

Learning with CNNs, Part I

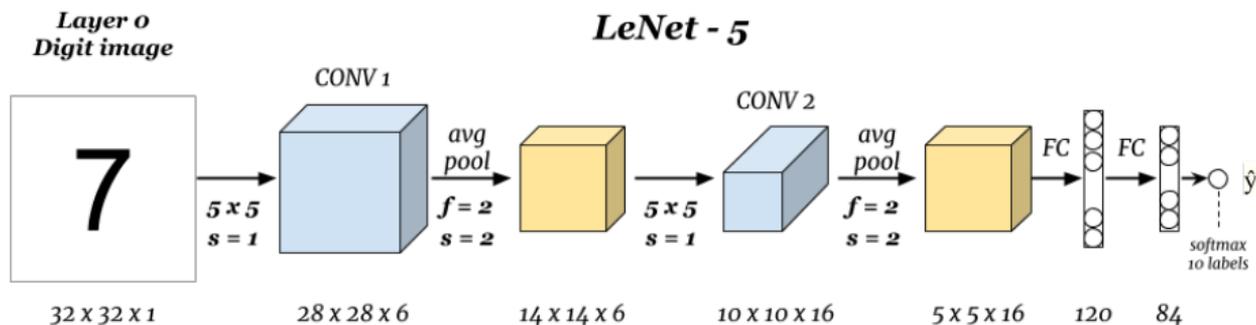
Tianxiang (Adam) Gao

School of Computing
DePaul University

Outline

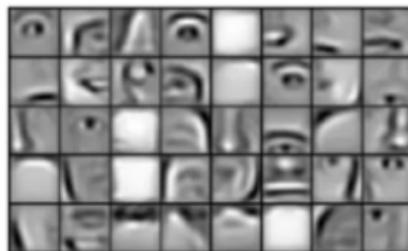
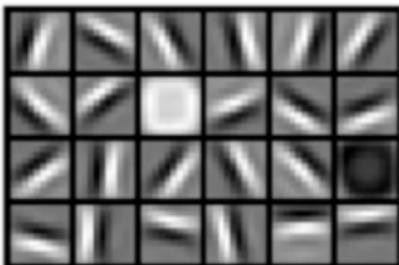
- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Semantic Segmentation
- 4 Face Recognition

Recap: Convolutional Neural Networks (CNNs)



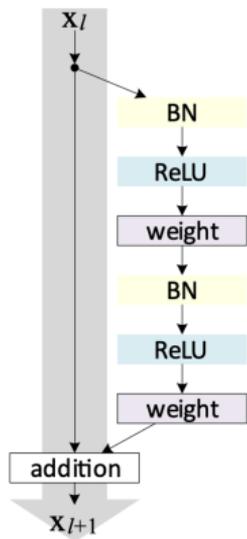
- **Convolution operation:** It slides a small **filter** over the input image, performing a **locally linear transformation** to produce a feature map that detects patterns.
- **Padding and stride:** Methods for controlling feature map size, preserving spatial dimensions, and improving *computational efficiency*.
- **Convolution over volumes and with multiple filters:** The input can be multi-channel, and so as the output with the use of multiple filters.
- **Weight Sharing and Sparsity:** Neurons in CNNs **share** weights across locations, with each output relying on a **small, localized** input region.

Recap: Hierarchical Feature Detection



- **Hierarchical Feature Detection:** Early layers capture basic features (like edges), which later layers combine into higher-level features.

Recap: Stabilizing CNN Training



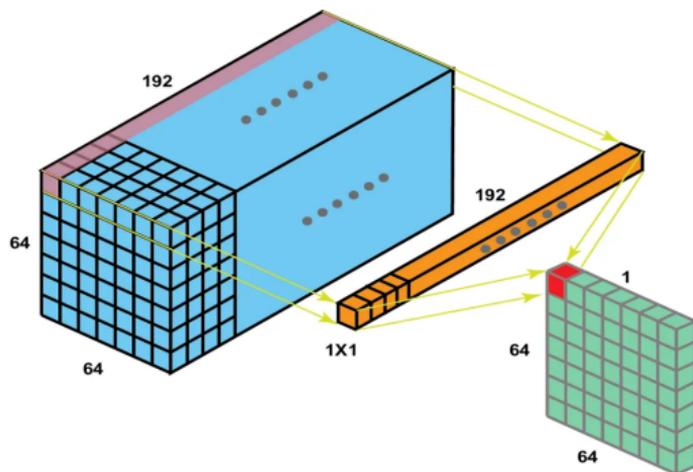
- **Pooling:** Reduces the spatial dimensions of feature maps by downsampling, typically using max or average pooling.
- **Batch Normalization:** Normalizes pre-activation values at hidden layers, mitigating internal covariate shift and accelerating training.
- **Skip Connections:** Create shortcuts by directly adding input to output, stabilizing information flow in deep neural networks.
- **Classic CNNs:** Spatial dimensions shrink while the number of channels increases as depth grows. Overparameterized and deeper CNNs are generally preferred.

Outline

- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Semantic Segmentation
- 4 Face Recognition

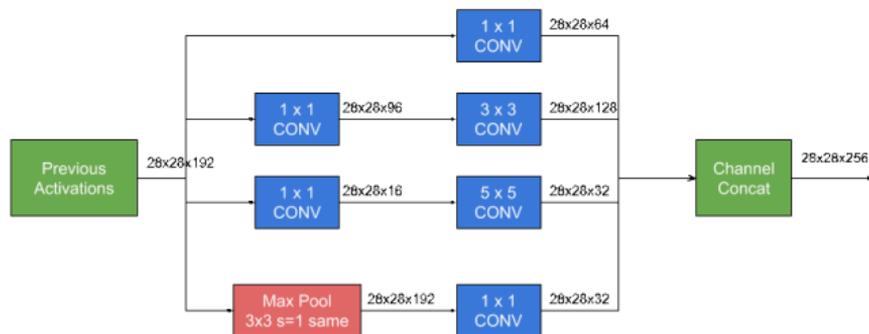
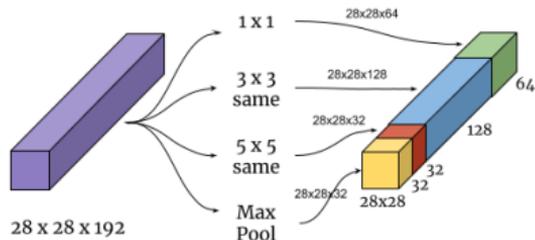
1 × 1 Convolutions

$$\underbrace{\begin{bmatrix} 3 & 0 & 1 \\ 1 & 5 & 8 \\ 2 & 7 & 2 \end{bmatrix}}_{\text{input image } 3 \times 3} * \underbrace{\begin{bmatrix} 2 \end{bmatrix}}_{\text{filter } 1 \times 1} = \underbrace{\begin{bmatrix} 6 & 0 & 2 \\ 2 & 10 & 16 \\ 4 & 14 & 4 \end{bmatrix}}_{\text{feature map } 3 \times 3}$$



- Recall that each channel in the input data serves as an **indicator map** for certain patterns detected by the filter.
- Each pixel, with **multiple channels**, represents a **set** of indicators that capture the distribution of **multiple patterns** within a local region of the original image.
- Applying a 1 × 1 convolution enables **cross-channel interaction** to detect complex patterns that integrate multiple underlying features.
- This approach is useful for combining lower-level features into **higher-level** representations.

Inception Networks



- Inception networks allow the model to freely select the best filter sizes within a layer, rather than manually choosing them, by providing multiple filter options.
- This approach can be computationally expensive (e.g., a 5x5 filter requires about 120 million operations).

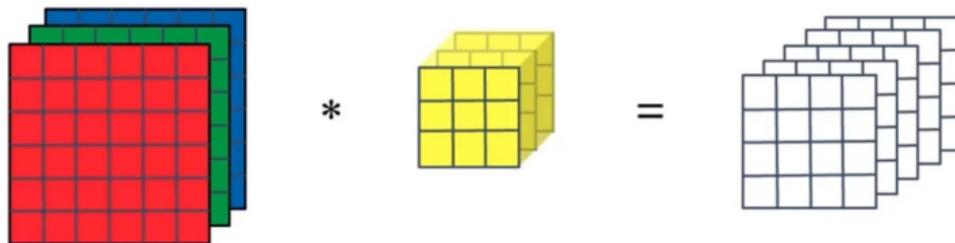
$$(28 \times 28 \times 32) \times (5 \times 5 \times 192) \approx 120 \text{ million}$$

- Applying 1 × 1 convolutions reduces computation by approximately 90%.

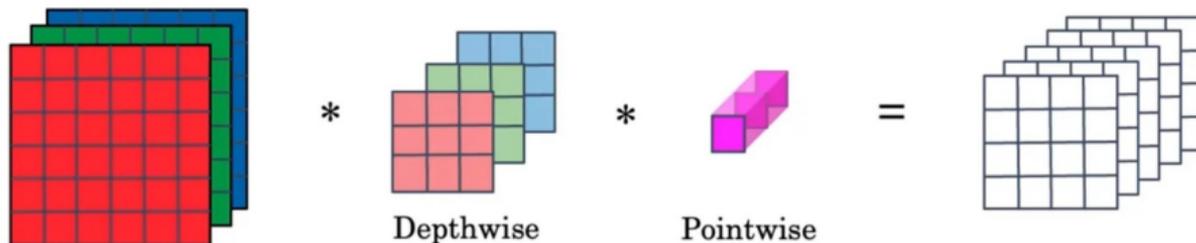
$$(28 \times 28 \times 16) \times (1 \times 1 \times 192) + (28 \times 28 \times 32) \times (5 \times 5 \times 16) \approx 2 \text{ million} + 10 \text{ million}$$

MobileNets

Normal Convolution



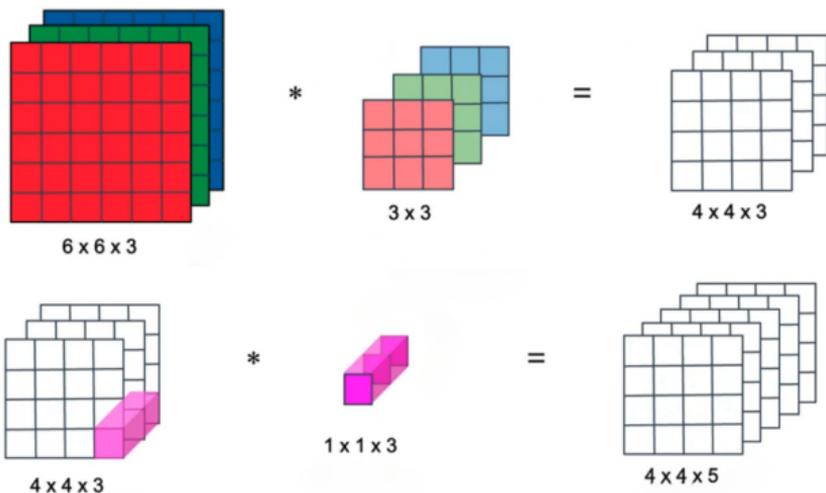
Depthwise Separable Convolution



The computational cost of a normal convolution is:

$$\underbrace{3 \times 3 \times 3}_{\# \text{ params}} \times \underbrace{4 \times 4}_{\# \text{ locations}} \times \underbrace{5}_{\# \text{ filters}} = \underbrace{2,160}_{\text{Computational cost}}$$

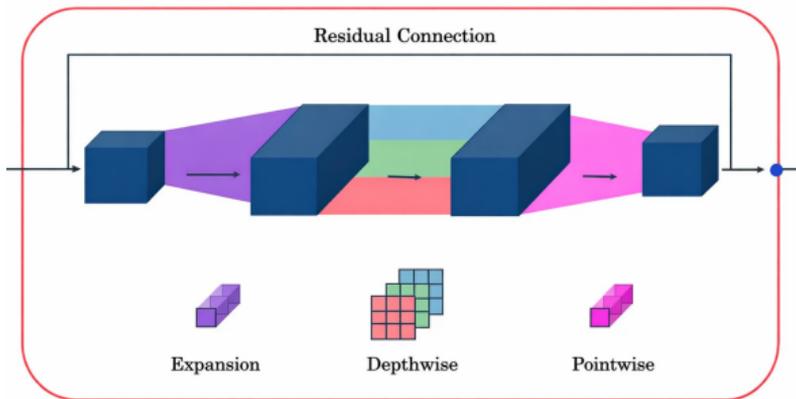
Depthwise Separable Convolution



- Each channel has an **independent** filter for the convolution operation.
- The result is then processed using a pointwise or 1×1 convolution.
- Computational cost ratio: $\frac{1}{C_{out}} + \frac{1}{f^2}$

$$672 = \underbrace{3 \times 3}_{\# \text{ params}} \times \underbrace{4 \times 4}_{\# \text{ locations}} \times \underbrace{3}_{\# \text{ filters}} + \underbrace{3}_{\# \text{ params}} \times \underbrace{4 \times 4}_{\# \text{ locations}} \times \underbrace{5}_{\# \text{ filters}}$$

MobileNetV2: Inverted Residual Blocks



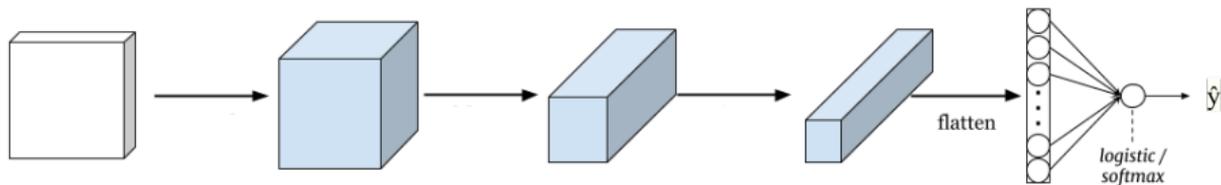
- The **expansion convolution** allows CNNs to learn richer functions by increasing the number of channels before depthwise convolution
- When deploying on resource-constrained devices, a **pointwise convolution** is used to project back to a smaller number of channels, reducing memory usage.
- VGG-16 \approx 528 MB, ResNet-152 \approx 230 MB, Inception \approx 85 MB, but MobileNet \approx **5-16 MB**.

Outline

- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Semantic Segmentation
- 4 Face Recognition

Transfer Learning

- **Definition:** Transfer learning leverages a **pre-trained** model on a large dataset to solve a new, related task on a smaller dataset.
- **Motivation:** Training deep networks from scratch demands significant data and computational resources. The initial layers in CNNs capture fundamental features that are transferable across various tasks, allowing the higher layers to focus on learning task-specific features.



- **Typical Workflow:**
 - ① **Pre-training:** A model is trained on a large dataset (e.g., ImageNet).
 - ② **Freezing Layers:** Early layers are frozen to retain their learned features.
 - ③ **Fine-tuning:** With a lower learning rate, higher layers are retrained on the target dataset, allowing for adaptation to new, task-specific features.
- Frozen early layers can be considered as **fixed** feature extractors, allowing activations to be precomputed, which reduces computation and speeds up training.
- For **larger** target datasets, freeze fewer layers to enable more flexibility in learning higher-level task-specific features during fine-tuning.

Data Augmentation

- **Definition:** Data augmentation involves generating new training samples by applying various transformations to existing data, helping improve model generalization.
- **Motivation:** Increases the diversity of the training dataset, which can reduce overfitting and improve the model's performance on unseen data.
- **Flips, rotations, and scaling:**



- **Random cropping:**



- **Color Adjustments:** Changes to colors or brightness



Mixup



- **Mixup** creates new training samples by blending pairs of images and their corresponding labels:

$$\tilde{\mathbf{x}} = (1 - \lambda)\mathbf{x}_1 + \lambda\mathbf{x}_2, \quad \text{and} \quad \tilde{\mathbf{y}} = (1 - \lambda)\mathbf{y}_1 + \lambda\mathbf{y}_2,$$

where $(\mathbf{x}_i, \mathbf{y}_i)$ are original training samples, and λ is the mixing factor drawn from a Beta distribution.

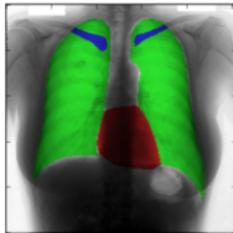
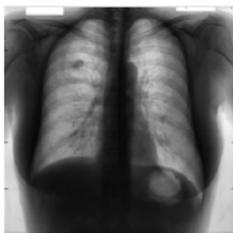
- This technique helps the model generalize by learning smoother transitions between classes.

Outline

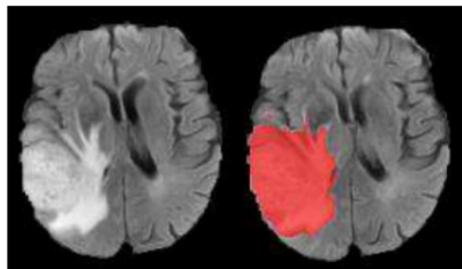
- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Semantic Segmentation**
- 4 Face Recognition

Semantic Segmentation with U-Net

Successful Applications:



Chest X-Ray



Brain MRI

Novikov, et al. "Fully convolutional architectures for multiclass segmentation in chest radiographs." IEEE Tran Med Img 2018
Dong, et al. "Automatic brain tumor detection and segmentation using U-Net based fully convolutional networks.", MIUA2017

Semantic Labeling

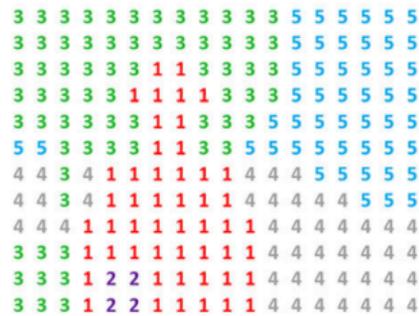
Per-Pixel Class Labeling:



Input



- 1: Person
- 2: Purse
- 3: Plants/Grass
- 4: Sidewalk
- 5: Building/Structures



Semantic Labels

- Assign a class label to every pixel in the image.
- Output is an image of the same dimensions as the input.

Deconvolution

Deconvolution: A deconvolution (or a **transpose convolution** or **up-sampling convolution**) is an operation that applies a filter to input data in a way that expands its spatial dimensions.

$$\underbrace{\begin{bmatrix} 3 & 0 \\ 1 & 5 \end{bmatrix}}_{\text{input } 2 \times 2} * \underbrace{\begin{bmatrix} 2 & 7 & 4 \\ 3 & 1 & 7 \\ 4 & 2 & 1 \end{bmatrix}}_{\text{filter } 3 \times 3} = \underbrace{\begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}}_{\text{feature map } 4 \times 4}$$

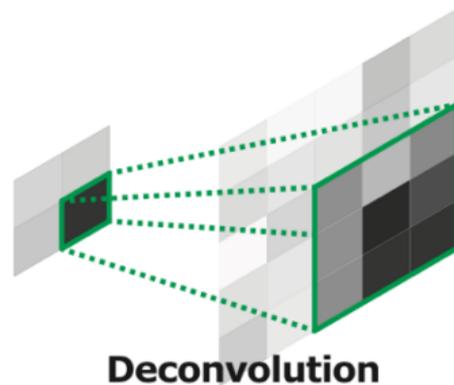
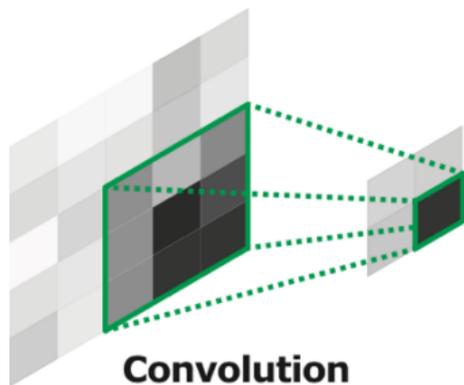
Deconvolution

Deconvolution: A deconvolution (or a **transpose convolution** or **up-sampling convolution**) is an operation that applies a filter to input data in a way that expands its spatial dimensions.

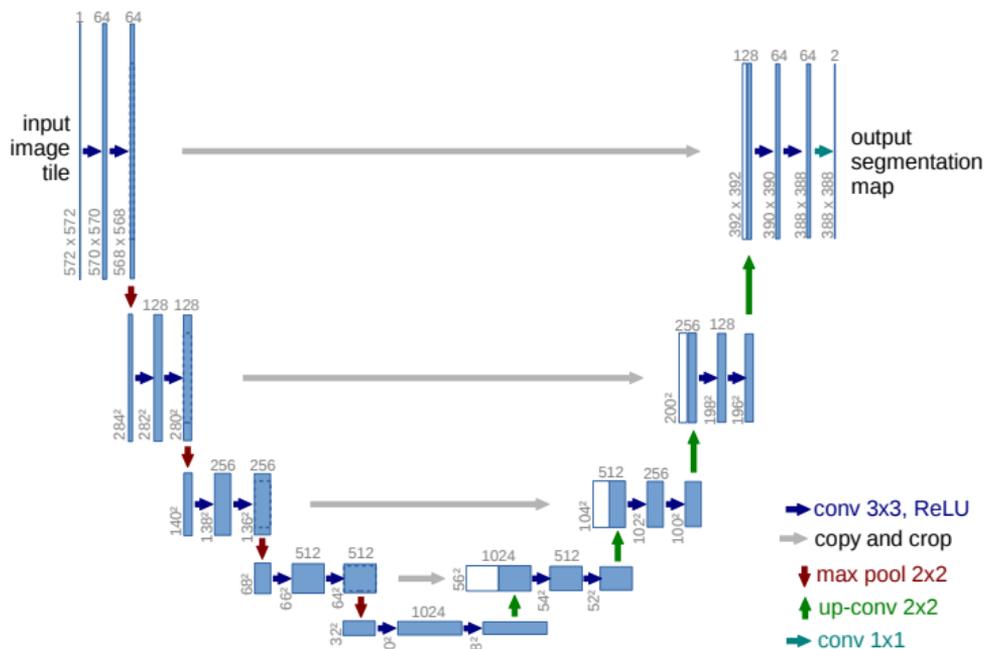
$$\underbrace{\begin{bmatrix} 3 & 0 \\ 1 & 5 \end{bmatrix}}_{\text{input } 2 \times 2} * \underbrace{\begin{bmatrix} 2 & 7 & 4 \\ 3 & 1 & 7 \\ 4 & 2 & 1 \end{bmatrix}}_{\text{filter } 3 \times 3} = \underbrace{\begin{bmatrix} 6 & 21 & 12 & 0 \\ 11 & 20 & 60 & 20 \\ 15 & 23 & 24 & 35 \\ 4 & 12 & 11 & 5 \end{bmatrix}}_{\text{feature map } 4 \times 4}$$

- The stride can be more than 1
- Padding is reversed by discarding boundary pixels.
- For overlaps, use averaging or summation.
- The filter represents patterns, with the input indicating where these patterns are detected.
- In unmax pooling, either duplicate pixels in the output or place the maximum value pixel while setting others to zero.

Illustration of Deconvolution



U-Net Architecture



- With skip connection, U-Net combines (or concatenates) **high-level abstract features** (from deeper layers) and **spatial details** (from earlier layers).

U-Net Output



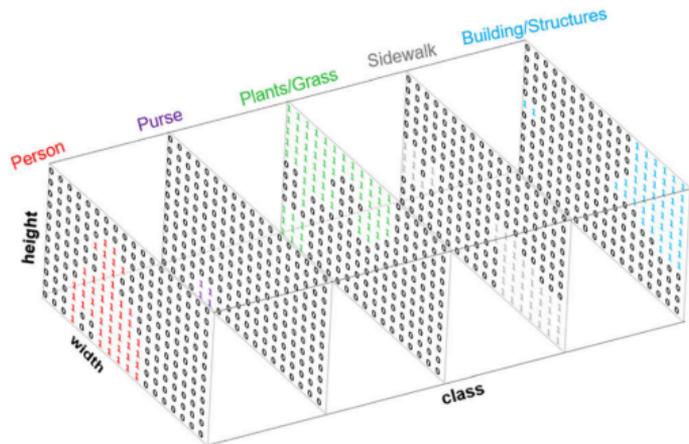
Input

segmented →

- 1: Person
- 2: Purse
- 3: Plants/Grass
- 4: Sidewalk
- 5: Building/Structures

3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	
3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	
3	3	3	3	3	1	1	3	3	3	3	5	5	5	5	
3	3	3	3	3	1	1	1	3	3	3	5	5	5	5	
3	3	3	3	3	1	1	3	3	3	5	5	5	5	5	
5	5	3	3	3	1	1	3	3	5	5	5	5	5	5	
4	4	3	4	1	1	1	1	1	4	4	4	5	5	5	
4	4	3	4	1	1	1	1	1	1	4	4	4	4	5	5
4	4	4	1	1	1	1	1	1	1	4	4	4	4	4	4
3	3	3	1	1	1	1	1	1	1	4	4	4	4	4	4
3	3	3	1	2	2	1	1	1	1	1	4	4	4	4	4
3	3	3	1	2	1	1	1	1	1	4	4	4	4	4	4

Semantic Labels



Outline

- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Semantic Segmentation
- 4 Face Recognition**

Verification vs. Recognition

Verification

- **Input:** An image
- **Output:** Confirms if the image matches the claimed person

Recognition

- **Database:** Contains K known identities
- **Input:** An image
- **Output:** Returns the person's ID if in the database; otherwise, not recognized

One-Shot Learning and Similarity Function

Problem: Traditional deep learning models require large amounts of labeled data, but face recognition needs to identify a person with only a few examples.

One-Shot Learning

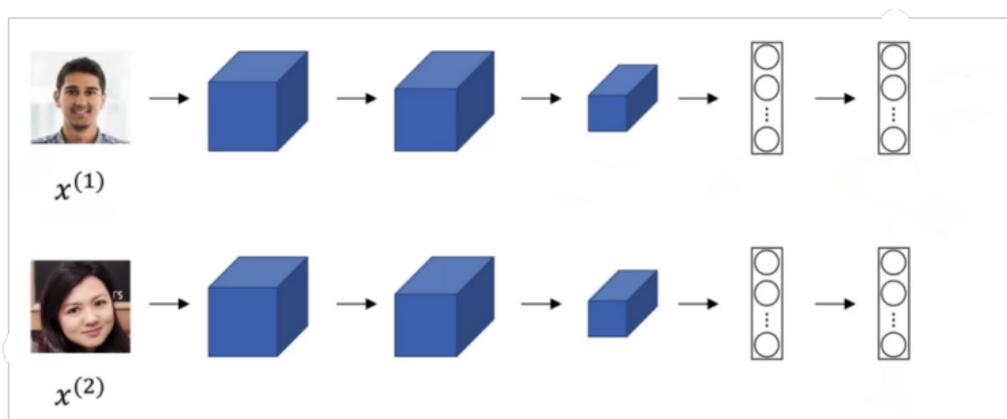
- Learns to recognize a person from a **single example** rather than relying on large datasets.
- Challenges with new identities that were not seen during training.
- Solution: Instead of classification, use a similarity function to compare images

Learning a Similarity Function

- Define a function $d(\mathbf{x}_1, \mathbf{x}_2)$ that measures the **difference** between two images
- If $d(\mathbf{x}_1, \mathbf{x}_2) \leq \text{threshold}$, classify them as the **same** person and return the identity; otherwise, reject.

Siamese Network

Definition: A Siamese network is a neural architecture consisting of **two identical** sub-networks that share weights and parameters.



- The ConvNet $f_{\theta}(x)$ represents the **embedding** (feature representation) of the input image x .
- The similarity between two images is measured using the distance between their **embeddings**:

$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \|\mathbf{f}_{\theta}(\mathbf{x}^{(i)}) - \mathbf{f}_{\theta}(\mathbf{x}^{(j)})\|$$

- The model learns parameters θ such that:
 - If $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ belong to the same person, $d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is small.
 - If they belong to different people, $d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is large.

Triplet Loss



Anchor



Positive



Negative

- Training uses image triplets: an **anchor** \mathbf{a} , a **positive** \mathbf{p} , and a **negative** \mathbf{n} .

- Objective:

$$d(\mathbf{a}, \mathbf{p}) + \alpha \leq d(\mathbf{a}, \mathbf{n})$$

where $\alpha > 0$ is the **margin**, preventing trivial solutions.

- Triplet loss:

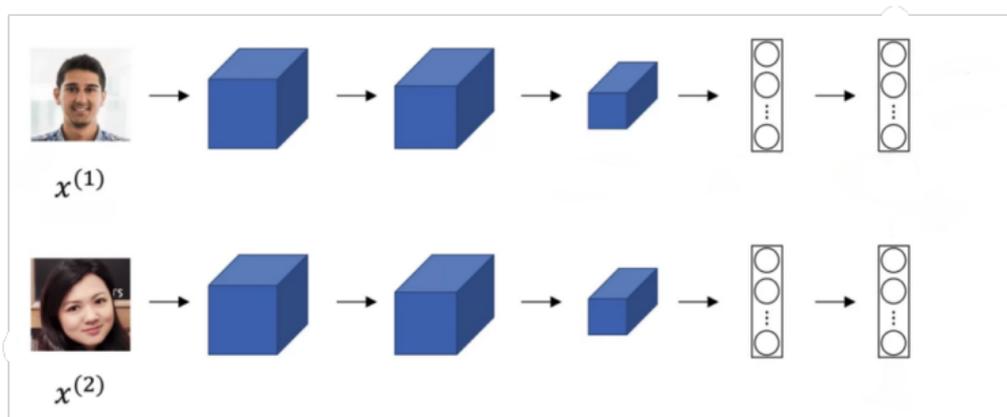
$$\ell(\mathbf{a}, \mathbf{p}, \mathbf{n}) = [d(\mathbf{a}, \mathbf{p}) + \alpha - d(\mathbf{a}, \mathbf{n})]_+,$$

where $[x]_+ = \max\{x, 0\}$.

$$\mathcal{L}(\theta) = \sum_{(\mathbf{a}, \mathbf{p}, \mathbf{n})} \ell(\mathbf{a}, \mathbf{p}, \mathbf{n})$$

- Requires multiple images per identity for effective learning.

Face Recognition via Binary Classification



- The ConvNet $f_{\theta}(x)$ extracts a feature **embedding** from the input image x .
- The similarity between two images is measured as the distance between their embeddings:

$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \|\mathbf{f}_{\theta}(\mathbf{x}^{(i)}) - \mathbf{f}_{\theta}(\mathbf{x}^{(j)})\|$$

- The distance $d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is passed through a **logistic activation** $\sigma(\cdot)$ to classify whether they belong to the same person:

$$\hat{y} = \sigma(d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}))$$