

Applications: Computer Vision

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Applications

1. Large-Scale Deep Learning
2. Computer Vision
3. Speech Recognition
4. Natural Language Processing
5. Other Applications

Topics in Computer Vision

- Overview
- Preprocessing
 - Contrast Normalization
 - Dataset Augmentation

Computer Vision and Deep Learning

- Computer Vision is one of the most active areas for deep learning research, since
 - Vision is a task effortless for humans but difficult for computers
- Standard benchmarks for deep learning algorithms are:
 - object recognition
 - OCR

Common tasks

- Small core of AI goals aimed at replicating human abilities
 - Object recognition
 - Detection of some form
 - Which object is present?
 - Annotating an image with bounding boxes around each object
 - Transcribing a sequence of symbols from image
 - Labeling each pixel with identity of object it belongs
 - Image synthesis
 - Because generative models are a guiding principle behind deep learning, large body of work on synthesis⁵

Preprocessing

- Some deep learning needs much preprocessing
- Computer vision requires little preprocessing
 - Pixel range
 - Images should be standardized, so pixels lie in same range $[0,1]$, $[-1,1]$, or $[0,255]$ etc
 - Picture size
 - Some architectures need a standard size. So images may need to be scaled
 - May not be needed with convolutional models which dynamically adjust size of pooling regions
 - Data set augmentation
 - Can be seen as a preprocessing step for training set

Training with large data sets

- Large data sets (Imagenet) & models (Alexnet)

- No preprocessing
- Learns invariances



ILSVRC 2012 tasks

1. Classification
2. Classification with localization
3. Fine-grained classification

- Alexnet for Imagenet has one preprocessor
 - Subtract mean across training examples of pixels
 - Dataset: ILSVRC subset of ImageNet: 1000 images in each of 1000 categories: 1.2m training, 50k validation, 150k testing
 - Architecture: CNN with 5 conv layers, max-pool layers, dropout layers, 3 fully connected layers.
 - Performance: top 5 error rate= 15.4% next was 26.2%

Contrast Normalization

- Image contrast can be safely removed
- Contrast refers to the magnitude of the difference between bright and dark pixels
- In deep learning different definition
 - Contrast = standard deviation of pixels
 - For image with r rows and c columns, and RGB image, contrast of entire image is

$$\sqrt{\frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 (X_{i,j,k} - \bar{\mathbf{x}})^2} \quad \text{where} \quad \bar{\mathbf{x}} = \frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 X_{i,j,k}$$

- When std dev is high, values differ more from mean

Global Contrast Normalization

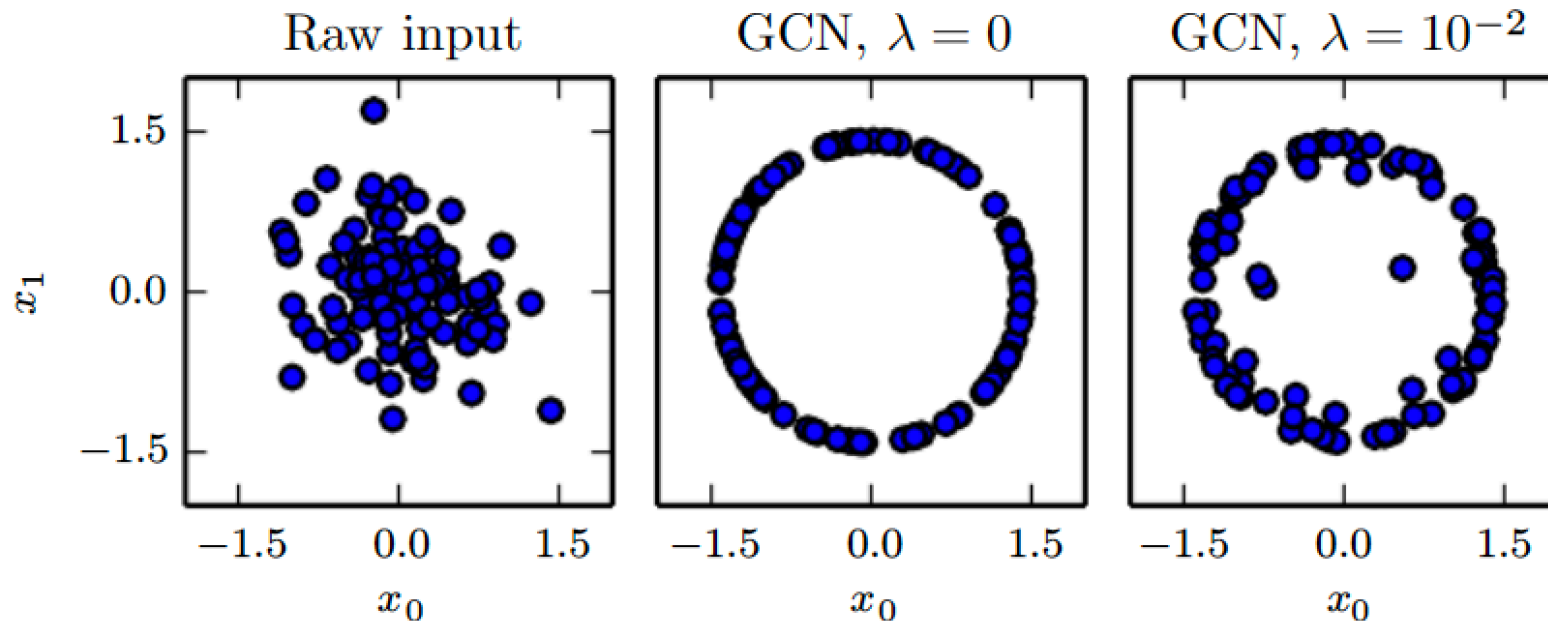
- Aims to prevent images from having varying amounts of contrast
- Subtract mean from each image, then rescale it so that std dev across pixels equals constant s
- Given an input image X , GCN produces an X'

$$X'_{i,j,k} = s \frac{X_{i,j,k} - \bar{X}}{\max \left\{ \epsilon, \sqrt{\lambda + \frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 (X_{i,j,k} - \bar{X})^2} \right\}}$$

- λ is a positive regularization term to bias the std deviation, the denominator is constrained to be at least ϵ

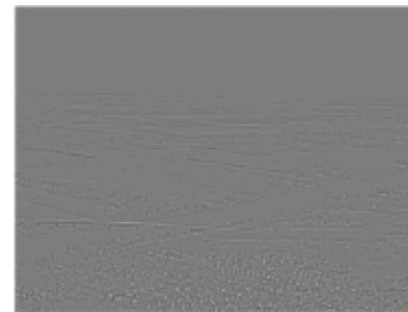
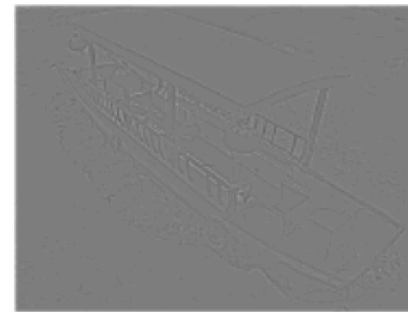
GCN maps examples onto sphere

- Raw input data may have any norm
- $\lambda=0$ maps all nonzero examples onto sphere
- $\lambda>0$ draws examples towards sphere but does not discard variations in norm



Local Contrast Normalization

- Contrast is normalized across each small window rather than entire image



Input image

GCN

LCN

Dataset Augmentation

- Increasing training set by adding modified training examples
 - with transformations that do not change the class
- Object recognition is helped because input may be transformed with many geometric operations
 - Classifiers benefit from random translations, rotations, flips of the input
- In specialized vision applications:
 - Perturbations of colors
 - Nonlinear geometric transformations of input