# Capsule Networks
## Dynamic Routing Between Capsules

Louis (Yiqing) Luo

July 6th, 2018

# Outline

# Outline

# Capsule Network - Intuition

- Drawbacks of CNN currently constructed
  1. Sub-sampling loses the precise spatial relationships between higher-level parts such as a nose and a mouth. The precise spatial relationships are needed for identity recognition
     - But overlapping the sub-sampling pools mitigates this.
  2. They cannot extrapolate their understanding of geometric relationships to radically new viewpoints.
  3. Discard (invariant) instead of disentangle (equivariant). Knowledge should be invariant, not activities inside networks.

# Capsule Network - Intuition

- Drawbacks of CNN currently constructed
  1. Sub-sampling loses the precise spatial relationships between higher-level parts such as a nose and a mouth. The precise spatial relationships are needed for identity recognition
     - But overlapping the sub-sampling pools mitigates this.
  2. They cannot extrapolate their understanding of geometric relationships to radically new viewpoints.
  3. Discard (invariant) instead of disentangle (equivariant). Knowledge should be invariant, not activities inside networks.



---

*Hinton: The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster.*
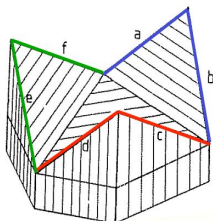
# Outline

# Better Representation of Images: Inverse Graphics

- Hinton claims that brain deconstruct hierarchy representation of visual information and tries to match patterns and relationships already learned.
    - representation of objects in brain does not depend on view angles.
- 3D graphics: relationships between 3D objects can be represented by poses
- Hinton: correct classification and recognition requires preservation of hierarchy pose relationship between objects

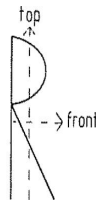# Better Representation of Images: Inverse Graphics

- Inverse Graphics is widely employed in computer vision.
- Uses hierarchical models in which spatial structure is modeled by matrices that represent the transformation from a coordinate frame embedded in the whole to a coordinate frame embedded in each part
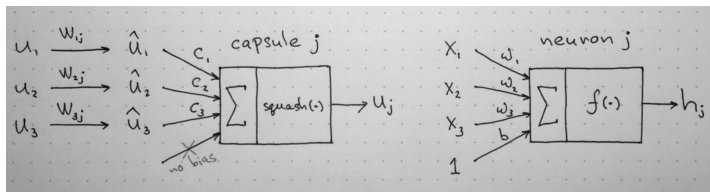- matrices are viewpoint invariant.



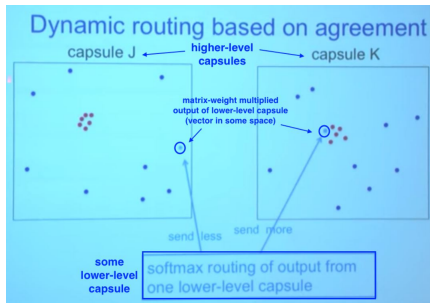(a)                              (b)

# Outline

# What are capsules

- each capsule will learn to recognize an implicitly defined visual entity over a limited domain of viewing conditions and deformations
- Outputs:
  1. probability that the entity is within its limited domain (length of the vector does note change)
  2. a set of instantiation parameters (eg. pose, lighting, possible deformations relative to implicitly defined canonical version of entities) (orientation of the vector changes)
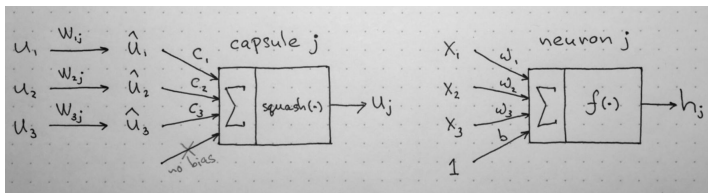
# Coincidence Filtering

- capsules receive multi-dimensional prediction vectors from capsules in layer below and look for tight predictions
- Outputs can also be described as:
  1. probability that the entity exist within its domain
  2. Center of gravity of the cluster = generalized pose
- good at filtering noise because coincidences rarely happens in high dimension space

# Outline

# Capsule vs Neurons



| Capsule vs. Traditional Neuron | | | |
|---|---|---|---|
| Input from low-level capsule/neuron | | vector($\mathbf{u}_i$) | scalar($x_i$) |
| Operation | Affine Transform | $\widehat{\mathbf{u}}_{j|i} = \mathbf{W}_{ij}\mathbf{u}_i$ | – |
| | Weighting | $\mathbf{s}_j = \sum_i c_{ij}\widehat{\mathbf{u}}_{j|i}$ | $a_j = \sum_i w_i x_i + b$ |
| | Sum | | |
| | Nonlinear Activation | $\mathbf{v}_j = \frac{\|\mathbf{s}_j\|^2}{1+\|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|}$ | $h_j = f(a_j)$ |
| Output | | vector($\mathbf{v}_j$) | scalar($h_j$) |

# Matrix Multiplication on Input Vectors

- inputs $u_1$, $u_2$, $u_3$ contains probability of features and their poses
- eg. for detecting a nose from a face:
  - $W_{ij}$ encodes relationship from $i_{th}$ feature (nose) to $j_{th}$ feature (face)

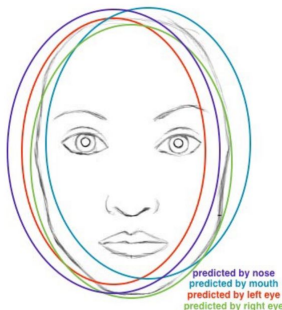$$\widehat{\mathbf{u}}_{j|i} = \mathbf{W}_{ij}\mathbf{u}_i$$

# Scalar Weighting

- measures the affinity between low level capsules and high level capsules
- $c_{ij}$ is adjusted each iteration by dynamic routing

$$\mathbf{s}_j = \sum_i c_{ij} \widehat{\mathbf{u}}_{j|i}$$

# Scalar Weighting - Dynamic Routing

- decides how to send output vector to higher level capsule j by changing $c_{ij}$
- important properties of $c_{ij}$
  1. $c_{ij}$ non-negative scalar
  2. $\Sigma_j c_{ij} = 1$ for all i
  3. # of weights for the vector $c_i$ = # of higher level capsules
  4. determined by iterative dynamic routing algorithm



predicted by nose
predicted by mouth
predicted by left eye
predicted by right eye

**Procedure 1** Routing algorithm.

1: **procedure** ROUTING($\hat{\mathbf{u}}_{j|i}, r, l$)
2:     for all capsule $i$ in layer $l$ and capsule $j$ in layer $(l+1)$: $b_{ij} \leftarrow 0$.
3:     **for** $r$ iterations **do**
4:         for all capsule $i$ in layer $l$: $\mathbf{c}_i \leftarrow \texttt{softmax}(\mathbf{b}_i)$          ▷ softmax computes Eq. 3
5:         for all capsule $j$ in layer $(l+1)$: $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}$
6:         for all capsule $j$ in layer $(l+1)$: $\mathbf{v}_j \leftarrow \texttt{squash}(\mathbf{s}_j)$          ▷ squash computes Eq. 1
7:         for all capsule $i$ in layer $l$ and capsule $j$ in layer $(l+1)$: $b_{ij} \leftarrow b_{ij} + \hat{\mathbf{u}}_{j|i}.\mathbf{v}_j$
    **return** $\mathbf{v}_j$

$$\mathbf{v}_j = \boxed{\frac{\|\mathbf{s}_j\|^2}{1+\|\mathbf{s}_j\|^2}} \boxed{\frac{\mathbf{s}_j}{\|\mathbf{s}_j\|}}$$

additional "squashing"     unit scaling

# Outline

# Properties of capsules

1. probability = invariant
   - does note change for possible appearances in manifold within limited domain covered by capsule

2. instantiation parameters = equivariant
   - change in viewing conditions, position etc. change correspondingly over the appearance manifold
     - face is centered around the nose
     - nose is 10 times smaller than the face
     - same orientation in space

# Outline

A Visual Representation of Capsule Connections in *Dynamic Routing Between Capsules*

# Loss Function



calculated for correct DigitCap                    calculated for incorrect DigitCaps

$$L_c = T_c \max(0, m^+ - ||\mathbf{v}_c||)^2 + \lambda(1 - T_c) \max(0, ||\mathbf{v}_c|| - m^-)^2$$

loss term for one DigitCap

L2 norm

0.5 constant used for numerical stability

L2 norm

1 when correct DigitCap, 0 when incorrect

zero loss when correct prediction with probability greater than 0.9, non-zero otherwise

1 when incorrect DigitCap, 0 when correct

zero loss when incorrect prediction with probability less than 0.1, non-zero otherwise

Note: correct DigitCap is one that matches training label, for each training example there will be 1 correct and 9 incorrect DigitCaps

# Outline

# Capsule Network Decoder Architecture



- reconstruction loss accounted for by sum of squares between original and predicted intensities

# Outline

- test errors on 28 x 28 MNIST with 60K and 10K training and testing datasets (with natural variance in skew, rotation, style, etc):

Figure 3: Sample MNIST test reconstructions of a CapsNet with 3 routing iterations. $(l, p, r)$ represents the label, the prediction and the reconstruction target respectively. The two rightmost columns show two reconstructions of a failure example and it explains how the model confuses a 5 and a 3 in this image. The other columns are from correct classifications and shows that model preserves many of the details while smoothing the noise.



Table 1: CapsNet classification test accuracy. The MNIST average and standard deviation results are reported from 3 trials.

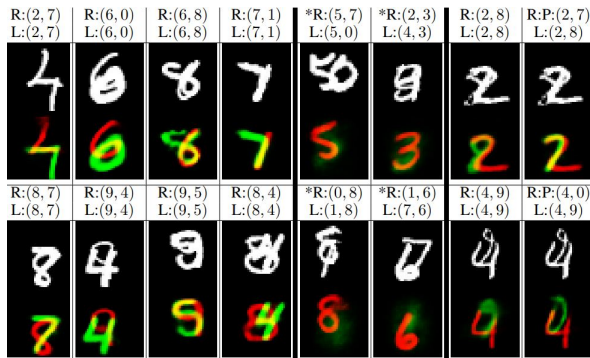| Method | Routing | Reconstruction | MNIST (%) | MultiMNIST (%) |
|---|---|---|---|---|
| Baseline | - | - | 0.39 | 8.1 |
| CapsNet | 1 | no | $0.34_{\pm 0.032}$ | - |
| CapsNet | 1 | yes | $0.29_{\pm 0.011}$ | 7.5 |
| CapsNet | 3 | no | $0.35_{\pm 0.036}$ | - |
| CapsNet | 3 | yes | $\mathbf{0.25}_{\pm 0.005}$ | **5.2** |

# Capsules in MNIST and Robustness to Affinity Transformations

- tested against traditional CNN with maxpooling and dropout
- both trained with MNIST digit placed randomly on a black background of 40 40 px
- tested on the affNIST
- traditional CNN: 99.23% accuracy on the expanded MNIST test set achieved 79% accuracy on the affNIST test set
- Capsule Network with similar parameters: similar accuracy (99.22%) on the expanded mnist test set only achieved 66% on the affNIST test set.

# Outline

# multiMNIST - Segmenting highly overlapping digits

- Dynamic routing can be viewed as a parallel attention mechanism that allows each capsule at one level to attend to some active capsules at the level below and to ignore others
- trained from scratch on the multiMNIST database (Lower L: predicted)

# multiMNIST - Results

- Network ability to reconstruct digits regardless of the overlap shows that each digit capsule can pick up the style and position from the votes it is receiving from PrimaryCapsules layer.

Figure 3: Sample MNIST test reconstructions of a CapsNet with 3 routing iterations. $(l, p, r)$ represents the label, the prediction and the reconstruction target respectively. The two rightmost columns show two reconstructions of a failure example and it explains how the model confuses a 5 and a 3 in this image. The other columns are from correct classifications and shows that model preserves many of the details while smoothing the noise.
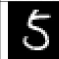


Table 1: CapsNet classification test accuracy. The MNIST average and standard deviation results are reported from 3 trials.

| Method | Routing | Reconstruction | MNIST (%) | MultiMNIST (%) |
|---|---|---|---|---|
| Baseline | - | - | 0.39 | 8.1 |
| CapsNet | 1 | no | $0.34_{\pm 0.032}$ | - |
| CapsNet | 1 | yes | $0.29_{\pm 0.011}$ | 7.5 |
| CapsNet | 3 | no | $0.35_{\pm 0.036}$ | - |
| CapsNet | 3 | yes | $\mathbf{0.25}_{\pm 0.005}$ | **5.2** |

# Outline

# CIFAR10

- CIFAR10 and achieved 10.6% error with an ensemble of 7 models each of which is trained with 3 routing iterations
  - standard CNN achieves 10.6% error
  - Each model has the same architecture as the simple model we used for MNIST except that there are three color channels and we used 64 different types of primary capsule
  - introduced a "none-of-the-above" category for the routing softmaxes, since we do not expect the final layer of ten capsules to explain everything in the image
- Drawback: tends to account for everything in the image so it does better when it can model the clutter than when it just uses an additional orphan category in the dynamic routing.
  - In CIFAR-10, the backgrounds are much too varied to model in a reasonable sized net which helps to account for the poorer performance.

# Thank you for Listening

# For Further Reading I

📕 Sara Sabour, Nicholas Frosst Geoffrey E. Hinton.
*Dynamic Routing Between Capsules.*
Google Brain, 2017.

📄 S. Someone.
On this and that.
*Journal of This and That*, 2(1):50–100, 2000.