



CS 540 Introduction to Artificial Intelligence

Deep Learning III

Fred Sala
University of Wisconsin-Madison

March 25, 2021

Announcements

- **Homeworks:**
 - HW7: Released. Grades for several HWs released soon.
- **Midterm:** grading in progress

Tuesday, March 23	Deep Learning II
Thursday, March 25	Deep Learning III
Tuesday, March 30	ML: Summary
Thursday, April 1	Games

- **Class roadmap:**

Artificial Intelligence



Outline

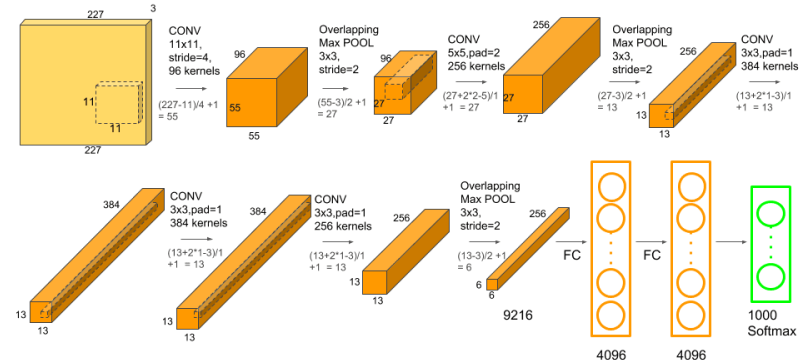
- CNNs with more layers: ResNets
 - Layer problems, residual connections, identity maps
- Data Augmentation & Regularization
 - Expanding the dataset, avoiding overfitting
- More Signal From our Data
 - Graph-structured data, graph neural networks

Last Time: CNNs

We talked about CNN components & architectures

- **Components:** convolutional layers, pooling layers (recall kernels, channels, strides, padding)
- **Architectures:** LeNet, AlexNet, VGG

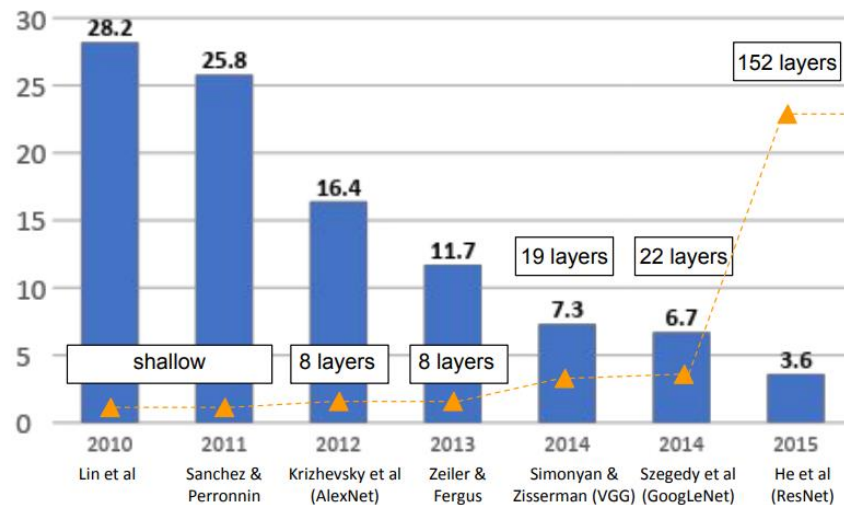
- Trend: bigger, deeper.



Credit: Mathworks

Evolution of CNNs

ImageNet competition (error rate)



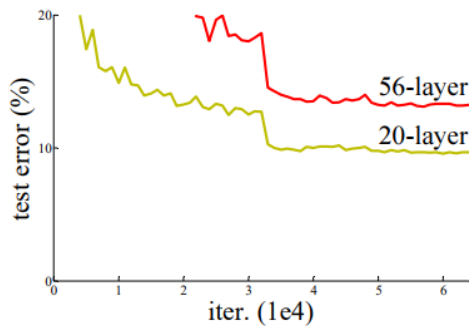
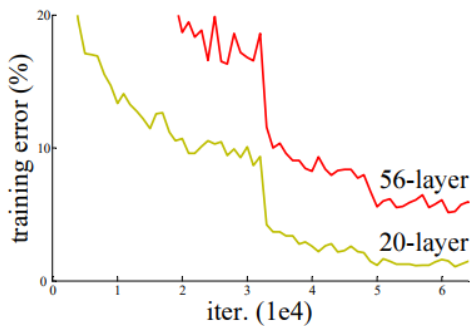
Credit: Stanford CS 231n

Simple Idea: Add More Layers

VGG: 19 layers. ResNet: 152 layers. **Add more layers...** sufficient?

- No! Some problems:
 - i) Vanishing gradients: more layers → more likely
 - ii) Instability: can't guarantee we learn **identity** maps

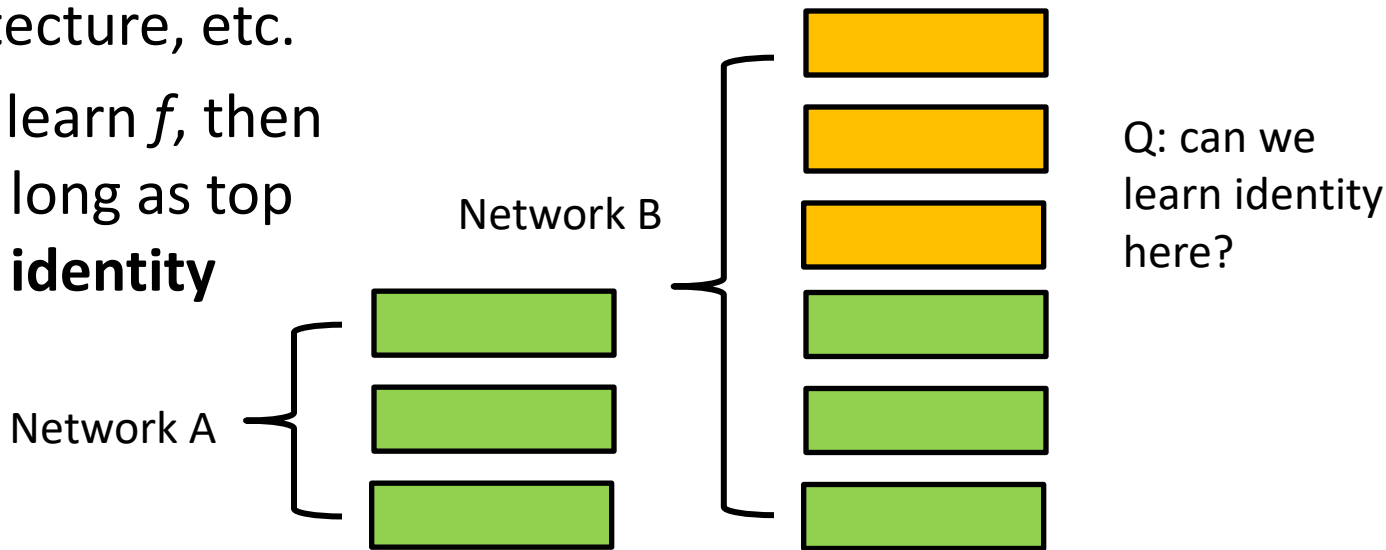
Reflected in training error:



Depth Issues & Learning Identity

Why would more layers result in **worse** performance?

- Same architecture, etc.
- If the A can learn f , then so can B, as long as top layers learn **identity**

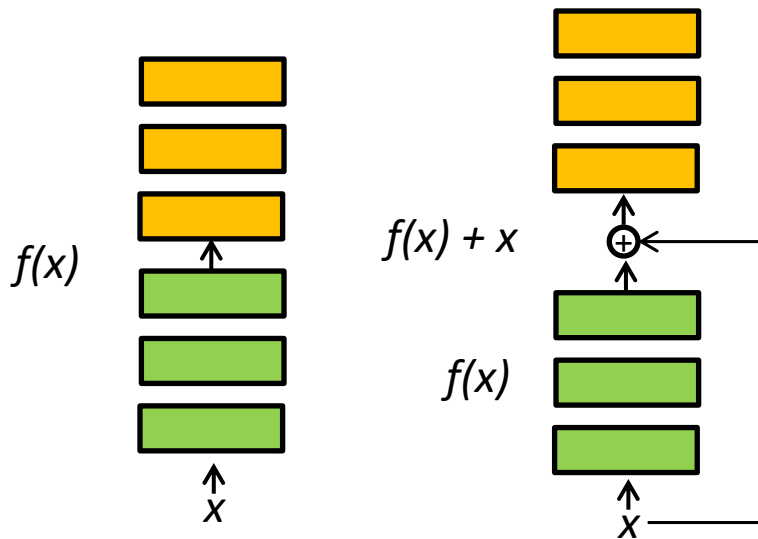


Idea: if layers can learn identity, **can't get worse**.

Residual Connections

Idea: Identity might be hard to learn, but zero is easy!

- Make all the weights tiny, produces zero for output
- Can easily transform learning identity to learning zero:



Left: Conventional layers block

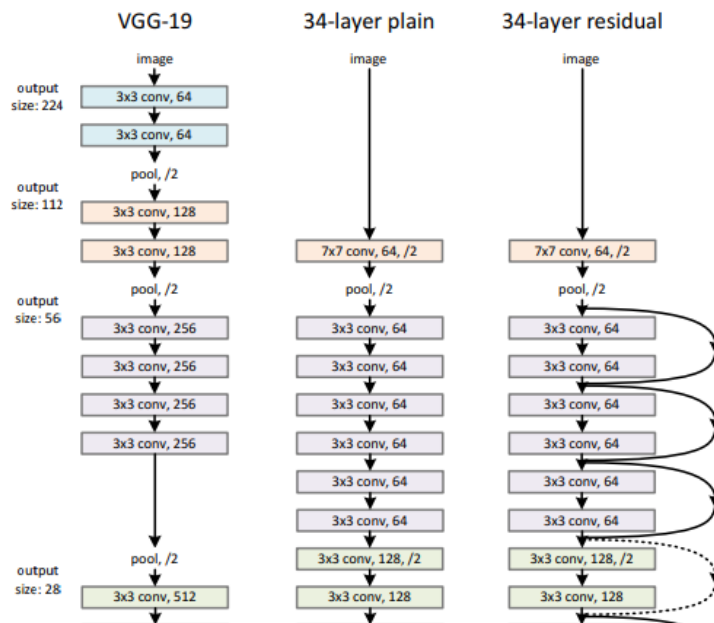
Right: **Residual** layer block

To learn identity $f(x) = x$, layers now need to learn $f(x) = 0 \rightarrow$ easier

ResNet Architecture

Idea: Residual (skip) connections help make learning easier

- Example architecture:
- Note: residual connections
 - Every two layers for ResNet34
- **Vastly better** performance
 - No additional parameters!
 - Records on many benchmarks



He et al: “Deep Residual Learning for Image Recognition”

A Bit More on ResNets

Idea: Residual (skip) connections help make learning easier

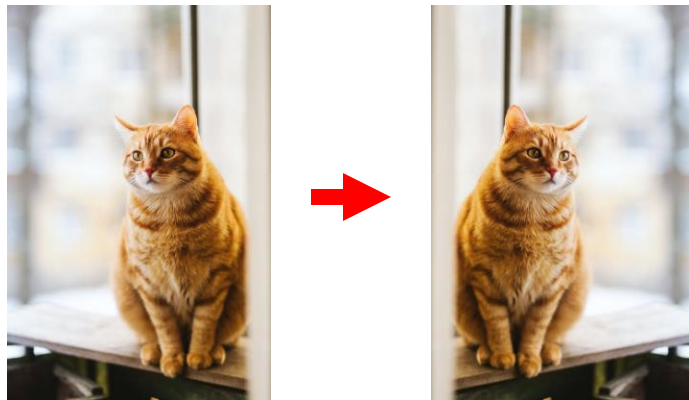
- Note: Can also analyze from **backpropagation** p.o.v
 - Residual connections add paths to computation graph
- Also uses **batch normalization**
 - Normalize the features at each layer to have same mean/variance
 - Common deep learning trick
- Highway networks: learn weights for residual connections

Ioffe and Szegedy: “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”

Data Concerns

What if we don't have a lot of data?

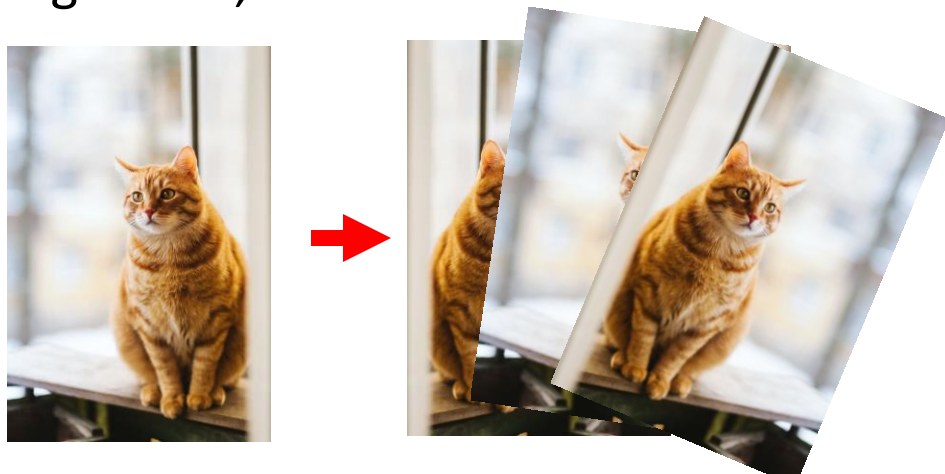
- We risk overfitting
- Avoiding overfitting: **regularization** methods
- Data augmentation: a classic way to regularize



Data Augmentation

Augmentation: transform + add new samples to dataset

- Transformations: based on domain
- Idea: build **invariances** into the model
 - **Ex:** if all images have same alignment, model learns to use it
- Keep the label the same!



Transformations

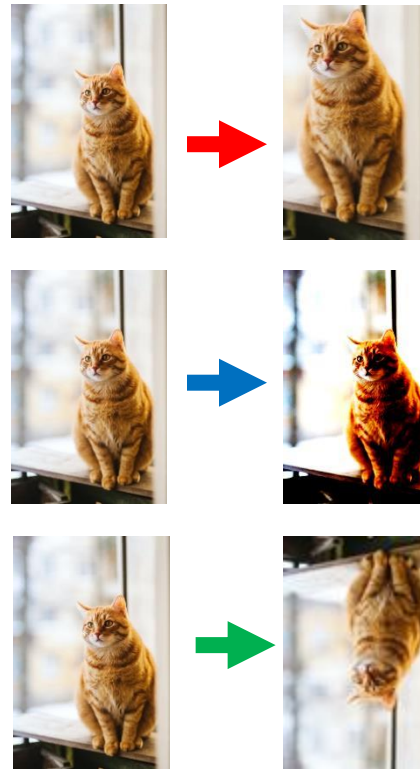
Examples of transformations for images

- **Crop** (and zoom)
- **Color** (change contrast/brightness)
- **Rotations+** (translate, stretch, shear, etc)

Many more possibilities. Combine as well!

Q: how to deal with this at **test time**?

- A: transform, test, average



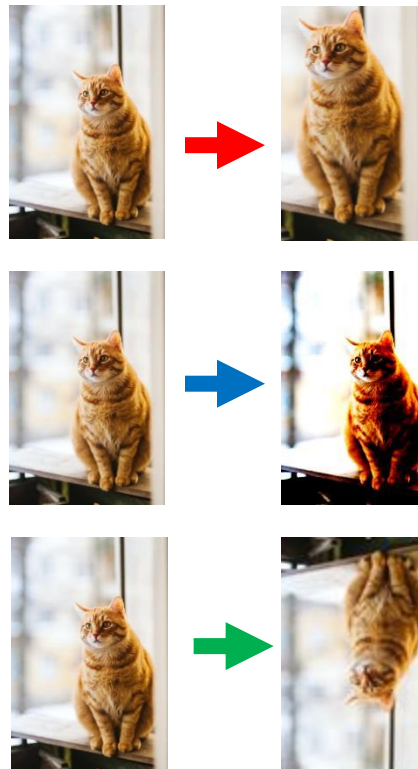
Combining & Automating Transformations

One way to automate the process:

- Apply every transformation and combinations
- **Downside:** most don't help...

Want a good policy, ie,     

- Active area of research: search for good policies
 1. **Ratner et al:** “Learning to Compose Domain-Specific Transformations for Data Augmentation”
 2. **Cubuk et al:** “AutoAugment: Learning Augmentation Strategies from Data”



Other Domains

Not just for image data. For example, on text:

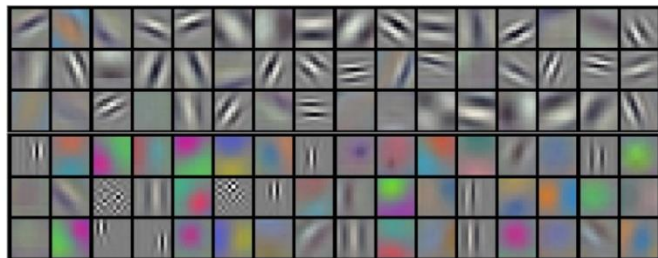
- Substitution
 - E.g., “It is a **great** day” → “It is a **wonderful** day”
 - Use a thesaurus for particular words
 - Or, use a model. Pre-trained word embeddings, language models
- Back-translation
 - “Given the low budget and production limitations, this movie is very good.”
→ “There are few budget items and production limitations to make this film a really good one”

Importance of Augmentation

Data augmentation is critical for top performance!

- You should use it!
- **AlexNet**: used (many papers re-used as well)
 - Random crops, rotations, flips. **2048x** expansion!
 - Color augmentation via PCA. **1% error rate reduction**

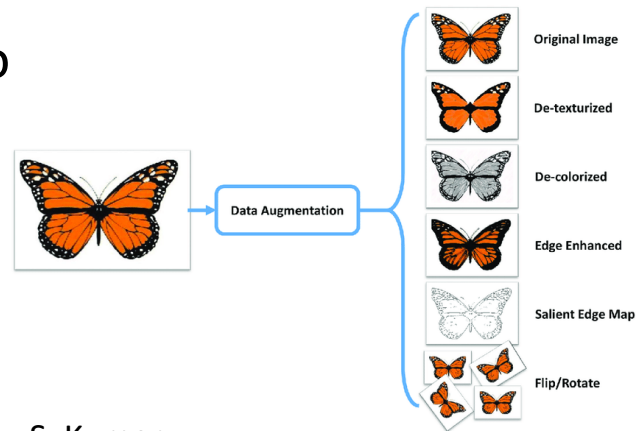
Krizhevsky et al: “ImageNet Classification with Deep Convolutional Neural Networks”



Other Forms of Regularization

Regularization has many interpretations

- **Goodfellow:** “any modification... to a learning algorithm that is intended to reduce its generalization error but not its training error.”
- A way of adding knowledge / side information
- Enforcing parsimony/simplicity



S. Kumar

Other Forms of Regularization

Classic regularizations

1. Modify loss functions

Ex: regularized least squares LR

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n (\theta_0 + x_i^T \theta - y_i)^2 + \lambda \|\theta\|_2^2$$

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i) + \lambda R(f_{\theta})$$

Standard
loss

Regularization
parameter

Regularizer

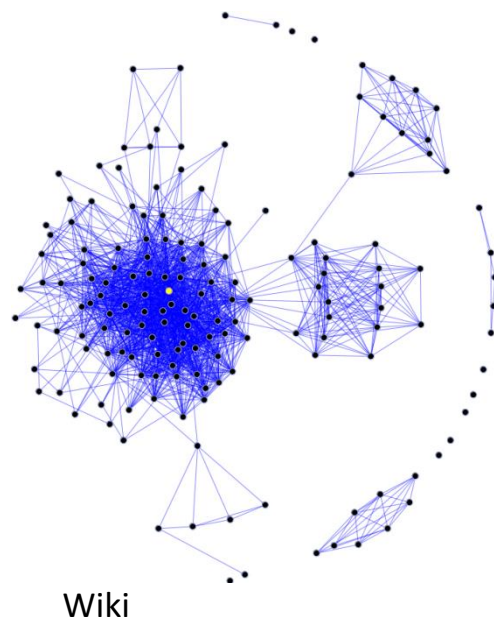
2. Modify architecture/training/data

a) Dropout, batch normalization, augmentation

Relationships in Data

So far, all of our data consists of points

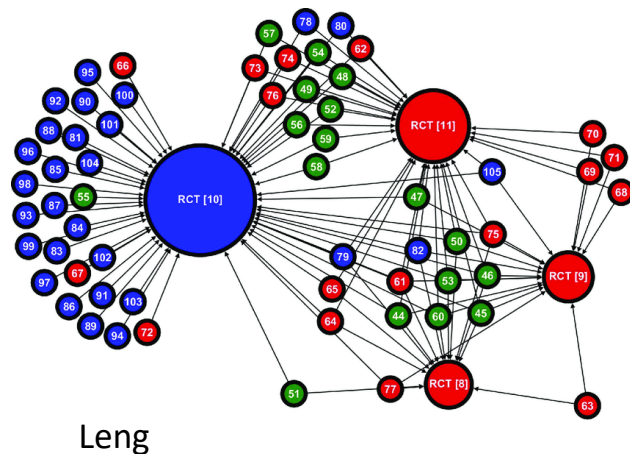
- Assume all are independent, “unrelated” in a sense $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
- Pretty common to have relationships between points
 - **Social networks**: individuals related by friendship
 - **Biology/chemistry**: bonds between compounds, molecules
 - **Citation networks**: Scientific papers cite each other



Signal from Relationships

Suppose we are classifying scientific papers

- **Features:** title, abstract, authors. **Labels:** math/science/eng.
- Could build a reasonable classifier with the above data
- **More signal** from relationships
 - Cite each other, more likely from the same field
 - Note: citations are not features; they're **links**
 - Need a new type of network to handle



Graph Neural Networks

Have: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n), G = (V, E)$

How should our new architecture look?

- Still want layers
 - linear transformation + non-linearity
- Now want to integrate neighbors
- Bottom: graph convolutional network

Hidden Layer Representation



$$H^{(\ell+1)} = \sigma(H^{(\ell)}W^{(\ell)})$$

Non-Linearity



Parameters



$$H^{(\ell+1)} = \sigma(A_G H^{(\ell)} W^{(\ell)})$$

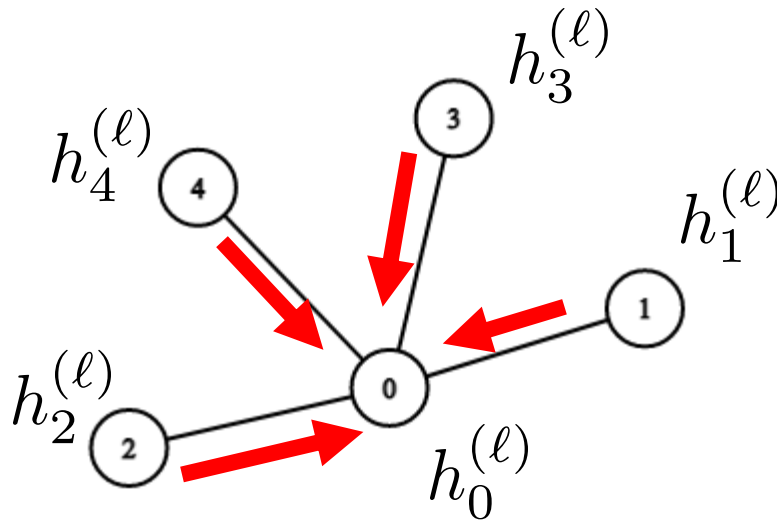


Graph Mixing

Graph Convolutional Networks

Let's examine the GCN architecture in more detail

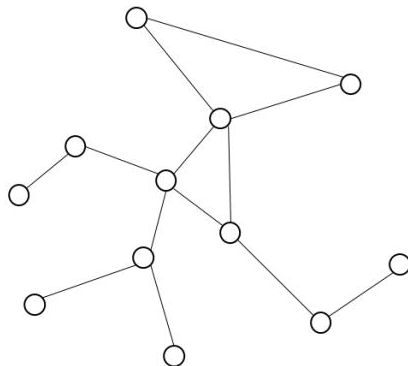
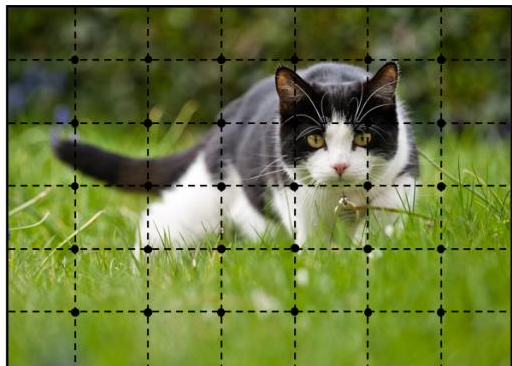
- Difference: “graph mixing” component
- At each layer, get representation at each node
- Combine node's representation with neighboring nodes
- “**Aggregate**” and “**Update**” rules



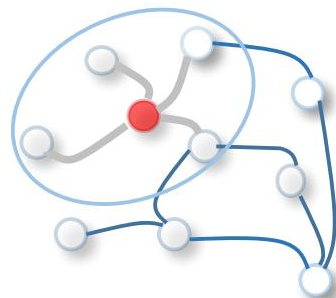
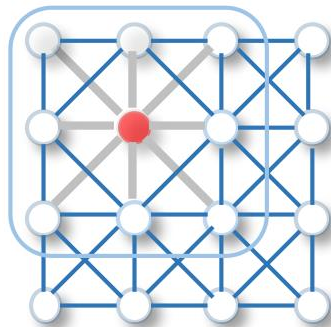
Graph Convolutional Networks

Note the resemblance to CNNs:

- Pixels: arranged as a very regular graph
- Want: more general configurations (less regular)



Wu et al, A Comprehensive Survey on Graph Neural Networks



Zhou et al, Graph Neural Networks: A Review of Methods and Applications

Summary

- Intro to deeper networks (resnets)
 - Dealing with problems by adding skip connections
- Intro to regularization
 - Data augmentation + other regularizers
- Basic graph neural networks



Acknowledgements: Inspired by materials by Fei-Fei Li, Ranjay Krishna, Danfei Xu (Stanford CS231n)