

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532L): **Multimodal Learning with Vision, Language and Sound**

Lecture 4: Convolutional Neural Networks (Part 1)



Course Logistics

- Azure credits have been distributed
- Instructions for using Azure are on Piazza
- Office hours (Shikib) on Wednesday (January 17th) 1-2pm
- Assignment 1 grading
- Assignment 2 will be out today (on CNNs)

- Start thinking about papers

- Start thinking about **project** and forming groups (proposal is in ~month)

Research **Projects**



Robert Herjavec

(Canadian businessman and investor) Shark Tank

Research **Projects**

Business

entrepreneur will take the best idea and run it into a ground."

"A good entrepreneur can take a mediocre idea and make it great, a bad



Robert Herjavec (Canadian businessman and investor) Shark Tank

Research **Projects**

Business

entrepreneur will take the best idea and run it into a ground."

Research

best project and run it into a ground." — Me

"A good entrepreneur can take a mediocre idea and make it great, a bad



Robert Herjavec (Canadian businessman and investor) Shark Tank

"A good idea can make a mediocre project great, a bad idea will take the

I want to **solve vision** / language / etc.



Marvin Minsky

- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)

 - Requires **forward thinking**, knowledge of the field
 - A sure way to get tenure
 - Difficult to do as the field matures

- This is excellent if X is something no one has done or thought about and is important (guaranteed success)



Rosland Picard, MIT — Affective Computing



- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)
- I think the right way to solve (or improve) X is Y

I want to **solve vision** / language / etc.

I want to do X (e.g., **image captioning**)

I think the **right way to solve** (or improve) X is Y

— More incremental; a lot of science is incremental ("standing on the shoulders of giants")

- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)
- I think the **right way to solve** (or improve) X is Y
 - More incremental; a lot of science is incremental ("standing on the shoulders of giants")
 - **Retrospective:** compare existing approaches see why they work what is missing (guaranteed success)
 - Perspective: come up with an idea or the insight that you truly believe and test it





- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)
- I think the **right way to solve** (or improve) X is Y
 - More incremental; a lot of science is incremental ("standing on the shoulders of giants")
 - **Retrospective:** compare existing approaches see why they work what is missing (guaranteed success)
 - **Perspective:** come up with an idea or the insight that you truly believe and test it
 - Requires through knowledge of the sub-field (lots of reading)





- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)
- think the right way to solve (or improve) X is Y
 - More incremental; a lot of science is incremental ("standing on the shoulders of giants")
 - **Retrospective:** compare existing approaches see why they work what is missing (guaranteed success)
 - **Perspective:** come up with an idea or the insight that you truly believe and test it
 - Requires through knowledge of the sub-field (lots of reading)
 - Requires strong intuition and high level (intuitive) thinking



Geoffrey Hinton







- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)
- think the **right way to solve** (or improve) X is Y
 - More incremental; a lot of science is incremental ("standing on the shoulders of giants")
 - **Retrospective:** compare existing approaches see why they work what is missing (guaranteed success)
 - **Perspective:** come up with an idea or the insight that you truly believe and test it
 - Requires through knowledge of the sub-field (lots of reading)
 - Requires strong intuition and high level (intuitive) thinking
 - Requires understanding of the mathematical tools and formulations to know what maybe possible



Geoffrey Hinton







- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)

think the **right way to solve** (or improve) X is Y

- More incremental; a lot of science is incremental ("standing on the shoulders of giants")
- **Retrospective:** compare existing approaches see why they work what is missing (guaranteed success)
- **Perspective:** come up with an idea or the insight that you truly believe and test it
- Requires through knowledge of the sub-field (lots of reading)
- Requires strong intuition and high level (intuitive) thinking
- Requires understanding of the mathematical tools and formulations to know what maybe possible
- Helps to bringing knowledge from other fields (field cross pollination)



Geoffrey Hinton







- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)
- I think the **right way to solve** (or improve) X is Y
- Mathematical formulation
- Implementation / engineering
- **Experimental** testing

On to todays lecture ...

Fully Connected Layer



Example: 200 x 200 image (small) x 40K hidden units

Fully Connected Layer



Example: 200 x 200 image (small) x 40K hidden units

= ~ 2 Billion parameters (for one layer!)



Fully Connected Layer



Example: 200 x 200 image (small) x 40K hidden units

= ~ 2 Billion parameters (for one layer!)

Spatial correlations are generally local

Waste of resources + we don't have enough data to train networks this large





Locally Connected Layer



Example: 200 x 200 image (small) x 40K hidden units

Filter size: 10 x 10

= ~ 4 Million parameters

Locally Connected Layer



Example: 200 x 200 image (small) x 40K hidden units

Filter size: 10 x 10

= ~ 4 Million parameters

Stationarity — statistics is similar at different locations



Example: 200 x 200 image (small) x 40K hidden units

Filter size: 10×10

= ~ 4 Million parameters

Share the same parameters across the locations (assuming input is stationary)

* slide adopted from Marc'Aurelio Renzato





Example: 200 x 200 image (small) x 40K hidden units

Filter size: 10×10

= ~ 4 Million parameters

= 100 parameters

Share the same parameters across the locations (assuming input is stationary)

* slide adopted from Marc'Aurelio Renzato










































$\begin{bmatrix} 0.11 & 0.11 & 0.11 \end{bmatrix}$ $\begin{bmatrix} 0.11 & 0.11 & 0.11 \\ 0.11 & 0.11 & 0.11 \end{bmatrix}$





Example: 200 x 200 image (small) x 40K hidden units

Filter size: 10 x 10

of filters: 20

Learn multiple filters



Example: 200 x 200 image (small) x 40K hidden units

Filter size: 10 x 10

of filters: 20

= 2000 parameters

Learn multiple filters

32 x 32 x 3 image (note the image preserves spatial structure)



3 depth

32 x 32 x 3 **image**





$5 \times 5 \times 3$ filter

Convolve the filter with the image (i.e., "slide over the image spatially, computing dot products")



Convolutional Layer 32 x 32 x 3 image





Filters always extend the full depth of the input volume

5 x 5 x 3 filter

Convolve the filter with the image (i.e., "slide over the image spatially, computing dot products"





32 x 32 x 3 **image**





1 number: the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

$$\mathbf{W}^T \mathbf{x} + b$$
, where $\mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$

32 x 32 x 3 **image**





1 number: the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

$$\mathbf{W}^T \mathbf{x} + b$$
, where $\mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$

How many **parameters** does the layer have?

32 x 32 x 3 **image**





1 number: the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

$$\mathbf{W}^T \mathbf{x} + b$$
, where $\mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$

How many **parameters** does the layer have? **76**

32 x 32 x 3 **image**





activation map

32 x 32 x 3 **image**





activation map





this results in the "new image" of size 28 x 28 x 6!



Convolutional Neural Network (ConvNet)







What filters do networks learn?



Layer 1





[Zeiler and Fergus, 2013]



What filters do networks learn?



[Zeiler and Fergus, 2013]



32 x 32 x 3 **image**





activation map



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

=> **5** x **5** output

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 2)

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 2)

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 2)

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 2)

=> **3 x 3 output**

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 3)

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 3)

7 height

Does not fit! **Cannot apply** 3 x 3 filter on 7 x 7 image with stride 3



N width



N x N input image (spatially) F x F filter

Output size: (N-F) / stride + 1

N height



N width



N x N input image (spatially) F x F filter

Output size: (N-F) / stride + 1

N height

Example: N = 7, F = 3

stride $1 = \frac{(7-3)}{1+1} = 5$ stride 2 = (7-3)/2 + 1 = 3stride 3 => (7-3)/3+1 = **2.33**



Convolutional Layer: Border padding

7 width



5 x 5 input image (spatially)3 x 3 filter(applied with stride 1)

pad with 1 pixel border

Output size: 7 x 7

7 height

Convolutional Layer: Border padding

7 width



5 x 5 input image (spatially) 3 x 3 filter (applied with **stride 3**)

pad with 1 pixel border

7 height
Convolutional Layer: Border padding

7 width



5 x 5 input image (spatially) 3 x 3 filter (applied with **stride 3**)

pad with 1 pixel border

7 height

Example: N = 7, F = 3

stride 1 => (9-3)/1+1 = 7stride 2 => (9-3)/2+1 = 4stride 3 => (9-3)/3+1 = 3



3 depth



28 height





28 width







With padding we can achieve no shrinking (32 -> 28 -> 24); shrinking quickly (which happens with larger filters) doesn't work well in practice





Convolutional Layer: 1x1 convolutions

56 x 56 x 64 **image**



Convolutional Neural Networks



VGG-16 Network



CNNs: Reminder Fully Connected Layers

Input

3072

(32 x 32 x 3 image -> stretches to 3072 x 1)





Convolutional Neural Networks



VGG-16 Network

Convolutional Neural Networks



VGG-16 Network

CNNs: Reminder Fully Connected Layers

Input

25,088

(7 x 7 x 512 image -> stretches to 25,088 x 1)



each neuron looks at the full input volume

Activation

4,096



CNNs: Reminder Fully Connected Layers



102,760,448 parameters!

* adopted from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

Activation

4,096



Convolutional Neural Networks



VGG-16 Network

Convolutional Neural Networks



VGG-16 Network



Let us assume the filter is an "eye" detector

How can we make detection spatially invariant (insensitive to position of the eye in the image)

* slide from Marc'Aurelio Renzato



Let us assume the filter is an "eye" detector

How can we make detection spatially invariant (insensitive to position of the eye in the image)

> By "pooling" (e.g., taking a max) response over a spatial locations we gain robustness to position variations



* slide from Marc'Aurelio Renzato

- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



e manageable and spatially invariant independently



- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



e manageable and spatially invariant independently



* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

How many **parameters**?

- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



e manageable and spatially invariant independently



Max **Pooling**

activation map





max pool with 2 x 2 filter and stride of 2

6 8 3 4

Average **Pooling**

activation map





avg pool with 2 x 2 filter and stride of 2

3.25 5.25 2 2

Pooling Layer Receptive Field

If convolutional filters have size KxK and stride 1, and pooling layer has pools of size PxP, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size: (P+K-1)x(P+K-1)



* slide from Marc'Aurelio Renzato

Pooling Layer Receptive Field

If convolutional filters have size KxK and stride 1, and pooling layer has pools of size PxP, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size: (P+K-1)x(P+K-1)



* slide from Marc'Aurelio Renzato

Convolutional Neural Networks



VGG-16 Network

Local Contrast Normalization Layer

ensures response is the same in both case (details omitted, no longer popular)





* images from Marc'Aurelio Renzato

Improving Single Model

Regularization

- L2, L1
- Dropout / Inverted Dropout
- Data augmentation

Improving Single Model

Regularization

- L2, L1
- Dropout / Inverted Dropout
- Data augmentation

L2 Regularization: Learn a more (dense) distributed representation

$R(\mathbf{W}) = ||\mathbf{W}|$

L1 Regularization: Learn a sparse representation (few non-zero wight elements)

 $R(\mathbf{W}) = ||\mathbf{W}|$

$$\|_2 = \sum_{i} \sum_{j} \mathbf{W}_{i,j}^2$$

$$\|_1 = \sum_i \sum_j |\mathbf{W}_{i,j}|$$

Improving Single Model

Regularization

- L2, L1
- Dropout / Inverted Dropout
- Data augmentation

L2 Regularization: Learn a more (dense) distributed representation

$R(\mathbf{W}) = ||\mathbf{W}|$

L1 Regularization: Learn a sparse representation

 $R(\mathbf{W}) = ||\mathbf{W}|$



Dropout

$$||_{2} = \sum_{i} \sum_{j} \mathbf{W}_{i,j}^{2}$$

n (few non-zero wight elements)

$$\|_1 = \sum_i \sum_j |\mathbf{W}_{i,j}|$$







Horizontal flips

Random crops & scales

Color Jitter

Horizontal flips

Random crops & scales





Color Jitter

Horizontal flips

Training: sample random crops and scales e.g., ResNet:

- 1. Pick random L in range [256, 480]
- 2. Resize training image, short size = L
- 3. Sample random 224x224 patch

Testing: average a fix set of crops e.g., ResNet:

Resize image to 5 scales (224, 256, 384, 480, 640) 2. For each image use 10 224x224 crops: 4 corners + center, + flips

Random crops & scales

Color Jitter



Horizontal flips

Random perturbations in contrast and brightness



Random crops & scales

Color Jitter


Regularization: Stochastic Depth

Effectively "dropout" but for layers

some layer (for each batch)



Huang et al., ECCV 2016]



Common "Wisdom": You need a lot of data to train a CNN



Common "Wisdom": You need a lot of data to train a CNN







Solution: Transfer learning — taking a model trained on the task that has lots of data and adopting it to the task that may not

Common "Wisdom": You need a lot of data to train a CNN



This strategy is PERVASIVE.





Solution: Transfer learning — taking a model trained on the task that has lots of data and adopting it to the task that may not

Train on ImageNet

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Imago

[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Train on ImageNet

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Why on **ImageNet**?

[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Train on **ImageNet**

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Why on **ImageNet**?

- Convenience, lots of **data**



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

- We know how to train these well

Train on **ImageNet**

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

- Why on **ImageNet**?
 - Convenience, lots of data
 - We know how to train these well



[Yosinski et al., NIPS 2014] Donahue et al., ICML 2014 Razavian et al., CVPR Workshop 2014

However, for some tasks we would need to start with something else (e.g., videos for optical flow)

Train on ImageNet

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Train on ImageNet

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image





[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Train on ImageNet

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Train on **ImageNet**

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Re-initialize and train



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Train on **ImageNet**

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Re-initialize and train



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Small dataset with C classes

Lower levels of the CNN are at task independent anyways

Train on **ImageNet**

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Re-initialize and train



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Larger dataset

Train on **ImageNet**

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Re-initialize and train



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Larger dataset

FC-1000 FC-1000 FC-4096 FC-4096 FC-4096 FC-4096 MaxPool MaxPool Conv-512 Conv-512 Conv-512 Conv-512 MaxPool MaxPool Conv-512 Conv-512 Conv-512 Conv-512 MaxPool MaxPool Freeze Conv-256 Conv-256 Conv-256 these Conv-256 MaxPool MaxPool layers Conv-128 Conv-128 Conv-128 Conv-128 MaxPool MaxPool Conv-64 Conv-64 Conv-64 Conv-64 Image Image

Train on **ImageNet**

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Re-initialize and train



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Larger dataset



* adopted from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

Freeze these layers

Train on **ImageNet**

Small dataset with C classes

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Re-initialize and train



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

Larger dataset



* adopted from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

Freeze these layers



[[]Yosinski et al., NIPS 2014]

Training: Train multiple independent models **Test:** Average their results

Training: Train multiple independent models **Test:** Average their results

~ 2% improved performance in practice

Training: Train multiple independent models **Test:** Average their results

~ 2% improved performance in practice

Alternative: Multiple snapshots of the single model during training!

Training: Train multiple independent models **Test:** Average their results

~ 2% improved performance in practice

Alternative: Multiple snapshots of the single model during training!

- **Improvement:** Instead of using the actual parameter vector, keep a moving average of the parameter vector and use that at test time (Polyak averaging)

CPU vs. GPU (Why do we need Azure?)



Data from https://github.com/jcjohnson/cnn-benchmarks

Frameworks: Super quick overview

1. Easily **build computational graphs**

2. Easily **compute gradients** in computational graphs

3. Run it all efficiently on a GPU (weap cuDNN, cuBLAS, etc.)



Frameworks: Super quick overview

Core DNN Frameworks

Caffe (UC Berkeley)

Caffe 2 (Facebook)

(Baidu)

Torch (NYU/Facebook)

PyTorch (Facebook)

CNTK (Microsoft)

Theano (U Montreal)

TensorFlow (Google)

Puddle

MXNet (Amazon)

Frameworks: Super quick overview

Core DNN Frameworks

Caffe (UC Berkeley) Caffe 2 (Facebook)

(Baidu)

Torch (NYU/Facebook)

PyTorch (Facebook)

Theano (U Montreal) **TensorFlow** (Google)

Puddle

CNTK (Microsoft)

MXNet (Amazon)

Wrapper Libraries

Keras TFLearn TensorLayer tf.layers **TF-Slim** tf.contrib.learn Pretty Tensor

Frameworks: PyTorch vs. TensorFlow

Dynamic vs. Static computational graphs

Frameworks: PyTorch vs. TensorFlow

Dynamic vs. **Static** computational graphs

With static graphs, framework can optimize the graph for you before it runs!



Frameworks: PyTorch vs. TensorFlow

Dynamic vs. Static computational graphs

Graph building and execution is intertwined. Graph can be different for every sample.

The cat ate a big rat

PyTorch: Three levels of abstraction

Tensor: Imperative ndarray, but runs on GPU

Variable: Node in a computational graph; stores data and gradients

Module: A neural network layer; may store state or learnable weights