

$$f \otimes h = \sum_k \sum_l f(k, l)h(k, l)$$

$f = \text{Image}$
 $h = \text{Kernel}$

Filtering, Convolution

Hue represents color in angle of HSV cone
0-60: red 60-120: yellow 120-180: green
80-240: cyan 240-300: blue 300-360: magenta
Filtering: Noise removal, feature extraction
Linear filtering: output is a linear combination of pixel values in some neighborhood

f

$$\begin{bmatrix} f_1 & f_2 & f_3 \\ f_4 & f_5 & f_6 \\ f_7 & f_8 & f_9 \end{bmatrix} \otimes \begin{bmatrix} h_1 & h_2 & h_3 \\ h_4 & h_5 & h_6 \\ h_7 & h_8 & h_9 \end{bmatrix}$$

$$f \otimes h = f_1h_1 + f_2h_2 + f_3h_3 + f_4h_4 + f_5h_5 + f_6h_6 + f_7h_7 + f_8h_8 + f_9h_9$$

Fourier Transform

$$\text{Definition} \quad F(u, v) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) e^{-j2\pi(ux+vy)} dx dy,$$

Image Frequency: obtained quantitatively using 2D Fourier Transform,

Fourier analysis -> decompose complex waves into simple constituent simple sin/cos waves
Low Pass Filtering: retains low spatial freq, remove high spatial freq components (noise, texture)
> generally image becomes blurry, bigger kernel -> more blurry
> denoising: low-pass filtering to remove signal comps with high spatial freq
> i.e. gaussian blur $f_x = ae^{-\frac{(x-b)^2}{2c^2}}$
> i.e. average blur (takes average of pixel values in kernel)
> i.e. non-linear filtering/median filtering -> takes median of kernel as output more robust against outliers (salt n pepper), computationally intense

High Pass Filtering: retains high spatial freq components

> Additive Gaussian Noise: $I(\text{distorted}) = I(\text{original}) + n$ (where n is norm dist)
generally more contrast but also more grainy
> Edge detection: detect significant change in intensity values.
provide shape info/sharpening
> i.e. prewitt operator to detect horizontal or vertical edges

(prewitt computes diff of neighbouring pixels to indicate likelihood of edges)

Detect vertical edge Detect horizontal edge

> i.e. Sobel operator: similar to prewitt but has larger weights for edges near centre
output edges are sharper than prewitt

Image Classification

Top-n accuracy: outputs k confidence for each of the k classes, generates top 1/2 etc. classes

eval: no. correct / no. images

Object localization: produce class label and bounding box for each input image

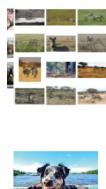
Error if label is incorrect or box has <50% overlap

Object Detection: model predicts sets of labels, bounding boxes and confidence scores.

Penalizes objects in image that are not annotated by model or >1 label for the same object

Challenges in image classification: Primitive data (arrays of values), Diff variation for every class (i.e. perspective, scale, clutter, occlusion, intra-class variations, deformation, lighting)

Data driven approach to solve: provide Big Data. Use learning algorithm.



Training/Learning (usually offline)

Training set: Images with known class information → Learn a model using some algorithm → A model for this specific classification problem and classifier

Testing/Evaluation (usually online)

Testing set (with label during evaluation, without label in an application) → Classifier algorithm → Predicted class information for this new image (compare with ground-truth during evaluation)

Linear Classifier: $L = f(x; W, b) = Wx + b$ OR $s = [Wb] [x]^T$

Given a test image x , produce the confidence score for each class using linear transformation (total: Kclasses) Higher confidence score for a class -> more likely to be the ground-truth class

Test image x : flatten to a Dx1 column vector, D is the image resolution times number of channel:

Input x : Dx1 Weight W : KxD

Bias b : Kx1 Score s : Kx1

Each row of W extracts features of a specific class

Loss Function: measures consistency of GT labels and score output for some W

softmax function (\exp , then normalize): \exp helps avoid small numbers -> rounding/floating point errors

$$\text{softmax}(f) = \frac{e^{f_m}}{\sum_{j=1}^K e^{f_j}} \quad \text{For the m-th class}$$

Cross-entropy loss (apply $-\log(\cdot)$) to only the ground-truth class):

$$L_i = -\log \frac{e^{f_{y_i}}}{\sum_{j=1}^K e^{f_j}} \quad \text{For the i-th training sample}$$

Softmax: Loss for N training samples (x_i, y_i) : $L = (1/N) * \sum(L_i)$

Output is normalized -> confidence across classes sum to 1

Cross-entropy Loss: 1. softmax (exp, normalize) 2. $-\log(\text{GT class})$

> Q: estimated prob (softmax) P: true prob (GT label)

Kullback-Liebler (KL) divergence: $D_{KL}(P||Q) = -\sum_i P(i) \log \frac{Q(i)}{P(i)} = \sum_i P(i)(\log \frac{1}{Q(i)} - \log \frac{1}{P(i)})$

Li = sum(max(0, s_i - s_j))

$D_{KL}(P||Q) = -\sum_x p(x) \log q(x) + \sum_x p(x) \log p(x)$

$\text{sigmoid}(df) = 1 / (1 + e^{-f})$

$\text{tanh} = (e^f z - e^{-f} z) / (e^f z + e^{-f} z)$

- saturation, vanishing gradient

- slow and difficult to train with grad desc

ReLU = max(0, z)

- reduce issue of vanishing grad

- sparse representation

CNN
> Image filtering is used as feature detection. Produces large outputs when patterns similar to the filter is detected in the image. Different filters used to detect different features

> Training images are used to find out what kernel kernels are useful

> Feature detection used in hierarchical manner to recognize object/high level concept

Convolutional Layer: provides local connectivity (receptive field, depth of input), shares params, uses filtering/convolution instead of matrix multiplication to detect features (edge, corner, faces)

> Assumption of image (locality of pixel dependencies, stationary of image statistics, translation invariance, Using the same set of filters for the whole image)

Pooling Layer: progressively reduce spatial size of feature map, reduce model params

> operates independently on each feature map

> can be overlap or non-overlap, avg or max pool, is translation invariant

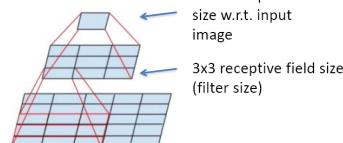
Receptive Field: Part of the image visible to a neuron. Neurons in higher layers can see large portions of the image

Stride = 1

Conv 2

Conv 1

Input layer (image)



$$\text{Output size of conv} = ((N-F)/\text{stride}) - 1$$

Inception Module

> use filters of diff size in the same layer, concat results as output
> 1x1 conv filter to reduce complexity (no. channels)
> increase width and depth of network
> all convolutions use ReLU, including 1x1 conv for reduction/projection

Inception module with dimensionality reduction

28x28x256 output



inception 3a

28x28x32

Challenges of going deeper:

Vanishing gradient: gradient is backpropagated to earlier layers, repeated multiplication may make the gradient very small
Difficult to learn an identity mapping in deep architecture

Deep Residual Learning: learn residual mapping instead of whole mapping

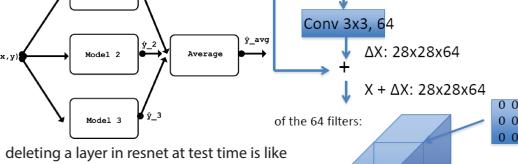
> optimal mapping is close to identity mapping, residual learning is a construction to ease the learning
Performs elem-wise addition channel by channel
Optimization in complex large scale system:
> to ensure going deeper doesn't degrade performance, introduce skip connections to facilitate learning of identity mapping
Ensemble Learning:
> Residue networks r like a collection of paths with many lengths. They sshow ensemble-like behaviour (don't strongly depend on each other) -> most gradients come from short paths (10-34 layers deep) to avoid vanish gradient
> Residual networks behave like ensembles of relatively shallow networks

Deep residual learning

- Learn the residual mapping instead of the entire mapping

Want ΔX to be small, so that the output feature maps would not cause dramatic change (degrade) in performance

- With skip connection, need $W \approx 0$
- Easy to learn with regularization on W



deleting a layer in resnet at test time is like zeroing half the paths. in feed-forward nets such as AlexNet or VGG, deleting a layer alters the only path from input to output.

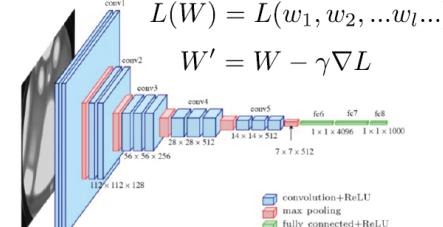
Error also increases smoothly when deleted several modules from a resnet or re-ordering the net by shuffling the modules. Resnet distribution of all path lengths follows a binom distribution. Gradient introduced appears to decay exponentially with no. modules

Gradient Descent: Method of learning W

Using the loss function $L(W) = L = 1/N * \sum(L_i + lr * R(W)) = L(W; (x1, y1), (x2, y2), (x3, y3) ...)$

> Grad Desc: Start from random W, iteratively improve W (reduce L(W))

- Update W by $W + \Delta W$, using the gradient



Gradient descent

$$\nabla L = [\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_l}]^T$$

$$w'_l = w_l - \gamma \frac{\partial L}{\partial w_l}$$

γ : step size (learning rate), a hyperparameter

Gradient is perpendicular (orthogonal) to the level curve. Steepest direction (steepest descent)

- Update W by $W + \Delta W$, using the gradient

$$w'_l = w_l - \gamma \frac{\partial L}{\partial w_l}$$

$$L = \frac{1}{N} \sum_i L_i + \lambda R(W)$$

$$\frac{\partial L}{\partial w_l} = \frac{1}{N} \sum_i \frac{\partial L_i}{\partial w_l} + \lambda \frac{\partial R(W)}{\partial w_l}$$

-Sum gradients for all (partial) training samples for one w_l

-Make one update of W once we have all the gradients

Mini-batch gradient descent / stochastic gradient descent: use small batch (64, 128, 256) for one update of W

$$\frac{\partial L}{\partial w_l} = \frac{1}{N_{batch}} \sum_i \frac{\partial L_i}{\partial w_l} + \lambda \frac{\partial R(W)}{\partial w_l}$$

Computing Gradient: 2 approaches

1. Compute $L(W)$ then differentiate > tedious & inflexible (if layers change, need to re-compute all)

2. Backprop - assume we have dl/dh

- 1. Grad descent (forward): $dl/dW(l) = dl/dh(l) \times dh(l)/dW(l)$
- 2. Backprop: $dl/dh(l-1) = dl/dh(l) \times dh(l)/dh(l-1)$

Modular design for backprop :

> Each module knows how to compute local gradients. Multiply by the global gradient from upstream

> Generalize to other modules, e.g. ReLU

> Obtain exactly same gradient as Approach 1


```

OpenCV and Numpy [15 points]
(a) Given the following python code,
>>> a = [1,3,4,5,2,10]
>>> b = []
>>> for i in a:
...     temp = (i*5) % 2
...     b.append(temp)
>>> print(b)
[1, 1, 0, 1, 0, 0]
Refactor the above python code using numpy and NO loops to achieve the same output (3 pts)
>>> import numpy as np
>>> a = [1,3,4,5,2,10]
### YOUR SOLUTION HERE
>>> a = np.array(a)
>>> b = (a*5) % 2
>>> b = list(b)
>>> print(b) # NOTE: b is a python list
[1, 1, 0, 1, 0, 0]
(b) As seen in Problem Set 1, we use the formula below to convert images to grayscale
0.299×R+0.587×G+0.114×B
Briefly explain why the implementation below of converting color images to grayscale is wrong. (2 pts)
>>> import cv2
>>> #args - image_path: String path of color image
>>> def rgb2gray(image_path):
    img = cv2.imread(image_path, cv2.IMREAD_COLOR)
    out = np.around(\n        img[:, :, 0]*0.299\n        + img[:, :, 1]*0.587\n        + img[:, :, 2]*0.114).astype(np.uint8)
    return out

```

OpenCV splices images in order of B,G,R instead of R,G,B hence the above indices for R and B are flipped.

(c) Given a 3-channel colour square image (10 pts):

```

>>> img = cv2.imread('sample.jpg', cv2.IMREAD_COLOR)
>>> img = cv2.resize(img, (50,50))
>>> H,W,_ = img.shape
>>> x=np.array([-1, 0, 1])
>>> M1 = cv2.getRotationMatrix2D((W/2,H/2),90,1)
>>> M2 = cv2.getRotationMatrix2D((W,H),90,1)

```

Choose the best matches between OpenCV code and numpy code that perform the same functions.

Briefly discuss the description/objective of the corresponding codes. If there is no match, write NIL

Numpy:

- (i) [img[:,2], img[:,1], img[:,0]]
- (ii) scale = np.sum(np.exp(-np.power(x,2)/2))
- (iii) np.transpose(img, [1, 0, 2])[:-1]
- (iv) np.pad(img, ((1,1), (1,1), (0, 0)), mode='constant', constant_values=0)

OpenCV:

- (1) cv2.warpAffine(img, M1, (H,W))
- (2) cv2.filter2D (img, -1, np.ones((5,5), np.float32)/25.0)
- (3) cv2.merge(cv2.split(img))
- (4) cv2.copyMakeBorder(img,1,1,1,1, cv2.BORDER_CONSTANT,value=(0,0,0))
- (5) cv2.getGaussianKernel(5, 1)
- (6) cv2.split(img)
- (7) cv2.warpAffine(img, M2, (W,H))

Answer:

| Numpy | OpenCV | Description |
|-------|--------|---|
| (i) | NIL | Get individual colour splits (Since img is read using cv2.imread() the img is already in BGR ie img[:,0] is blue hence (i) should be [img[:,0], img[:,1], img[:,2]]) |
| (ii) | NIL | Create gaussian kernel of size (3,1) |
| (iii) | 1 | Rotate image 90 degrees anticlockwise |
| (iv) | 4 | Add one row and column of zero pad around width and height image of shape |

| Answer: | | |
|---------|--------|---|
| Numpy | OpenCV | Description |
| (i) | NIL | Get individual colour splits (Since img is read using cv2.imread() the img is already in BGR ie img[:,0] is blue hence (i) should be [img[:,0], img[:,1], img[:,2]]) |
| (ii) | NIL | Create gaussian kernel of size (3,1) |
| (iii) | 1 | Rotate image 90 degrees anticlockwise |
| (iv) | 4 | Add one row and column of zero padding around width and height image of shape |

Calculating no. params in a CNN

$W_c = K^2 \times C \times N$ W_c = Number of weights of the Conv Layer.

$B_c = N$ B_c = Number of biases of the Conv Layer.

$P_c = W_c + B_c$ P_c = Number of parameters of the Conv Layer.

There are no params in pooling layers

Calculating Params in FC layers

1. If FC is connected to Conv layer

$W_{cf} = O^2 \times N \times F$ W_{cf} = Number of weights of a FC Layer which is connected to a Conv Layer.

$B_{cf} = F$ B_{cf} = Number of biases of a FC Layer which is connected to a Conv Layer.

$P_{cf} = W_{cf} + B_{cf}$ O = Size (width) of the output image of the previous Conv Layer.

F = Number of neurons in the FC Layer.

2. If FC is connected to FC layer

$W_{ff} = F_{-1} \times F$ W_{ff} = Number of weights of a FC Layer which is connected to an FC Layer.

$B_{ff} = F$ B_{ff} = Number of biases of a FC Layer which is connected to an FC Layer.

$P_{ff} = W_{ff} + B_{ff}$ P_{ff} = Number of parameters of a FC Layer which is connected to an FC Layer.

F = Number of neurons in the FC Layer.

F_{-1} = Number of neurons in the previous FC Layer.