

- Consider the output at a certain layer of a convolutional neural network
- After this layer we could apply a maxpooling layer

- Or a 1×1 convolution
- Or a 3×3 convolution



- Consider the output at a certain layer of a convolutional neural network
- After this layer we could apply a maxpooling layer

▲ロ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ● ● ● ● ● ●

- Or a 1×1 convolution
- Or a 3×3 convolution
- Or a 5×5 convolution



- Consider the output at a certain layer of a convolutional neural network
- After this layer we could apply a maxpooling layer
- Or a 1×1 convolution
- Or a 3×3 convolution
- Or a 5×5 convolution
- Question: Why choose between these options (convolution, maxpooling, filter sizes)?

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで



- Consider the output at a certain layer of a convolutional neural network
- After this layer we could apply a maxpooling layer
- Or a 1×1 convolution
- Or a 3×3 convolution
- $\bullet~{\rm Or}$ a 5×5 convolution
- Question: Why choose between these options (convolution, maxpooling, filter sizes)?
- Idea: Why not apply all of them at the same time and then concatenate the feature maps?

・ロト ・周ト ・ヨト ・ヨト ・ヨー



• Well this naive idea could result in a large number of computations

◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ ─臣 ─の�?

Mitesh M. Khapra CS7015 (Deep Learning) : Lecture 11



- Well this naive idea could result in a large number of computations
- If P = 0 & S = 1 then convolving a $W \times H \times D$ input with a $F \times F \times D$ filter results in a (W - F + 1)(H - F + 1) sized output

▲ロ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ● ● ● ● ● ●



- Well this naive idea could result in a large number of computations
- If P = 0 & S = 1 then convolving a $W \times H \times D$ input with a $F \times F \times D$ filter results in a (W - F + 1)(H - F + 1) sized output
- Each element of the output requires $O(F \times F \times D)$ computations



- Well this naive idea could result in a large number of computations
- If P = 0 & S = 1 then convolving a $W \times H \times D$ input with a $F \times F \times D$ filter results in a (W - F + 1)(H - F + 1) sized output
- Each element of the output requires $O(F \times F \times D)$ computations
- Can we reduce the number of computations?

• Yes, by using 1×1 convolutions

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ / 圖 / のへで

58/68

Mitesh M. Khapra CS7015 (Deep Learning) : Lecture 11



- Yes, by using 1×1 convolutions
- Huh?? What does a 1 × 1 convolution do ?

▲□▶ ▲御▶ ▲臣▶ ▲臣▶ ―臣 …のへで

Mitesh M. Khapra CS7015 (Deep Learning) : Lecture 11



- Yes, by using 1×1 convolutions
- Huh?? What does a 1 × 1 convolution do ?

• It aggregates along the depth



- Yes, by using 1×1 convolutions
- Huh?? What does a 1 × 1 convolution do ?
- It aggregates along the depth
- So convolving a $D \times W \times H$ input with $D_1 \ 1 \times 1 \ (D_1 < D)$ filters will result in a $D_1 \times W \times H$ output (S = 1, P = 0)



- Yes, by using 1×1 convolutions
- Huh?? What does a 1 × 1 convolution do ?
- It aggregates along the depth
- So convolving a $D \times W \times H$ input with $D_1 \ 1 \times 1 \ (D_1 < D)$ filters will result in a $D_1 \times W \times H$ output (S = 1, P = 0)
- If $D_1 < D$ then this effectively reduces the dimension of the input and hence the computations



- $\bullet\,$ Yes, by using 1×1 convolutions
- Huh?? What does a 1 × 1 convolution do ?
- It aggregates along the depth
- So convolving a $D \times W \times H$ input with $D_1 \ 1 \times 1 \ (D_1 < D)$ filters will result in a $D_1 \times W \times H$ output (S = 1, P = 0)
- If $D_1 < D$ then this effectively reduces the dimension of the input and hence the computations
- Specifically instead of $O(F \times F \times D)$ we will need $O(F \times F \times D_1)$ computations

지원에 지원에 지원에 있는 것



- Yes, by using 1×1 convolutions
- Huh?? What does a 1 × 1 convolution do ?
- It aggregates along the depth
- So convolving a D×W×H input with D₁ 1×1 (D₁ < D) filters will result in a D₁×W×H output (S = 1, P = 0)
- If $D_1 < D$ then this effectively reduces the dimension of the input and hence the computations
- Specifically instead of $O(F \times F \times D)$ we will need $O(F \times F \times D_1)$ computations
- We could then apply subsequent 3×3 , 5×5 filter on this reduced output

▲ロト ▲母ト ▲ヨト ▲ヨト ヨー のくぐ



 But we might want to use different dimensionality reductions before the 3 × 3 and 5 × 5 filters

イロト (個) (日) (日) ヨー つくで

Mitesh M. Khapra CS7015 (Deep Learning) : Lecture 11



- But we might want to use different dimensionality reductions before the 3 × 3 and 5 × 5 filters
- So we can use D_1 and D_2 1 × 1 filters before the 3 × 3 and 5 × 5 filters respectively

Mitesh M. Khapra CS7015 (Deep Learning) : Lecture 11



- But we might want to use different dimensionality reductions before the 3 × 3 and 5 × 5 filters
- So we can use D_1 and D_2 1 × 1 filters before the 3 × 3 and 5 × 5 filters respectively
- We can then add the maxpooling layer followed by dimensionality reduction



- But we might want to use different dimensionality reductions before the 3 × 3 and 5 × 5 filters
- So we can use D_1 and D_2 1 × 1 filters before the 3 × 3 and 5 × 5 filters respectively
- We can then add the maxpooling layer followed by dimensionality reduction
- $\bullet\,$ And a new set of 1×1 convolutions



- But we might want to use different dimensionality reductions before the 3 × 3 and 5 × 5 filters
- So we can use D_1 and D_2 1 × 1 filters before the 3 × 3 and 5 × 5 filters respectively
- We can then add the maxpooling layer followed by dimensionality reduction
- $\bullet\,$ And a new set of 1×1 convolutions
- And finally we concatenate all these layers



- But we might want to use different dimensionality reductions before the 3 × 3 and 5 × 5 filters
- So we can use D_1 and D_2 1 × 1 filters before the 3 × 3 and 5 × 5 filters respectively
- We can then add the maxpooling layer followed by dimensionality reduction
- $\bullet\,$ And a new set of 1×1 convolutions
- And finally we concatenate all these layers
- This is called the **Inception module**



- But we might want to use different dimensionality reductions before the 3 × 3 and 5 × 5 filters
- So we can use D_1 and D_2 1 × 1 filters before the 3 × 3 and 5 × 5 filters respectively
- We can then add the maxpooling layer followed by dimensionality reduction
- $\bullet\,$ And a new set of 1×1 convolutions
- And finally we concatenate all these layers
- This is called the **Inception module**
- We will now see **GoogLeNet** which contains many such inception modules





◆□ → < 畳 → < Ξ → < Ξ → Ξ < つ Q ○ 60/68</p>











▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで





▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで





▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで







< □ > < ■ > < ■ > < ■ > < ■ > < ■ > < ■ > < ■ > ○ < ⊙ 60/68



▲□ → ▲圖 → ▲ 差 → 差 - のへで 60/68




▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで



▲□ → ▲圖 → ▲ 差 → 差 - のへで 60/68









▲□▶ ▲□▶ ▲国▶ ▲国▶ - 国 - のへで

60/68

• **Important Trick:** Got rid of the fully connected layer

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ / 圖 / のへで

61/68



- Important Trick: Got rid of the fully connected layer
- Notice that output of the last layer is $7 \times 7 \times 1024$ dimensional

・ロト ・ 日 ・ モ ト ・ 日 ・ ・ つ へ つ ・



- Important Trick: Got rid of the fully connected layer
- Notice that output of the last layer is $7 \times 7 \times 1024$ dimensional
- What if we were to add a fully connected layer with 1000 nodes (for 1000 classes) on top of this



- Important Trick: Got rid of the fully connected layer
- Notice that output of the last layer is $7 \times 7 \times 1024$ dimensional
- What if we were to add a fully connected layer with 1000 nodes (for 1000 classes) on top of this
- We would have $7 \times 7 \times 1024 \times 1000 = 49M \ parameters$



- Important Trick: Got rid of the fully connected layer
- Notice that output of the last layer is $7 \times 7 \times 1024$ dimensional
- What if we were to add a fully connected layer with 1000 nodes (for 1000 classes) on top of this
- We would have $7 \times 7 \times 1024 \times 1000 = 49M \ parameters$
- Instead they use an average pooling of size 7×7 on each of the 1024 feature maps



- Important Trick: Got rid of the fully connected layer
- Notice that output of the last layer is $7 \times 7 \times 1024$ dimensional
- What if we were to add a fully connected layer with 1000 nodes (for 1000 classes) on top of this
- We would have $7 \times 7 \times 1024 \times 1000 = 49M \ parameters$
- Instead they use an average pooling of size 7×7 on each of the 1024 feature maps
- This results in a 1024 dimensional output



- Important Trick: Got rid of the fully connected layer
- Notice that output of the last layer is $7 \times 7 \times 1024$ dimensional
- What if we were to add a fully connected layer with 1000 nodes (for 1000 classes) on top of this
- We would have $7 \times 7 \times 1024 \times 1000 = 49M \ parameters$
- Instead they use an average pooling of size 7 × 7 on each of the 1024 feature maps
- This results in a 1024 dimensional output
- Significantly reduces the number of



• $12 \times$ less parameters than AlexNet

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

61/68



• $12 \times$ less parameters than AlexNet

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ ― 臣 … のへぐ

• $2 \times$ more computations

• GoogLeNet

• ResNet

Mitesh M. Khapra CS7015 (Deep Learning) : Lecture 11

▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへぐ

62/68



• Suppose we have been able to train a shallow neural network well



- Suppose we have been able to train a shallow neural network well
- Now suppose we construct a deeper network which has few more layers (in orange)



- Suppose we have been able to train a shallow neural network well
- Now suppose we construct a deeper network which has few more layers (in orange)
- Intuitively, if the shallow network works well then the deep network should also work well by simply learning to compute identity functions in the new layers



- Suppose we have been able to train a shallow neural network well
- Now suppose we construct a deeper network which has few more layers (in orange)
- Intuitively, if the shallow network works well then the deep network should also work well by simply learning to compute identity functions in the new layers
- Essentially, the solution space of a shallow neural network is a subset of the solution space of a deep neural network

イロト (雪) (日) (日) (日) (日)

• But in practice it is observed that this doesn't happen

イロト 不得下 イヨト イヨト





- But in practice it is observed that this doesn't happen
- Notice that the deep layers have a higher error rate on the test set

イロト イポト イヨト イヨト

э.



• Consider any two stacked layers in a CNN

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで

65/68



- Consider any two stacked layers in a CNN
- The two layers are essentially learning some function of the input

▲□▶ ▲御▶ ▲臣▶ ▲臣▶ ―臣 …のへで



- Consider any two stacked layers in a CNN
- The two layers are essentially learning some function of the input
- What if we enable it to learn only a residual function of the input?



• Why would this help?

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ / 圖 / のへで

66/68



- Why would this help?
- Remember our argument that a deeper version of a shallow network would do just fine by learning identity transformations in the new layers



- Why would this help?
- Remember our argument that a deeper version of a shallow network would do just fine by learning identity transformations in the new layers
- This identity connection from the input allows a ResNet to retain a copy of the input



- Why would this help?
- Remember our argument that a deeper version of a shallow network would do just fine by learning identity transformations in the new layers
- This identity connection from the input allows a ResNet to retain a copy of the input
- Using this idea they were able to train really deep networks

▲ロト ▲母ト ▲ヨト ▲ヨト ヨー のくぐ



 1^{st} place in all five main tracks

• ImageNet Classification: "Ultradeep" 152-layer nets

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

ResNet, 152 layers

LILL'LLLL

1^{st} place in all five main tracks

- ImageNet Classification: "Ultradeep" 152-layer nets
- ImageNet Detection: 16% better than the 2nd best system

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

ResNet, 152 layers

1^{st} place in all five main tracks

- ImageNet Classification: "Ultradeep" 152-layer nets
- ImageNet Detection: 16% better than the 2nd best system
- ImageNet Localization: 27% better than the 2nd best system

ResNet, 152 layers

ResNet, 152 layers

1^{st} place in all five main tracks

- ImageNet Classification: "Ultradeep" 152-layer nets
- ImageNet Detection: 16% better than the 2nd best system
- ImageNet Localization: 27% better than the 2nd best system
- COCO Detection: 11% better than the 2nd best system

ResNet, 152 layers

1^{st} place in all five main tracks

- ImageNet Classification: "Ultradeep" 152-layer nets
- ImageNet Detection: 16% better than the 2nd best system
- ImageNet Localization: 27% better than the 2nd best system
- COCO Detection: 11% better than the 2nd best system
- **COCO Segmentation:** 12% better than the 2nd best system



ResNet, 152 layers

Mitesh M. Khapra CS7015 (Deep Learning) : Lecture 11

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ○ ○ ○ ○



Bag of tricks

- Batch Normalizaton after every CONV layer
- Xavier/2 initialization from [He et al]

ResNet, 152 layers
UDUCUUUUU

Bag of tricks

- Batch Normalizaton after every CONV layer
- Xavier/2 initialization from [He et al]

・ロト ・ 日 ・ モ ト ・ 日 ・ ・ つ へ つ ・

• SGD + Momentum(0.9)

ResNet, 152 layers

Bag of tricks

- Batch Normalizaton after every CONV layer
- Xavier/2 initialization from [He et al]
- SGD + Momentum(0.9)
- Learning rate:0.1, divided by 10 when validation error plateaus

・ロト ・ 一下・ ・ ヨト

ResNet, 152 layers

ResNet, 152 layers

Bag of tricks

- Batch Normalizaton after every CONV layer
- Xavier/2 initialization from [He et al]
- SGD + Momentum(0.9)
- Learning rate:0.1, divided by 10 when validation error plateaus

• Mini-batch size 256

ResNet, 152 layers

Bag of tricks

- Batch Normalizaton after every CONV layer
- Xavier/2 initialization from [He et al]
- SGD + Momentum(0.9)
- Learning rate:0.1, divided by 10 when validation error plateaus

- Mini-batch size 256
- Weight decay of 1e-5

ResNet, 152 layers

Bag of tricks

- Batch Normalizaton after every CONV layer
- Xavier/2 initialization from [He et al]
- SGD + Momentum(0.9)
- Learning rate:0.1, divided by 10 when validation error plateaus

- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used