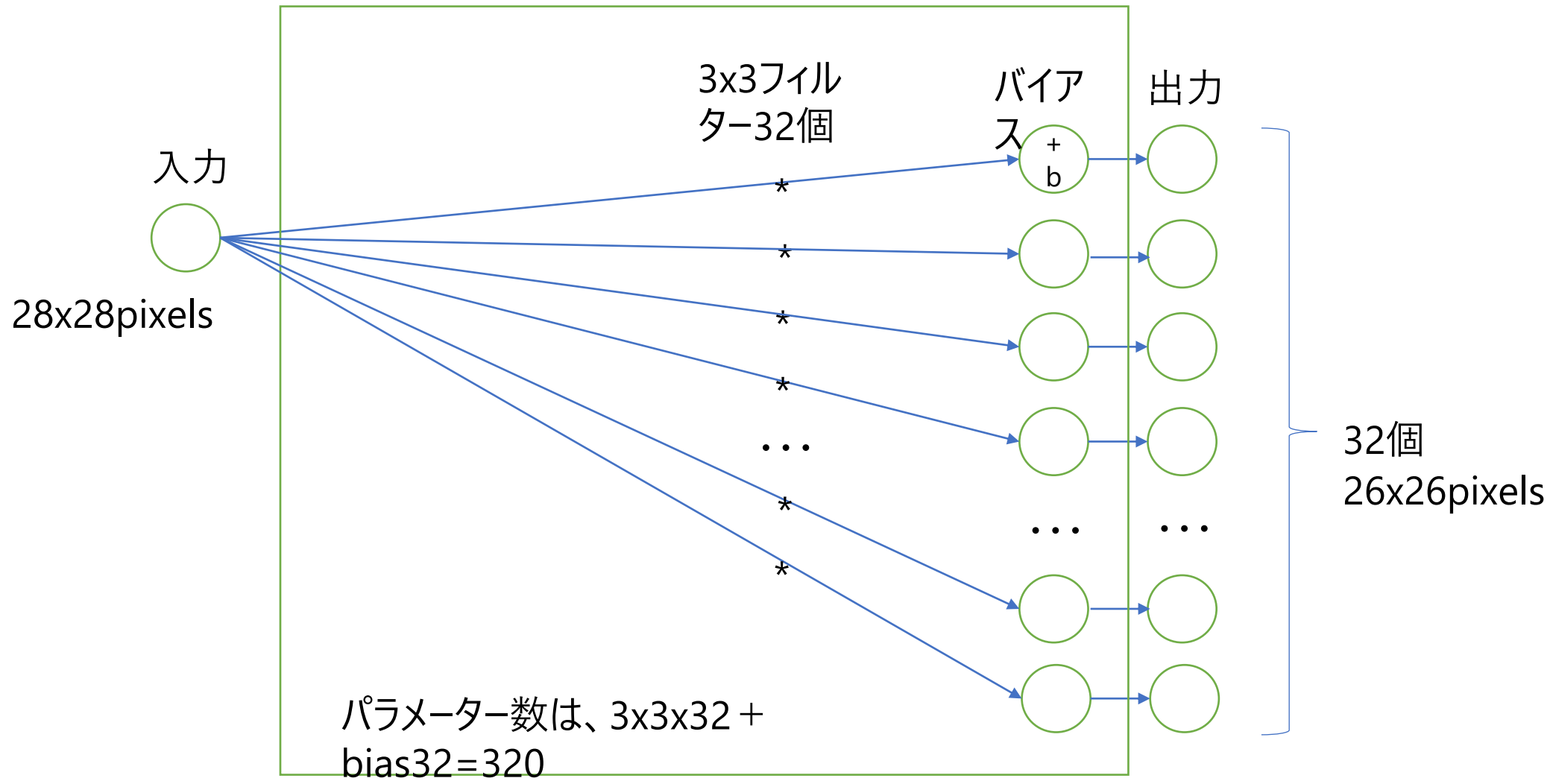
確率に変換

Netクラス\_\_Init\_\_

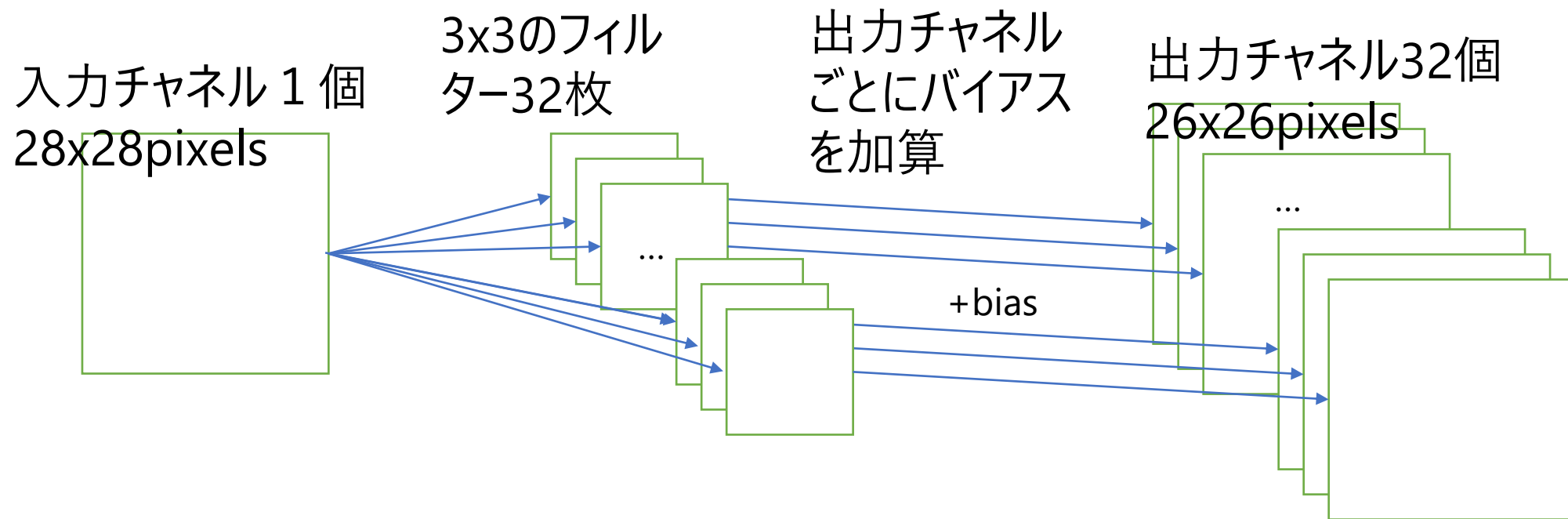
`nn.Conv2d(1, 32, 3, 1)`



# conv1



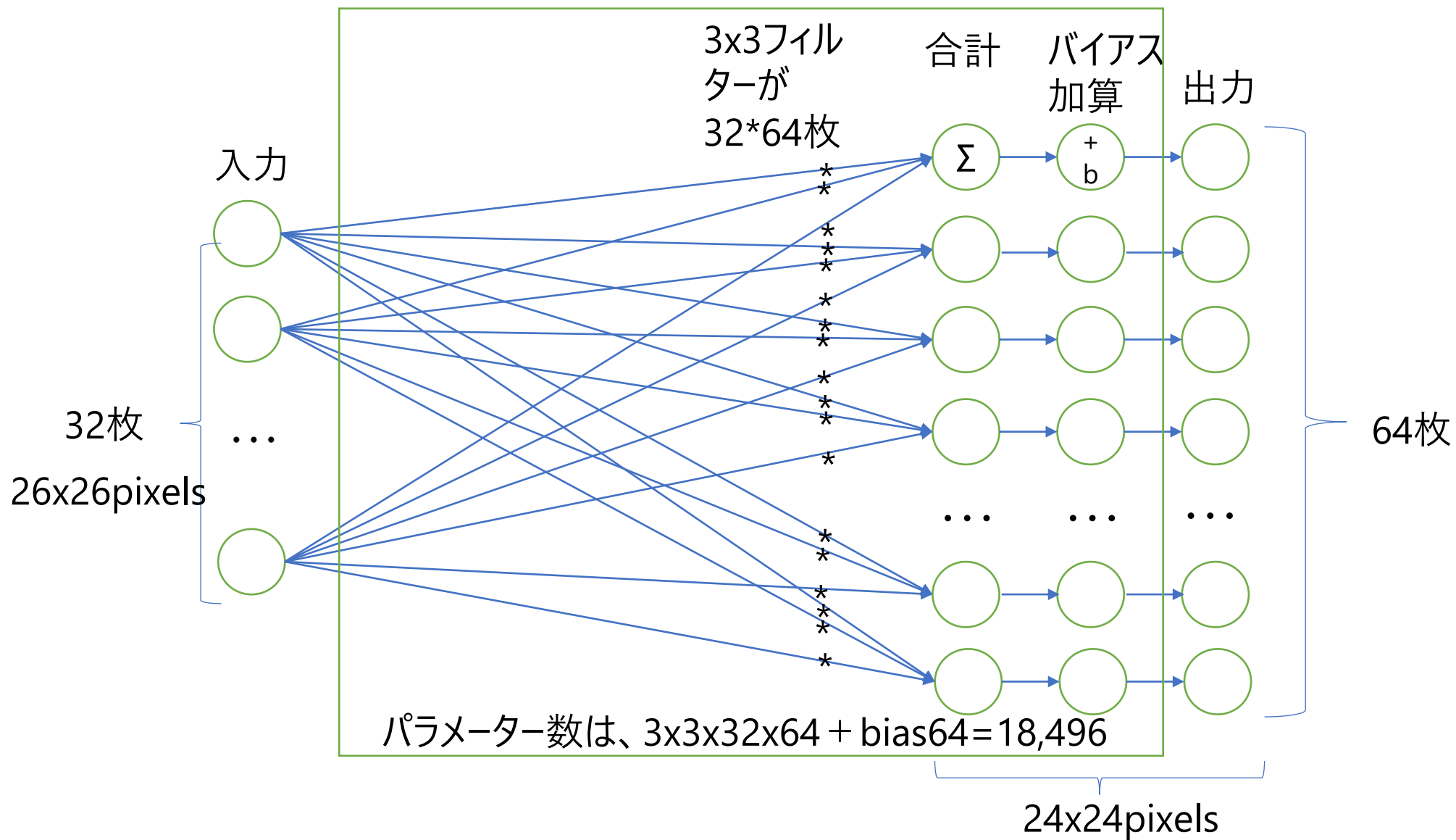
# conv1



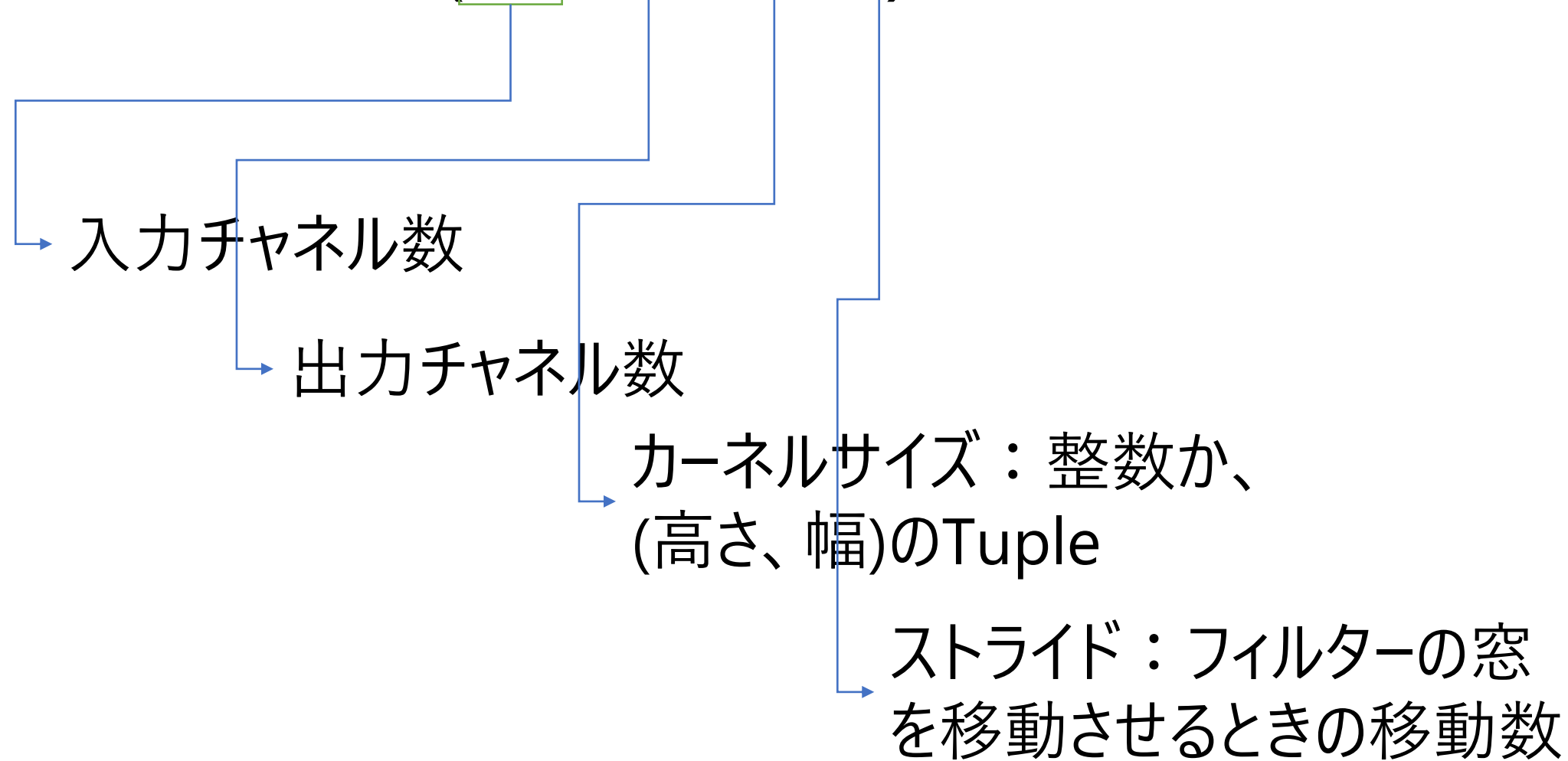
パラメータ数は、 $3 \times 3 \times 32 + \text{bias}32 = 320$



# conv2



```
nn.Conv2d(32, 64, 3, 1)
```



# conv2

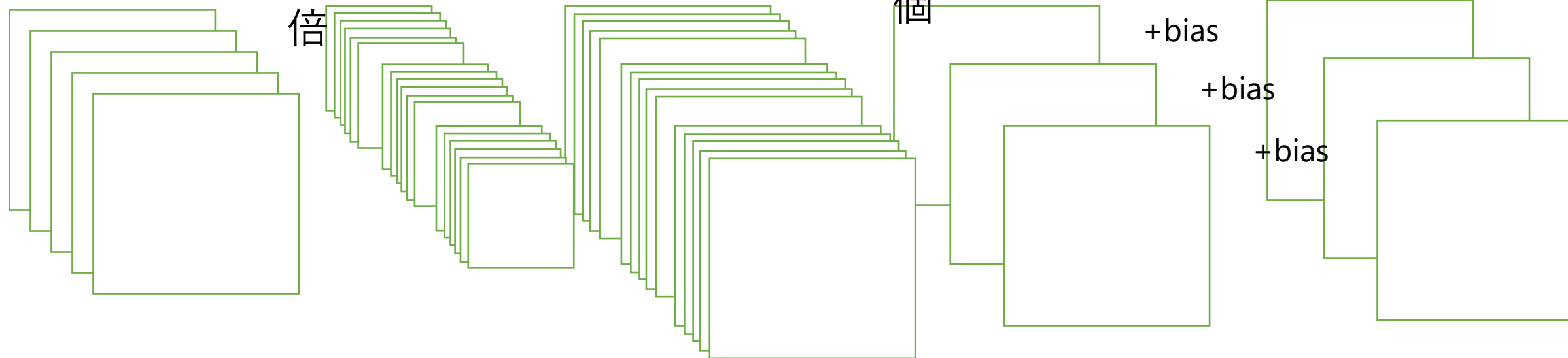
入力32個  
26x26pixels

フィルター  
32個を64  
倍

32x64個  
の結果

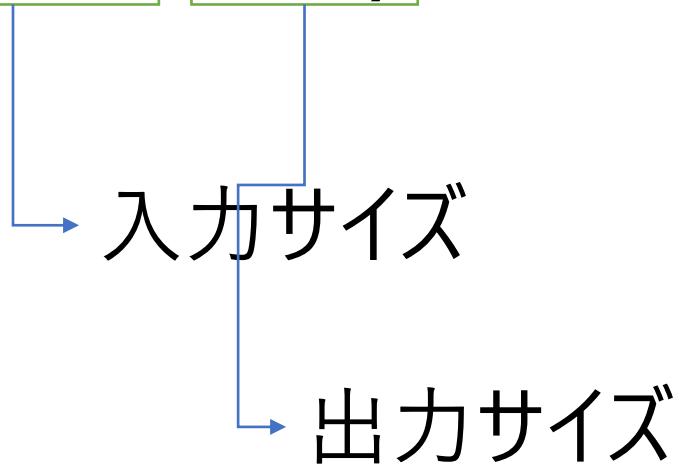
入力32個分の  
結果合計が64  
個

出力64個  
24x24pixels

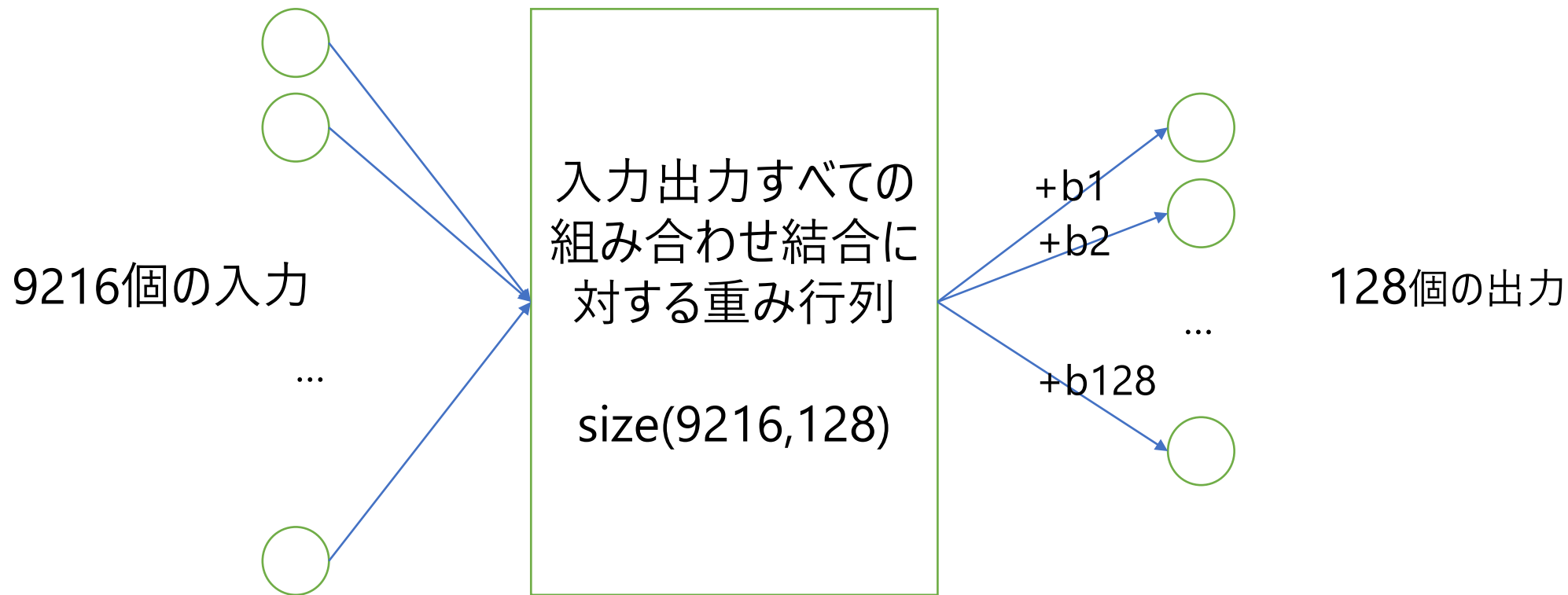


パラメーター数は、 $3 \times 3 \times 32 \times 64 + \text{bias} 64 = 18,496$

nn.Linear(9216, 128)



# fc1



実際には、ミニバッチの数分、まとめて計算され、  
入力も出力も、1次元ベクトルではなく、ミニバッチを行方向に収めた行列となる

nn.Linear(128, 10)

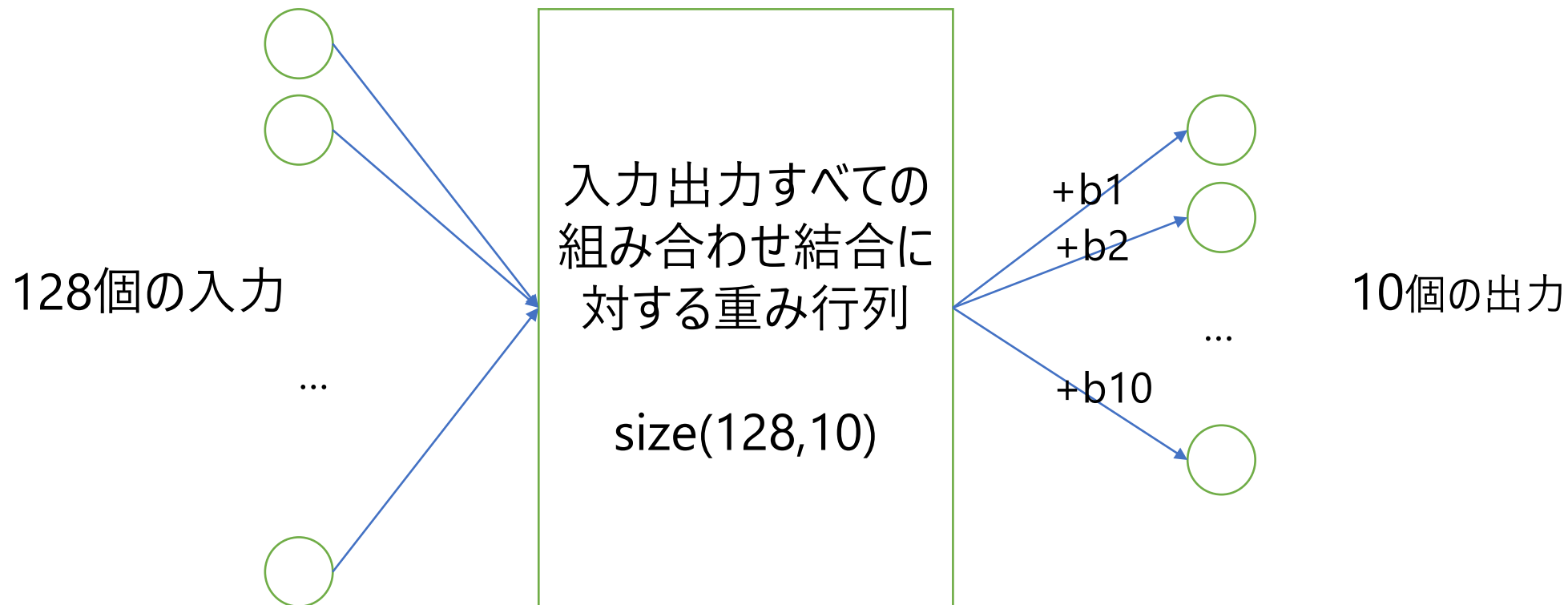


入力サイズ



出力サイズ

# fc2

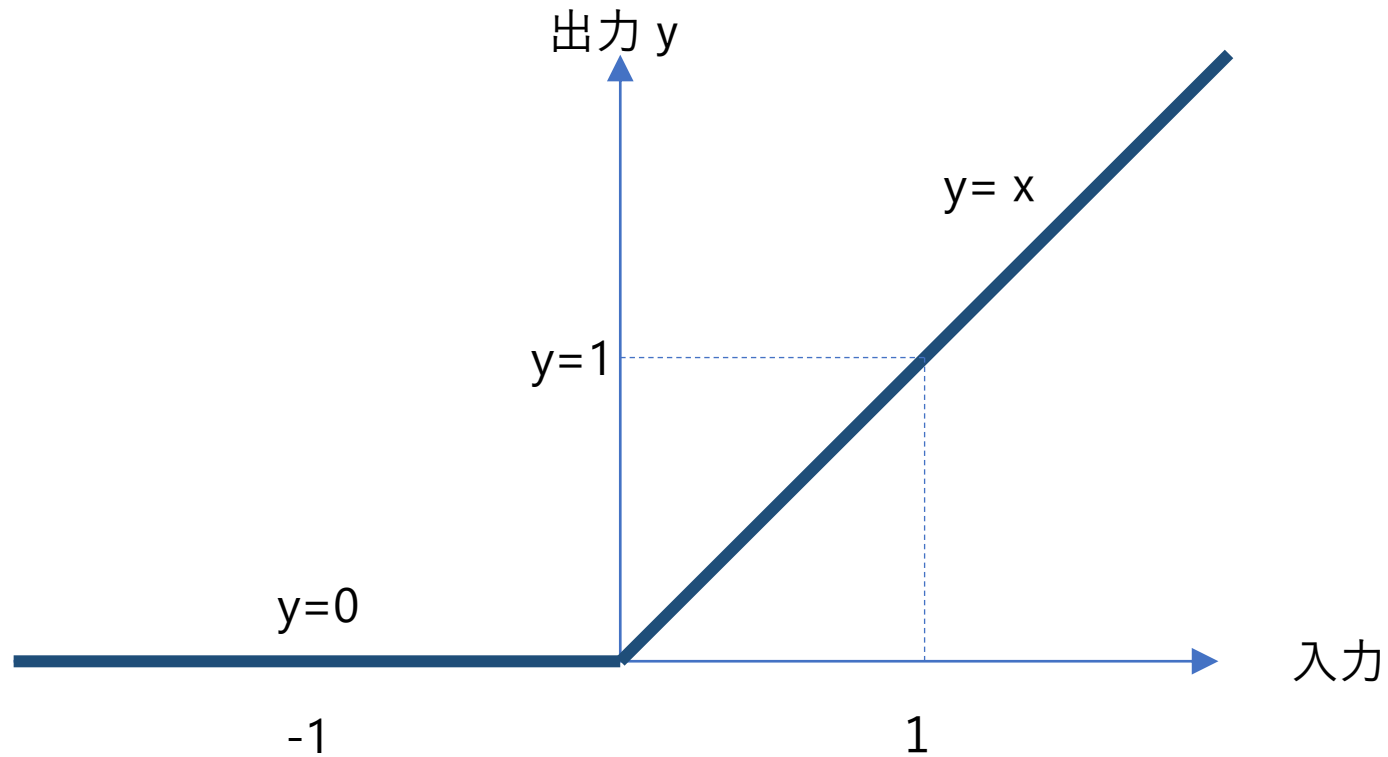


実際には、ミニバッチの数分、まとめて計算され、  
入力も出力も、1次元ベクトルではなく、ミニバッチを行方向に収めた行列となる

Netクラス forward



# 活性化関数ReLU(Rectified Linear Unit)



max\_pool2D(x, 2)

入力：（ミニバッチ数、チャンネル数、高さ、幅）の  
sizeを持つテンソル

カーネルサイズ：整数

省略可能引数にストライド（ウィンドウの移動幅）があり、  
ストライドのデフォルトはカーネルサイズ

MaxPooling：広い範囲の情報を集約する。  
データを小さくして計算量を小さくする。

1	2	3	4
6	7	8	9
1	2	3	4
6	7	8	9

入力



7	9
7	9

出力

# torch.flatten(x, 1)

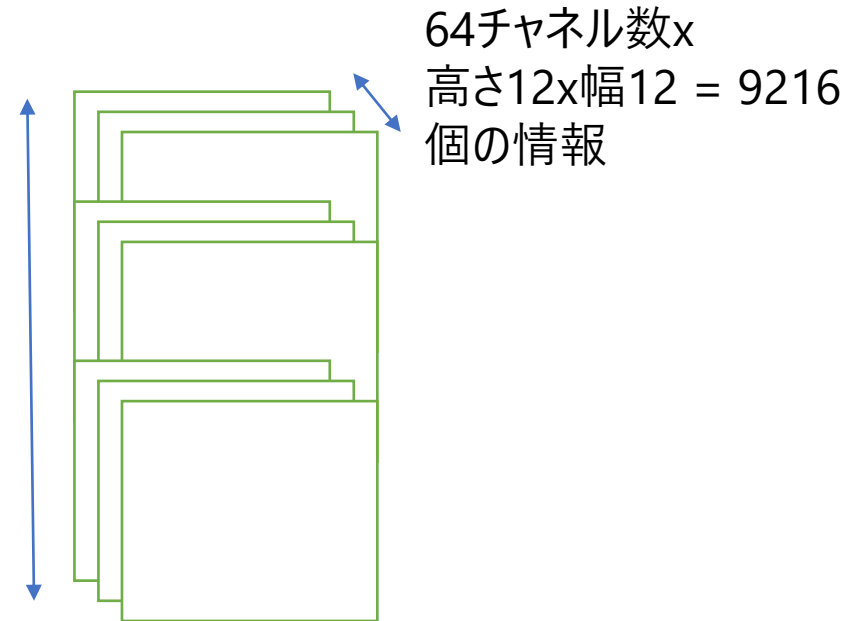
入力テンソル：ここでは、  
sizeは (ミニバッチ数、チャンネル数、高さ、幅)

平板化する最初の次元 1 -> size が  
(ミニバッチ数、平板化されたデータ) となる。

tensor([[...],[...],..., [...]])

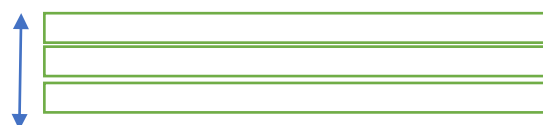
ミニバッチの数だけ

ミニバッチ数分



9216要素のベクトル

ミニバッチ数分

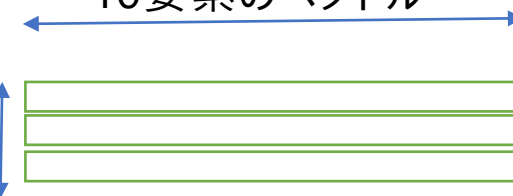


F.softmax(x, dim=1)

入力

ミニバッチ数分

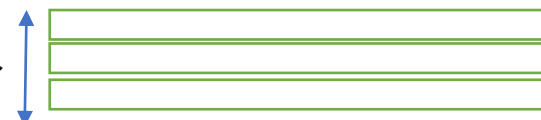
10要素のベクトル



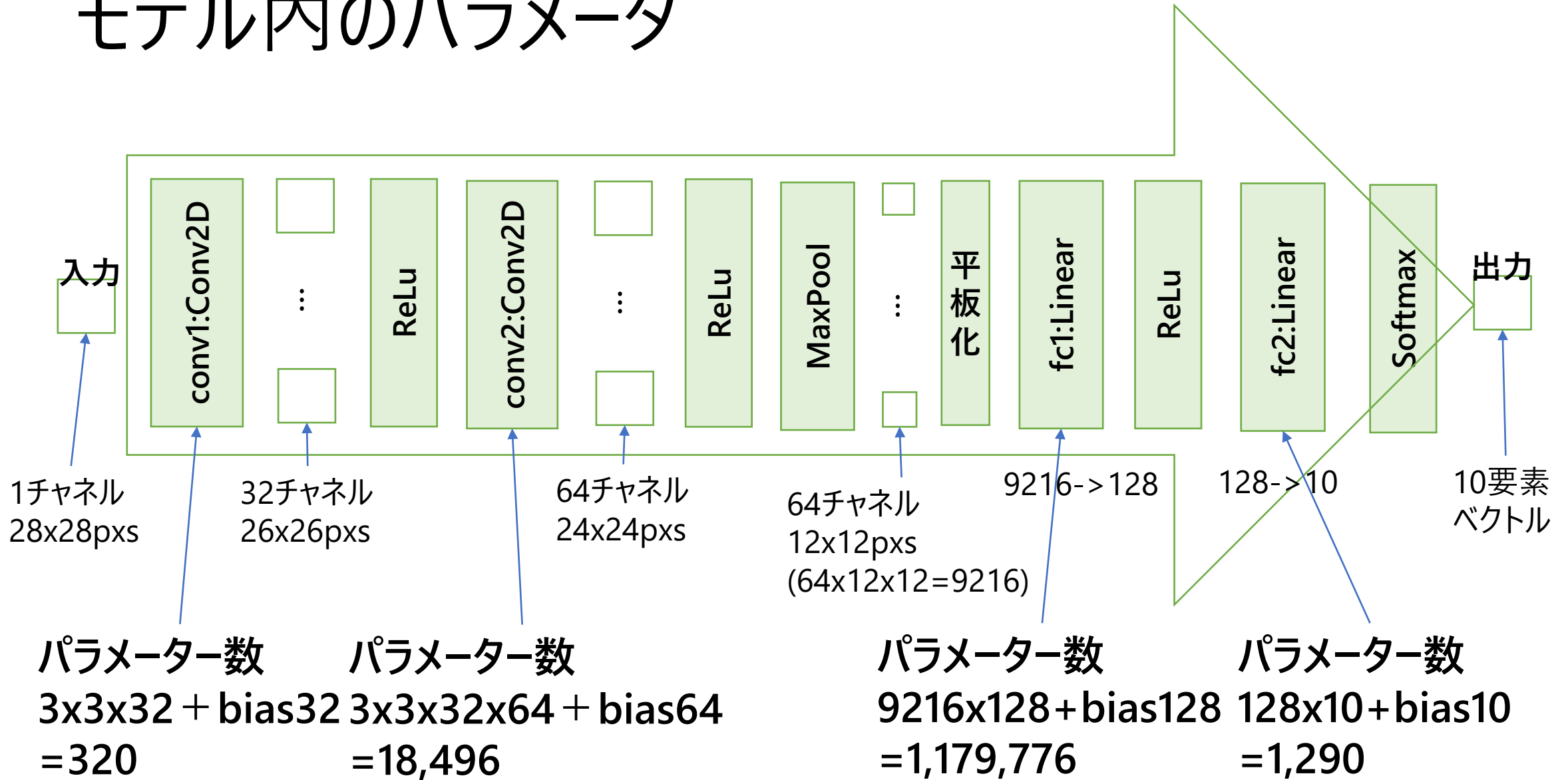
softmaxをとる次元：0がミニバッチ、1が10要素のベクトル

Softmax適用

ミニバッチ数分



# モデル内のパラメータ



main

# データセット、データローダー

- データセット
  - 入力データを、Pytorch Tensor形式に変換したり、値を正規化したり、など変換してから、送る。
- データローダー
  - ミニバッチ単位にまとめて、モデルに送る。



```
optim.SGD(model.parameters(), lr=0.1)
```

モデルのパラメータ

学習レート

SGD (Stochastic Gradient Decent) 確率的勾配降下法。ほかにも。

勾配降下法なのだが、入力データはシャッフルしてあってランダムなので、確率的。

SGD以外にいろんなOptimizerが準備されている。

StepLR(optimizer, step\_size=1, gamma=0.7)

エポック何個で学習具合を低減させるか

学習レートの逓減率

エポックが変わったときに、段階的に、一律に学習レートを低減させる。  
ほかに低減させるやり方のOptionが何個か準備されている。

train, test

# ミニバッチ

$$(x_1, x_2, x_3) \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{pmatrix} = (x_1w_{11} + x_2w_{21} + x_3w_{31}, x_1w_{12} + x_2w_{22} + x_3w_{32})$$



データを何個分かまとめる

ミニバッチ  
3個

$$\left\{ \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{pmatrix} \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{pmatrix} = \begin{pmatrix} x_{11}w_{11} + x_{12}w_{21} + x_{13}w_{31} & x_{11}w_{12} + x_{12}w_{22} + x_{13}w_{32} \\ x_{21}w_{11} + x_{22}w_{21} + x_{23}w_{31} & x_{21}w_{12} + x_{22}w_{22} + x_{23}w_{32} \\ x_{31}w_{11} + x_{32}w_{21} + x_{33}w_{31} & x_{31}w_{12} + x_{32}w_{22} + x_{33}w_{32} \end{pmatrix} \right.$$

GPUの能力活用

*outputs = model(inputs)* が、forward呼び出しに化ける訳

Pytorchは、Pythonの呼び出し可能オブジェクトの機能を多用している

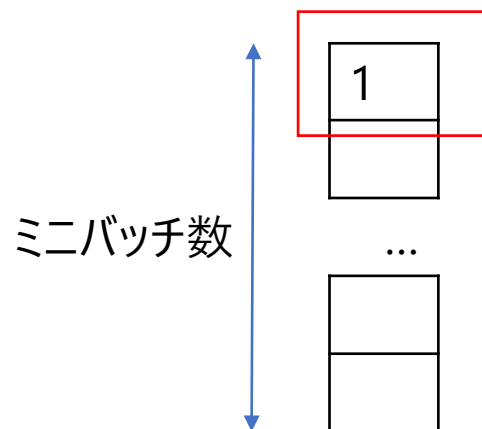
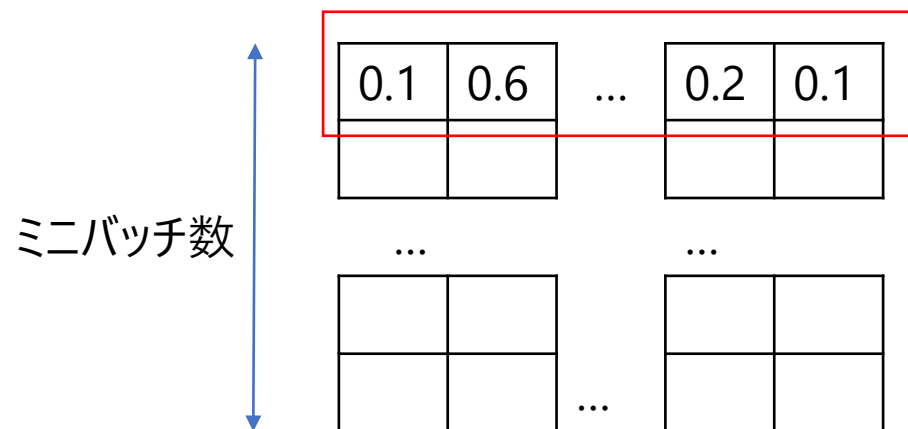
- object. call (self[, args...])
  - Called when the instance is "called" as a function; if this method is defined, x(arg1, arg2, ...) roughly translates to type(x). call (x, arg1, ...).
  - オブジェクト() と書くと、クラス.\_\_call\_\_(インスタンス、引数...) に化ける
- Pytorchのコード
  - moduleクラスの call メソッドは、中で forward 呼び出しを行っている。

# F.cross\_entropy(output, target)

検査するデータ

正解ラベル 1 が立っているインデックス 0,...,9

10要素、各数字である確率



赤の場合、出力は、 $-\log(0.6)$

発展：softmaxの結果0..1の値でロスを計算するより、0..1の各値の対数として $-\infty$ ..0にしてから、ロス計算に回したほうが、分解能が高い。

- Net
  - forward
    - ...
    - `output = F.log_softmax(x, dim=1)`
    - ...
- Train/Test
  - ...
  - `loss = F.nll_loss(output, target)`
  - ...

発展：計算結果が毎回変わる要素を減らすため、乱数発生器の初期値を固定してもいい。

- main
  - ...
  - `torch.manual_seed(args.seed)`
  - ...



# 発展：過学習しないような対策

- Net

- `__init__`

- ...

- `self.dropout1 = nn.Dropout2d(0.25)`

- `self.dropout2 = nn.Dropout2d(0.5)`

- ...

- `forward`

- ...

- `x = self.dropout1( F.max_pool2d( F.relu(self.conv2(x)) )`

- ...

- `x = self.dropout2( F.relu(self.fc1(x)) )`

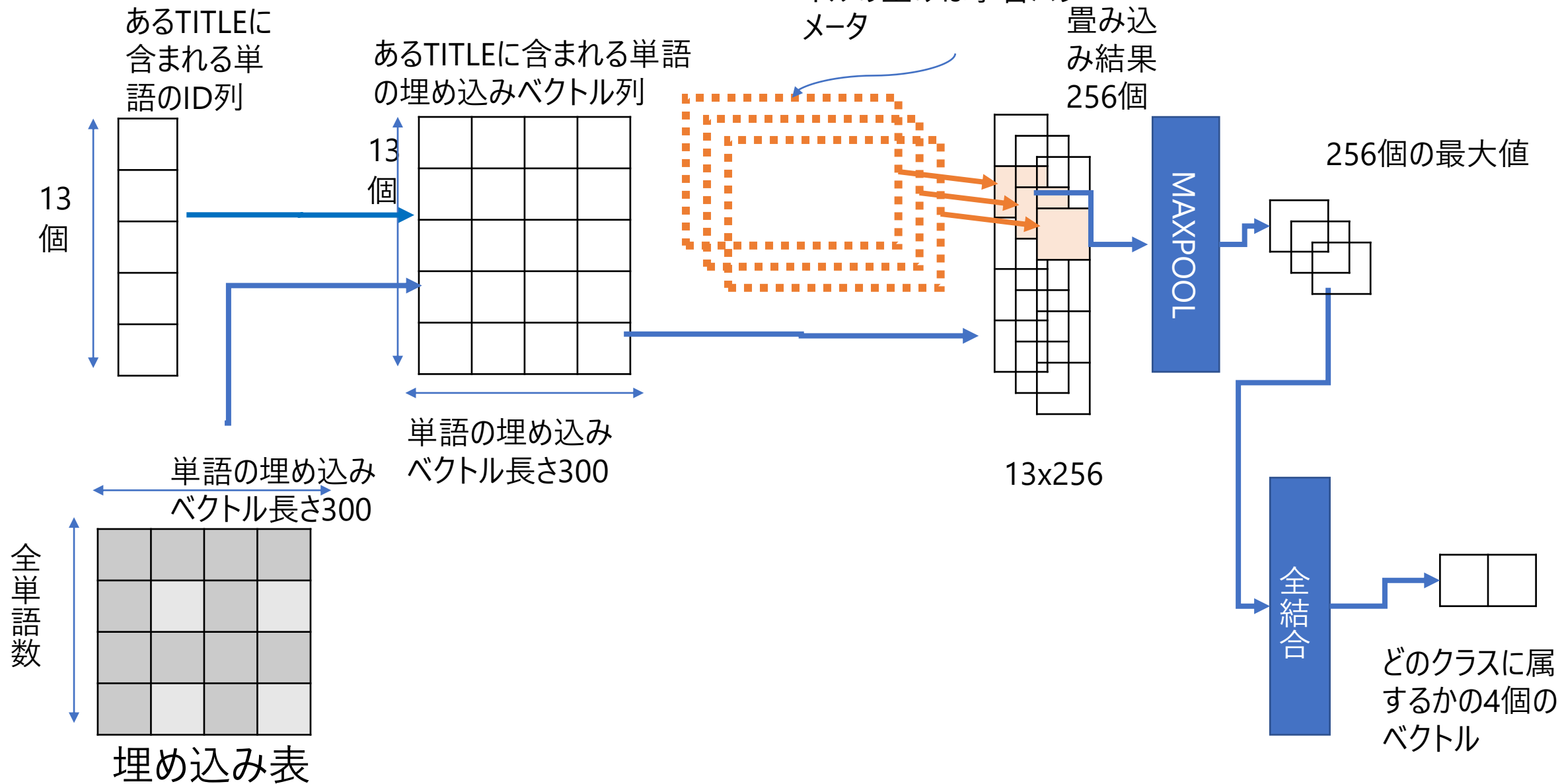
- モデルにランダムに欠落ノードを作ることによって、常時接続した単一のモデルにはない汎用性が仕込まれるらしい。経験的な知見でしかない。

# 100本ノック第9章課題80,86

- [「100本ノック」の9章の課題](#) の80(データ準備と86/87がCNNに関する課題)です。
- 「NLP、CNNRNNTransformer.ipynb」というノートをコピーして下さい。80、86回答例コードを完成し、実行ログを残して下さい。

# 86 (ミニバッチサイズ = 1)

3x埋め込みベクトルサイズのカーネル256個、カーネルの重みは学習パラメータ



# 確認クイズ

- 確認クイズをやってください。